ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

УПРАВЛЕНИЕ, ВЫЧИСЛИТЕЛЬНАЯ ТЕХНИКА И ИНФОРМАТИКА

TOMSK STATE UNIVERSITY JOURNAL OF CONTROL AND COMPUTER SCIENCE

Научный журнал

2020 № 50

Зарегистрирован в Федеральной службе по надзору в сфере массовых коммуникаций, связи и охраны культурного наследия (свидетельство о регистрации ПИ № ФС 77-29497 от 27 сентября 2007 г.)

Подписной индекс в объединённом каталоге «Пресса России» 44031

Журнал включен в «Перечень рецензируемых научных изданий, в которых должны быть опубликованы основные научные результаты диссертаций на соискание ученой степени кандидата наук, на соискание ученой степени доктора наук», Высшей аттестационной комиссии

Founder - Tomsk State University

EDITORIAL BOARD

Alexander Gortsev - Editor-in-Chief, Doctor of Sc., Prof., Head of the Applied Mathematics Department Tomsk State University. Tel: +73822529599 Valery Smagin - Deputy Editor-in-Chief, Doctor of Sc., Prof. of the Applied Mathematics Department Tomsk State University. Tel: +73822529599 Lyudmila Nezhelskaya - Executive Editor, Doctor. of Sc., Prof. of the Applied Mathematics Department Tomsk State University.

E-mail: vestnik_uvti@mail.tsu.ru

Sergey Vorobeychikov - Doctor of Sc., Prof. of the System Analysis and Mathematical Modeling Department Tomsk State University

Vladimir Vishnevsky – Doctor of Sc., Prof. Head of the laboratory Institute of Control Sciences of Russian Academy of Sciences (Moscow, Russia).

Gennady Koshkin - Doctor of Sc., Prof. of the System Analysis and Mathematical Modeling Department Tomsk State University

Yury Kostyuk - Doctor of Sc., Prof. of the Theoretical Informatics Departmen Tomsk State University

Anjela Matrosova - Doctor of Sc., Prof. of the Programming Department Tomsk State University

Anatoly Nazarov- Doctor of Sc., Prof., Head of the Probability Theory and Mathematical Statistics Department Tomsk State University

Konstantin Samouylov- Doctor of Sc., Prof., Head of the Applied Probability and Informatics Department RUDN University (Moscow, Russia)

Eugene Semenkin – Doctor of Sc., Prof. System Analysis and Operations Research Department Reshetnev Siberian State University of Science and Technology (Krasnoyarsk, Russia)

Sergey Sushchenko – Doctor of Sc., Prof., Head of the Applied of Information Department Tomsk State University

Mais Farkhadov – Doctor of Sc., Head of the laboratory Institute of Control Sciences of Russian Academy of Sciences (Moscow, Russia).

Gurami Tsitsiashvili - Doctor of Sc., Prof., Chief researcher Institute for Applied Mathematics Far Eastern Branch of RAS, Prof. Far Eastern Federal University (Vladivostok, Russia)

Editorial address:

Institute of Applied Mathematics and Computer Science, unit of Applied Mathematics

National Research Tomsk State University 36 Lenina Avenue, Tomsk, 634050 Telephone / fax: +73822529599

E-mail: vestnik_uvti@mail.tsu.ru

EDITORIAL COUNCIL

PhD, Prof. University VII Paris, France Vladimir Dombrovskii Doctor of Sc., Prof.

Ana Rosa Cavalli

Tomsk State University Russia

Alexander Dudin Doctor of Sc., Prof.

Belarusian State University Minsk, Republic Belorussia

Enco Orsingher PhD, Prof. University of Rome Italy

Paolo Prinetto

Prof. Politechnic Institute Torino, Italy

Gilbert Saporta PhD, Prof.

Pierre and Marie Curie University, Paris, France Raimund Ubar Doctor of Sc., Prof. University of Technology

Tallinn, Estonia Reindert Nobel

Doctor of Sc., Associate Prof. Vrije University of Amsterdam

Netherlands

Nina Yevtushenko Doctor of Sc., Prof. Ivannikov V.P. ISP RAS Moscow, Russia Yervant Zorian

PhD, Fellow & Chief Architect, Synopsys, Mountain View, CA, USA

Учредитель – Томский государственный университет

РЕДАКЦИОННАЯ КОЛЛЕГИЯ

Горцев Александр Михайлович – гл. редактор, проф., д-р техн. наук, зав. кафедрой прикладной математики ТГУ. Тел. +73822529599

Смагин Валерий Иванович – зам. гл. редактора, проф., д-р техн. наук, проф. кафедры прикладной математики ТГУ. Тел. +73822529599

Нежельская Людмила Алексеевна – ответственный секретарь, доц., д-р физ.мат. наук, проф. кафедры прикладной математики ТГУ E-mail: vestnik_uvti@mail.tsu.ru

Воробейчиков Сергей Эрикович – д-р физ.-мат. наук, проф. кафедры системного анализа и математического моделирования ТГУ

Вишневский Владимир Миронович - проф., д-р техн. наук, зав. лабораторией Института проблем управления РАН (г. Москва)

Кошкин Геннадий Михайлович – проф., д-р физ.-мат. наук, проф. кафедры системного анализа и математического моделирования ТГУ

Костюк Юрий Леонидович - проф., д-р техн. наук, проф. кафедры теоретической информатики ТГУ

Матросова Анжела Юрьевна – проф., д-р техн. наук, проф. кафедры программирования ТГУ **Назаров Анатолий Андреевич** – проф., д-р техн. наук, зав. кафедрой теории

вероятностей и математической статистики ТГУ

Самуйлов Константин Евгеньевич – проф., д-р техн. наук, зав. кафедрой прикладной информатики и теории вероятностей РУДН (г. Москва)

Семенкин Евгений Станиславович – проф., д-р техн. наук, проф. каф. системного анализа и исследования операций, СибГУ им. акад. М.Ф. Решетнева (г. Красноярск)

Сущенко Сергей Петрович – проф., д-р техн. наук, зав. кафедрой прикладной информатики ТГУ

Фархадов Маис Паша Оглы – д-р техн. наук, зав. лабораторией Института проблем управления РАН (г. Москва)

Цициашвили Гурами Шалвович – проф., д-р физ.-мат. наук, гл. науч. сотр. Института прикладной математики ДВО РАН, проф. ДВФУ (г. Владивосток)

Адрес редакции и издателя: 634050, Томск, пр. Ленина, 36

Национальный исследовательский Томский государственный университет, Институт прикладной математики и компьютерных наук,

отделение прикладной математики Телефон / факс: +73822529599 E-mail: vestnik_uvti@mail.tsu.ru

РЕДАКЦИОННЫЙ СОВЕТ

Ана Роза Кавалли д-р философии, проф

Университет VII, Париж, Франция

Владимир Домбровский д-р техн. наук, проф ТГУ, Томск, Россия

Александр Дудин

д-р физ.-мат. наук, проф. БГУ, Минск,

Республика Беларусь

Енцо Орзингер д-р философии, проф.

Римский университет,

Италия

Паоло Принетто

проф. Политехнический институт, Турин, Италия

Университет им. Пьера и Марии, Кюри, Париж, Франция Раймонд Убар

д-р, проф.

Жильберт Сапорта

д-р философии, проф.

Технологический университет, Таллинн, Эстония Рейндерт Нобель

д-р. доцент Свободный университет, Амстердам, Нидерланды

Нина Евтушенко

д-р техн. наук, проф. ИСП РАН им. Иванникова В.П.,

Москва, Россия Ервант Зориан

д-р философии, гл. науч. сотр. фирмы «Синопсис», США

JOURNAL INFO

Tomsk State University Journal of Control and Computer Science is an independent peer-reviewed research journal that welcomes submissions from across the world.

Tomsk State University Journal of Control and Computer Science is issued four times per year, and can be subscribed to in the Russian Press Joint Catalogue (Subscription Index 44031

The publication in the journal is free of charge and may be in Russian or in English. The topics of the journal are the following:

- control of dynamical systems,
- mathematical modeling,
- data processing,
- informatics and programming,
- discrete function and automation,
- designing and diagnostics of computer systems.

Rules of registration articles are given in a site:

http://journals.tsu.ru/informatics/

ISSN 2311-2085 (Online), ISSN 1998-8605 (Print).

О ЖУРНАЛЕ

Журнал «Вестник Томского государственного университета. Управление, вычислительная техника и информатика» выходит ежеквартально и распространяется по подписке

Статьи публикуются на русском и английском языках.

Тематика публикаций журнала:

- управление динамическими системами.
- математическое моделирование, • обработка информации,
- информатика и программирование,
- дискретные функции и автоматы,
- проектирование и диагностика вычислительных систем.

Журнал входит в систему Российского Индекса Научного Цитирования (РИНЦ).

Правила оформления статей приведены на сайте:

http://journals.tsu.ru/informatics

ISSN 2311-2085 (Online), ISSN 1998-8605 (Print).

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020

Управление, вычислительная техника и информатика

№ 50

СОДЕРЖАНИЕ СОЛТЕПТЅ

УПРАВЛЕНИЕ ДИНАМИЧЕСКИМИ СИСТЕМАМИ		CONTROL OF DYNAMICAL SYSTEMS			
Пашинская Т.Ю., Домбровский В.В. Стратегии прогнозирующего управления инвестиционным портфелем на финансовом рынке со скрытым переключением режимов		Pashinskaya T.Yu., Dombrovskii V.V. Predictive control strategies for investment portfolio in the financial market with hidden regime switching			
ОБРАБОТКА ИНФОРМАЦИИ		DATA PROCESSING			
Батраева И.А., Нарцев А.Д., Лезгян А.С. Использование анализа семантической близости слов при решении задачи определения		Batraeva I.A., Nartsev A.D., Lezgyan A.S. Using the analysis of semantic proximity of words in solving the problem			
жанровой принадлежности текстов методами глубокого обучения	14	of determining the genre of texts within deep learning	14		
Воробьев А.В., Воробьева Г.Р. Подход к повышению производительности программных процессов обработки и		Vorobev A.V., Vorobeva G.R. Approach to improving the performance of software processes for processing and			
хранения больших объемов геомагнитных данных Копать Д.Я., Матальцкий М.А. Анализ ожидаемых доходов в открытых	23	storing large volumes of geomagnetic data	23		
марковских сетях с различными особенностями	31	Markov networks with various features	31		
Оценивание современной стоимости <i>п</i> -летней ренты для смешанного страхования жизни Rouban A.I., Mikhalev A.S.	39	Estimation of present value of n-year life annuity for endowment insurance	39		
The global optimization method with selective averaging of the discrete decision variables	47	The global optimization method with selective averaging of the discrete decision variables	47		
с останавливающейся интенсивностью входного потока	56	with the staying intensity of the input flow	56		
Active parametrical identification of stochastic linear continuous-discrete systems based on the experiment	61	Chubich V.M., Filippova E.V. Active parametrical identification of stochastic linear continuous-discrete systems based on the experiment	<i>c</i> 1		
design in the presence of abnormal observations ИНФОРМАТИКА И ПРОГРАММИРОВАНИЕ	61	design in the presence of abnormal observations INFORMATICS AND PROGRAMMING	61		
Золоторевич Л.А.		Zolotorevich L.A.			
Аппаратная защита цифровых устройств	69	Hardware protection of digital devices Isaeva O.S., Kulaysov N.V., Isaev S.V. Method of structural and graphical analysis	69		
интеллектуальной имитационной модели	79	and verification of intellectual simulation model	79		
Численные исследования пропускной способности транспортного протокола с механизмом прямой коррекции ошибок в межсегментном пространстве	89	Numerical studies of transport protocol throughput with forward error correction mechanism in intersegment space	89		
Пазников А.А. Распределенная очередь с ослабленной семантикой выполнения операций в модели удаленного	07	Paznikov A.A. Distributed relaxed queue in remote memory access model designing and diagnostics	07		
доступа к памяти		of computer systems			
Солдатов А.И., Матросова А.Ю., Ким О.Х., Солдатов А.А., Костина М.А. Программируемая коммутационная среда		Soldatov A.I., Matrosova A.Yu., Kim O.H., Soldatov A.A., Kostina M.A. Programmable switching area			
ОБЗОРЫ		REVIEWS	117		
Петухова Н.В., Фархадов М.П., Качалов Д.Л. Разгрузка и консолидация вычислительных ресурсов	122	Petukhova N.V., Farkhadov M.P., Kachalov D.L. Unloading and consolidation of computing resources	122		
в среде туманных и граничных вычислений	130	in the environment of fog and boundary computing	123		

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020

Управление, вычислительная техника и информатика

№ 50

УПРАВЛЕНИЕ ДИНАМИЧЕСКИМИ СИСТЕМАМИ

УДК 519.865.5

DOI: 10.17223/19988605/50/1

Т.Ю. Пашинская, В.В. Домбровский

СТРАТЕГИИ ПРОГНОЗИРУЮЩЕГО УПРАВЛЕНИЯ ИНВЕСТИЦИОННЫМ ПОРТФЕЛЕМ НА ФИНАНСОВОМ РЫНКЕ СО СКРЫТЫМ ПЕРЕКЛЮЧЕНИЕМ РЕЖИМОВ

Рассматривается задача управления инвестиционным портфелем на финансовом рынке с переключением режимов с учетом явных ограничений на объемы вложений и займов и транзакционных издержек. Предполагается, что параметры финансовых активов изменяются в соответствии с эволюцией дискретной скрытой марковской цепи. Для оценки параметров используется адаптивный ЕМ-алгоритм. Представлены результаты численного моделирования с использованием реальных данных российского фондового рынка.

Ключевые слова: инвестиционный портфель; скрытая марковская цепь; прогнозирующее управление; ограничения.

Задача управления инвестиционным портфелем (ИП) является одной из ключевых в финансовой инженерии. Финансовые временные ряды представляют собой нестационарные динамические стохастические системы с высокой волатильностью и скачкообразными изменениями. В связи с этим для описания динамики ИП широко используются модели с марковскими скачками.

Задаче управления ИП на финансовом рынке с марковским переключением режимов посвящены работы [1–7]. В этих работах предполагается, что цепь Маркова является наблюдаемой. Однако на практике при управлении реальным ИП состояние цепи, как правило, не доступно прямому наблюдению.

В работах [8, 9] рассматривается задача управления ИП на скачкообразном рынке со скрытой сменой режимов цепи. В частности, работа [8] посвящена задаче управления по критерию «meanvariance». Оценки параметров модели скрытой цепи Маркова получены с использованием ЕМ-алгоритма. Оптимизационная задача сводится к решению уравнений Гамильтона—Якоби—Беллмана. В работе [9] исследуется задача оптимизации ИП по критерию «mean-variance» с учетом квадратичных транзакционных издержек и ограничений. Для решения задачи используется метод управления с прогнозирующей моделью (Model Predictive Control).

В данной работе рассматривается динамическая задача управления ИП на финансовом рынке с переключением режимов с учетом явных ограничений на объемы вложений и займов и транзакционных издержек. Задача управления ИП формулируется как динамическая задача слежения со скользящим горизонтом инвестирования за эталонным портфелем, имеющим заданную доходность. Предполагается, что параметры финансовых активов изменяются в соответствии с эволюцией дискретной скрытой марковской цепи. Для оценки параметров используется адаптивный ЕМ-алгоритм, предложенный в работе [10]. Представлены результаты численного моделирования с использованием реальных данных российского фондового рынка.

1. Описание модели ИП и определение оптимальной стратегии управления

Рассмотрим ИП, состоящий из n рисковых вложений и безрискового финансового актива (например, банковский счет или надежные облигации). Допускаются также возможность займа по

безрисковой ставке и участие в операциях «продажи без покрытия». Управление портфелем осуществляется путем перераспределения капитала между различными видами инвестиций посредством банковского счета [11].

Пусть $x_i(k)$ $(i=\overline{1,n})$ — объем вложений в i-й рисковый актив в момент времени k; $x_{n+1}(k) \geq 0$ — объем вложений в безрисковый актив; $u_i^+(k) \geq 0$ — объем капитала, переведенного с банковского счета в i-й рисковый актив в k-м периоде; $u_i^-(k) \geq 0$ — объем капитала, переведенного с i-го рискового актива на банковский счет. Если $x_i(k) < 0$ $(i=\overline{1,n})$, то это означает участие в операции «продажа без покрытия» на сумму $|x_i(k)|$.

Допускается также возможность займа по безрисковой ставке. Объем займа безрискового актива равен $x_{n+2}(k) \ge 0$; v(k) – объем заемного капитала, перераспределяемого между банковским и кредитным счетами в k-м периоде: v(k) > 0 означает заем в размере v(k), v(k) < 0 означает возврат кредита в размере |v(k)|; $r_1(k+1)$ – ставка доходности безрискового актива за период (k,k+1], $r_2(k+1)$ – ставка займа безрискового актива за период (k,k+1], $r_1(k+1) < r_2(k+1)$.

Динамика вложений в рисковый актив *i*-го вида $x_i(k)$ ($i = \overline{1,n}$) удовлетворяет уравнению

$$x_i(k+1) = [1 + \eta_i(k+1)][x_i(k) + u_i^+(k) - u_i^-(k)], \tag{1}$$

где $\eta_i(k+1)$ — ставка доходности i-го рисково актива за период времени [k, k+1], определяемая по формуле $\eta_i(k+1) = (Z_i(k+1) - Z_i(k)) / Z_i(k)$, $Z_i(k)$ — рыночная цена i-го рискового актива в момент времени k (наблюдаемая величина).

Предполагается, что транзакционные издержки при покупке и продаже рисковых активов удерживаются из банковского счета (безрискового вложения), динамика которого имеет вид:

$$x_{n+1}(k+1) = [1 + r_1(k+1)][x_{n+1}(k) + v(k) - (1 + \lambda^+) \sum_{i=1}^n u_i^+(k) + (1 - \lambda^-) \sum_{i=1}^n u_i^-(k)],$$
 (2)

где λ^+ — доля капитала $u_i^+(k)$, идущая на уплату транзакционных издержек при покупке рискового актива i-го вида, а λ^- — доля капитала $u_i^-(k)$, идущая на уплату издержек при продаже рискового актива i-го вида.

Динамика кредитного счета описывается уравнением

$$x_{n+2}(k+1) = [1 + r_2(k+1)][x_{n+2}(k) + v(k)].$$
(3)

Поскольку $x_{n+1}(k+1) \ge 0$, $x_{n+2}(k+1) \ge 0$, то справедливы неравенства

$$x_{n+1}(k) + v(k) - (1+\lambda^{+}) \sum_{i=1}^{n} u_{i}^{+}(k) + (1-\lambda^{-}) \sum_{i=1}^{n} u_{i}^{-}(k) \ge 0,$$
(4)

$$x_{n+2}(k) + v(k) \ge 0.$$
 (5)

Будем полагать, что объем операций «продажа без покрытия» по активу i-го вида ограничен величиной $d_i(k) \geq 0$, следовательно, справедливо неравенство:

$$x_i(k) + u_i^+(k) - u_i^-(k) \ge -d_i(k), (i = \overline{1, n}),$$
 (6)

если «продажи без покрытия» запрещены, то $d_i(k) = 0$. Объем заемных средств также ограничен величиной $d_0(k) \ge 0$, следовательно,

$$x_{n+2}(k) + v(k) \le d_0(k). \tag{7}$$

Величины $d_i(k)$ (i=0,...,n) часто зависят от величины общего капитала ИП V(k), что можно учесть, положив $d_i(k)=\gamma_i V(k)$, где $\gamma_i>0$ – постоянный коэффициент.

Капитал инвестиционного портфеля V(k) описывается уравнением

$$V(k) = \sum_{i=1}^{n+1} x_i(k) - x_{n+2}(k).$$
 (8)

Будем полагать, что эволюция доходностей рисковых активов $\eta_i(k)$ (i = 1, ..., n) описывается разностной аппроксимацией уравнений геометрического (экономического) броуновского движения с параметрами, зависящими от состояния цепи Маркова [1, 7]:

$$\eta_i \left[\theta(k), k \right] = \mu_i \left[\theta(k), k \right] + \sum_{i=1}^n \sigma_{ij} \left[\theta(k), k \right] w_j(k), \tag{9}$$

где $\mu_i[\theta(k), k]$ – ожидаемая доходность i-го рискового вложения; $\sigma[\theta(k), k] = {\sigma_{ij}[\theta(k), k]}_{i,j=1,...,n}$ – матрица волатильностей; $\{w_j(k); k=0, 1, ...; j=1, ..., n\}$ – независимые между собой дискретные белые шумы с нулевым средним и единичной дисперсией; $\theta(k) = [\delta(\alpha(k), 1), ..., \delta(\alpha(k), \nu)]^T$, $\delta(\alpha(k), j)$ – функция Кронекера $(j=1, 2, ..., \nu)$; $\alpha(k)$ – однородная дискретная цепь Маркова, принимающая значения из конечного множества $\{1, 2, ..., \nu\}$, с матрицей переходных вероятностей

$$P = [P_{ij}], (i, j \in \{1, 2, ..., v\}), P_{ji} = P\{\alpha(k+1) = j | \alpha(k) = i\}, \sum_{i=1}^{v} P_{ji} = 1,$$

и начальным распределением $p_i = P\left\{\alpha(0) = i\right\}, i = \overline{1, \nu}, \sum\limits_{i=1}^{\nu} p_i = 1.$

Последовательности $w_j(k)$ и $\alpha(k)$ независимы. Марковская цепь $\alpha(k)$ определяет состояние (режим) рынка, например рынок в состоянии высокой или низкой волатильности.

Ожидаемые доходности и волатильности принимают одно из возможных значений из заданного набора в зависимости от состояния цепи Маркова:

$$\mu_{i}[\theta(k),k] \in \left\{\mu_{i}^{(1)},...,\mu_{i}^{(v)}\right\}, \sigma[\theta(k),k] \in \left\{\sigma^{(1)},...,\sigma^{(v)}\right\}, \sigma^{(l)} = \left\{\sigma_{ij}^{(l)}\right\}, \ \left(i,j=\overline{1,n}\right), \ \left(l=\overline{1,v}\right).$$

С учетом (9), уравнение (1) примет вид:

$$x_{i}(k+1) = [1 + \mu_{i} [\theta(k+1), k+1] + \sum_{i=1}^{n} \sigma_{ij} [\theta(k+1), k+1] w_{j}(k+1)] [x_{i}(k) + u_{i}^{+}(k) - u_{i}^{-}(k)].$$
 (10)

Введем обозначения: $x(k) = [x_1(k), x_2(k), ..., x_{n+2}(k)]^{\mathrm{T}}$ — вектор, определяющий состояние портфеля в момент времени k; $u(k) = \begin{bmatrix} v(k) & u_1^+(k) & ... & u_n^+(k) & u_1^-(k) & ... & u_n^-(k) \end{bmatrix}^{\mathrm{T}}$ — вектор управляющих переменных. Тогда с учетом (2), (3), (10), эволюция капитала ИП может быть представлена в виде разностного уравнения [7]:

$$x(k+1) = \left[A_0[\theta(k+1), k+1] + \sum_{i=1}^{n} A_j[\theta(k+1), k+1] w_j(k+1)\right] x(k) + \left[B_0[\theta(k+1), k+1] + \sum_{i=1}^{n} B_j[\theta(k+1), k+1] w_j(k+1)\right] u(k),$$
(11)

где

$$A_0[\theta(k), k] = \operatorname{diag} \left\{ b_0[\theta(k), k], 1 + r_1(k), 1 + r_2(k) \right\}, \quad A_j[\theta(k), k] = \operatorname{diag} \left\{ \sigma_{1j}[\theta(k), k], ..., \sigma_{nj}[\theta(k), k], 0, 0 \right\},$$

$$B_0[\theta(k),k] = \begin{bmatrix} \bar{\mathbf{0}}_n^T & b_0[\theta(k),k] & -b_0[\theta(k),k] \\ 1 + r_1(k) & -(1+\lambda^+)b_1(k) & (1-\lambda^-)b_1(k) \\ 1 + r_2(k) & \bar{\mathbf{0}}_n & \bar{\mathbf{0}}_n \end{bmatrix},$$

$$B_{j}[\theta(k),k] = \begin{bmatrix} \bar{0}_{n}^{T} & \bar{b}_{j}[\theta(k),k] & -\bar{b}_{j}[\theta(k),k] \\ 0 & \bar{0}_{n} & \bar{0}_{n} \\ 0 & \bar{0}_{n} & \bar{0}_{n} \end{bmatrix},$$

$$\begin{split} b_0[\theta(k),k] &= \mathrm{diag} \left\{ 1 + \mu_1[\theta(k),k], ..., 1 + \mu_n[\theta(k),k] \right\}, \ b_1(k) = [1+r_1(k)] 1_n, \\ \bar{b}_j[\theta(k),k] &= \mathrm{diag} \left\{ \sigma_{1,j}[\theta(k),k], ..., \sigma_{n,j}[\theta(k),k] \right\}, \ j = \overline{1,n}, \ \bar{0}_n = [0,...,0]_n, \ 1_n = [1,...,1]_n. \end{split}$$

Ограничения $u_i^+(k) \ge 0$, $u_i^-(k) \ge 0$ и (4)–(7) могут быть записаны в матричном виде:

$$D(k) \le S(k)u(k),\tag{12}$$

где

$$S(k) = \begin{bmatrix} \overline{0}_{n}^{\mathrm{T}} & I_{n} & 0_{n} \\ \overline{0}_{n}^{\mathrm{T}} & 0_{n} & I_{n} \\ \overline{0}_{n}^{\mathrm{T}} & I_{n} & -I_{n} \\ 1 & -(1+\lambda^{+})1_{n} & (1-\lambda^{-})1_{n} \\ 1 & \overline{0}_{n} & \overline{0}_{n} \\ -1 & \overline{0}_{n} & \overline{0}_{n} \end{bmatrix}, \quad D(k) = \begin{bmatrix} \overline{0}_{n}^{\mathrm{T}} \\ \overline{0}_{n}^{\mathrm{T}} \\ \overline{X}(k) \\ -x_{n+1}(k) \\ -x_{n+2}(k) \\ x_{n+2}(k) - d_{0}(k) \end{bmatrix}, \quad \overline{X} = \begin{bmatrix} -x_{1}(k) - d_{1}(k) \\ \dots \\ -x_{n}(k) - d_{n}(k) \end{bmatrix}.$$

Будем определять стратегию управления ИП путем перераспределения капитала между различными видами инвестиций так, чтобы капитал реального портфеля с минимально возможными отклонениями следовал капиталу некоторого определяемого инвестором эталонного портфеля с желаемой доходностью μ_0 , эволюция которого описывается уравнением

$$V^{0}(k+1) = [1 + \mu_{0}]V^{0}(k), V^{0}(0) = V(0).$$
(13)

Критерий качества управления со скользящим горизонтом инвестирования имеет вид:

$$J(k+m|k) = \sum_{i=1}^{m} E \left\{ \rho_{1}(k+i) \left[V(k+i|k) - V^{0}(k+i) \right]^{2} - \rho_{2}(k+i) \left[V(k+i|k) - V^{0}(k+i) \right] + u^{T}(k+i-1|k) R(k+i-1)u(k+i-1|k) \middle| V(k), \theta(k) \right\},$$
(14)

где m – горизонт прогноза, k – текущий момент времени; V(k+i|k)=cx(k+i|k), $c=[1,...,1,-1]_{n+2}$, – прогнозное значение капитала ИП согласно уравнению динамики (11); $u(k+i|k)=[v(k+i|k), u_1^+(k+i|k), ..., u_n^+(k+i|k), u_1^-(k+i|k), ..., u_n^-(k+i|k)]^{\mathrm{T}}$ – вектор прогнозирующих управлений; $\rho_1(k+i) \geq 0$, $\rho_2(k+i) \geq 0$ – весовые коэффициенты (скалярные величины); R(k+i) > 0 – положительно определенная симметричная матрица размерности $(2n+1) \times (2n+1)$.

Критерий (14) может быть записан в виде:

$$J(k+m|k) = \sum_{i=1}^{m} E\left\{x^{\mathrm{T}}(k+i)R_{1}(k+i)x(k+i) - -R_{2}(k+i)x(k+i) + u^{\mathrm{T}}(k+i-1|k)R(k+i-1)u(k+i-1|k)|x(k), \theta(k)\right\},$$
(15)

где $R_1(k+i) = c^{\mathrm{T}}c$ и $R_2(k+i) = [2\rho_1(k+i)V^0(k+i) + \rho_2(k+i)]c$.

Решение данной задачи управления ИП дается следующей теоремой.

Теорема. Пусть капитал ИП описывается уравнением (11) при ограничениях (12). Стратегия прогнозирующего управления u(k+i|k) (i=0,1,...,m-1) со скользящим горизонтом m, минимизирующая критерий (15), при ограничениях (12) на каждом шаге k определяется уравнением

$$u(k) = \begin{bmatrix} I_{2n+1} & 0_{2n+1} & \dots & 0_{2n+1} \end{bmatrix} U(k),$$

где I_{2n+1} — единичная матрица размерности 2n+1, 0_{2n+1} — квадратная нулевая матрица размерности 2n+1; $U(k) = [u^{\mathrm{T}}(k\,|\,k),...,u^{\mathrm{T}}(k+m-1\,|\,k)]^{\mathrm{T}}$ — последовательность прогнозирующих управлений, которая определяется из решения задачи квадратичного программирования с критерием вида:

$$Y(k+m|k) = [2x^{T}(k)G(k) - F(k)]U(k) + U^{T}(k)H(k)U(k),$$

при ограничениях

$$U_{\min}(k) \le \overline{S}(k)U(k) \le U_{\max}(k),$$

 $H(k) = \{H_{ts}(k)\}, G(k) = \{G_t(k)\}, F(k) = \{F_t(k)\}\ (s,t=1,m)$ – блочные матрицы, блоки которых удовлетворяют уравнениям:

$$\begin{split} H_{tt}(k) &= \sum_{i_{t}=1}^{V} \sum_{j=0}^{n} \left(B_{j}^{(i_{t})}(k+t)\right)^{\mathrm{T}} Q^{(i_{t})}(k) B_{j}^{(i_{t})}(k+t) + R(k+t-1), \\ H_{ts}(k) &= \sum_{i_{t}=1}^{V} \left(B_{0}^{(i_{t})}(k+t)\right)^{\mathrm{T}} \sum_{i_{t+1}=1}^{V} \dots \sum_{i_{s-1}=1}^{V} \left(A_{0}^{(i_{t+1})}(k+t+1)\right)^{\mathrm{T}} \dots \left(A_{0}^{(i_{s-1})}(k+s-1)\right)^{\mathrm{T}} \times \\ &\times \sum_{j=0}^{n} \sum_{i_{s}=1}^{V} \left(A_{j}^{(i_{s})}(k+s)\right)^{\mathrm{T}} Q^{(i_{t},\dots,i_{s})}(k) B_{j}^{(i_{s})}(k+s), s > t, \\ H_{ts}(k) &= \left(H_{st}(k)\right)^{\mathrm{T}}, s < t, \\ G_{t}(k) &= \sum_{i_{1}=1}^{V} \dots \sum_{i_{t-1}=1}^{V} \left(A_{0}^{(i_{1})}(k+1)\right)^{\mathrm{T}} \dots \left(A_{0}^{(i_{t-1})}(k+t-1)\right)^{\mathrm{T}} \sum_{j=0}^{n} \sum_{i_{t}=1}^{V} \left(A_{j}^{(i_{t})}(k+t)\right)^{\mathrm{T}} Q^{(i_{1},i_{2},\dots,i_{t})}(k) B_{j}^{(i_{t})}(k+t), \\ F_{t}(k) &= \sum_{i_{t}=1}^{V} Q_{2}^{(i_{t})}(k) B_{0}^{(i_{t})}(k+t). \end{split}$$

Последовательности матриц $Q^{(i_t,...,i_s)}(k), Q_2^{(i_t,...,i_s)}(k), (s,t=\overline{1,m})$ определяется уравнениями:

$$Q^{(i_{t},\dots,i_{s})}(k) = \Theta^{(i_{t},\dots,i_{s})}(k)R_{1}(k+s) + \sum_{j=0}^{n} \sum_{i_{s+1}=1}^{\nu} \left(A_{j}^{(i_{s+1})}(k+s+1)\right)^{T} Q^{(i_{t},\dots,i_{s+1})}(k)A_{j}^{(i_{s+1})}(k+s+1), t = \overline{1,m-2}, t < s < m,$$

$$Q^{(i_{t})}(k) = e_{i_{t}} P^{t} \Theta(k)R_{1}(k+t) + \sum_{j=0}^{n} \sum_{i_{t+1}=1}^{\nu} \left(A_{j}^{(i_{t+1})}(k+t+1)\right)^{T} Q^{(i_{t},i_{t+1})}(k)A_{j}^{(i_{t+1})}(k+t+1), t = \overline{1,m-1},$$

$$Q^{(i_{t},\dots,i_{s})}_{2}(k) = R_{2}(k+s)\Theta^{(i_{t},\dots,i_{s})}(k) + \sum_{i_{s+1}=1}^{\nu} Q_{2}^{(i_{t},\dots,i_{s+1})}(k)A_{0}^{(i_{s+1})}(k+s+1), t = \overline{1,m-2}, t < s < m,$$

$$Q^{(i_{t})}_{2}(k) = R_{2}(k+t)e_{i_{t}} P^{t} \Theta(k) + \sum_{i_{s+1}=1}^{\nu} Q_{2}^{(i_{t},i_{t+1})}(k)A_{0}^{(i_{t+1})}(k+t+1), t = \overline{1,m-1},$$

$$Q^{(i_{t},\dots,i_{s})}_{2}(k) = R_{2}(k+t)e_{i_{t}} P^{t} \Theta(k) + \sum_{i_{s+1}=1}^{\nu} Q_{2}^{(i_{t},i_{t+1})}(k)A_{0}^{(i_{t+1})}(k+t+1), t = \overline{1,m-1},$$

с начальными условиями:

$$Q^{(i_m)}(k) = e_i \ P^m \Theta(k) R_1(k+m), Q^{(i_1, \dots, i_m)}(k) = \Theta^{(i_1, \dots, i_m)}(k) R_1(k+m), t = \overline{1, m-1}, \tag{17}$$

$$Q_2^{(i_m)}(k) = e_{i_m} P^m \Theta(k) R_2(k+m), Q_2^{(i_1, \dots, i_m)}(k) = \Theta^{(i_1, \dots, i_m)}(k) R_2(k+m), t = \overline{1, m-1},$$
(18)

$$\Theta^{(i_t, \dots, i_s)}(k) = P_{i_s, i_{s-1}} P_{i_{s-1}, i_{s+1}} \dots P_{i_{t+1}, i_t} \theta_{i_t}(k+t \mid k), t = \overline{1, m-1}, s > t,$$
(19)

где $\theta_{i_t}(k+t\,|\,k)$ – компонента вектора прогноза состояния цепи Маркова

$$\Theta(k+t \mid k) = E\{\Theta(k+t) \mid \Theta(k)\} = P^{t}\Theta(k), e_{i,} = [0,...,0,1,0,...,0]_{1 \times v}, i_{t} = \overline{1,v}, t = \overline{1,m}.$$

Методика доказательства теоремы основана на результатах, приведенных в работе [7].

2. Адаптивный алгоритм фильтрации марковской цепи

При определении оптимальной стратегии прогнозирующего управления предполагалось, что состояние марковской цепи $\alpha(k)$ в момент времени k доступно наблюдению. Однако на практике при управлении реальным ИП состояние цепи Маркова не доступно прямому наблюдению.

Для оценки параметров модели со скрытыми марковскими переключениями будем использовать адаптивный ЕМ-алгоритм, предложенный в работе [10]. В дальнейшем будем предполагать, что вектор доходностей рисковых активов подчиняется условному многомерному нормальному распределению с параметрами, зависящими от состояния:

$$\eta(k) | \alpha(k) \sim N(\mu[\alpha(k)], \sigma[\alpha(k)])$$
.

Это означает, что в динамике доходностей (9) величины $w_j(k)$ подчиняются стандартному нормальному распределению. Параметрами, подлежащими оценке, являются векторы ожидаемых доход-

ностей $\mu^{(1)},...,\mu^{(v)}$, матрицы волатильностей рисковых активов $\sigma^{(1)},...,\sigma^{(v)}$ в каждом состоянии цепи и матрица переходных вероятностей P, а также состояние цепи в момент времени k $\theta(k)$.

Обозначим l(k) — совместную вероятность появления последовательности $Y_k = \{\eta(1), ..., \eta(k)\}$ и нахождения цепи в состоянии i в момент времени k:

$$l^{(i)}(k) = f(\alpha(k) = i, Y_k).$$

Обозначим f_k – вектор плотностей распределения доходностей $\eta(k)$ в каждом состоянии цепи:

$$f_k = [f^{(1)}(\eta(k)), ..., f^{(v)}(\eta(k))],$$

$$f^{(i)}(\eta(k)) = f(\eta(k) \mid \alpha(k) = i) = \frac{1}{\left(2\pi\right)^{n/2} \left|\sigma^{(i)}\right|^{1/2}} \exp\left[-\frac{(\eta(k) - \mu^{(i)})^T \left(\sigma^{(i)}\right)^{-1} (\eta(k) - \mu^{(i)})}{2}\right], i = \overline{1, \nu}.$$

Оценки параметров модели скрытой цепи Маркова пересчитываются на каждом шаге $k=2,\,3,\,...,\,T$, с появлением нового наблюдения вектора доходностей рисковых активов. Пошагово алгоритм оценки имеет вид:

1. Задаются начальные значения вероятностей перехода $P_{ji}(1)$, начальное распределение $p_i(1)$ и значения параметров нормального распределения $\mu^{(i)}(1)$, $\sigma^{(i)}(1)$, $\sigma^{(i)}(1)$, ($i=\overline{1,v}$) (здесь индекс в скобках (1) определяет номер итерации алгоритма). Начальные значения величин $l^{(i)}(1)$ вычисляются по формуле:

$$l^{(i)}(1) = p_i(1) f^{(i)}(\eta(1)), i = \overline{1, \nu}.$$

2. На каждом шаге k=2,3,...,T величины $l^{(i)}(k)$ пересчитываются, суммируя вероятности всех возможных путей, которые ведут в новое состояние j, по формуле

$$l^{(j)}(k) = \sum_{i=1}^{\nu} l^{(i)}(k-1)\hat{P}_{ji}(k-1)f^{(j)}(\eta(k)), j = \overline{1,\nu}.$$

Вероятности фильтрации равны

$$\hat{\xi}_{i,k|k} = P\{\alpha(k) = i \mid Y_k\} = \frac{f(\alpha(k) = i, Y_k)}{f(Y_k)} = \frac{l^{(i)}(k)}{1_{v}l(k)}, 1_{v} = [1, ..., 1]_{v}, i = \overline{1, v}.$$

Оцененные совместные апостериорные вероятности равны

$$\zeta_{ij,k|k} = P\left\{\alpha(k-1) = i, \alpha(k) = j \mid Y_k\right\} = \frac{l^{(i)}(k-1)\hat{P}_{ji}(k-1)f^{(j)}(\eta(k))}{1...l(k)}, i, j = \overline{1, \nu}.$$

Оцененные апостериорные вероятности перехода равны

$$\hat{P}_{ji}(k) = \frac{\sum_{\tau=2}^{k-1} \hat{\xi}_{i,\tau|\tau}}{\sum_{\tau=2}^{k} \hat{\xi}_{i,\tau|\tau}} \hat{P}_{ji}(k-1) + \frac{\zeta_{ij,k|k}}{\sum_{\tau=2}^{k} \hat{\xi}_{i,\tau|\tau}}, i, j = \overline{1,\nu}.$$
(20)

Оцененные параметры рисковых финансовых активов равны

$$\mu^{(i)}(k) = \frac{\sum_{\tau=1}^{k-1} \hat{\xi}_{i,\tau|\tau}}{\sum_{\tau=1}^{k} \hat{\xi}_{i,\tau|\tau}} \mu^{(i)}(k-1) + \frac{\hat{\xi}_{i,k|k} \eta(k)}{\sum_{\tau=1}^{k} \hat{\xi}_{i,\tau|\tau}},$$
(21)

$$\sigma^{(i)}(k) = \frac{\sum_{\tau=1}^{k-1} \hat{\xi}_{i,\tau|\tau}}{\sum_{\tau=1}^{k} \hat{\xi}_{i,\tau|\tau}} \sigma^{(i)}(k-1) + \frac{\hat{\xi}_{i,k|k} \left(\eta(k) - \mu^{(i)}(k) \right) \left(\eta(k) - \mu^{(i)}(k) \right)^{T}}{\sum_{\tau=1}^{k} \hat{\xi}_{i,\tau|\tau}}, i = \overline{1, \nu}.$$
(22)

В выражениях (16)–(19) в качестве оценки состояния цепи Маркова в момент времени k будем использовать сглаживающие вероятности $\theta(k) = \hat{\xi}_{k|k}, \hat{\xi}_{k|k} = [\hat{\xi}_{1,k|k},...,\hat{\xi}_{v,k|k}]^T$.

3. Численное моделирование

В данном разделе приводятся результаты численного моделирования с использованием реальных данных российского фондового рынка. Для моделирования использовались цены закрытия наиболее ликвидных акций, торгующихся на Московской бирже, а именно: ПАО «Сбербанк России» (SBER), ПАО «Газпром» (GAZP), ПАО «Газпром нефть» (SIBN), ПАО «ГМК "Норильский никель"» (GMKN), ПАО «Банк ВТБ» (VTBR), ПАО «ЛУКОЙЛ» (LKOH), ПАО «НК "Роснефть"» (ROSN). Данные взяты с www.finam.ru. Рассматривались все возможные комбинации портфелей, состоящих из n=5 рисковых активов и одного безрискового актива с доходностью $r_1=0,00001$ (0,001% в день). Допускалось привлечение заемных средств по ставке $r_2=0,0001$ (0,01% в день). Предполагалось, что в начальный момент времени весь капитал инвестирован в безрисковый актив, следовательно

$$x_i(0) = x_{n+2}(0) = 0, (i = \overline{1, n}), x_{n+1}(0) = V(0) = V^0(0) = 1.$$

Управление портфелем осуществлялось в каждый торговый день. Весовая матрица полагалась равной $R(k+i)={\rm diag}\{10^{-3},\ldots,10^{-3}\}$ для всех k,i. Объем заемных средств ограничивался величиной $d_0(k)=3V(k)$. Операции «продажи без покрытия» запрещены, т.е. $d_i(k)=0$, ($i=\overline{1,n}$). Транзакционные издержки составляли $\lambda^+=\lambda^-=0.0006$.

Предполагалось, что финансовый рынок может находиться в двух состояниях (v=2). Оценка параметров скрытой цепи Маркова производилась согласно алгоритму, описанному в разделе 2. При этом для оценки состояния скрытой цепи и матрицы переходных вероятностей использовалась подгруппа из двух акций, входящих в портфель. Векторы средних значений и матрицы волатильностей для разных состояний цепи оценивались для портфеля в целом. Данный подход обусловлен тем, что применение полного портфеля для оценки состояния цепи и матрицы переходных вероятностей приводит к частым переключениям состояний цепи и, как следствие, к низкому качеству управления [9]. Полученные на каждом шаге k параметры модели скрытой цепи Маркова далее использовались для определения оптимальной стратегии прогнозирующего управления. Горизонт прогноза m=5.

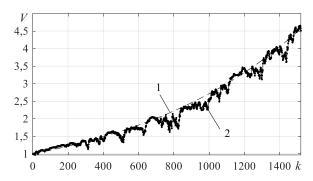


Рис 1. Динамика капиталов эталонного ИП (линия I) и управляемого ИП (линия 2) Fig. 1. Control portfolio value (line I) and reference portfolio value (line 2)

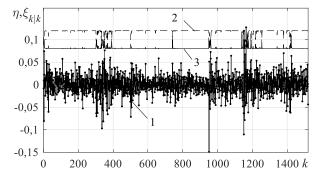


Рис. 2. Динамика доходности акции SBER (линия I (η)) и сглаживающие вероятности (линия 2 ($\xi_{1,k|k}$) — состояние 1, линия 3 ($\xi_{2,k|k}$) — состояние 2)

Fig. 2. Daily return of SBER (line 1 (η)) and smoothed probabilities (line 2 ($\xi_{1,kk}$) – state 1, line 3 ($\xi_{2,kk}$) – state 2)

Численно была реализована стратегия слежения за эталонным инвестиционным портфелем с доходностью $\mu_0 = 0{,}001$ (0,1% в день), $\rho_1(k+i) = 1$, $\rho_2(k+i) = 0{,}02$. Капитал реального управляемого ИП вычислялся по формулам (1)–(3) и (8), где использовались реальные наблюдаемые значения доходностей в момент времени k+1.

Далее приводятся типичные результаты моделирования. Портфель составлен из рисковых активов LKOH, GAZP, SBER, ROSN, GMKN. Период инвестирования: 27.05.2010-17.06.2016 (1 520 торговых дней). Для оценки состояния рыночного режима и матрицы переходных вероятностей использовались акции SBER и GMKN. На рис. 1 показана динамика капиталов эталонного портфеля $V^0(k)$ и управляемых портфелей V(k). Рисунок 2 иллюстрирует динамику доходности акции SBER и сглаживающие вероятности состояний цепи Маркова, приведенные к значениям доходностей для наглядного отображения на графике.

Рис. 1 показывает, что капитал реального портфеля следует капиталу эталонного портфеля.

Заключение

В данной работе предложен метод управления ИП с прогнозирующей моделью на финансовом рынке переключением режимов в соответствии со скрытой цепью Маркова с учетом явных ограничений на объемы вложений и займов и транзакционных издержек. Для оценки параметров используется адаптивный ЕМ-алгоритм. Результаты численного моделирования с использованием реальных данных демонстрируют эффективность предложенной стратегии управления.

ЛИТЕРАТУРА

- 1. Costa O.L.V., Araujo M.V. A generalized multi-period portfolio optimization with Markov switching parameters // Automatica. 2008. V. 44, No. 10. P. 2487–2497.
- Гальперин В.А., Домбровский В.В., Федосов Е.Н. Динамическое управление инвестиционным портфелем на диффузионно-скачкообразном финансовом рынке с переключающимися режимами // Автоматика и телемеханика. 2005. № 5. С. 175–189.
- 3. Bäuerle N., Rieder U. Portfolio optimization with Markov-modulated stock prices and interest rates // IEEE Transactions on Automatic Control. 2004. V. 49, No. 3. P. 442–447.
- 4. Sotomayor L.R., Cadenillas A. Explicit Solutions of Consumption-investment Problems in Financial Markets with Regime-switching // Mathematical Finance. 2009. V. 19, No. 2. P. 251–279.
- 5. Wu H. Mean-variance portfolio selection with a stochastic cash flow in a Markov-switching Jump-Diffusion Market // J. Optim. Theory Appl. 2013. V. 158. P. 918–934.
- 6. Levy M., Kaplanski G. Portfolio selection in two-regime world // European J. of Operational Research. 2015. V. 241. P. 514–524.
- 7. Dombrovskii V., Pashinskaya T. Design of model predictive control for constrained Markov jump linear systems with multiplicative noises and online portfolio selection // Int. J. Robust Nonlinear Control. 2020. V. 30, No. 3. P. 1050–1070.
- 8. Ishijima H., Uchida M. Log Mean-Variance Portfolio Selection Under Regime Switching // Asia-Pacific Financial Markets. 2011. V. 18, No. 2. P. 213–229.
- 9. Nystrup P., Boyd S., Lindström E., Madsen H. Multi-period portfolio selection with drawdown control // Ann. Oper. Res. 2018. V. 282. P. 245–271.
- Stenger B., Ramesh V., Paragios N., Coetzee F., Buhmann J.M. Topology Free Hidden Markov Models: Application to Background Modeling // Proc. of the 8th IEEE Int. Conf. on Computer Vision. 2001. V. 1. P. 294–301.
- 11. Dombrovskii V.V., Dombrovskii D.V., Lyashenko E.A. Investment portfolio optimization with transaction costs and constraints using model predictive control // Proc. of the 8th Russian-Korean Int. Symposium on Science and Technology. KORUS. Tomsk, Russia: IEEE, 2004. P. 202–205.

Поступила в редакцию 23 мая 2019 г.

Pashinskaya T.Y., Dombrovskii V.V. (2020) PREDICTIVE CONTROL STRATEGIES FOR INVESTMENT PORTFOLIO IN THE FINANCIAL MARKET WITH HIDDEN REGIME SWITCHING. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 4–13

DOI: 10.17223/19988605/50/1

Consider an investment portfolio consisting of n risky assets and one risk-free asset (e.g., a bank account). Let $x_i(k)$ $(i = \overline{1,n})$ denote the amount of the wealth invested in the ith risky asset; $x_{n+1}(k) \ge 0$ is the amount invested in a risk-free asset; $u_i^+(k) \ge 0$ is the amount of money by which an investor buys the ith risky asset; $u_i^-(k) \ge 0$ is the amount of money by which an investor sells the ith risky asset. The volume of borrowing of a risk-free asset is equal to $x_{n+2}(k) \ge 0$; v(k) is the amount of borrowing that transferred from borrowing account to bank account; $r_1(k+1)$ is the riskless lending rate over time period (k, k+1], $r_2(k+1)$ is the riskless borrowing rate.

The *i*th stock holding $x_i(k)$ $(i = \overline{1,n})$ satisfies the following stochastic difference equation:

$$x_i(k+1) = [1 + \eta_i(k+1)][x_i(k) + u_i^+(k) - u_i^-(k)],$$

where $\eta_i(k+1)$ is the return of the *i*th risky asset. The dynamics of the bank account is given by:

$$x_{n+1}(k+1) = [1 + r_1(k+1)][x_{n+1}(k) + v(k) - (1+\lambda^+) \sum_{i=1}^n u_i^+(k) + (1-\lambda^-) \sum_{i=1}^n u_i^-(k)],$$

where λ^+ , λ^- are fractions of the amount transacted on purchase $u_i^+(k)$ and sell $u_i^-(k)$ of the *i*th stock, respectively. The evolution of the borrowing account is the following: $x_{n+2}(k+1) = [1+r_2(k+1)][x_{n+2}(k)+v(k)]$. The wealth process satisfies V(k) = cx(k), $c = [1, ..., 1, -1]_{n+2}, x(k) = [x_1(k), ..., x_{n+2}(k)]^T$.

The following constraints are taken into account:

$$x_{n+1}(k) + v(k) - (1 + \lambda^{+}) \sum_{i=1}^{n} u_{i}^{+}(k) + (1 - \lambda^{-}) \sum_{i=1}^{n} u_{i}^{-}(k) \ge 0, \quad x_{n+2}(k) + v(k) \ge 0,$$
 (1)

$$x_i(k) + u_i^+(k) - u_i^-(k) \ge -d_i(k); \quad x_{n+2}(k) + v(k) \le d_0(k); u_i^+(k) \ge 0, u_i^-(k) \ge 0, (i = \overline{1, n}). \tag{2}$$

The evolution of the risky assets returns $\eta_i(k)$ is described by the equation:

$$\eta_i \left[\theta(k), k \right] = \mu_i \left[\theta(k), k \right] + \sum_{j=1}^n \sigma_{ij} \left[\theta(k), k \right] w_j(k),$$

where $\mu_i[\theta(k), k]$ is the expected return; $\sigma[\theta(k), k] = {\{\sigma_{ij}[\theta(k), k]\}_{i,j=1,...,n}}$ is the volatility matrix; $\{w_j(k); j=1, ..., n\}$ are independent noises with zero mean and unit variance; $\theta(k) = [\delta(\alpha(k), 1), ..., \delta(\alpha(k), \nu)]^T$, $\delta(\alpha(k), j)$ is a Kronecker function $(j = 1, 2, ..., \nu)$; $\alpha(k) \in \{1, 2, ..., \nu\}$ is a discrete-time Markov chain.

Our objective is to control the investment portfolio by tracking a deterministic portfolio with a desired return μ_0 , which evolution is described by the equation:

$$V^{0}(k+1) = [1 + \mu_{0}]V^{0}(k), V^{0}(0) = V(0).$$

In this paper, we design portfolio control strategies subject to constraints (1)–(2) under the quadratic performance criterion with receding horizon m:

$$\begin{split} J(k+m\,|\,k) &= \sum_{i=1}^{m} E\Big\{\rho_{1}(k+i)[V(k+i\,|\,k) - V^{0}(k+i)]^{2} - \\ &-\rho_{2}(k+i)[V(k+i\,|\,k) - V^{0}(k+i)] + u^{\mathrm{T}}(k+i-1\,|\,k)R(k+i-1)u(k+i-1\,|\,k)\Big|V(k), \theta(k)\Big\}, \end{split}$$

where $u(k+i|k) = [v(k+i|k), u_1^+(k+i|k), ..., u_n^+(k+i|k), u_1^-(k+i|k), ..., u_n^-(k+i|k)]^T$ is the predictive control vector; $\rho_1(k+i) \ge 0$, $\rho_2(k+i) \ge 0$ are the weight coefficients (scalar values); R(k+i) > 0 is a symmetric weight matrix of dimension $(2n+1) \times (2n+1)$.

We assume that the state of the Markov chain $\theta(k)$ is not observed. To estimate the parameters of the hidden Markov model, the on-line adaptive EM-algorithm is applied. We present the numerical modelling results based on the real data from the Russian stock exchange.

Keywords: investment portfolio; hidden Markov chain; model predictive control; constraints.

PASHINSKAYA Tatiana Yurievna (Candidate of Physics and Mathematics, Associate Professor, National Research Tomsk State University, Tomsk, Russian Federation).

E-mail: tani4kin@mail.ru

DOMBROVSKII Vladimir Valentinovich (Doctor of Technical Sciences, Professor, National Research Tomsk State University, Tomsk, Russian Federation).

E-mail: dombrovs@ef.tsu.ru

REFERENCES

- 1. Costa, O.L.V. & Araujo, M.V. (2008) A generalized multi-period portfolio optimization with Markov switching parameters. *Automatica*. 44(10). pp. 2487–2497. DOI: 10.1016/j.automatica.2008.02.014
- 2. Galperin, V.A., Dombrovsky, V.V. & Fedosov, E.N. (2005) Dynamic control of the investment portfolio in the jump-diffusion financial market with regime-switching. *Automation and Remote Control*. 66(5). pp. 837–850. DOI: 10.1007/s10513-005-0127-9

- 3. Bäuerle, N. & Rieder, U. (2004) Portfolio optimization with Markov-modulated stock prices and interest rates. *IEEE Transactions on Automatic Control*. 49(3). pp. 442–447. DOI: 10.1109/TAC.2004.824471
- 4. Sotomayor, L.R. & Cadenillas, A. (2009) Explicit Solutions of Consumption-investment Problems in Financial Markets with Regime-switching. *Mathematical Finance*. 19(2). pp. 251–279. DOI: 10.1111/j.1467-9965.2009.00366.x
- 5. Wu, H. (2013) Mean-variance portfolio selection with a stochastic cash flow in a Markov-switching Jump-Diffusion Market. *Journal of Optimization Theory and Application*. 158. pp. 918–934. DOI: 10.1007/s10957-013-0292-x
- 6. Levy, M. & Kaplanski, G. (2015) Portfolio selection in two–regime world. *European Journal of Operational Research*. 241. pp. 514–524. DOI: 10.1016/j.ejor.2014.10.012
- Dombrovskii, V. & Pashinskaya, T. (2020) Design of model predictive control for constrained Markov jump linear systems with multiplicative noises and online portfolio selection. *International Journal of Robust Nonlinear Control*. 30(3). pp. 1050–1070. DOI: 10.1002/rnc.4807
- 8. Ishijima, H. & Uchida, M. (2011) Log Mean-Variance Portfolio Selection Under Regime Switching. *Asia-Pacific Financial Markets*. 18(2). pp. 213–229. DOI: 10.1007/s10690-010-9132-2
- Nystrup, P., Boyd, S., Lindström, E. & Madsen, H. (2018) Multi-period portfolio selection with drawdown control. Annals of Operations Research. 282(2). pp. 1–27. DOI: 10.1007/s10479-018-2947-3
- Stenger, B., Ramesh, V., Paragios, N., Coetzee, F. & Buhmann, J.M. (2001) Topology Free Hidden Markov Models: Application to Background Modeling. *Proc. to the 8th IEEE International Conference on Computer Vision*. 1. pp. 294–301. DOI: 10.1109/ICCV.2001.937532
- 11. Dombrovskii, V.V., Dombrovskii, D.V. & Lyashenko, E.A. (2004) Investment portfolio optimization with transaction costs and constraints using model predictive control. *IEEE: Proc. of the 8th Russian-Korean Int. Symposium on Science and Technology. KORUS.* Tomsk, Russia. pp. 202–205.

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020

Управление, вычислительная техника и информатика

№ 50

ОБРАБОТКА ИНФОРМАЦИИ

УДК 004.855

DOI: 10.17223/19988605/50/2

И.А. Батраева, А.Д. Нарцев, А.С. Лезгян

ИСПОЛЬЗОВАНИЕ АНАЛИЗА СЕМАНТИЧЕСКОЙ БЛИЗОСТИ СЛОВ ПРИ РЕШЕНИИ ЗАДАЧИ ОПРЕДЕЛЕНИЯ ЖАНРОВОЙ ПРИНАДЛЕЖНОСТИ ТЕКСТОВ МЕТОДАМИ ГЛУБОКОГО ОБУЧЕНИЯ

Рассматриваются вопросы применения сверточных нейронных сетей для анализа текстов с точки зрения определения их жанровой принадлежности. Описана разработанная архитектура сверточной нейронной сети с использованием векторного представления слов на основе модели word2vec, приведены результаты экспериментов по обучению сети.

Ключевые слова: машинное обучение; сверточные нейронные сети; модель word2vec; интеллектуальный анализ текстов.

Автоматизация извлечения различной информации из текстов стала одной из основных проблем, связанных с информационным поиском. Так как тексты чаще всего либо являются слабо структурированными, либо вообще не обладают структурой с точки зрения решаемой задачи, то особо важным стало направление интеллектуального анализа текстов, включающее в себя методы классификации и анализа текстов на основе алгоритмов машинного обучения.

В частности, одной из задач анализа текстов является тематическая классификация, которая позволяет определить принадлежность текста к определенной группе тем. Особенно актуальна такая классификация для решения задач корпусной лингвистики, так как в большинстве существующих на сегодняшний день корпусов деление по темам и жанрам выполняется вручную или исходя из тематики источников текста [1]. Особенностью классификации текстов языковых корпусов является то, что для них важна, скорее, литературная классификация по темам (война, история, фантастика, сказки и т.д.) и жанрам (песни, стихи, повествование и т.п.). В данной работе рассматривается решение задачи классификации текстов для использования в языковых корпусах.

Более формально задача жанровой классификации может быть сформулирована так: даны текст на естественном языке и множество возможных жанров, к которым он может принадлежать. Требуется определить жанр текста. Если текст относится к нескольким жанрам одновременно, то определить основной жанр.

В последнее время наибольшую популярность для решения задач классификации приобрели глубокие нейронные сети, так как они позволяют достичь наивысшей точности среди всех известных моделей машинного обучения. В частности, сверточные нейронные сети совершили прорыв в классификации изображений. В настоящее время они успешно справляются и с некоторыми задачами автоматической обработки текстов. Более того, как утверждается в некоторых исследованиях [2–5], сверточные сети подходят для этого даже лучше рекуррентных нейронных сетей, которые чаще всего используются для анализа текстовых последовательностей [6]. С другой стороны, использование сверточных сетей для классификации текстов мало исследовано. Поэтому исследование применения сверточных нейронных сетей для задачи классификации текстов в качестве альтернативы рекуррентным нейронным сетям представляет практический интерес.

Для решения поставленной задачи требуется получить способ представления данных в виде, пригодном для обработки сверточной нейронной сетью. Например, в виде матрицы вещественных чисел. Наиболее распространенным является способ отображения каждого слова в многомерное векторное пространство. В рамках данной работы векторные представления слов строились на основе модели word2vec [7].

Таким образом, поставленная задача решалась сверточной нейронной сетью, на вход которой подавались векторные представления слов, полученные при обучении модели word2vec.

1. Предварительная обработка данных

На вход модели должен подаваться заранее обработанный корпус текстов. Предварительная обработка состоит из следующих этапов:

- Удаление всех знаков препинания, чисел и слов «нецелевых» языков (не предназначенных для обработки моделью).
- Разбиение текста на предложения. Для этого был выбран пакет библиотек Natural Language Toolkit (NLTK). Данная библиотека применяет регулярные выражения, а также некоторые алгоритмы машинного обучения для обработки естественного языка. Базовая версия NLTK не поддерживает разбиение русскоязычных текстов на предложения, поэтому использовалась модификация, расширяющая функционал библиотеки [8, 9].
- Удаление «стоп-слов» слов, не несущих определенной смысловой нагрузки, но при этом затрудняющих обработку исходного текста. Обычно для каждой специфической задачи применяется свой словарь стоп-слов, однако для нашей задачи достаточно стандартного словаря, содержащего буквы, частицы, предлоги, союзы, местоимения, числительные. Установлено, что удаление стоп-слов из тренировочного набора значительно снижает вычислительную стоимость, а также повышает точность модели.
- К корпусу текстов применяется стемминг или лемматизация. Это позволяет сократить размер словаря и искать семантически близкие слова, а не разные формы одного слова. Стемминг это поиск основы слова, причем не обязательно совпадающей с корнем. Он имеет высокую скорость работы, но наиболее эффективен для английского языка, так как в нем для нахождения основы слова обычно достаточно удалить окончание. Для русского языка стемминг малоэффективен, поэтому применяется более ресурсоемкий алгоритм лемматизации. Лемматизация это процесс приведения слова к начальной форме. В данной работе лемматизация осуществлялась морфологическим анализатором MyStem [10, 11].
- Дополнение предложений до одинаковой длины с использованием нейтрального слова, так как сверточные нейронные сети способны обрабатывать только последовательности одинаковой длины.

2. Построение векторного представления слов (модель word2vec)

Как уже было сказано, на начальном этапе необходимо перевести слова естественного языка в форму, пригодную для анализа сверточной нейронной сетью. Для этого лучше всего подходит векторное представление слов. Кроме того, среди всех моделей выберем ту, которая наиболее точно отражает реальные взаимосвязи между словами, а именно семантическую близость. Отметим, что модель не должна быть слишком требовательной к вычислительным ресурсам, чтобы было возможно совершать обучение сети на достаточно больших объемах данных.

Для выявления семантических связей между словами воспользуемся предположением лингвистики – дистрибутивной гипотезой: лингвистические единицы, встречающиеся в схожих контекстах, имеют близкие значения.

Во многих моделях обработки текстов входные данные кодируются унарным кодом (one-hot encoding) – вектором, размерность которого равна мощности словаря. Элемент, соответствующий

номеру слова в словаре, равен единице, а остальные элементы равны нулю. Однако у этого метода есть ряд существенных недостатков:

- словари естественных языков могут быть достаточно объемными и исчисляться десятками и сотнями тысяч слов; следовательно, если каждое слово кодировать таким вектором, объем данных становится слишком большим;
- при таком способе кодирования теряется связь между словами: все слова считаются разными и никак не связанными между собой.

В силу вышесказанного, one-hot encoding не подходит для анализа семантической близости слов. Поэтому для данной задачи воспользуемся другим способом кодирования – распределенным представлением слов.

Распределенное (или векторное) представление слов – это способ представления слов в виде векторов евклидова пространства, размерность которого обычно равна нескольким сотням. Основная идея заключается в том, что геометрические отношения между точками евклидова пространства будут соответствовать семантическим отношениям между словами. Например, слова, представленные двумя близко расположенными точками векторного пространства, будут, скорее всего, синонимами или просто тесно связанными по смыслу словами. Семантическая близость слов вычисляется как расстояние между векторами, для чего используется так называемая косинусная мера [12].

В 2013 г. группой исследователей Google под руководством Томаша Миколова была разработана нейросетевая модель для анализа семантики естественных языков, названная word2vec. В ее основу легли идея распределенного представления слов и дистрибутивная гипотеза, позволяющая рассматривать тексты с точки зрения статистики.

Word2vec включает в себя две различные архитектуры – CBOW (Continuous Bag of Words – непрерывный мешок слов) и Skip-gram. CBOW пытается предсказать слово, исходя из текущего контекста, а Skip-gram, наоборот, пытается предсказать контекст по текущему слову. Для реализации модели была выбрана архитектура Skip-gram, которая, несмотря на меньшую скорость обучения, лучше работает с редкими словами.

Предварительно обработанный текст можно подавать на вход модели, после чего будут выполнены следующие действия:

- считывается корпус текстов и рассчитывается, сколько раз в нем встретилось каждое слово;
- из этих слов формируется словарь, который сортируется по частоте слов; также из словаря для сокращения его размера удаляются редкие слова;
- модель идет по субпредложению (обычно предложение исходного текста или абзац) окном определенного размера; под размером окна понимается максимальная длина между текущим словом и словом, которое предсказывается. Оптимальный размер окна – 10 слов;
- к данным, находящимся в текущем окне, применяется нейронная сеть прямого распространения с линейной функцией активации скрытого слоя и функцией активации softmax для выходного слоя.

Из всего вышесказанного ясно, что матрицы, задающие скрытый и выходной слои, получаются чрезвычайно большими. Это делает обучение сети долгим процессом. Поэтому используются различные оптимизации, которые позволяют существенно снизить временные и вычислительные затраты, незначительно потеряв в точности. Одной из таких модификаций является субсемплирование. Дело в том, что в больших корпусах некоторые слова могут встречаться сотни миллионов раз. Такие слова зачастую несут меньшую информационную ценность, чем редкие слова. Чтобы избежать дисбаланса между редкими и часто встречающимися словами, используется простой подход: каждое слово отбрасывается с вероятностью, зависящей от частоты вхождения этого слова в текст.

В качестве модели word2vec была выбрана реализация из библиотеки Gensim [13]. Гиперпараметры модели:

- размерность векторного пространства 300;
- размер сканирующего окна 10;

- константа в формуле субсемплирования 0,00001;
- количество эпох 5.

В качестве обучающих данных был выбрал корпус русскоязычных текстов Максима Мошкова [14]. Он содержит более 25 тыс. книг общим объемом примерно 450 млн слов. Сформированный словарь содержал около 1,3 млн слов.

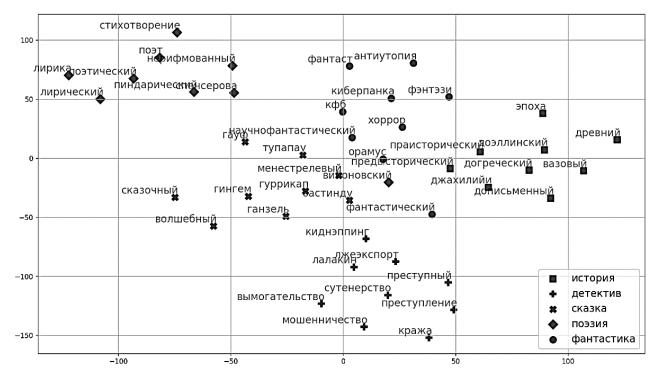


Рис. 1. Диаграмма кластеров слов различных жанров Fig. 1. Diagram of clusters of words of various genres

Результат обучения модели можно видеть на рис. 1, где видно четкое выделение кластеров, пересекающихся лишь по тем словам, которые действительно могут быть отнесены к нескольким группам одновременно.

3. Сверточная нейронная сеть

Сверточная нейронная сеть — это специальная архитектура нейронной сети, основной принцип которой заключается в том, что обработка некоторой области данных осуществляется независимо от расположения этой области. Применительно к задачам обработки естественного языка сверточные сети позволяют анализировать семантику слов в зависимости от их контекста, так как в большинстве случаев, достаточно рассмотреть сравнительно небольшой фрагмент текста. Рассмотренный подход может быть ошибочным для некоторых слов, лексическое значение которых правильно определяется на основе литературного произведения целиком или значительной его части, однако доля таких слов в тексте, как правило, незначительна.

Входные данные представляют собой матрицу, размерности которой равны количеству предложений в обучающей выборке и максимальной длине предложений (при этом каждое слово заменено своим векторным представлением). Заметим, что при таком представлении данных имеет смысл осуществлять свертку только по одному измерению – по ширине, поэтому сверточные фильтры будут одномерными.

За основу архитектуры сети была взята конфигурация, предложенная в работе [15]. На основе анализа экспериментальных результатов были выявлены некоторые недостатки, которые были устранены следующими модификациями:

- для борьбы с переобучением был добавлен дополнительный слой Dropout (на каждом этапе обучения некоторые нейроны исключаются из рассмотрения, что в некотором смысле приводит к рассмотрению новой конфигурации сети и препятствует чрезмерной адаптации нейронов друг к другу);
- проблема внутреннего сдвига переменных (возникает при использовании мини-батчей при обучении глубоких нейронных сетей) решалась применением нормализации;
- сформирован сверточный блок с применением фильтров разных размеров, что привело к увеличению точности классификации;
 - увеличено число полносвязных слоев;
- архитектура сети была доработана для осуществления классификации на произвольное количество классов: функция активации последнего слоя была заменена на softmax, что позволяет интерпретировать выход сети как вектор вероятностей принадлежности текста каждому классу.

На основе результатов многочисленных экспериментов были выбраны следующие гиперпараметры модели:

- слои Dropout: вероятности 0,5 (для входа блоки свертки) и 0,8 (для выхода блока свертки);
- слои свертки: размеры одномерных фильтров 3, 5, 8; количество фильтров 10; функция активации Relu;
 - слои субдискретизации: функция субдискретизации взятие максимума;
 - полносвязный слой: число нейронов 50, функция активации Relu;
- выходной (полносвязный) слой: число нейронов равно количеству классов (в нашем случае 5), функция активации Softmax;
 - размер одного мини-батча: 64.

Схематично разработанная архитектура представлена на рис. 2.

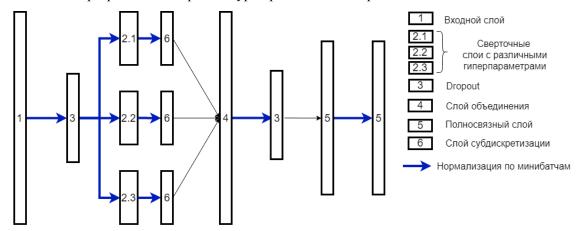


Рис. 2. Архитектура разработанной сети Fig. 2. Architecture of the developed network

В качестве алгоритма обучения был выбран адаптивный алгоритм градиентного спуска – Adam. Он основан на следующей идее: шаг изменения должен быть меньше у тех параметров, которые в большей степени варьируют в данных, и больше у тех, которые меньше изменяются на различных примерах. Как показывает практика, такой метод обучения работает эффективнее и сходится к правильным весам быстрее, чем стохастический градиентный спуск. Несмотря на свои преимущества, адаптивные варианты градиентного спуска не решают проблему переобучения. Поэтому необходимо следить за качеством обобщающей способности модели [16].

4. Результаты обучения

Для того чтобы экспериментально проверить эффективность работы построенной модели, было выбрано пять классов: история, детективы, детская литература, поэзия и песни, фантастика и фэнте-

зи. Среди данных классов наиболее специфичным для распознавания является класс «поэзия и песни». Дело в том, что из-за применения процесса лемматизации на этапе предварительной обработки данных, текст теряет рифму и стихотворный размер. Кроме того, стихотворные произведения обычно обладают сравнительно небольшой длиной. Это затрудняет распознавание текстов данного класса на основе семантики. Как уже было сказано, на вход сети подаются дополненные до одинаковой длины предложения. Поэтому сеть может использовать информацию о количестве добавленных нейтральных слов для классификации не только по семантике, но и по длине предложений.

Сеть обучалась 100 эпох: точность на тренировочной выборке составила 78,64%, точность на тестовой выборке – более 73,12%. График зависимости ошибки на тренировочных и тестовых данных от количества эпох приведен на рис. 3.

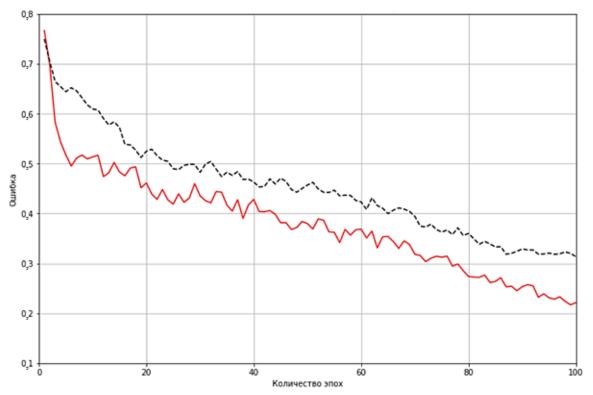


Рис. 3. Зависимость ошибки на тренировочных и тестовых данных от количества эпох (сплошная линия – ошибка на тренировочных данных, пунктирная – ошибка на тестовых данных) Fig. 3. Correlation of an error and the number of epochs for training and testing data (solid line – an error on training sample, dotted line – an error on testing sample)

Как видно из графика, в течение последних 15 эпох ошибка на валидационной выборке существенно не менялась, тогда как ошибка на тестовой выборке продолжала снижаться. Чтобы предотвратить переобучение, тренировка модели была остановлена. Каждые 5 эпох производилось сохранение весов, что позволило в качестве итоговой модели выбрать сеть с минимальной ошибкой на тестовой выборке.

Задача классификации текстов по темам или жанрам решалась во многих исследованиях. Наибольший интерес представляет работа [17], в которой классификация велась по темам. В работе использовались различные модели машинного обучения, в частности сверточные (полученная точность – 70,46%) и рекуррентные (точность – 72,12%) нейронные сети, метод опорных векторов (точность – 70,22%). Следует отметить, что нами была достигнута более высокая точность классификации по сравнению с аналогичными архитектурами сверточных нейронных сетей, а также рекуррентными сетями, которые зачастую показывают наивысшие результаты при анализе текстовых последовательностей.

Рассмотрим работу сети на некоторых примерах (таблица).

Примеры работы модели

Название жанра	Произведение	Вероятность принадлежности произведения каждому классу						
		История	Детективы	Детская	Поэзия	Фантастика и		
				литература	и песни	фэнтези		
История	«Петр I»	0,4229827	0,20211824	0,10318473	0,09724073	0,17456163		
	А.Н. Толстой							
Детективы	«Собака Баскервилей»	0,17708729	0,4044393	0,13971927	0,07812358	0,20063058		
	Артур Конан Дойл							
Детская	«Малыш и Карлсон»	0.11668574	0,13317753	0,38233585	0.10727942	0.26052145		
литература	Астрид Линдгрен	0,11006574	0,13317733	0,5025555	0,10727942	0,20032143		
Поэзия	«Руслан и Людмила»	0,16811042	0,12634562	0,16529457	0,37694713	0,16330227		
и песни	А.С. Пушкин				0,37034713	0,10330227		
Поэзия	«Привет, Андрей»	0.05886933	0.09415500	0.05023616	0,73858399	0,05815552		
и песни	И.Ю. Николаев	0,03880933	0,07413300	0,03023010	0,73030399	0,03013332		
Фантастика	«Гарри Поттер»	0,10806894	0,10925354	0,36183408	0,02905480	0,39178870		
и фэнтези	Дж.К. Роулинг							

Видно, что вероятность принадлежности произведения определенному жанру вполне коррелирует с литературным пониманием этого текста. Действительно, поэма «Руслан и Людмила», в первую очередь рассматривается как стихотворное произведение, однако содержит элементы исторического рассказа и фэнтези. «Собака Баскервилей», например, со значительной вероятностью относится к классу «Фантастика и фэнтези», что в некоторой степени объясняется высокой степенью мистицизма в данном произведении. Вероятность того, что «Гарри Поттер» относится к классу «Фантастика и фэнтези», близка к вероятности класса «Детская литература», что также соответствует нашим представлениям.

Заключение

Таким образом, для решения поставленной задачи была разработана архитектура сверточной нейронной сети, на вход которой подавались векторные представления слов, полученных на основе модели word2vec. Предложенная модель сверточной нейронной сети является корректной и достаточно точно отражает литературные представления о жанре текстов. Поэтому данная модель может быть применена для автоматизации обработки текстов в корпусной лингвистике.

ЛИТЕРАТУРА

- 1. Батраева И.А., Крючкова А.А. Разработка программного обеспечения диалектологических корпусов // Компьютерные науки и информационные технологии : материалы Междунар. науч. конф. Саратов : Наука, 2018. С. 45–49.
- Bai S., Kolter J.Z., Koltun V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. 2018. arXiv preprint arXiv: 1803.01271.
- 3. Conneau A., Schwenk H., Barrault L., LeCun Y. Very deep convolutional networks for text classification. 2017. arXiv preprint arXiv: 1606.01781.
- 4. Zhang X., Zhao J., LeCun Y. Character-level Convolutional Networks for Text Classification. 2016. arXiv preprint arXiv: 1509.01626.
- 5. Yin W., Kann K., Yu M., Schütze H. Comparative study of CNN and RNN for natural language processing. 2017. arXiv preprint arXiv: 1702.01923.
- 6. Yogatama D., Dyer Chr., Ling W., Blunsom Ph. Generative and discriminative text classification with recurrent neural networks. 2017. arXiv preprint arXiv: 1703.01898.
- 7. Rong Xin. Word2vec parameter learning explained. 2014. arXiv preprint arXiv: 1411.2738.
- 8. NLTK 3.4 documentation. URL: http://www.nltk.org/ (accessed: 08.02.2019).
- 9. Train NLTK punkt tokenizers. URL: https://github.com/mhq/train_punkt (accessed: 08.02.2019).
- 10. Яндекс технология MyStem. URL: https://tech.yandex.ru/mystem/ (дата обращения: 08.02.2019).
- 11. Segalovich I. A fast morphological algorithm with unknown word guessing induced by a dictionary for a web search engine // Proc. of the International Conf. on Machine Learning. Models, Technologies and Applications. MLMTA'03, June 23–26, 2003, Las Vegas, Nevada, USA. P. 1–8.
- 12. Mikolov T., Sutskever I., Chen K., Corrado G., Dean J. Distributed representations of words and phrases and their compositionality. 2013. arXiv preprint arXiv: 1310.4546

- 13. Gensim documentation. URL: https://radimrehurek.com/gensim/tutorial.html (accessed: 10.02.2019).
- 14. Библиотека Максима Мошкова. URL: http://lib.ru (дата обращения: 20.01.2019).
- 15. Kim Y. Convolutional neural networks for sentence classification // Proc. of the 2014 Conf. on Empirical Methods in Natural Language Processing (EMNLP 2014). 2014. P. 1746–1751.
- 16. Флах П. Машинное обучение. Наука и искусство построения алгоритмов, который извлекают знания из данных. М.: ДМК Пресс, 2015. 402 с.
- 17. Kamran K., Donald E., Mojtaba H., Kiana J.M., Matthew S., Laura E. HDLTex: Hierarchical Deep Learning for Text Classification. 2017. arXiv preprint arXiv:1709.08267.

Поступила в редакцию 2 июня 2019 г.

Batraeva I.A., Nartsev A.D., Lezgyan A.S. (2020) USING THE ANALYSIS OF SEMANTIC PROXIMITY OF WORDS IN SOLVING THE PROBLEM OF DETERMINING THE GENRE OF TEXTS WITHIN DEEP LEARNING. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 14–22

DOI: 10.17223/19988605/50/2

The relevant objective in the processing of text corpora is the classification of texts by topics and genres. Usually this work is done manually, so processing large text corpora is an extremely long process. Moreover, an unambiguous classification is not always possible: in most cases, the same text can be attributed to several topics and genres, with only one of them being the principal one. Therefore, the full automation of the classification process or limiting the choice of a researcher to the list of the most likely topics and genres is of practical interest.

To solve the problem, the authors propose to use convolutional neural networks, which, on the one hand, are efficient in classifications, and, on the other hand, are not used and studied properly for text recognition.

To present the data in a form suitable for processing by a convolutional neural network, the word2vec model was chosen. This model allows us to conduct vector representations of words that reflect their semantic proximity. To implement the word2vec model, the Skip-gram architecture was chosen, which, despite the slow learning rate, works well with rare words.

Based on the results of numerous experiments, the most optimal model hyperparameters were selected. The output of a trained model is the probability of attribution of a work to each class. Based on the analysis of the obtained results, we can conclude that the proposed model of the convolutional neural network is correct and fairly accurately reflects the literary perception of the genre.

Keywords: machine learning; convolutional neural networks; word2vec model; text natural language processing.

BATRAEVA Inna Aleksandrovna (Candidate of Physics and Mathematics, Head of the Department of Programming Technologies, Saratov State University, Saratov, Russian Federation).

E-mail: BatraevaIA@info.sgu.ru

NARTSEV Andrey Dmitrievich (Saratov State University, Saratov, Russian Federation).

E-mail: narcev.andrey@gmail.com

LEZGYAN Artem Sarkisovich (Saratov State University, Saratov, Russian Federation).

E-mail: lezgyan@yandex.ru

REFERENCES

- 1. Batraeva, I.A. & Kryuchkova, A.A. (2018) Razrabotka programmnogo obespecheniya dialektologicheskikh korpusov [Developing software for dialect corpora]. In: Tverdokhlebov, V. (ed.) *Komp'yuternye nauki i informatsionnye tekhnologii* [Computer Science and Information Technologies]. Saratov: Saratov State University. pp. 45–49.
- 2. Bai, S., Kolter, J.Z. & Koltun, V. (2018) An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. arXiv preprint arXiv: 1803.01271
- 3. Conneau, A., Schwenk, H., Barrault, L. & LeCun, Y. (2017) Very deep convolutional networks for text classification. arXiv preprint arXiv: 1606.01781
- 4. Zhang, X., Zhao, J. & LeCun, Y. (2016) Character-level Convolutional Networks for Text Classification. arXiv preprint arXiv: 1509.01626
- 5. Yin, W., Kann, K., Yu, M. & Schütze, H. (2017) Comparative study of CNN and RNN for natural language processing. arXiv preprint arXiv: 1702.01923
- 6. Yogatama, D., Dyer, Chr., Ling, W. & Blunsom, Ph. (2017) Generative and discriminative text classification with recurrent neural networks. arXiv preprint arXiv: 1703.01898
- 7. Rong, Xin. (2014) Word2vec parameter learning explained. arXiv preprint arXiv: 1411.2738
- 8. NLTK.org. (n.d.) NLTK 3.4 documentation. [Online] Available from: http://www.nltk.org/ (Accessed: 8th April 2019).

- 9. Github.com. (n.d.) *Train NLTK punkt tokenizers*. [Online] Available from: https://github.com/mhq/train_punkt (Accessed: 8th April 2019).
- 10. Yandex.ru. (n.d.) Yandeks tekhnologiya MyStem [MyStem Yandex Technology]. [Online] Available from: https://tech.yandex.ru/mystem/ (Accessed: 8th April 2019).
- 11. Segalovich, I. (2003) A fast morphological algorithm with unknown word guessing induced by a dictionary for a web search engine. *Models, Technologies and Applications. MLMTA'03*. Proc. of the International Conference on Machine Learning. June 23–26, 2003. Las Vegas, Nevada, USA. pp. 1–8.
- 12. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. (2013) Distributed representations of words and phrases and their compositionality. arXiv preprint arXiv: 1310.4546
- 13. Radimrehurek.com. (n.d.) *Gensim documentation*. [Online] Available from: https://radimrehurek.com/gensim/tutorial.html (Accessed: 8th April 2019).
- 14. Lib.ru. (n.d.) Biblioteka Maksima Moshkova [Maxim Moshkov's Library]. [Online] Available from: http://lib.ru (Accessed: 8th April 2017).
- 15. Kim, Y. (2014) Convolutional neural networks for sentence classification. *Proc. of the 2014 Conf. on Empirical Methods in Natural Language Processing (EMNLP 2014)*. pp. 1746–1751.
- 16. Flach, P. (2015) Mashinnoe obuchenie. Nauka i iskusstvo postroeniya algoritmov, kotoryy izvlekayut znaniya iz dannykh [Machine Learning: The Art and Science of Algorithms That Make Sense of Data]. Translated from English. Moscow: DMK Press.
- 17. Kamran, K., Donald, E., Mojtaba, H., Kiana, J., Matthew, S. & Laura, E. (2017) *HDLTex: Hierarchical Deep Learning for Text Classification*. arXiv preprint arXiv:1709.08267

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

УДК 004.6

DOI: 10.17223/19988605/50/3

А.В. Воробьев, Г.Р. Воробьева

ПОДХОД К ПОВЫШЕНИЮ ПРОИЗВОДИТЕЛЬНОСТИ ПРОГРАММНЫХ ПРОЦЕССОВ ОБРАБОТКИ И ХРАНЕНИЯ БОЛЬШИХ ОБЪЕМОВ ГЕОМАГНИТНЫХ ДАННЫХ

Работа выполнена при поддержке гранта РФФИ № 20-07-00011-а.

Обсуждаются вопросы повышения вычислительной скорости процессов аналитической обработки больших объемов геомагнитных данных, являющихся результатом непрерывного наблюдения за параметрами геомагнитного поля распределенными магнитными станциями и обсерваториями. Предложена гибридная архитектура, сочетающая особенности реляционной, иерархической и колоночной моделей данных, использующая правила ссылочной целостности и POSIX-структуру адресации компонентов. Проводится анализ эффективности предложенного подхода на основе оценки вычислительных затрат на хранение и обработку геомагнитных данных.

Ключевые слова: геомагнитные данные; реактивность программного обеспечения; аналитическая обработка; большие данные.

Одним из основных источников знаний о характере и закономерностях пространственновременного распределения параметров магнитного поля Земли и его вариаций являются геомагнитные данные, регистрируемые магнитными станциями и обсерваториями в режиме реального времени. При этом специализированное программное обеспечение для хранения и обработки геомагнитных данных к настоящему времени не разработано, а анализ данных выполняется отдельными исследователями посредством загрузки результатов наблюдений, хранящихся в репозиториях геомагнитных данных, техническое сопровождение которых осуществляется мировыми и региональными центрами геомагнитных данных [1. С. 390; 2. С. 2].

Общепринятым способом представления геомагнитных данных является формат IAGA2002, развиваемый Международной ассоциацией геомагнетизма и аэрономии [3. С. 5]. В структуре документа выделены: служебный заголовок, экспликация геомагнитных данных, значения параметров геомагнитного поля с соответствующими временными метками. Значения параметров и их временные метки заданы в ASCI-кодировке и разделены равным числом пробелов. Такое описание данных обеспечивает возможность использования формата для представления значений на длительном временном интервале — от нескольких секунд до многих месяцев.

Посуточное распределение результатов наблюдений параметров геомагнитного поля и его вариаций по отдельным файлам, низкоскоростные протоколы передачи данных, отсутствие вебсервисов и АРІ – далеко не полный перечень проблем, с которыми сталкивается разработчик программных средств для обработки геомагнитных данных формата IAGA2002. При этом наибольшую сложность с технической точки зрения представляет производительность программного продукта. Кроме того, локальное сохранение загруженных из репозиториев геомагнитных данных сопряжено с существенными затратами дискового пространства: например, годовой архив минутных значений результатов наблюдений параметров геомагнитного поля и его вариаций занимает в среднем объем в 40 МБ. На сегодняшний день в общей сложности доступны результаты более чем десятилетних наблюдений почти 300 магнитных станций и обсерваторий, что пропорционально увеличивает такие аппаратные затраты. Вместе с тем технические возможности научных организаций, занимающихся

исследованиями геомагнитного поля и его вариаций, зачастую ограничены, что не позволяет хранить подобные архивы наблюдений полностью и тем более выполнять их масштабную аналитическую обработку и визуализацию. Большие объемы геомагнитных данных и производительность программных средств их обработки напрямую связаны: к примеру, выполнение однопредикатного запроса к годовому архиву геомагнитных наблюдений одной магнитной обсерватории занимает в среднем 70 с при условии локального размещения обрабатываемых данных. Очевидно, что увеличение объемов обрабатываемых данных и сложности запросов к ним, а также использование, например, низкоскоростных протоколов для обращения к удаленным репозиториям в разы снизит производительность программного обеспечения.

Еще одна проблема связана с избыточностью формата IAGA2002. Обилие служебных символов, многократное повторение крайне редко изменяемых метаданных магнитных станций и обсерваторий в каждом суточном файле с результатами наблюдений приводит к тому, что объем полезной информации в IAGA2002-документе составляет менее 30% от его общего объема. При этом большинство разрабатываемых в научных организациях программных средств и систем зачастую ориентированы на использование устаревших технологий, не предназначенных для обработки данных такого большого объема.

Указанные проблемы приводят к необходимости совершенствования формата представления геомагнитных данных для обеспечения возможности создания высокопроизводительных программных средств их обработки и визуализации. Для решения поставленной задачи в настоящей работе предлагается новый гибридный формат долговременного хранения геомагнитных данных, представленный совокупностью трех взаимосвязанных компонент и отличающийся тем, что использует правила ссылочной целостности для объединения реляционной, иерархической и колончатой моделей данных, применяемых для описания метаданных и геомагнитных данных, а также реализует комбинацию текстового и бинарного форматов представления информации с целью повышения реактивности программных средств аналитической обработки геомагнитных данных, с одной стороны, и сокращения затрат требуемого объема физической памяти – с другой. Предлагаемый формат используется для представления данных в гибридном хранилище в составе предложенного авторами единого пространства геомагнитных данных [1. С. 395].

Результаты проведенных сравнительных экспериментов показали, что предложенный формат обеспечивает существенное повышение производительности вычислений, проводимых применительно к наборам разнородных геомагнитных данных, а также позволяет значительно сократить вычислительные затраты, связанные с их физическим хранением.

1. Структура описания метаданных

Служебный заголовок геомагнитных данных содержит признаковое описание магнитной обсерватории / станции, крайне редко изменяется и повторяется в каждом файле со значениями параметров геомагнитного поля, зарегистрированных обсерваторией / станцией. Очевидным шагом оптимизации формата представления геомагнитных данных является устранение избыточности служебного заголовка. Для этого предлагается отделить служебный заголовок и объединить метаданные всех магнитных станций и обсерваторий.

Метаданные магнитной станции / обсерватории, представленные множеством разноформатных объектов и их признаков, могут быть описаны посредством реляционной модели, заданной несколькими сущностями (рис. 1). Родительские сущности представляют собой обобщенные справочники параметров обсерватории, а каждый экземпляр дочерней описывает определенную станцию / обсерваторию посредством набора значений атрибутов. Сущности заданы в нормальной форме Бойса—Кодда и связаны друг с другом отношением типа «один-ко-многим».

Сущность «Observatory» предназначена для представления обобщенных данных о магнитной обсерватории / вариационной станции. Идентификатором каждого ее экземпляра выступает трех-

значный IAGA-код (поле «IAGAcode», текстовый формат, фиксированная размерность в 3 символа), который присваивается каждой станции / обсерватории, зарегистрированной в магнитной сети (независимо от ее принадлежности научной организации). Официальное название обсерватории, представленное в ее технической документации, задается в поле «Name» (текстовый формат, динамическая размерность). Для представления геодезических координат магнитной станции / обсерватории, таких как широта, долгота и высота над уровнем моря, использованы поля «Geodetic Longitude», «Geodetic Latitude», «Elevation» соответственно. Кроме того, в поле «Digital Sampling» задается значение скорости сбора данных с цифровых устройств или оцифровки аналогового сигнала в магнитной обсерватории (число одинарной точности). Также в поле «Data Interval Type» (текстовый формат, фиксированная размерность в 1 символ) предусмотрено хранение данных о временном интервале публикации геомагнитных данных (мгновенные регистрируемые значения или средние значения для интервалов от 1 с).

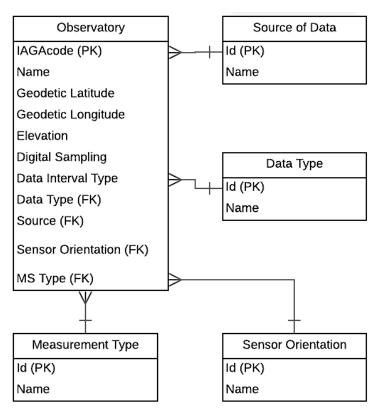


Рис. 1. Реляционная модель для описания метаданных Fig. 1. Relational Model for Metadata Description

Остальные сущности модели являются независимыми и содержат справочную информацию, используемую при описании магнитных станций / обсерваторий. Так, экземпляры сущности «Меаsurement Туре» (поле «Name», текстовый формат, фиксированная размерность в 4 символа) указывают на наименования регистрируемых станцией параметров геомагнитного поля (допустимые значения: DHIF, DHZF и XYZF). В сущности «Data Type» (поле «Name», текстовый формат, фиксированная размерность в 1 символ) указываются допустимые типы геомагнитных данных (временный (Р), окончательный (D), квази-окончательный (Q) или вариационный (V)). Физическая ориентация приборов наблюдения задается в сущности «Sensor Orientation», а курирующая станцию / обсерваторию научная организация — в сущности «Source of Data».

Единый доступ к данным об обсерваториях / станциях позволяет оперативно сформировать набор метаданных по соответствующему IAGA-коду, при этом отсутствует физическое дублирование хранимых данных, присутствующее в применяемом в настоящее время формате представления геомагнитных данных. При этом выделение метаданных магнитных станций позволяет на 80% сокра-

тить затраты памяти, требуемой для физического хранения геомагнитных данных, зарегистрированных обсерваторией за год.

2. Структура описания каталогов данных

Результаты геомагнитных наблюдений физически размещены в иерархической системе директорий, в большинстве решений доступной по протоколу FTP. Структура директорий такова, что корневым элементом является суррогатный каталог с именем, например, магнитной сети, далее он декомпозируется на директории, соответствующие календарным годам наблюдений, каждая из которых делится на каталоги для хранения результатов измерений по месяцам. Такая иерархическая архитектура базируется на принципах построения POSIX-систем с использованием соответствующей адресации.

Древовидная файловая структура может быть описана посредством иерархии элементов формата разметки XML (Extensible Markup Language), где корнем является суррогатный элемент с именем станции / обсерватории, а дочерними по отношению к нему — одноуровневые элементы, соответствующие календарным годам наблюдений. При этом все наблюдения должны быть агрегированы в директорию, где каждой станции / обсерватории соответствует XML-файл с геомагнитными данными. В результате входными параметрами для получения данных являются код магнитной станции / обсерватории и искомый год регистрации наблюдений за параметрами геомагнитного поля и его вариаций. На программном уровне формирование запроса выполняется последовательным применением операций работы с файлами и XPath-запроса непосредственно в теле XML-документа. Централизованное размещение всех геомагнитных данных одной станции / обсерватории позволит существенно повысить производительность программных запросов к ним, поскольку считывание файла и обращение к нему осуществляются единожды, а все последующие действия выполняются со сформированным на его основе виртуальным объектом.

3. Структура описания геомагнитных данных

Постоянно растущий объем геомагнитных данных снижает целесообразность применения текстового формата их хранения в плане как затрат физической памяти, так и производительности выполняемых при этом вычислений. Так, обработка однопредикатного запроса к годовым геомагнитным данным в условиях применения персонального компьютера со средней производительностью (процессор с частотой 1,6 ГГц, 2 ядра, оперативная память 4 Гб, скорость интернет-соединения 342,7 Мбит/с) занимает около 8 с, что существенно превышает общепринятое (с точки зрения эргономики программного обеспечения) время отклика, составляющее 3 с. Отметим, что параметры сетевого соединения здесь имеют принципиальное значение, поскольку в соответствии с концепцией единого пространства геомагнитных данных [1. С. 398] результаты геомагнитных измерений хранятся на сервере, обращение к которому осуществляется по протоколу HTTP(s).

Предварительно целесообразно отметить ряд параметров, которые представляются избыточными с точки зрения необходимости их физического хранения. Прежде всего к ним относится порядковый номер дня в году — параметр, который может быть оперативно вычислен с помощью библиотечных функций на основании календарной даты. Физическое хранение даты и времени регистрации параметра геомагнитного поля в каждой строке суточного файла наблюдений неэффективно, но эта проблема решается применением формата XML в описании геомагнитных данных обсерватории (поэтому в качестве временной метки выбран не порядковый номер дня в году, а дата, что обеспечивает уникальность элемента в составе описания магнитной станции). Остальные параметры, заданные в структуре геомагнитных данных, представляют собой непосредственно результаты измерений, заданные в формате разделенной пробелами строки.

Особенность аналитической обработки геомагнитных данных связана с тем, что наибольшая вычислительная нагрузка приходится на большие выборки записей, зачастую с группированием и агрегированием. При этом количество операций записи не так велико, а добавление новых записей

обычно осуществляется крупными блоками. Небольшое количество столбцов, громоздкие и частые операции выборок, редкие и крупные обновления данных – признаки, указывающие на целесообразность организации хранения геомагнитных данных с помощью колоночных СУБД (имеется в виду именно модель данных, поскольку на физическом уровне колоночное представление обычно используется в архитектуре хранилищ данных). Такие СУБД обеспечивают высокую скорость и гибкость выполнения сложных запросов при сохранении преимуществ использования структурированного языка SQL, а также соответствуют обязательным требованиям ACID.

Колоночная организация хранения геомагнитных данных позволит существенно повысить производительность операций их обработки. Это связано в первую очередь с тем, что при построчной записи чтение с диска происходит более линейно. Более предсказуемое чтение файла при построчной записи позволяет операционной системе эффективнее использовать дисковый кэш.

На сегодняшний день широкое распространение получил колоночно-ориентированный формат представления данных Арасhe Parquet, отличительной особенностью которого является возможность программного управления механизмом сжатия данных в столбцах. Кроме того, Parquet реализован с использованием алгоритма измельчения и сборки записей, вмещающих сложные структуры данных, которые также можно использовать для их хранения.

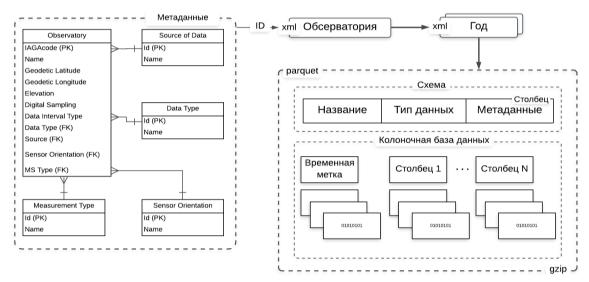
Еще одним важным преимуществом Parquet является его бинарный формат, обеспечивающий хранение данных в том виде, в котором они представляются компьютеру в процессе работы программы. Поэтому при чтении файла не выполняются дополнительные преобразования, что существенно повышает скорость работы с данными, что и требуется для повышения производительности программной обработки геомагнитных данных. Колоночный Parquet на программном уровне позволяет не считывать все данные при выполнении запросов, извлекая только значения определенных столбцов, что также повышает производительность обработки данных. Сжатие по столбцам позволяет существенно сэкономить место при физическом хранении геомагнитных данных.

Набор геомагнитных данных в формате Parquet представлен двумя разделами. Первый из них является схемой документа и содержит описание структурных и параметрических ограничений представления данных: определяются состав столбцов, их наименования и последовательность, алгоритм сжатия и пр. Второй компонент представляет собой геомагнитные данные — значения параметров геомагнитного поля и его вариаций. Для хранения данных выделено 5 столбцов (колонок): один под временную метку, а оставшиеся — под три компонента и полный вектор геомагнитного поля соответственно. Для упрощения структуры в документе выделена только одна страница, а все столбцы образуют одну группу. Поэтому ко всем составляющим Parquet-документа применен один и тот же алгоритм сжатия (в нашем случае — gzip).

4. Интеграция компонент гибридного формата хранения геомагнитных данных

В общем виде хранение геомагнитных данных подразумевает смешанную логическую и физическую интеграцию предложенных выше компонент (рис. 2). Образуется иерархия структур данных, корневым элементом которой выступает реляционная структура с метаданными магнитных станций и обсерваторий. Результаты геомагнитных наблюдений физически размещаются в едином каталоге, в котором каждой обсерватории выделен ХМL-документ с именем, содержащим IAGA-код. В составе ХМL-документа каждый соответствующий году наблюдений элемент содержит блок CDATA, в котором размещается набор геомагнитных данных в бинарном формате Parquet. При этом анализатор запросов, предусмотренный в архитектуре единого пространства геомагнитных данных [1. С. 398], обеспечивает проверку ссылочной целостности как по заданному IAGA-коду, так и по указанным временным меткам.

Взаимодействие с хранилищем данных осуществляется строго в соответствии с иерархической структурой. По IAGA-коду из реляционной структуры выгружаются метаданные.



Puc. 2. Гибридная архитектура представления геомагнитных данных Fig. 2. Hybrid architecture of geomagnetic data presentation

Далее тот же код используется для обращения к XML-файлу станции / обсерватории, а оттуда посредством XPath-запроса выбирается секция CDATA с искомыми геомагнитными данными. При необходимости выполняются фильтрация, группирование и агрегирование результатов наблюдений с использованием языка запросов SQL.

5. Экспериментальные исследования

Оценка эффективности предложенного гибридного формата хранения геомагнитных данных выполнена на основании сравнительного анализа распространенных форматов данных (рис. 3). По результатам исследования распространенных форматов и архитектур данных [4. С. 18] отобраны следующие: IAGA2002 (он же CSV) [3. С. 5]; реляционная база данных (RDB, relational database, на примере СУБД MS SQL Server 2017); XML [5. С. 1147]; JSON [6. С. 7]; AVRO [7. С. 267]; HDF5 [8; 9. С. 393]; neo4j [10. С. 232; 11. С. 11].

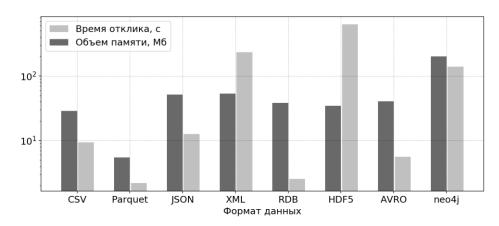


Рис. 3. Результаты сравнительного анализа форматов для хранения геомагнитных данных Fig. 3. Results of comparative analysis of geomagnetic data formats

Критериями оценки эффективности гибридного формата хранения геомагнитных данных определены реактивность программной обработки данных и объем требуемого для их размещения дискового пространства. Выбор первого из критериев связан с тем, что существующие технологии аналитической обработки геомагнитных данных недостаточно эффективны в плане затрат вычислительных ресурсов на выполнение операций, а также времени на сбор и интеграцию данных на этапе их пред-

варительной обработки. Требуется оценить, насколько предлагаемый формат позволит повысить производительность выполнения операций обработки данных. Второй критерий оценки эффективности является в большей степени вспомогательным, поскольку в современных условиях развития технологий облачных хранилищ проблема занимаемого данными объема дискового пространства теряет свою остроту. Однако в большинстве случаев при анализе изменения значений параметров геомагнитного поля и его вариаций исследователи прибегают к аккумулированию всех необходимых данных на персональном компьютере, что требует больших объемов дискового пространства.

Исследование эффективности гибридного формата выполнено на примере получения выборки из годового архива минутных наблюдений станции с IAGA-кодом BOX за период 01–06.03.2018. Тем самым имеет место двухпредикатный запрос, выполнение которого предполагает обращение к суточному архиву геомагнитных данных по IAGA-коду станции (BOX), формирование набора данных, а выборку данных из соответствующих секций CDATA.

Экспериментальные исследования показали, что минимальное время отклика программного сценария обработки геомагнитных данных достигается при использовании для их хранения формата Parquet (2,2 c), что примерно в 4,3 раза меньше, чем для формата IAGA2002/CSV.

Согласно результатам исследований, применение предложенного формата для хранения геомагнитных данных позволяет минимизировать требования к объему дискового пространства. Так, по сравнению с форматом IAGA2002/CSV, для хранения годового архива геомагнитных наблюдений одной станции требуется примерно в 5,2 раза меньше объема дискового пространства.

Заключение

В результате проведенных исследований предложен гибридный формат хранения геомагнитных данных, который отличается тем, что использует правила ссылочной целостности для объединения реляционной, иерархической и колоночной моделей данных, применяемых для описания метаданных и геомагнитных данных, а также использует комбинацию текстового и бинарного форматов представления информации с целью повышения реактивности программных средств аналитической обработки геомагнитных данных, с одной стороны, и сокращения затрат требуемого объема физической памяти – с другой.

ЛИТЕРАТУРА

- 1. Воробьев А.В., Воробьева Г.Р., Юсупова Н.И. Концепция единого пространства геомагнитных данных // Тр. СПИИРАН. 2019. Т. 18, № 2. С. 390–415.
- 2. Geomagnetic Observations and Models / ed. by M. Mandea, M. Korte. Dordrecht: Springer, 2011. P. 149–181. (IAGA Special Sopron Book Series 5). https://link.springer.com/book/10.1007/978-90-481-9858-0 (accessed: 22.05.2019).
- 3. Intermagnet technical reference manual. Version 4.6 / ed. by Benoît St-Louis. Edinburgh, 2012. 92 p. https://www.intermagnet.org/publications/intermag_4-6.pdf (accessed: 22.05.2019).
- 4. Carrera D., Rosales J., Blanco G.A.T. Optimizing Binary Serialization with an Independent Data Definition Format // Int. J. of Computer Applications. 2018. V. 180, No. 28. P. 15–18.
- 5. Yahui Y. Impact data-exchange based on XML // Proc. 7th Int. Conf. Computer Science & Education (ICCSE). 2012. P. 1147–1149.
- 6. Peng D., Cao L., Xu W. Using JSON for Data bn exchanging in Web Service Applications // J. of Computational Information System. 2011. V. 7 (16). P. 5883–5890.
- 7. Plase D., Niedrite L., Taranovs R. Comparison of HDFS compact data formats: Avro Versus Parquet // Mokslas-Lietuvos ateitis. 2017. No. 9. P. 267–276.
- 8. HDF5. URL: https://www.hdfgroup.org/HDF5/ (accessed: 22.05.2019).
- 9. Emeakaroha V., Healy P. et al. Analysis of Data Interchange Formats for Interoperable and Efficient Data Communication in Clouds // Proc. of the 2013 IEEE/ACM 6th Int. Conf. on Utility and Cloud Computing. P. 393–398.
- 10. Femy P.F.M., Reshma K.R., Surekha S.M. Outcome analysis using Neo4j graph database // Int. J. on Cybernetics & Informatics (IJCI). 2016. V. 5, No. 2. P. 229–236.
- 11. Angles R., Gutierrez C. Survey of graph database models // ACM Computing Surveys. 2008. V. 40, No. 1. P. 1–39.

Поступила в редакцию 19 июля 2019 г.

Vorobev A.V., Vorobeva G.R. (2020) APPROACH TO IMPROVING THE PERFORMANCE OF SOFTWARE PROCESSES FOR PROCESSING AND STORING LARGE VOLUMES OF GEOMAGNETIC DATA. *Vestnik Tomskogo gosudarstvennogo universiteta*. *Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 23–30

DOI: 10.17223/19988605/50/3

The issues of increasing the computational speed of software processes for the analytical processing of large volumes of geomagnetic data, which are the result of continuous monitoring of the parameters of the geomagnetic field by a great number of distributed ground magnetic stations and observatories, are discussed. A comparative review of the existing geomagnetic data architecture (presented in the framework of the specified IAGA-2002 format provided by International Association of Geomagnetism and Aeronomy), as well as popular data formats is given, and arguments are presented in favor of the need to improve the approach to organizing the results of geomagnetic observations.

To solve this problem, a new hybrid format for long-term storage of geomagnetic data is presented, represented by a set of three interrelated components and characterized in that it uses the rules of referential integrity to combine relational, hierarchical and columnar data models used to describe metadata and geomagnetic data, and also sets POSIX-component addressing structure and implements a combination of textual and binary formats for presenting information. The main purpose of the proposed architecture is to increase the reactivity of software tools for analytic processing of geomagnetic data, on the one hand, and reducing the cost of the required amount of physical memory, on the other hand.

The results of the comparison of the proposed hybrid format for presenting geomagnetic data with the existing approach to describing geomagnetic observation data (IAGA-2002), as well as other common formats for presenting large volumes of structured and semi-structured data (XML, JSON, Avro, etc.) are presented. In this case, the criteria for evaluating the effectiveness of a hybrid format for storing geomagnetic data determined the reactivity of software data processing and the amount of required disk space for their placement. The results of the experiment showed that the proposed format provides a significant increase in computing performance (about 4 times), conducted in relation to sets of heterogeneous geomagnetic data, and also significantly reduces the computational costs associated with their physical storage (approximately 5 times).

Keywords: geomagnetic data; software reactivity; analytical processing; big data.

VOROBEV Andrei Vladimirovich (Candidate of Technical Sciences, Associate Professor, Ufa State Aviation Technical University, Ufa, Russian Federation).

E-mail: geomagnet@list.ru

VOROBEVA Gulnara Ravilevna (Candidate of Technical Sciences, Associate Professor, Ufa State Aviation Technical University, Ufa, Russian Federation).

E-mail: gulnara.vorobeva@gmail.com

REFERENCES

- 1. Vorobev, A.V., Vorobeva, G.R. & Yusupova, N.I. (2019). Conception of geomagnetic data integrated space. Tr. SPIIRAN *SPIIRAS Proceedings*. 18(2). pp. 390–415. DOI: 10.15622/sp.18.2.390-415
- 2. Mandea, M. & Korte, M. (eds) (2011) Geomagnetic Observations and Models. Dordrecht: Springer. pp. 149-181.
- 3. Trigg, D.F. & coles, R.L. (ed.) (2012) *Intermagnet Technical Reference Manual*. 4.6. Edinburgh: [s.n.]. [Online] Available from: https://www.intermagnet.org/publications/intermag_4-6.pdf (Accessed: 22nd May 2019).
- 4. Carrera, D., Rosales, J. & Blanco, G.A.T. (2018) Optimizing Binary Serialization with an Independent Data Definition Format. *International Journal of Computer Applications*. 180(28). pp. 15–18. DOI: 10.5120/ijca2018916670
- Yahui, Y. (2012) Impact data-exchange based on XML. Proc. 7 th Int. Conf. Computer Science & Education (ICCSE). pp. 1147–1149. DOI: 10.1109/ICCSE.2012.6295268
- Peng, D., Cao, L. & Xu, W. (2011) Using JSON for Data bn exchanging in Web Service Applications. *Journal of Computational Information System*. 7(16). pp. 5883–5890.
- 7. Plase, D., Niedrite, L. & Taranovs R. (2017) Comparison of HDFS compact data formats: Avro Versus Parquet. *Mokslas Lietuvos ateitis*. 9. pp. 267–276.
- 8. HDF5. (n.d.) [Online] Availablef rom: https://www.hdfgroup.org/HDF5/ (Accessed: 22nd May 2019).
- 9. Emeakaroha, V., Healy, P. et al. (2013) Analysis of Data Interchange Formats for Interoperable and Efficient Data Communication in Clouds. *Proc. of the 2013 IEEE/ACM 6th Int. Conf. on Utility and Cloud Computing*. pp. 393–398. DOI: 10.1109/UCC.2013.79
- 10. Femy, P.F.M., Reshma, K.R. & Surekha, S.M. (2016) Outcome analysis using Neo4j graph database. *International Journal on Cybernetics & Informatics (IJCI)*. 5(2). 2016. pp. 229–236. DOI: 10.5121/ijci.2016.5225. 229
- 11. Angles, R. & Gutierrez, C. (2008) Survey of graph database models. *ACM Computing Surveys*. 40(1). pp. 1–39. DOI: 10.1145/1322432.1322433

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

УДК 519.872

DOI: 10.17223/19988605/50/4

Д.Я. Копать, М.А. Маталыцкий

АНАЛИЗ ОЖИДАЕМЫХ ДОХОДОВ В ОТКРЫТЫХ МАРКОВСКИХ СЕТЯХ С РАЗЛИЧНЫМИ ОСОБЕННОСТЯМИ

Проведено исследование системы разностно-дифференциальных уравнений, которой удовлетворяют ожидаемые доходы открытых марковских сетей массового обслуживания с различными особенностями. Число состояний сети в этом случае и число уравнений в данной системе бесконечны. Потоки поступающих в сеть заявок являются простейшими и независимыми, времена обслуживания заявок распределены по экспоненциальным законам. Доходы от переходов между состояниями сети являются детерминированными функциями, зависящими от ее состояния и времени, а доходы систем в единицу времени, когда они не меняют своих состояний, зависят только от этих состояний. Для нахождения ожидаемых доходов систем сети предложен модифицированный метод последовательных приближений, совмещенный с методом рядов.

Ключевые слова: система разностно-дифференциальных уравнений; открытая сеть массового обслуживания; метод последовательных приближений.

Сети массового обслуживания (CeMO) с доходами в нестационарном режиме изучались в работе [1]. Заявка при переходе из одной СМО в другую приносит последней некоторый доход, а доход первой СМО уменьшается на эту величину. При этом доходы от переходов между состояниями сетей зависели от их состояний и времени или являлись случайными величинами (CB) с заданными моментами первого и второго порядков. В статьях [1–3] приведены результаты по анализу, оптимизации и выбору оптимальных стратегий управления в марковских сетях с доходами, описаны различные применения их в качестве стохастических моделей прогнозирования ожидаемых доходов в информационно-телекоммуникационных системах и сетях (ИТСС), в страховых компаниях, логистических транспортных системах, производственных системах и других объектах. Как известно, функционирование любой марковской СеМО можно описать при помощи цепей Маркова с непрерывным временем и, как правило, с большим или счетным числом состояний. В простейшем случае марковские цепи с небольшим числом состояний и доходами от переходов между состояниями, являющимися константами, были рассмотрены в монографии [4].

Следует отметить, что марковские сети с положительными и отрицательными заявками были исследованы Е. Gelenbe в статьях [5–9] как модели поведения компьютерных вирусов в ИТСС и называются ныне G-сетями. Нахождение нестационарных вероятностей состояний марковской G-сети с сигналами и групповым удалением заявок модифицированным методом последовательных приближений, совмещенным с методом рядов, изложено в [10].

В последние годы большое внимание было уделено исследованию марковских сетей с доходами и различными особенностями: с ограниченным временем ожидания заявок и ненадежными СМО [11]. В [12] рассматривалась марковская G-сеть с доходами в случае, когда доходы от переходов между состояниями могут зависеть от ее состояний и времени. Для сети, допускающей мультипликативное представление для совместного стационарного распределения вероятностей состояний, для ожидаемых доходов систем сети выведена система разностно-дифференциальных уравнений (РДУ), для решения которой также предложено использовать метод последовательных приближений, совмещенный с методом рядов. В данной статье эти результаты обобщены на случай других марковских сетей с заявками многих классов. Марковские СеМО являются математическими моделями различных

реальных объектов, которые обычно функционируют на каком-то конечном промежутке времени, например [0; T].

Требуется найти ожидаемые (средние) доходы систем сети, полученные модифицированным методом последовательных приближений, совмещенным с методом рядов, для марковских сетей с различными особенностями.

1. Система РДУ для ожидаемых доходов

На основании ранее полученных результатов [2, 10–12] было замечено, что в общем случае систему РДУ для ожидаемых доходов открытой марковской сети, в которой могут присутствовать положительные и отрицательные заявки и сигналы различных классов, системы обслуживания могут подвергаться поломкам, заявки могут быть «нетерпеливыми» и с иными различными особенностями, можно записать в общем случае в виде:

$$\frac{d\vec{V}\left(\vec{d},\vec{k},\vec{l},t\right)}{dt} = -\Lambda\left(\vec{d},\vec{k},\vec{l}\right)\vec{V}\left(\vec{d},\vec{k},\vec{l},t\right) + \sum_{i^*,j^*=0}^{n} \sum_{\alpha,\beta,\gamma,\theta,\eta=0}^{\Psi r} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^*j^*\alpha m\beta b\gamma\theta\eta}\left(\vec{d},\vec{k},\vec{l}\right) \times \\
\times \vec{V}\left(\vec{d}+I_{i^*}-I_{j^*},\vec{k}+\tilde{I}_{\alpha}+m\tilde{I}_{\beta}-b\tilde{I}_{\gamma},\vec{l}+\tilde{I}_{\theta}-\tilde{I}_{\eta},t\right) + \vec{E}\left(\vec{d},\vec{k},\vec{l}\right), \tag{1}$$

где \tilde{I}_{α} — вектор размерности Ψr , состоящий из нулей, за исключением компоненты с номером α , которая равна 1, Ψr — некоторое целое положительное число, r — число типов заявок, I_{α} — вектор размерности n, состоящий из нулей, за исключением компоненты с номером α , которая равна 1, I_{ic} — вектор размерности nr, состоящий из нулей, за исключением компоненты с номером (i-1)r+c, которая равна 1, \vec{d} — вектор размерности n с компонентами d_i , где d_i — количество исправных линий обслуживания в i-й СМО, \vec{k} — вектор размерности Ψr с компонентами k_{ic} , где k_{ic} — количество положительных заявок типа c в i-той СМО, \vec{l} — вектор размерности Ψr с компонентами l_{ic} , где l_{ic} — количество сигналов типа c в i-й СМО, i = \vec{l} , \vec{n} , c = \vec{l} , \vec{n} . Здесь $\vec{V}^T(\vec{d},\vec{k},\vec{l},t)$ = $=(v_1(\vec{d},\vec{k},\vec{l},t),v_2(\vec{d},\vec{k},\vec{l},t),...,v_n(\vec{d},\vec{k},\vec{l},t))$, где $v_i(\vec{d},\vec{k},\vec{l},t)$ — ожидаемый (средний) доход, который получает i-я СМО за время t, если в начальный момент времени сеть находится в состоянии $(\vec{d},\vec{k},\vec{l})$, $\Lambda(\vec{d},\vec{k},\vec{l})$, Θ_{i,j' алярьую $(\vec{d},\vec{k},\vec{l})$, $\vec{E}(\vec{d},\vec{k},\vec{l})$ — некоторые функции, различные для каждой сети обслуживания, $\vec{E}^T(\vec{d},\vec{k},\vec{l}) = (E_1(\vec{d},\vec{k},\vec{l}), E_2(\vec{d},\vec{k},\vec{l}), ..., E_n(\vec{d},\vec{k},\vec{l}))$.

Предположим, что ряд $\sum_{\vec{i},\vec{j}=1}^{n} \sum_{\alpha,\beta,\gamma,\theta,\eta=0}^{\Psi_r} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{\vec{i}^*j^*\alpha m\beta b\gamma\theta\eta} \left(\vec{d},\vec{k},\vec{l}\right)$ сходится. Ранее в работах [1, 10–12] это было доказано для конкретных марковских сетей.

Из системы (1) следует:

$$\vec{V}(\vec{d},\vec{k},\vec{l},t) = e^{-\Lambda(\vec{d},\vec{k},\vec{l},t)} \left(\vec{V}(\vec{d},\vec{k},\vec{l},0) + \int_{0}^{t} e^{\Lambda(\vec{d},\vec{k},\vec{l},t)x} \times \left(\sum_{i^*,j^*=1}^{n} \sum_{\alpha,\beta,\gamma,\theta,\eta=0}^{nr} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^*j^*am\beta b\gamma\theta\eta} (\vec{d},\vec{k},\vec{l}) \right) \left(\sum_{i^*,j^*=1}^{n} \sum_{\alpha,\beta,\gamma,\theta,\eta=0}^{nr} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^*j^*am\beta b\gamma\theta\eta} (\vec{d},\vec{k},\vec{l}) \times \vec{V}(\vec{d}+I_{i^*}-I_{j^*},\vec{k}+I_{\alpha}+mI_{\beta}-bI_{\gamma},\vec{l}+I_{\theta}-I_{\eta},x) \right) dx \right) + \\
+ \frac{\vec{E}(\vec{d},k,\vec{l})}{\Lambda(\vec{d},\vec{k},\vec{l})} \left[1 - e^{-\Lambda(\vec{d},\vec{k},\vec{l})t} \right]. \tag{2}$$

Обозначим через $\vec{V}_q(\vec{d},\vec{k},\vec{l},t)$ приближение $\vec{V}(\vec{d},\vec{k},\vec{l},t)$ на q-й итерации, а $\vec{V}_{q+1}(\vec{d},\vec{k},\vec{l},t)$ – решение системы (1), полученное методом последовательных приближений. Тогда из (2) вытекает, что

$$\vec{V}_{q+1}(\vec{d}, \vec{k}, \vec{l}, t) = e^{-\Lambda(\vec{d}, \vec{k}, \vec{l})t} \left(\vec{V}(\vec{d}, \vec{k}, \vec{l}, 0) + \int_{0}^{t} e^{\Lambda(\vec{d}, \vec{k}, \vec{l})x} \left(\sum_{i^{*}, j^{*} = 1}^{n} \alpha_{\beta, \gamma, \theta, \eta = 0} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^{*}j^{*}\alpha m\beta b\gamma \theta \eta} (\vec{d}, \vec{k}, \vec{l}) \times \right) \times \vec{V}_{q}(\vec{d} + I_{i^{*}} - I_{j^{*}}, \vec{k} + I_{\alpha} + mI_{\beta} - bI_{\gamma}, \vec{l} + I_{\theta} - I_{\eta}, x) dx + \frac{\vec{E}(\vec{d}, k, \vec{l})}{\Lambda(\vec{d}, \vec{k}, \vec{l})} \left[1 - e^{-\Lambda(\vec{d}, \vec{k}, \vec{l})t} \right].$$

$$(3)$$

2. Нахождение ожидаемых доходов методом последовательных приближений

Аналогично [10] можно показать, что последовательность $\{\vec{V}_q(\vec{d},\vec{k},\vec{l},t)\}$, q=0,1,2,..., построенная по схеме (3), при любом ограниченном по t нулевом приближении $\vec{V}_0(\vec{d},\vec{k},\vec{l},t)$ сходится при $q\to\infty$ к единственному решению системы (1), а каждое последовательное приближение с течением времени сходится к стационарному решению системы (1), которое удовлетворяет соотношению

$$\Lambda\left(\vec{d},\vec{k},\vec{l}\right)\vec{V}\left(\vec{d},\vec{k},\vec{l}\right) = \sum_{i^*,j^*=1}^{n} \sum_{\alpha,\beta,\gamma,\theta,\eta=1}^{\Psi_r} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^*j^*\alpha m\beta b\gamma\theta\eta} \left(\vec{d},\vec{k},\vec{l}\right) \times \\
\times \vec{V}\left(\vec{d}+I_{i^*}-I_{j^*},\vec{k}+I_{\alpha}+mI_{\beta}-bI_{\gamma},\vec{l}+I_{\theta}-I_{\eta}\right) + \vec{E}\left(\vec{d},\vec{k},\vec{l}\right).$$
(4)

Кроме того, справедливо следующее утверждение.

Теорема. Любое приближение $\vec{V}_q(\vec{d},\vec{k},\vec{l},t)$, $q \ge 1$, представимо в виде сходящегося степенного ряда

$$\vec{V}_{q}(\vec{d}, \vec{k}, \vec{l}, t) = \sum_{l=0}^{\infty} \vec{g}_{ql}^{+-} (\vec{d}, \vec{k}, \vec{l}) t^{l},$$
 (5)

коэффициенты которого удовлетворяют рекуррентным соотношениям

$$\vec{g}_{q+1l}^{+-}(\vec{d},\vec{k},\vec{l}) = \frac{-\Lambda(\vec{d},\vec{k},\vec{l})^{l}}{l!} \left\{ \vec{V}(\vec{d},\vec{k},\vec{l},0) + \sum_{u=0}^{l-1} \frac{(-1)^{u+1}u!}{\Lambda(\vec{d},\vec{k},\vec{l})^{u+1}} \vec{G}_{qu}^{+-}(\vec{d},\vec{k},\vec{l}) \right\}, l \ge 0,$$

$$\vec{g}_{q0}^{+-}(\vec{d},\vec{k},\vec{l}) = \vec{V}(\vec{d},\vec{k},\vec{l},0), \ \vec{g}_{0l}^{+}(\vec{d},\vec{k},\vec{l}) = \vec{V}(\vec{d},\vec{k},\vec{l},0) \delta_{l0},$$

$$\vec{G}_{ql}^{+-}(\vec{d},\vec{k},\vec{l}) = \sum_{i^{*},j^{*}=1}^{n} \sum_{\alpha,\beta,\gamma,\theta,\eta=0}^{\Psi_{r}} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^{*}j^{*}\alpha m\beta b\gamma\theta\eta} (\vec{d},\vec{k},\vec{l}) \times$$

$$\times \vec{g}_{ql}(\vec{d}+I_{i^{*}}-I_{j^{*}},\vec{k}+I_{\alpha}+mI_{\beta}-bI_{\gamma},\vec{l}+I_{\theta}-I_{\eta}).$$
(6)

Доказательство. Докажем, что коэффициенты степенного ряда (5) удовлетворяют рекуррентным соотношениям (6). Подставим последовательные приближения (6) в (3). Тогда, учитывая, что

$$e^{-\Lambda(\vec{d},\vec{k},\vec{l})t} \int_{0}^{t} e^{\Lambda(\vec{d},\vec{k},\vec{l})x} x^{l} dx = \left[\frac{1}{\Lambda(\vec{d},\vec{k},\vec{l})} \right]_{j=l+1}^{l+1} \frac{\left[-\Lambda(\vec{d},\vec{k},\vec{l}) \right]^{l}}{j!}, l = 0,1,2,...,$$

получим

$$\begin{split} \sum_{l=0}^{\infty} \vec{g}_{ql}^{\; +-} \Big(\vec{d}, \vec{k}, \vec{l} \Big) t^l &= e^{-\Lambda \left(\vec{d}, \vec{k}, \vec{l} \right) t} \vec{V} \Big(\vec{d}, \vec{k}, \vec{l}, 0 \Big) + \sum_{i^*, j^* = 1}^{n} \sum_{\alpha, \beta, \gamma, \theta, \eta = 0}^{\Psi_r} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{i^*j^* \alpha m \beta b \gamma \theta \eta} \Big(\vec{d}, \vec{k}, \vec{l} \Big) \times \\ &\times \vec{g}_{q-1l} \Big(\vec{d} + I_{i^*} - I_{j^*}, \vec{k} + I_{\alpha} + m I_{\beta} - b I_{\gamma}, \vec{l} + I_{\theta} - I_{\eta} \Big). \end{split}$$

Используя обозначения (6), этот ряд можно переписать в виде:

$$\sum_{l=0}^{\infty} \vec{g}_{ql}^{+-} (\vec{d}, \vec{k}, \vec{l}) t^{l} = e^{-\Lambda(\vec{d}, \vec{k}, \vec{l})t} \vec{V} (\vec{d}, \vec{k}, \vec{l}, 0) + \sum_{l=0}^{\infty} \vec{G}_{ql}^{+-} (\vec{d}, \vec{k}, \vec{l}) \left[\frac{1}{\Lambda(\vec{d}, \vec{k}, \vec{l})} \right]^{l+1} l! \sum_{u=l+1}^{\infty} \frac{\left[-\Lambda(\vec{d}, \vec{k}, \vec{l}) \right]^{u}}{u!} t^{u}.$$

Поменяв местами индексы суммирования и разлагая $e^{-\Lambda(\vec{d},\vec{k},\vec{l}\,)t}$ в ряд по степеням t, будем иметь

$$\sum_{l=0}^{\infty} \vec{g}_{ql}^{+-} (\vec{d}, \vec{k}, \vec{l}) t^{l} = \sum_{l=0}^{\infty} \frac{\left[-\Lambda(\vec{d}, \vec{k}, \vec{l}) \right]^{l}}{l!} \left\{ \vec{V}(\vec{d}, \vec{k}, \vec{l}, 0) + \sum_{u=0}^{l-1} \frac{(-1)^{u+1} u!}{\left[\Lambda(\vec{d}, \vec{k}, \vec{l}) \right]^{u+1}} G_{qu}^{+-} (\vec{d}, \vec{k}, \vec{l}) \right\} t^{l}.$$
 (7)

Приравнивая в левой и правой частях выражения (7) коэффициенты при t^l , получим соотношения (6) для коэффициентов ряда (5). Доказательство того, что радиус сходимости ряда (6) равен $+\infty$, можно провести, используя формулу Коши–Адамара, аналогично [12].

3. Анализ сети с разнотипными заявками и сигналами

Рассмотрим, например, G-сеть с разнотипными положительными заявками и сигналами. В сеть из внешней среды поступают простейший поток обычных (положительных) заявок с интенсивностью λ^+ и дополнительный поток сигналов, который также является простейшим с интенсивностью $\lambda^{(1)}$, $i=\overline{1,n}$. Все поступающие потоки независимы. Каждая положительная заявка входного потока независимо от других заявок направляется в i-ю CMO как заявка типа c с вероятностью p_{0ic}^+ , $\sum\limits_{i=1}^{n}\sum\limits_{c=1}^{r}p_{0ic}^+=1$. Если линия обслуживания в i-й CMO свободна, то заявка поступает на обслуживание, иначе она становится в очередь. Положительная заявка при переходе из одной CMO в другую приносит ей некоторый доход, а доход первой CMO уменьшается, соответственно, на эту величину. Длительности обслуживания положительных заявок c-го типа в i-й CMO распределены по экспоненциальному закону с параметром μ_{ic} , $i=\overline{1,n}$, $c=\overline{1,r}$. Будем считать, что заявки на обслуживание из очереди выбираются случайно, т.е. если в i-й CMO находится k_{is} заявок класса s, то вероятность того, что на обслуживании в ней будет заявка класса c, равна $\frac{k_{ic}}{r}$, $i=\overline{1,n}$, $c=\overline{1,r}$.

Сигнал входного потока независимо от других сигналов направляется в i-ю СМО как сигнал типа c с вероятностью p_{0ic}^- , $\sum_{i=1}^n \sum_{c=1}^r p_{0ic}^- = 1$. Сигнал типа c, поступающий в СМО, в которой нет положительных заявок данного типа, не оказывает никакого влияния и сразу исчезает из нее. В противном случае могут произойти следующие события: поступающий сигнал мгновенно перемещает положительную заявку типа c из системы i-й СМО в j-ю СМО как заявку типа s с вероятностью q_{icjs} , в этом случае сигнал называют триггером, или с вероятностью $q_{in0} = 1 - \sum_{j=1}^n \sum_{s=1}^r q_{icjs}$ сигнал срабатывает как отрицательная заявка и уничтожает в i-й СМО положительную заявку типа c. Таким образом, отрицательная заявка является частным случаем сигнала, когда $q_{icjs} = 0$, $q_{ic0} = 1$. После окончания обслуживания положительной заявки типа c в i-й СМО она направляется в j-ю СМО с вероятностью p_{icjs}^+ опять как положительная заявка типа s, а с вероятностью p_{icjs}^- как сигнал типа s, и с вероятностью p_{icjs}^- опять как положительная заявка типа s, а с вероятностью p_{icjs}^- как сигнал типа s, и с вероятностью p_{icjs}^- опять как положительная заявка типа s, а с вероятностью p_{icjs}^- как сигнал типа s, и с вероятностью

Под состоянием i-й СМО в момент времени t будем понимать вектор $\vec{k}_i(t) = (k_{i1}(t), k_{i2}(t), ..., k_{ir}(t))$, где $k_{ic}(t)$ – число положительных заявок типа c в i-й СМО в момент времени t, а под состоянием сети – вектор $\vec{k}(t) = (\vec{k}, t) = (\vec{k}_1(t), ..., \vec{k}_n(t))$, который образует цепь Маркова со счетным числом состояний и непрерывным временем.

Введем в рассмотрение вектор I_{ic} , состоящий из нулей, за исключением компоненты с номером r(i-1)+c, которая равна единице; $M_{ic}\left(\vec{k}_i\right)=\mu_{ic}\frac{k_{ic}+1}{\displaystyle\sum_{s^*=1}^r k_{is^*}+1}$; $v_i\left(\vec{k}_i,t\right)$ – ожидаемый доход, который

получает i-я СМО за время t, если в начальный момент времени сеть находится в состоянии \vec{k} ; u(x) – единичная функция Хевисайда, $u(x) = \begin{cases} 1, x > 0, \\ 0, x \le 0, \end{cases}$. Наша цепь Маркова может осуществлять следующие переходы в состояние (\vec{k} , t) за время Δt :

- 1) из состояния $(\vec{k}-I_{js},t)$, $j\neq i, s\neq c$, в этом случае в j-ю СМО за время Δt поступит положительная заявка типа s с вероятностью $\lambda^+p_{0js}^+u(k_{js})\Delta t+o(\Delta t)$, $j=\overline{1,n}, s=\overline{1,r}$; доход системы S_i в этом случае составит $r_i(\vec{k})\Delta t+v_i(\vec{k}-I_{js},t)$; если i=j, s=c, то доход системы S_i составит $r_{0ic}(\vec{k}-I_{ic})+v_i(\vec{k}-I_{ic},t)$, где $r_{0ic}(\vec{k}-I_{ic})$, доход i-й системы от данного перехода, $r_i(\vec{k})$ доход системы в единицу времени за пребывание в состоянии \vec{k} ;
- 2) из состояния $(\vec{k} + I_{js}, t)$, $j \neq i$, $s \neq c$, в данном случае в j-ю СМО за время Δt поступит сигнал типа s, который сработает как отрицательная заявка и уничтожит в ней положительную заявку своего типа, или после завершения обслуживания положительная заявка типа s уйдёт из сети, или переходит в m-ю СМО как сигнал типа l, но не обнаружит там положительных заявок данного типа с вероятностью $(\lambda^{(1)}p_{0js}^-q_{js0} + M_{js}(\vec{k}_j)p_{js0}^- + M_{js}(\vec{k}_j)p_{jsml}^-(1-u(k_{ml}))\Delta t + o(\Delta t), j = \overline{1,n}, s = \overline{1,r};$ доход системы S_i в этом случае составит $r_i(\vec{k})\Delta t + v_i(\vec{k} + I_{js}, t)$, если i = j, s = c, то доход системы S_i составит $-R_{ic0}(\vec{k} + I_{ic}) + v_i(\vec{k} + I_{ic}, t)$, где $R_{ic0}(\vec{k} + I_{ic})$ доход i-й системы от данного перехода;
- 3) из состояния $\left(\vec{k}+I_{ml}-I_{dh},\,t\right)$, $m\neq i,\,l\neq c,\,d\neq j,\,h\neq s$, при этом положительная заявка типа l после обслуживания в m-й СМО перейдет в d-ю СМО в качестве положительной заявки типа h, или из внешней среды в m-ю СМО поступит сигнал типа l, который действует как триггер и сразу переместит положительную заявку из системы S_m в систему S_d в качестве положительной заявки типа h; вероятность такого события равна $\left(M_{ml}\left(\vec{k}_m\right)p_{mldh}^+ + \lambda^{(1)}p_{0ml}^-q_{mldh}\right)u\left(k_{dh}\right)\Delta t + o\left(\Delta t\right),\,\,d,m=\overline{1,n};$ $l,h=\overline{1,r};$ доход системы S_i в этом случае составит $r_i(\vec{k})\Delta t + v_i(\vec{k}+I_{ml}-I_{dh},t);$ если $m=j,\,l=s,\,i=d,\,c=h,\,$ то доход S_i составит $-r_{jsic}(\vec{k}+I_{js}-I_{ic})+v_i(\vec{k}+I_{js}-I_{ic},t);$ если $k=i,\,l=c,\,j=d,\,s=h,\,$ то доход S_i составит $r_{icjs}(\vec{k}-I_{js}+I_{ic})+v_i(\vec{k}-I_{js}+I_{ic},t);$
- 4) из состояния $(\vec{k} + I_{ml} + I_{dh}, t)$, при этом после окончания обслуживания положительной заявки типа l в СМО S_m она направится в СМО S_d в качестве сигнала типа d, который сработает в ней как отрицательная заявка типа d, уничтожит в S_d положительную заявку своего типа; вероятность такого события равна $M_{ml}(\vec{k}_m)p_{mldh}^-q_{dh0}\Delta t + o(\Delta t)$; доход системы S_i в этом случае составит $r_i(\vec{k})\Delta t + v_i(\vec{k} + I_{ml} + I_{dh}, t)$; если m = j, l = s, I = d, c = h, то доход S_i составит $r_{icjs}(\vec{k} + I_{js} + I_{ic}) + v_i(\vec{k} + I_{js} + I_{ic}, t)$;

5) из состояния $(\vec{k} + I_{ml} + I_{dh} - I_{\alpha\beta}, t)$, в этом случае после окончания обслуживания заявки типа m в СМО S_m она направится в СМО S_j как сигнал типа s, который мгновенно переместит положительную заявку своего типа из системы S_j в систему S_m как положительную заявку типа l; вероятность такого события равна $M_{ml}(\vec{k}_m)p_{mldh}^-q_{dh\alpha\beta}u(k_{dh})\Delta t + o(\Delta t)$; доход системы S_i в этом случае составит $r_i(\vec{k} + I_{ml} + I_{dh} - I_{\alpha\beta}, t)\Delta t + v_i(\vec{k} + I_{ml} + I_{dh} - I_{\alpha\beta}, t)$; при m=i, l=c доход составит $-r_{dh\alpha}(\vec{k} + I_{ic} + I_{dh} - I_{\alpha\beta}, t) + v_i(\vec{k} + I_{ic} + I_{dh} - I_{\alpha\beta}, t)$; при $\alpha=i$, $\beta=c$ он будет равен $-r_{dh\alpha}(\vec{k} + I_{ml} + I_{dh} - I_{ic}, t) + v_i(\vec{k} + I_{ml} + I_{dh} - I_{ic}, t)$; иначе $r_{dh\alpha}(\vec{k} + I_{ml} + I_{dh} - I_{ic}, t) + v_i(\vec{k} + I_{ml} + I_{dh} - I_{a\beta}, t)$; иначе $r_{dh\alpha}(\vec{k} + I_{ml} + I_{dh} - I_{ic}, t) + v_i(\vec{k} + I_{ml} + I_{dh} - I_{a\beta}, t)$;

6) из состояния (\vec{k},t) , при этом в каждую СМО $S_i,\ i=\overline{1,n}$, не поступают ни положительные заявки любых типов, ни сигналы, или сигналы при поступлении не будут обнаруживать заявок своего типа и в этих СМО за время Δt не обслужилось ни одной заявки; вероятность такого события равна $1-\sum\limits_{i=1}^{n}\sum\limits_{c=1}^{r}\left[\lambda^{+}p_{0ic}^{\ +}+\lambda^{-}p_{0ic}^{\ -}+\mu_{ic}\right]\Delta t+o(\Delta t),\ i=\overline{1,n}$; доход системы S_i в этом случае составит $r_i(\vec{k})\Delta t+v_i(\vec{k},t)$.

Используя формулу полной вероятности, поделив обе ее части на Δt и переходя к пределу при $\Delta t \to 0$, получим, что ожидаемые доходы систем рассматриваемой в данном случае сети удовлетворяют системе РДУ (1), где:

$$\begin{split} &\Lambda\left(\vec{d},\vec{k},\vec{l}\right) = \Lambda\left(\vec{k}\right) = \sum_{i=1}^{n} \sum_{c=1}^{r} \left[\lambda_{0ic}^{+} + \lambda_{0ic}^{(1)} + \mu_{ic}\right], \\ &\Theta_{i'j'amplipipin}\left(\vec{d},\vec{k},\vec{l}\right) = \Theta_{i'j'amplipipin}\left(\vec{k}\right) = \delta_{i'j'}\delta_{\bar{d}\bar{i}_k}\delta_{a((i-1)r+c)}\delta_{nj}\delta_{i\bar{0}}\left(\delta_{m0}\delta_{bi}\delta_{\gamma((j-1)r+s)} \times \right) \\ &\times \delta_{0j}\delta_{0s}[\lambda^{(1)}p_{0ic}^{-}q_{ic0} + \mu_{ic}\frac{k_{ic}+1}{\sum_{s'=1}^{r}k_{is'}+1}p_{ic0} + \mu_{ic}\frac{k_{ic}+1}{\sum_{s'=1}^{r}k_{is'}+1}\sum_{m=1}^{r}\sum_{l=1}^{r}p_{ioml}^{-}\left(1-u\left(k_{ml}\right)\right)] + \\ &+ \delta_{0i}\delta_{0c}\lambda^{+}p_{0js}^{+}u\left(k_{js}\right) + \left[\mu_{ic}\frac{k_{ic}+1}{\sum_{l'=1}^{r}k_{il'}+1}p_{icjs}^{+} + \lambda^{(1)}p_{0ic}q_{icjs}\right]u\left(k_{js}\right) + \\ &+ \delta_{ml}\delta_{b0}\delta_{\beta((j-1)r+s)}\mu_{ic}\frac{k_{ic}+1}{\sum_{l'=1}^{r}k_{ic}+1}p_{icjs}^{-}\delta_{0m}\delta_{0l} + \delta_{ml}\delta_{b1}\delta_{\beta((j-1)r+s)}\delta_{\gamma((m-1)r+l)}\mu_{ic}\frac{k_{ic}+1}{\sum_{l=1}^{r}k_{ic}+1}p_{icjs}^{-}q_{joml}u\left(k_{js}\right), \\ &E_{i}\left(\vec{d},\vec{k},\vec{l}\right) = r_{i}\left(\vec{k}\right) + \lambda^{+}p_{0ic}^{+}u\left(k_{ic}\right)r_{0ic}\left(\vec{k}-I_{ic}\right) - \\ &- \left[\lambda^{(1)}p_{0js}^{-}q_{js0} + \mu_{js}\frac{k_{js}+1}{\sum_{s'=1}^{r}k_{js'}+1}p_{js0} + \mu_{js}\frac{k_{js}+1}{\sum_{s'=1}^{r}k_{js'}+1}p_{joml}^{-}\left(1-u\left(k_{ml}\right)\right)\right]R_{ic0}\left(\vec{k}-I_{ic},t\right) + \\ &+ \sum_{j\neq i}^{n}\sum_{c=1}^{r}\left(\mu_{ic}\frac{k_{ic}+1}{\sum_{l'=1}^{r}k_{il'}+1}p_{icjs}^{+} + \lambda^{(1)}p_{0ic}q_{icjs}\right)u\left(k_{js}\right)r_{icjs}\left(\vec{k}-I_{ic}+I_{js},t\right) - \\ &- \sum_{j\neq i}^{n}\sum_{c=1}^{r}\left(\mu_{js}\frac{k_{js}+1}{\sum_{l'=1}^{r}k_{jl'}+1}p_{jsic}^{+} + \lambda^{(1)}p_{0js}q_{jsic}\right)u\left(k_{js}\right)r_{icjs}\left(\vec{k}-I_{js}+I_{ic},t\right) + \\ &+ \sum_{j\neq i}^{n}\sum_{c=1}^{r}\left(\mu_{js}\frac{k_{js}+1}{\sum_{l'=1}^{r}k_{jl'}+1}p_{jsic}^{+} + \lambda^{(1)}p_{0js}q_{jsic}\right)u\left(k_{js}\right)r_{icjs}\left(\vec{k}-I_{js}+I_{jc},t\right) + \\ &+ \sum_{j\neq i}^{n}\sum_{c=1}^{n}\left$$

$$\begin{split} &+\sum_{d,\alpha=1}^{n}\sum_{h,\beta=1}^{r}\mu_{ic}\frac{k_{ic}+1}{\sum_{l=1}^{r}k_{il}+1}p_{icdh}^{-}q_{dh\alpha\beta}u(k_{dh})r_{id\alpha}\left(\vec{k}+I_{ic}+I_{dh}-I_{\alpha\beta},\ t\right)+\\ &+\sum_{m,\alpha=1}^{n}\sum_{l,\beta=1}^{r}\mu_{ml}\frac{k_{ml}+1}{\sum_{l=1}^{r}k_{ml}+1}p_{mlic}^{-}q_{ic\alpha\beta}u(k_{dh})r_{mi\alpha}\left(\vec{k}+I_{ml}+I_{ic}-I_{\alpha\beta},\ t\right)+\\ &+\sum_{m,d=1}^{n}\sum_{l,h=1}^{r}\mu_{ml}\frac{k_{ml}+1}{\sum_{l=1}^{r}k_{ml}+1}p_{mldh}^{-}q_{dhic}u(k_{dh})r_{mdi}\left(\vec{k}+I_{ml}+I_{dh}-I_{ic},\ t\right). \end{split}$$

Заключение

Проведено исследование в нестационарном режиме открытых марковских CeMO с различными особенностями. Рассмотрена обобщенная система РДУ для ожидаемых доходов в системах сети, состоящая из счетного числа таких уравнений. Когда доходы от переходов между состояниями сети зависят только от ее состояний, для решения системы предложено применить метод последовательных приближений, совмещенный с методом рядов. Исследованы свойства последовательных приближений.

ЛИТЕРАТУРА

- 1. Маталыцкий М.А. О некоторых результатах анализа и оптимизации марковских сетей с доходами и их применении // Автоматика и телемеханика. 2009. № 10. Р. 97–113.
- 2. Паньков А.В. Анализ стохастической модели расходов на содержание гибкого вычислительного кластера // Современные математические методы анализа и оптимизации информационно-телекоммуникационных сетей: материалы 20-й междунар. науч. конф., Минск, 26–29 янв. 2009 г. / ред. А.Н. Дудин [и др.]. Минск: РИВШ, 2009. Вып. 20. С. 184–188.
- 3. Gelenbe E. G-networks: a unifying model for neural and queueing networks // Annals of Operations Research. 1994. V. 48. P. 433-461.
- 4. Ховард Р. Динамическое программирование и марковские процессы. М.: Сов. радио, 1964. 109 с.
- 5. Gelenbe E. Product form queueing networks with negative and positive customers // J. of Applied Probability. 1991. V. 28. P. 656–663.
- Gelenbe E., Labed A. G-networks with multiple classes of signals and positive customers // European J. of Operational Research. 1998. V. 108. P. 293–305.
- 7. Gelenbe E. G-networks with signals and batch removal // Probability in the Engineering and Informational Sciences. 1993. V. 7. P. 335–342.
- 8. Fourneaua J.M., Gelenbe E., Surosc R. G-networks with multiple classes of negative and positive customers // Theoretical Computer Science. 1996. V. 155, is. 1. P. 141–156.
- 9. Gelenbe E. G-Networks: Multiple Classes of Positive Customers, Signals, and Product Form // Results Performance Evaluation of Complex Systems: Techniques and Tools. Springer, 2002. P. 1–16. (Lecture Notes in Computer Science. V. 2459)
- 10. Matalytski M. Finding non-stationary state probability of G-networks with signal and customers batch removal // Probability in the Engineering and Informational Sciences. 2017. V. 31, No. 4. P. 346–412.
- 11. Статкевич С.Э. Анализ НМ-сети с ненадёжными системами обслуживания и случайными доходами от переходов между её состояниями // Вестник Гродненского государственного университета. Сер. 2. 2010. № 3. С. 40–52.
- 12. Маталыцкий М.А., Науменко В.В. Стохастические сети с нестандартными перемещениями заявок. Гродно : Гродненский гос. ун-т, 2016. 347 с.

Поступила в редакцию 5 июля 2019 г.

Matalytski M.A., Kopat D.Ya. (2020) ANALYSIS OF EXPECTED REVENUES IN OPEN MARKOV NETWORKS WITH VARIOUS FEATURES. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 31–38

DOI: 10.17223/19988605/50/4

Investigation at non-stationary regime open Markov QN with various features has been carried out. The generalized system of DDE for expected revenues in network systems, consisting from countable of number equations has been considered. When the revenues from transitions between network states can depend only on these states, a successive approximation method in combination with a series method has been proposed for the decision system. The properties of successive approximations are investigated. The considered system has the form:

$$\begin{split} \frac{d\vec{V}\left(\vec{d},\vec{k},\vec{l}\,,t\right)}{dt} &= -\Lambda\left(\vec{d}\,,\vec{k}\,,\vec{l}\,\right) \vec{V}\left(\vec{d}\,,\vec{k}\,,\vec{l}\,,t\right) + \sum_{\vec{i}^{\,\prime}\,,\vec{j}^{\,\prime}=0}^{\,\prime} \sum_{\alpha,\beta,\gamma,\theta,\eta=0}^{\,\Psi r} \sum_{m=0}^{\infty} \sum_{b=0}^{1} \Theta_{\vec{i}^{\,\prime}\,\vec{j}^{\,\prime}\alpha m\beta b\gamma\theta\eta}\left(\vec{d}\,,\vec{k}\,,\vec{l}\,\right) \times \\ &\times \vec{V}\left(\vec{d}\,+I_{\vec{i}^{\,\prime}}-I_{\vec{j}^{\,\prime}}\,,\vec{k}\,+\tilde{I}_{\alpha}+m\tilde{I}_{\beta}-b\tilde{I}_{\gamma},\vec{l}\,+\tilde{I}_{\theta}-\tilde{I}_{\eta},t\right) + \vec{E}\left(\vec{d}\,,\vec{k}\,,\vec{l}\,,t\right), \end{split}$$

where $\vec{V}^{\mathrm{T}}\left(\vec{d},\vec{k},\vec{l},t\right) = \left(v_1\left(\vec{d},\vec{k},\vec{l},t\right),v_2\left(\vec{d},\vec{k},\vec{l},t\right),\dots,v_n\left(\vec{d},\vec{k},\vec{l},t\right)\right)$, $v_i\left(\vec{d},\vec{k},\vec{l},t\right)$ are the expected revenue obtained by the i-th QS in time t, \tilde{I}_{α} is a vector of dimension Ψr , consisting of zeros, with the exception of the component with number α , which is 1, Ψr is a positive number, r is the number of customer types, I_{α} is a vector of dimension n, consisting of zeros, with the exception of the component with the number α , which is 1, \vec{d} is a vector of dimension n, consisting of components d_i , where d_i is a number of serviceable channels in the i-th QS at the moment t, \vec{k} is a vector of dimension Ψr , consisting of components k_{ic} , where k_{ic} is a number of positive customers of type c in queue i at the moment t, \vec{l} is a vector of dimension Ψr , consisting of components l_{ic} , where l_{ic} is a number of signals of type c in the queue i at the moment t, $\Lambda\left(\vec{d},\vec{k},\vec{l}\right)$, $\Theta_{\vec{l},\vec{l}' \text{ cump}br_i\Theta_i}\left(\vec{d},\vec{k},\vec{l}\right)$, $\vec{E}\left(\vec{d},\vec{k},\vec{l}\right)$ are some functions that are different for each service network.

The form of successive approximations and their properties have been investigated in the paper.

Keywords: the system of difference-differential equations (DDE); open queueing network; the method of successive approximations.

MATALYTSKI Mihail Alekseevich (Doctor of Physics and Mathematics, Professor, Grodno State University, Belarus). E-mail: m.matalytski@gmail.com

KOPAT Dmitry Yaroslavovich (Post-graduate Student, Grodno State University, Belarus).

E-mail: dk80395@mail.ru

REFERENCES

- 1. Matalytski, M. (2009) On some results in analysis and optimization of Markov networks with incomes and their application. *Automatic and Remote Control*. 70(10). pp. 1683–1697. DOI: 10.1134/S0005117909100075
- 2. Pankov, A.V. (2009) [Analysis of the stochastic cost model for maintaining a flexible computing cluster]. Sovremennye matematicheskie metody analiza i optimizatsii informatsionno-telekommunikatsionnykh setey [Modern Mathematical Methods for the Analysis and Optimization of Information and Telecommunication Networks]. Proc. of the 20th International Conference. Minsk. Janruary 26–29, 2009. pp. 184–188.
- Gelenbe, E. (1994) G-networks: a unifying model for neural and queueing networks. Annals of Operations Research. 48. pp. 433–461. DOI: 10.1007/BF02033314
- Howard, R. (1964) Dinamicheskoe programmirovanie i markovskie protsessy [Dynamic programming and Markov process].
 Translated from English by V.V. Rykov. Moscow: Sovetskoe Radio.
- 5. Gelenbe, E. (1991) Product form queueing networks with negative and positive customers. *Journal of Applied Probability*. 28. pp. 656–663. DOI: 10.2307/3214499
- 6. Gelenbe, E. & Labed, A. (1998) G-networks with multiple classes of signals and positive customers. *European Journal of Operational Research*. 108. pp. 293–305. DOI: 10.1016/S0377-2217(97)00371-8
- 7. Gelenbe, E. (1993) G-networks with signals and batch removal. *Probability in the Engineering and Informational Sciences*. 7. pp. 335–342. DOI: 10.1017/S0269964800002953
- 8. Fourneaua, J.M., Gelenbe, E. & Surosc, R. (1996) G-networks with multiple classes of negative and positive customers. *Theoretical Computer Science*. 155(1). pp. 141–156. DOI: 10.1016/0304-3975(95)00018-6
- Gelenbe, E. (2002) G-Networks: Multiple classes of positive customers, signals, and product form. In: Calzarossa, M.C. & Tucci, S. (eds) Results performance evaluation of complex systems: Techniques and Tools. Springer-Verlag Berlin Heidelberg. pp. 1–16. DOI: 10.1007/3-540-45798-4
- 10. Matalytski, M. (2017) Finding non-stationary state probability of G-networks with signal and customers batch removal. *Probability in the Engineering and Informational Sciences*. 31(4). pp.346–412. DOI: 10.1017/S0269964817000109
- 11. Statkevich, S. (2010) Analysis of HM-networks with unreliable systems and random revenues for transitions between states. *Vestnik GrSU*. 1(3). pp. 40–52.
- 12. Matalytski, M. & Naumenko, V. (2016) Stochastic network with nonstandart moving customers. Grodno: GrSU.

Управление, вычислительная техника и информатика

№ 50

УДК 369:519.2

DOI: 10.17223/19988605/50/5

О.В. Губина, Г.М. Кошкин

ОЦЕНИВАНИЕ СОВРЕМЕННОЙ СТОИМОСТИ *n*-ЛЕТНЕЙ РЕНТЫ ДЛЯ СМЕШАННОГО СТРАХОВАНИЯ ЖИЗНИ

Рассматривается задача оценивания *п*-летней ренты для смешанного страхования жизни, которое часто предлагается страховыми компаниями. Находятся главная часть асимптотической среднеквадратической ошибки и порядок смещения оценки ренты, доказывается ее асимптотическая нормальность.

Ключевые слова: смешанное страхование жизни; *п*-летняя рента; непараметрическая оценка; среднеквадратическая ошибка; асимптотическая нормальность.

Суть смешанного страхования жизни, или n-летнего страхования на дожитие, заключается в следующем. Человек заключает договор страхования на n лет. Выплата по договору производится либо в момент смерти застрахованного бенефициарию, если застрахованный умер в течение n лет, либо в момент окончания срока действия договора, если застрахованный дожил до конца этого срока. Этот вид договора выполняет функции как страхования, так и накопления средств, тем самым являясь наиболее привлекательным для клиента.

В страховую компанию обращаются люди, достигшие определенного возраста x лет, поэтому все случайные события (страховые случаи), связанные с этим человеком, имеют условный характер. Для человека в возрасте x лет целесообразнее использовать не продолжительность жизни X, а остаточное время жизни T(x) = X - x. Согласно [1–3] остаточное время жизни T(x) имеет функцию распределения

$$F_{x}(t) = P(T(x) \le t) = \frac{S(x) - S(x+t)}{S(x)}$$

и плотность

$$f_{x}(t) = \frac{d}{dt}F_{x}(t) = -\frac{d}{dt}S_{x}(t) = \frac{f(x+t)}{S(x)}, \quad 0 \le t < \infty,$$

где S(x) – функция выживания, f(u) = -S'(u) – плотность распределения продолжительности жизни X.

Определим для смешанного страхования жизни современную величину страховой выплаты z:

$$z = \begin{cases} e^{-\delta T(x)}, & T(x) \le n, \\ e^{-\delta n}, & T(x) > n, \end{cases}$$
 (1)

где δ обозначает банковскую процентную ставку. В данном случае величина z, определяемая выражением (1), показывает настоящую долю будущей страховой выплаты, принимаемой за условную единицу. Чем больше срок страхования, тем меньше выплаты застрахованного за счет использования банковской процентной ставки.

В качестве n-летней пожизненной ренты для смешанного страхования по аналогии с [3] и формулами (1) и (2) из статьи [4] получаем

$$\overline{a}_{x:\overline{n}|} = \frac{1 - \frac{1}{S(x)} \int_{0}^{n} e^{-\delta t} f(x+t) dt - \frac{e^{-\delta n} S(x+n)}{S(x)}}{\delta}.$$
 (2)

С помощью замены переменных преобразуем интеграл в (2):

$$\int_{0}^{n} e^{-\delta t} f(x+t)dt = e^{\delta x} \int_{x}^{x+n} e^{-\delta t} dF(t) = \Phi_{n}(x,\delta),$$

$$\int_{x}^{x+n} e^{-\delta t} dF(t) = J_{n}(x,\delta).$$

Тогда формула (2) принимает вид

$$\overline{a}_{x,\overline{n}|} = \frac{1}{\delta} \left(1 - \frac{e^{\delta x}}{S(x)} \int_{x}^{x+n} e^{-\delta t} dF(t) - \frac{e^{-\delta n} S(x+n)}{S(x)} \right), \tag{3}$$

или

$$\overline{a}_{x:\overline{n}|} = \frac{1}{\delta} \left(1 - \left(\frac{\Phi_n(x,\delta)}{S(x)} + \frac{e^{-\delta n}S(x+n)}{S(x)} \right) \right). \tag{4}$$

Далее будут использоваться как формула (3), так и формула (4).

1. Синтез оценки

Пусть имеется случайная выборка $X_1, ..., X_N$ продолжительности жизни X, по которой необходимо оценить ренту (3).

Воспользуемся вместо неизвестных F(x) и S(x) их непараметрическими оценками: эмпирическими функциями распределения $F_N(x) = \frac{1}{N} \sum\limits_{i=1}^N I(X_i \leq x)$ и выживания $S_N(x) = \frac{1}{N} \sum\limits_{i=1}^N I(X_i > x)$, где I(A) – индикатор события A. Подставив $F_N(x)$ и $S_N(x)$ в выражения для смешанной ренты (3) или (4), получим следующую оценку подстановки:

$$\overline{a}_{x,\overline{n}|}^{N} = \frac{1}{\delta} \left(1 - \left(\frac{e^{\delta x}}{S_{N}(x) \cdot N} \sum_{i=1}^{N} \exp(-\delta X_{i}) \mathbf{I}(x < X_{i} \le (x+n)) + \frac{e^{-\delta n} S_{N}(x+n)}{S_{N}(x)} \right) \right) = \frac{1}{\delta} \left(1 - \left(\frac{e^{\delta x} J_{n,N}(x,\delta)}{S_{N}(x)} + \frac{e^{-\delta n} S_{N}(x+n)}{S_{N}(x)} \right) \right) = \frac{1}{\delta} \left(1 - \left(\frac{\Phi_{n,N}(x,\delta)}{S_{N}(x)} + \frac{e^{-\delta n} S_{N}(x+n)}{S_{N}(x)} \right) \right).$$
(5)

Отметим, что в оценке (5) вместо эмпирических функций $F_N(x)$ и $S_N(x)$ можно воспользоваться их гладкими модификациями [5–19].

2. Свойства оценки *п*-летней ренты

Найдем сначала главную часть асимптотической среднеквадратической ошибки (СКО) и порядок смещения оценки (5). Для этого нам понадобится теорема 1 из [20], которую ниже сформулируем в виде леммы.

Введем следующие обозначения согласно [20]: $t_N = \left(t_{1N}, t_{2N}, ..., t_{sN}\right)^{\mathrm{T}} - s$ -мерная векторная статистика с компонентами $t_{jN} = t_{jN}(x) = t_{jN}(x; X_1, ..., X_N), \ j = \overline{1, s}, \ x \in R^a, \ R^a - \alpha$ -мерное евклидово пространство. Пусть $\left\{d_N\right\}$ — последовательность положительных чисел, таких что $\lim_{n \to \infty} d_N = \infty$; функция $H(t): R^s \to R^1$, где $t = t(x) = \left(t_1(x), ..., t_s(x)\right)^{\mathrm{T}}$ является s-мерной ограниченной вектор-функцией; $N_s(\mu; \sigma)$ есть s-мерная нормально распределенная случайная величина с вектором средних $\mu = \mu(x) = \left(\mu_1, ..., \mu_s\right)^{\mathrm{T}}$ и ковариационной матрицей $\sigma = \sigma(x); \ \nabla H(t) = \left(H_1(t), ..., H_s(t)\right)^{\mathrm{T}}$, где $H_j(t) = \frac{\partial H(z)}{\partial z_j}\Big|_{z=t}$,

 $j = \overline{1, s}; \Rightarrow$ — символ сходимости по распределению; ||x|| — евклидова норма вектора x; \Re — множество натуральных чисел.

Определение 1. Функция $H(t): R^s \to R^1$ и последовательность $\{H(t_N)\}$ принадлежат классу $N_{v,s}(t;\gamma)$, если:

1) существует є-окрестность

$$\sigma = \left\{ z : |z_i - t_i| < \varepsilon, i = \overline{1, s} \right\},\,$$

в которой функция H(z) и все ее частные производные вплоть до порядка v непрерывны и ограничены;

2) для всевозможных значений величин $X_1,...,X_N$ последовательность $\{H(t_N)\}$ мажорируется числовой последовательностью $C_0d_N^\gamma$, такой что $d_N\uparrow\infty$ при $N\to\infty,\ 0\le\gamma<\infty.$

Лемма. Пусть:

1) H(z), $\{H(t_N)\} \in \mathbb{N}_{2,s}(t;\gamma)$;

2)
$$E ||t_N - t||^i = O(d_N^{-i/2}), i \in \Re.$$

Тогда для любых $k \in \Re$

$$\left| E \left[H(t_N) - H(t) \right]^k - E \left[\nabla H(t) \cdot (t_N - t) \right]^k \right| = o \left(d_N^{-(k+1)/2} \right). \tag{6}$$

При k=1 получаем главную часть смещения оценки $H(t_n)$, а при k=2 – ее СКО.

Теорема 1. Если S(x) > 0, S(x+n) > 0, S(t) — непрерывна в точках x и x+n, то

1) для смещения оценки ренты (5) выполняется следующее соотношение:

$$E\left|\overline{a}_{x:\overline{n}|}^{N}-\overline{a}_{x:\overline{n}|}\right|=o\left(N^{-1}\right);$$

2) СКО оценки (5) задается выражением

$$u^{2}\left(\overline{a}_{x:\overline{n}|}^{N}\right) = E\left(\overline{a}_{x:\overline{n}|}^{N} - \overline{a}_{x:\overline{n}|}\right)^{2} = \frac{\sigma\left(\overline{a}_{x:\overline{n}|}\right)}{N} + o\left(N^{-3/2}\right),\tag{7}$$

где $\sigma(\overline{a}_{x:\overline{n}|})$ определяется по формуле (8).

Доказательство. Для оценки $\overline{a}^N_{x:\overline{n}|}$ в обозначениях леммы имеем:

$$\begin{split} t_N &= \left(t_{1N}, t_{2N}, t_{3N}\right)^{\mathrm{T}} = \left(\Phi_{n,N}(x,\delta), S_N(x), S_N(x+n)\right)^{\mathrm{T}};\\ d_N &= N; \quad t = \left(t_1, t_2, t_3\right)^{\mathrm{T}} = \left(\Phi_n(x,\delta), S(x), S(x+n)\right)^{\mathrm{T}};\\ H(t) &= \frac{1}{\delta} \left(1 - \frac{t_1 + e^{-\delta n}t_3}{t_2}\right) = \frac{1}{\delta} \left(1 - \frac{\Phi_n(x,\delta) + e^{-\delta n}S(x+n)}{S(x)}\right) = \overline{a}_{x:\overline{n}|};\\ H(t_N) &= \frac{1}{\delta} \left(1 - \frac{\Phi_{n,N}(x,\delta) + e^{-\delta n}S_N(x+n)}{S_N(x)}\right) = \overline{a}_{x:\overline{n}|};\\ \nabla H(t) &= \left(H_1(t), H_2(t), H_3(t)\right)^{\mathrm{T}} = \left(\frac{1}{\delta S(x)}, -\frac{\Phi_n(x,\delta) - e^{-\delta n}S(x+n)}{\delta S^2(x)}, -\frac{e^{-\delta n}}{\delta S(x)}\right)^{\mathrm{T}} \neq 0. \end{split}$$

Последовательность $\left\{H(t_N)\right\}$ удовлетворяет условию 1 леммы с константами $C_0=\frac{1}{\delta}\Big(1+e^{-\delta n}\Big),$ $\gamma=0$. Действительно,

$$\begin{split} |H(t_N)| &= \frac{1}{\delta} \left| 1 - \frac{\Phi_{n,N}(x,\delta) + e^{-\delta n} S_N(x+n)}{S_N(x)} \right| \leq \frac{1}{\delta} \left(1 + \frac{\Phi_{n,N}(x,\delta) + e^{-\delta n} S_N(x+n)}{S_N(x)} \right) \leq \\ &\leq \frac{1}{\delta} \left(1 + \frac{e^{\delta x} \sum\limits_{i=1}^N \exp(-\delta X_i) I\left(x < X_i \leq (x+n)\right) + e^{-\delta n} \sum\limits_{i=1}^N I\left(X_i > (x+n)\right)}{\sum\limits_{i=1}^N I\left(X_i > x\right)} \right) \leq \\ &\leq \frac{1}{\delta} \left(1 + \frac{e^{\delta x} e^{-\delta x} \sum\limits_{i=1}^N I\left(x < X_i \leq (x+n)\right) + e^{-\delta n} \sum\limits_{i=1}^N I\left(X_i > (x+n)\right)}{\sum\limits_{i=1}^N I\left(X_i > x\right)} \right) \leq \frac{1}{\delta} \left(1 + e^{-\delta n} \right). \end{split}$$

Функция H(t) удовлетворяет условию 1, так как $t_2 = S(x) > 0$. Также эта функция удовлетворяет условию 2 согласно лемме 3.1 [21], так как для всех $i \in \Re$ выполняются следующие неравенства:

$$E\left\{e^{i\delta x}e^{-i\delta X}I^{i}\left(x < X \le (x+n)\right)\right\} \le e^{i\delta x}e^{-i\delta x}\left[S(x) - S(x+n)\right] = S(x) - S(x+n) \le 1,$$

$$E\left\{I^{i}(X > x)\right\} = S(x) \le 1, \quad E\left\{I^{i}\left(X > (x+n)\right)\right\} = S(x+n) \le 1$$

Отметим, что $S_N(x)$ является несмещенной оценкой S(x), а $J_{n,N}(x,\delta)$ — несмещенной оценкой функционала $J_n(x,\delta)$. Известно, что отношение двух несмещенных оценок может иметь смещение. Нахождение смещения отношения, как правило, является сложной задачей и требует использования результатов работы [20]. Найдем порядок смещения оценки. Так как $E(t_N-t)=0$, то

$$\left| E\left(\overline{a}_{x:\overline{n}|}^{N} - \overline{a}_{x:\overline{n}|}\right) - E\left[\nabla H(t)(t_{N} - t)\right] \right| = \left| E\left(\overline{a}_{x:\overline{n}|}^{N} - \overline{a}_{x:\overline{n}|}\right) \right| = o\left(N^{-1}\right).$$

Для оценки $J_{n,N}(\delta)$ вычислим дисперсию:

$$\begin{split} DJ_{n,N}(x,\delta) &= D\left\{\frac{1}{N}\sum_{i=1}^{N}I(x < X_i \le (x+n))e^{-\delta X_i}\right\} = \frac{1}{N^2}\sum_{i=1}^{N}D\left\{I(x < X_i \le (x+n))e^{-\delta X_i}\right\} = \\ &= \frac{1}{N}\left(\int\limits_{0}^{\infty}I(x < X_i \le (x+n))e^{-2\delta X_i}dF(X_i) - J_n^2(x,\delta)\right) = \frac{1}{N}\left(J_n(x,2\delta) - J_n^2(x,\delta)\right). \end{split}$$

Теперь, учитывая что $\Phi_n(x,\delta) = e^{\delta x} J_n(x,\delta)$, найдем компоненты ковариационной матрицы трехмерной статистики t_N :

$$\begin{split} \sigma_{11} &= ND\{\Phi_{n,N}(x,\delta)\} = \Phi_n(x,2\delta) - \Phi_n^2(x,\delta); \quad \sigma_{22} = ND\{S_N(x)\} = S(x)(1-S(x)); \\ \sigma_{33} &= ND\{S_N(x+n)\} = S(x+n)(1-S(x+n)); \quad \sigma_{12} = \sigma_{21} = N \operatorname{cov}\left(S_N(x), \Phi_{n,N}(x,\delta)\right) = \\ &= N\left(E\{S_N(x) \cdot \Phi_{n,N}(x,\delta)\} - E\{S_N(x)\}E\{\Phi_{n,N}(x,\delta)\}\right) = (1-S(x))\Phi_n(x,\delta); \\ \sigma_{13} &= \sigma_{31} = N \operatorname{cov}\left(S_N(x+n), \Phi_{n,N}(x,\delta)\right) = \\ &= N\left(E\{S_N(x+n)\Phi_{n,N}(x,\delta)\} - E\{S_N(x+n)\}E\{\Phi_{n,N}(x,\delta)\}\right) = (1-S(x+n))\Phi_n(x,\delta); \\ \sigma_{23} &= \sigma_{32} = N \operatorname{cov}\left(S_N(x), S_N(x+n)\right) = (1-S(x))S(x+n). \end{split}$$

Используя предыдущий результат о смещении и найденную ковариационную матрицу, получаем СКО оценки:

$$u^{2}\left(\overline{a}_{x:\overline{n}|}^{N}\right) = E\left[\nabla H(t)(t_{N}-t)\right]^{2} + O\left(N^{-3/2}\right) = \frac{\sigma(\overline{a}_{x:\overline{n}|})}{N} + O\left(N^{-3/2}\right),$$

где

$$\sigma(\overline{a}_{x,\overline{n}|}) = \sum_{p=1}^{3} \sum_{j=1}^{3} H_{j}(t) \sigma_{jp} H_{p}(t) = H_{1}^{2}(t) \sigma_{11} + H_{2}^{2}(t) \sigma_{22} + H_{3}^{2}(t) \sigma_{33} + 2H_{1}(t) H_{2}(t) \sigma_{12} + H_{2}^{2}(t) H_{3}(t) \sigma_{13} + 2H_{2}(t) H_{3}(t) \sigma_{23} = \frac{\Phi(x, 2\delta)}{\delta S^{2}(x)} - \frac{\Phi^{2}(x, \delta)}{\delta S^{2}(x)} + \frac{\Phi(x, \delta) e^{-\delta n} S(x+n)}{\delta S^{2}(x)} - \frac{e^{-2\delta n} S^{2}(x+n)}{\delta S^{3}(x)} - \frac{2e^{-\delta n} S(x+n)}{\delta S^{2}(x)}.$$
(8)

Теорема доказана.

Для нахождения предельного распределения оценки (5) нам понадобятся две теоремы.

Теорема 2 (центральная предельная теорема в многомерном случае) [22]. Пусть $t_1, t_2, ..., t_N, ...$ – последовательность независимых одинаково распределенных *s*-мерных векторов, $E\{t_s\}=0$,

$$\sigma(x) = E\{t_s^{\mathrm{T}}t_s^{\mathrm{}}\}, \;\; S_N = \sum_{s=1}^N t_s^{\mathrm{}}$$
 . Тогда при $N \to \infty$

$$\frac{S_N}{\sqrt{N}} \Rightarrow N_s (0, \sigma(x)).$$

Теорема 3 (асимптотическая нормальность $H(t_N)$) [23]. Пусть:

1)
$$\sqrt{d_N} \cdot t_N \Rightarrow N_s \{\mu, \sigma(x)\};$$

2) функция H(z) дифференцируема в точке μ , $\nabla H(\mu) \neq 0$.

Тогда

$$\sqrt{d_N}(H(t_N) - H(\mu)) \Rightarrow N_1 \left\{ \sum_{j=1}^s H_j(\mu) \mu_j, \sum_{p=1}^s \sum_{j=1}^s H_j(\mu) \sigma_{jp} H_p(\mu) \right\}.$$

Теорема 4 (асимптотическая нормальность оценки (5)). В условиях теоремы 1

$$\sqrt{n}(\overline{a}_{x:\overline{n}|}^{N} - \overline{a}_{x:\overline{n}|}) \Rightarrow N_{1}(0,\sigma(\overline{a}_{x:\overline{n}|})).$$

Доказательство. Так как $t_N = \left(t_{1N}, t_{2N}, t_{3N}\right)^{\mathrm{T}} = \left(\Phi_{n,N}(x,\delta), S_N(x), S_N(x+n)\right)^{\mathrm{T}}$, то в обозначениях теоремы 3 имеем: $s=3, \ \sigma(x) = \sigma(\overline{a}_{x\overline{n}|})$. Таким образом,

$$\sqrt{N}\left\{\Phi_{n,N}(x,\delta) - \Phi_n(x,\delta), S_N(x) - S(x), S_N(x+n) - S(x+n)\right\} \Rightarrow N_3\left(0,\sigma(\overline{a}_{x\overline{n}|})\right),$$

где
$$\sigma(\overline{a}_{x:\overline{n}|}) = \begin{bmatrix} \sigma_{11}\sigma_{12}\sigma_{13} \\ \sigma_{21}\sigma_{22}\sigma_{23} \\ \sigma_{31}\sigma_{32}\sigma_{33} \end{bmatrix}.$$

Функция H(z) дифференцируема в точке t и $\nabla H(t) \neq 0$. Следовательно, выполнены все условия теоремы 3, и для оценки ренты (5) получаем

$$\sqrt{n}(\overline{a}_{x:\overline{n}|}^{N} - \overline{a}_{x:\overline{n}|}) \Rightarrow N_{1}(0,\sigma(\overline{a}_{x:\overline{n}|})).$$

Теорема доказана.

Заключение

В статье построены оценки *п*-летней ренты для смешанного страхования жизни, которое часто предлагается страховыми компаниями. Найдены главная часть асимптотической СКО и порядок смещения оценки ренты, доказывается ее асимптотическая нормальность. Рассмотренный подход к оцениванию индивидуальной смешанной ренты может быть распространен также и на смешанные ренты, связанные с коллективным страхованием [24, 25].

ЛИТЕРАТУРА

- 1. Bowers N., Gerber H., Hickman J., Jones D., Nesbitt C. Actuarial mathematics. Itasca: Society of Actuaries, 1986. 624 p.
- 2. Gerber H. Life insurance mathematics. 3rd ed. New York: Springer-Verlag, 1997. 118 p.
- 3. Фалин Г.И. Математические основы теории страхования жизни и пенсионных схем. М.: Анкил, 2002. 262 с.
- 4. Губина О.В., Кошкин Г.М. Оценивание современной стоимости непрерывной *п*-летней временной пожизненной ренты // Известия высших учебных заведений. Физика. 2015. Т. 58, № 11/2. С. 235–241.
- 5. Nadaraya E.A. Some new estimates of distribution function // Theory of Probability and its Applications. 1964. V. 9, No. 3. P. 497–500.
- Azzalini A. A note on the estimation of a distribution function and quantiles by a kernel method // Biometrika. 1981. V. 68, No. 1. P. 326–328.
- 7. Reiss R.-D. Nonparametric estimation of smooth distribution functions // Scand. J. Statist. 1981. V. 8. P. 116-119.
- 8. Falk M. Relative efficiency and deficiency of kernel type estimators of smooth distribution functions // Statist. Neerlandica. 1983. V. 37. P. 73–83.
- 9. Swanepoel J.W.H. Mean integrated squared error properties and optimal kernels when estimating a distribution function // Comm. Statist. Theory Methods. 1988. V. 17, No. 11. P. 3785–3799.
- 10. Jones M.C. The performance of kernel density functions in kernel distribution function estimation // Statist. Probab. Lett. 1990. V. 9. P. 129–132.
- 11. Shirahata S., Chu I.S. Integrated squared error of kernel-type estimator of distribution function // Ann. Inst. Statist. Math. 1992. V. 44, No. 3, P. 579–591.
- 12. Sarda P. Smoothing parameter selection for smooth distribution functions // J. Statist. Plann. Inference Inf. 1993. V. 35. P. 65–75.
- 13. Altman N., Leger C. Bandwidth selection for kernel distribution function estimation // J. Statist. Plann. Inference. 1995. V. 46. P. 195–214.
- 14. Bowman A., Hall P., Prvan T. Trust bandwidth selection for the smoothing of distribution functions // Biometrika. 1998. V. 85, No. 4. P. 799–808.
- 15. Chu I.S. Bootstrap smoothing parameter selection for distribution function estimation // Math. Japon. 1995. V. 41, No. 1. P. 189-197.
- 16. Shao Y., Xiang X. Some extensions of the asymptotics of a kernel estimator of a distribution function // Statist. Probab. Lett. 1997. V. 34. P. 301–308.
- 17. Una-Alvarez J., Gonzalez-Manteiga W., Cadarso-Suarez C. Kernel distribution function estimation under the Koziol-Green model // J. Statist. Plann. Inference. 2000. V. 87. P. 199–219.
- 18. Кошкин Г.М. Гладкое рекуррентное оценивание функции надежности // Известия высших учебных заведений. Физика. 2015. Т. 58, № 7. С. 128–134.
- 19. Fuks I., Koshkin G. Smooth Recurrent Estimation of Multivariate Reliability Function // Proc. The Int. Conference on Information and Digital Technologies 2015 (IDT 2015), 7–9 July 2015, Zilina, Slovakia. P. 84–89.
- 20. Кошкин Г.М. Моменты отклонений оценки подстановки и ее кусочно-гладких аппроксимаций // Сибирский математический журнал. 1999. Т. 40, № 3. С. 604–618.
- 21. Ибрагимов И.А., Хасьминский Р.З. Асимптотическая теория оценивания. М.: Наука, 1979. 528 с.
- 22. Боровков А.А. Теория вероятностей. М.: Наука, 1986. 432 с.
- 23. Кошкин Г.М. Асимптотические свойства функций от статистик и их применения к непараметрическому оцениванию // Автоматика и телемеханика. 1990. № 3. С. 82–97.
- 24. Губина О.В., Кошкин Г.М. Оценивание коллективной ренты статуса совместной жизни // Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2016. № 2 (35). С. 30–36.
- 25. Кошкин Г.М., Губина О.В. Оценивание коллективной ренты статуса выживания последнего // Известия высших учебных заведений. Физика. 2016. Т. 59, № 8/2. С. 57–60.

Поступила в редакцию 10 июля 2019 г.

Gubina O.V., Koshkin G.M. (2020) ESTIMATION OF PRESENT VALUE OF n-YEAR LIFE ANNUITY FOR ENDOWMENT INSURANCE. *Vestnik Tomskogo Gosudarstvennogo Universiteta*. *Upravlenie, vychislitelnaja tehnica i informatika* [Tomsk State University Journal of Control and Computer Science]. 50, pp. 39–46

DOI: 10.17223/19988605/50/5

In the paper, we constructed a nonparametric estimator of the n-year life annuity for the endowment insurance and studied its asymptotic properties.

For the n-year endowment life insurance, define the present value of the insurance payment z as follows:

$$z = \begin{cases} e^{-\delta T(x)}, & T(x) \le n, \\ e^{-\delta n}, & T(x) > n, \end{cases}$$

where δ denotes a force of interest, x is the age of an individual, X is its lifetime, T(x) = X - x is its future lifetime. In this case, the value of z shows the present share of future insurance payments taken as some unit. The longer the insurance period the lower the payment of the insured using bank interest rates.

As the *n*-year life annuity for the endowment insurance we take

$$\overline{a}_{x:\overline{n}|} = \frac{1}{\delta} \left(1 - \frac{e^{\delta x}}{S(x)} \int_{x}^{x+n} e^{-\delta t} dF(t) - \frac{e^{-\delta n} S(x+n)}{S(x)} \right),$$

where S(x) = P(X > x) is a survival function, $F(x) = P(X \le x) = 1 - S(x)$ is a distribution function.

Assume that we have a random sample $X_1, ..., X_N$ of N individuals' lifetimes. Using the empirical survival function $S_N(x) = \frac{1}{N} \sum_{i=1}^N \mathrm{I}(X_i > x)$, where I(A) is the indicator of an event A, obtain the following estimator of $\overline{a}_{x,\overline{n}|}$:

$$\overline{a}_{x,\overline{n}|}^{N} = \frac{1}{\delta} \left(1 + \frac{e^{\delta x} \sum_{i=1}^{N} \exp(-\delta X_{i}) I\left(x < X_{i} \le (x+n)\right) + e^{-\delta n} \sum_{i=1}^{N} I\left(X_{i} > (x+n)\right)}{\sum_{i=1}^{N} I\left(X_{i} > x\right)} \right).$$

We found the principal term of the asymptotic mean square error of this estimator and proved its asymptotic normality.

Keywords: endowment life insurance; n-year life annuity; nonparametric estimation; mean squared error; asymptotic normality.

GUBINA Oxana Viktorovna (Post-graduate Student, National Research Tomsk State University, Tomsk, Russian Federation). E-mail: gov7@mail.ru

KOSHKIN Gennady Mikhailovich (Doctor of Physics and Mathematics, Professor, National Research Tomsk State University, Tomsk, Russian Federation).

E-mail: kgm@mail.tsu.ru

REFERENCES

- 1. Bowers, N., Gerber, H., Hickman, J., Jones, D. & Nesbitt, C. (1986) Actuarial mathematics. Itasca: Society of Actuaries.
- 2. Gerber, H. (1997) Life insurance mathematics. 3rd ed. New York: Springer-Verlag.
- 3. Falin, G.I. (2002) *Matematicheskie osnovy teorii strakhovaniya zhizni i pensionnykh skhem* [Mathematical Foundations of the Theory of Life Insurance and Pension Schemes]. Moscow: Ankil.
- 4. Gubina, O.V. & Koshkin, G.M. (2015) Estimation of the actuarial present value of the continuous *n*-year time life annuity. *Russian Physics Journal*. 58(11/2). pp. 235–241. DOI: 10.17223/19988605/30/5
- Nadaraya, E.A. (1964) Some new estimates of distribution function. Theory of Probability and its Applications. 9(3). pp. 497–500.
 DOI: 10.1137/1109069
- 6. Azzalini, A. (1981) A note on the estimation of a distribution function and quantiles by a kernel method. *Biometrika*. 68(1). pp. 326–328. DOI: 10.1093/biomet/68.1.326
- 7. Reiss, R.-D. (1981) Nonparametric estimation of smooth distribution functions. *Scandinavian Journal of Statistics*. 8. pp.116–119. DOI: 10.1007/BF02613619
- 8. Falk, M. (1983) Relative efficiency and deficiency of kernel type estimators of smooth distribution functions. *Statist. Neerlandica*. 37. pp. 73-83. DOI: 10.1111/j.1467-9574.1983.tb00802.x
- 9. Swanepoel, J.W.H. (1988) Mean integrated squared error properties and optimal kernels when estimating a distribution function. *Comm. Statist. Theory Methods.* 17(11). pp. 3785–3799. DOI: 10.1080/03610928808829835
- 10. Jones, M.C. (1990) The performance of kernel density functions in kernel distribution function estimation. *Statistics and Probability Letters*. 9. pp. 129–132. DOI: 10.1016/0167-7152(92)90006-Q
- 11. Shirahata, S. & Chu, I.S. (1992) Integrated squared error of kernel-type estimator of distribution function. *Annual Inst. Statist. Math.* 44(3). pp. 579–591. DOI: 10.1007/BF00050707
- 12. Sarda, P. (1993) Smoothing parameter selection for smooth distribution functions. *Journal of Statistical Planning and Inference*. 35. pp. 65–75. DOI: 10.1016/0378-3758(93)90068-H
- 13. Altman, N. & Leger, C. (1995) Bandwidth selection for kernel distribution function estimation. *Journal of Statistical Planning and Inference*. 46. pp. 195–214. DOI: 10.1016/0378-3758(94)00102-2
- 14. Bowman, A., Hall, P. & Prvan, T. (1998) Trust bandwidth selection for the smoothing of distribution functions. *Biometrika*. 85(4). pp. 799–808. DOI: 10.1093/biomet/85.4.799
- 15. Chu, I.S. (1995) Bootstrap smoothing parameter selection for distribution function estimation. *Math. Japon.* 41(1). pp 189–197.
- Shao, Y. & Xiang, X. (1997) Some extensions of the asymptotics of a kernel estimator of a distribution function. Statistics and Probability Letters. 34. pp. 301–308. DOI: 10.1016/S0167-7152(96)00194-0
- Una-Alvarez, J., Gonzalez-Manteiga, W. & Cadarso-Suarez, C. (2000) Kernel distribution function estimation under the Koziol-Green model. *Journal of Statistical Planning and Inference*. 87. pp. 199–219.
- 18. Koshkin, G.M. (2015) Smooth Recurrent Estimators of the Reliability Functions. *Russian Physics Journal*. 58(7). pp. 1018–1025. DOI: 10.1007/s11182-015-0603-9

- Fuks, I. & Koshkin, G. (2015) Smooth recurrent estimation of multivariate reliability function. *Proc. of the Int. Conference on Information and Digital Technologies* 2015. IDT 2015. Zilina, Slovakia. July 7–9, 2015. pp. 84–89. DOI 10.1109/DT.2015.7222955
- 20. Koshkin, G.M. (1999) Deviation moments of the substitution estimator and its piecewise smooth approximations. *Sibirskiy matematicheskiy zhurnal Siberian Mathematical Journal*. 40(3). pp. 515–527. DOI: 10.1007/BF02679759
- 21. Ibragimov, I.A. & Hasminskii, R.Z. (1981) Statistical Estimation: Asymptotic Theory. Berlin; New York: Springer.
- 22. Borovkov, A.A. (1986) Probability Theory. Moscow: Nauka.
- 23. Koshkin, G.M. (1990) Asymptotic properties of functions of statistics and their application to nonparametric estimation. *Automation and Remote Control*. 51(3). pp. 345–357.
- 24. Gubina, O.V. & Koshkin, G.M. (2016) Collective annuity estimation of joint-life status. *Vestnik Tomskogo gosudarstvennogo universiteta*. *Upravlenie, vychislitel'naya tekhnika i informatika Tomsk State University Journal of Control and Computer Science*. 2(35). pp. 30–36. DOI: 10.17223/19988605/35/3
- 25. Koshkin, G.M. & Gubina, O.V. (2016) Estimation of collective annuity of the last-survivor status. *Russian Physics Journal*. 59(8/2). pp. 57–60.

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

УДК 517.977

DOI: 10.17223/19988605/50/6

A.I. Rouban, A.S. Mikhalev

THE GLOBAL OPTIMIZATION METHOD WITH SELECTIVE AVERAGING OF THE DISCRETE DECISION VARIABLES

In the paper, the functional of selective averaging of discrete decision variables is proposed. The positive selectivity coefficient is entered into a positive decreasing kernel of functional and with growth of selectivity coefficient the mean gives optimum values (in a limit) of decision discrete variables in a problem of global optimization. Based on the estimate of the selective averaging functional, a basic global optimization algorithm is synthesized on a set of discrete variables with ordered possible values under inequality constraints. The basis is a computational scheme for optimizing continuous variables and its transformation for optimization with respect to discrete variables. On a test example the high convergence rate and a noise stability of base algorithm are shown. Simulations have shown that the estimate of the probability of making a true decision reaches unit.

Keywords: global optimization; discrete variable; selective averaging of decision variables; multiextreme function; constraints of inequality type.

The problem of search of a global extremum of objective functions on admissible set of decision variables (continuous, discrete or continuous-discrete) belongs to the very complex class [1–27]. The specificity of global optimization is caused by multiextremal view of the objective functions, their discontinuous nature, very different sensitivity of objective functions in relation to decision variables, the discrete nature of part or all the variables, the presence of noises, the presence of a considerable quantity of decision variables and constraints of inequalities type and equalities type. If objective functions are not calculated, they are measured in points of admissible set of decision variables.

The majority of algorithms is oriented only on cases of continuous decision variables and on elemental constraints of inequalities type. The comparative analysis of base algorithm of a method with selective averaging of continuous decision variables [6, 7, 10, 14, 24] in relation to existing heuristic algorithms showed its advantages on the rate of convergence, a noise stability and total computing complexity. In [27] the variant of base algorithm of a solution of a problem of global optimization on set of continuous-discrete variables is proposed. It is rational to develop this effective method of global optimization on the set of discrete variables.

In this paper for a solution of global optimization problems on a set of discrete variables the approach based on the selective averaging of decision variables with adaptive reorganization of an admissible set of trial movements is proposed.

It is constructed the functional of selective averaging of discrete decision variables. The selectivity coefficient is entered into kernel of functional. With an increase in the core selectivity coefficient, it becomes possible for the functional to distinguish the positions of the global minimum. It is shown that when selectivity coefficient of kernel tends to infinity the averaging gives optimal value of decision discrete variables.

The computing scheme of base algorithm of global optimization at continuous variables is transformed to the similar scheme at discrete variables. The base algorithm of global optimization on a set of discrete variables with ordered possible values at the presence of inequalities constraints is synthesized.

Trial and working steps are also separated in time. Before performance of each working step the series of calculations of the minimized function in the sampling points is carried out. Based on this information at the fixed selectivity coefficient of kernel, the selective averaging of decision variables will be executed numerically. For each discrete variable the continuous auxiliary non-dimension variable which contains

numbers of possible values of a discrete variable and the identical the adjoining each other subintervals covering these numbers is put in compliance. Due to one-to-one compliance the transition in the sampling points from continuous variables to discrete possible values is carried out. The same (but already reverse) transition occurs for received averaged values of decision variables.

Also, the adaptive reorganization of sizes of a set of possible trial movements is carried out and functions of inequalities constraints are considered.

The example shows the high convergence rate and noise stability of the base algorithm. It also provides near-to-one the estimate of the probability of obtaining a true solution.

1. Statement of problem

The problem of search of a global minimum of objective function f(y) on a set of discrete variables with ordered possible values at the presence of inequalities constraints is solving:

$$f(y) = \text{globmin}, \varphi_i(y) \le 0, j = \overline{1,m},$$
 (1)

where $y = (y_1, ..., y_h)$ is vector h of discrete variables. Each discrete variable y_t has r_t possible ordered values $y_{t,1}, ..., y_{t,r_t}$.

Inequalities constraints select (narrow) admissible set of possible values in which search of minimum of global minimum is carried out. It's required to define the position y_{\min} of global minimum of objective function f(y) on limited set of change of its variables.

Function f(y) is multiextreme and can be distorted by noises. Functions of constraints can be non-convex. Search of extremum is carried out on basis only measurements or calculations of specified functions f(y), $\varphi_j(y)$, $j = \overline{1,m}$ in the selected sampling points which satisfy to the inequalities constraints: $\varphi_j(y) \le 0$, $j = \overline{1,m}$.

We assume that a global minimum of function f(y) on an admissible set of points with discrete values of variables is unique.

2. The selective averaging of discrete decision variables

The selective averaging of decision continuous variables [6, 7, 10, 14, 24] is a mathematical expectation with special probability density function of these variables. Probability density function at a decreasing kernel (with increase of its normalized argument from 0 to 1) allows to approach with growth of selectivity coefficient of a kernel to the specified average value of decision variables i.e. to the true position of global minimum. This theoretical result gave to chance of construct the structure of a base numerical algorithm, which successfully applied at the presence of inequalities constraints. Due to the expansion of possibilities of this algorithm, the algorithms of single-objective global optimization at the constraints such as inequalities and equalities are synthesized. The algorithms of the solution of other extreme problems are obtained: multi-objective and minimax global optimization, search of the main minima of multiextremal functions [14].

Let's transfer idea of selective averaging [14] on solution the problem of global optimization (1) on set of possible values of the discrete decision variables which satisfy the inequalities constraints (1).

At first we will consider the one-dimensional version of optimization problem (1). Discrete variable y has r possible ordered values $(y_1, ..., y_r)$.

We enter the notation: f_{\min} , f_{\max} are smallest and largest values of minimized function on an admissible set of possible values; y_{\min} is admissible value of discrete variable at which function f(y) is reaching the minimum value: $f(y_{\min}) = f_{\min}$.

The averaged value (mathematical expectation) of decision variable is equal to value:

$$\overline{[y]}_s = \sum_{\mu=1}^r \vec{p}_s(g_\mu) y_\mu , \qquad (2)$$

where $\vec{p}_s(g_{\mu}) = p_s(g_{\mu}) / \sum_{k=1}^r p_s(g_k)$ is normalized (on 1) positive decreasing kernel (analog of probability

of a random event): $\sum_{\mu=1}^{r} \overline{p}_{s}(g_{\mu}) = 1$; kernels $p_{s}(g_{\mu})$ and $\overline{p}_{s}(g_{\mu})$ lie in interval [0; 1]; s is selectivity coeffi-

cient of kernel $(1 \le s)$, at the decision it's enough to set of it from integer sequence: 1, 2, 3, ...; $g_{\mu} = (f(y_{\mu}) - f_{\min})/(f_{\max} - f_{\min})$ is arguments of kernels which also lie in interval [0; 1]: $0 \le g_{\mu} \le 1$, $\mu = \overline{1, r}$.

In optimal point $f(y_{\min}) = f_{\min}$ and $g_{\min} = 0$, in point of the maximum value $f(y_{\max}) = f_{\max}$ and $g_{\max} = 1$. The presence in the argument of kernel of values f_{\min} and f_{\max} allows cover of all range of change of the optimized function and to pass to non-dimension variable. Due to of it the subsequent numerical algorithms become more universal. They independent of the units of measure of minimized function. And also the rate of convergence of algorithms and accuracy tracking of position of extremum increases. The absence of information about f_{\min} and f_{\max} is compensated by their estimates which calculated on each working step based on measurements (calculations) of optimized function in trial points.

Let's stop on a type of kernel $p_s(g)$ because normalized kernel $\vec{p}_s(g)$ repeats its form. It's convenient to represent the $p_s(g)$ in the form of raising to the degree s rather simple decreasing kernel: $p_s(g) = [p(g)]^s$. Possible degree kernels p(g): linear p(g) = 1 - g, parabolic $p(g) = 1 - g^2$, cubic $p(g) = 1 - g^3$ and etc. The example of other type of kernels is an exponential kernel $p(g) = e^{-g}$, a hyperbolic kernel $p(g) = g^{-1}$. Further we will consider only degree kernels.

With growth of selectivity coefficient s all components of discrete function $p_s(g_l)$, $l=\overline{1,r}$, approaches to zero, except $p_s(g_{\min}=0)=1$ and $p_s(g_{\max}=1)=0$. Then normalized discrete function tends to Kroneker's function [28]. It's equal to 1 in a point $g_{\min}=0$ and equal to 0 in other points. At a result in the right part of procedure of averaging (2) Kroneker's function «prick out» optimal value of decision variable y_{\min} , i.e. $\overline{[y]}_s \xrightarrow[s \to \infty]{} y_{\min}$.

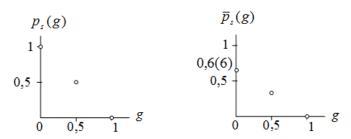


Fig. 1. The values of linear kernel $p_s(g)$ and normalized kernel $\overline{p}_s(g)$ at three possible values of a discrete variable and selectivity coefficient s=1

On Fig. 1 the simplest case of the linear kernel $p_s(g)$ and its normalized analog $\overline{p}_s(g)$ when discrete variable has three possible values: y_{\min}, y_2, y_{\max} . Let's say in the point y_2 the value g is equal to 0.5. With growth of selectivity coefficient s the kernel $p_s(0.5) = 1/2^s$ and normalized kernel $\overline{p}_s(0.5) = (1/2^s)/(1+(1/2^s)) = 1/(1+2^s)$ tends to zero. Respectively the $\overline{p}_s(0) = 1/(1+(1/2^s))$ tends to 1 and function $\overline{p}_s(g)$ tends to Kroneker's function with the special point $g_{\min} = 0$. With increase the number of possible values of discrete variable y, all noted limit regularities are preserved.

At increase the number of discrete variables $y = (y_1, ..., y_h)$ (see (1)) all noted properties also take place. Each discrete variable y_t has the r_t possible ordered values $y_{t,1}, ..., y_{t,r_t}$. All admissible points have the h coordinates and satisfy to inequalities constraints. One of these points $y_{\min} = (y_{1,\min}, ..., y_{h,\min})$ corresponds to a global minimum $f(y_{\min}) = f_{\min}$ and $g_{\min} = 0$.

For any single-valued function $\phi(y)$ the selective averaging has the form:

$$\overline{[\phi(y)]}_s = \sum_{\mu=1}^R \overline{p}_s(g_\mu) \phi(y_\mu).$$

Here: y_{μ} is admissible point with one of combinations of possible values of its coordinates, μ is number of this point, $\phi(y) \equiv \phi(y_1, ..., y_h)$, $\phi(y_{\mu}) \equiv \phi(y_{1,\mu}, ..., y_{h,\mu})$; R is number of admissible points: (combinations of possible values of discrete variables) and each point which satisfy to inequalities constraints.

Limit result

$$\overline{[\phi(y)]}_s \xrightarrow[s \to \infty]{} \phi(y_{\min}) \equiv \phi(y_{1,\min}, ..., y_{h,\min})$$

is the same as for one-dimensional case.

We take the h functions $\phi_t(y) = y_t$, $t = \overline{1, h}$ and we receive limit values for all coordinates

$$\overline{[y_t]}_s \xrightarrow[s \to \infty]{} y_{t,\min}, t = \overline{1, h},$$

where $[\overline{y_t}]_s = \sum_{\mu=1}^R \overline{p}_s(g_{\mu}) y_{t,\mu}$, $t = \overline{1, h}$, $y_{t,\mu}$ is coordinate with number t in point with number μ . This is coordinate-wise averaging. It's written in vector form:

$$\overline{[y]}_s \xrightarrow[s \to \infty]{} y_{\min} \equiv (y_{1,\min}, ..., y_{h,\min}),$$

where $\overline{[y]}_s = \sum_{\mu=1}^R \vec{p}_s(g_\mu) y_\mu$, y_μ is admissible point with number μ at one of combinations of possible values of its variables.

By search of a maximum of function f the arguments of kernels are calculated on another: $g_{\mu} = (f_{\text{max}} - f(y_{\mu}))/(f_{\text{max}} - f_{\text{min}})$ and the point y_{max} corresponds to global maximum value: $f(y_{\text{max}}) = f_{\text{max}}$ and $g_{\text{max}} = 0$.

3. Basic algorithm of optimization

The method is based on separation in time of trial and working steps, uniform distribution of sampling points on admissible set of possible values of discrete decision variables, numerical selective averaging (calculation of mathematical expectation) of decision variables by results of calculated (or measured) values of optimized function in sampling points. Adaptive reorganization of the sizes of the set of sampling points is carried out also on each working step.

The basic algorithm of global optimization based on selective averaging of decision continuous variables allows implementation on a set of discrete and continuous-discrete variables. In [27] the specified algorithm was generalized on a case of continuous and discrete variables with ordered possible values. Based on the scheme presented in [27], in this work the algorithm for a case only of discrete variables is constructed.

The solution of problem of optimization with discrete variables is based on transition from each discrete variable y_t to the corresponding auxiliary continuous variable \overline{x}_t . From possible values $y_{t,1}, \ldots, y_{t,r_t}$ of a discrete variable y_t transition to their numbers are carries out. Calculations on each working step are conducted for the number $\overline{x}_{t,N}^{(i)}$ possible values of discrete variables y_t . Averaged values of variables (estimates of mathematical expectation) and the sizes (also averaged) the variation sets of variables are calculated.

For receiving of n sampling points $y^{(i)}$, $i = \overline{1, n}$ consistently by uniform distribution are generated on a «rectangular» set of a points and from them only those which satisfy inequalities constraints are left.

After the l-th working step the generation of sampling points is carried out equally:

$$\overline{x}_{t}^{(i)} = \overline{\overline{x}}_{t}^{l} + \overline{\Delta x}_{t}^{l} u_{\overline{x}_{t}}^{(i)}, \ u_{\overline{x}_{t}}^{(i)} \in [-1; 1], t = \overline{1, h}, \ i = \overline{1, h},$$
(3)

In a formula (3) $u_{\bar{x}_i}^{(i)}$ is elements of pseudo-random sequence of uniform distributed of the continuous random value. This uniform law of distribution causes identical probabilities of emergence of possible values (and their numbers) all discrete variables.

In the received sampling points value of minimized function $f^{(i)} \equiv f(y^{(i)}), i = \overline{1,n}$ is calculated. Further, the position of the minimum is specified.

New value \bar{x}^{l+1} on auxiliary variables and the sizes $\Delta \bar{x}^{l+1}$ of «rectangular» set of trial movements are calculated by the following formulas:

$$\overline{x}_{t}^{l+1} = \sum_{i=1}^{n} \overline{x}_{t,N}^{(i)} \overline{p}_{s,\min}^{(i)}, t = \overline{1,h},$$

$$\overline{p}_{s,\min}^{(i)} = \frac{p_{s}(g_{\min}^{(i)})}{\sum_{j=1}^{n} p_{s}(g_{\min}^{(j)})}, g_{\min}^{(i)} = \frac{f^{(i)} - \widehat{f}_{\min}}{\widehat{f}_{\max} - \widehat{f}_{\min}}, \Delta \overline{x}_{t}^{l+1} = \gamma_{q} \left(\sum_{i=1}^{n} |\overline{x}_{t,N}^{(i)} - \overline{\overline{x}}_{t}^{l}|^{q} \overline{p}_{s,\min}^{(i)}\right)^{1/q}, t = \overline{1,h},$$
(4)

 $l = 0, 1, 2, ...; [0 < \gamma_q, q \in \{1, 2, ...\}, 0 < s] \,.$

Here: $\widehat{f}_{\max} = \max\{f^{(i)}, i = \overline{1,n}\}, \ \widehat{f}_{\min} = \min\{f^{(i)}, i = \overline{1,n}\}, \ p_s(\cdot)$ is positive kernel, s is selectivity coefficient of kernel. The positive kernels $\overline{p}_{s,\min}^{(i)}$ normalized on 1 on a system of n sampling points:

 $\sum_{i=1}^{n} \overline{p}_{s,\min}^{(i)} = 1.$ The argument of kernel is non-dimension variable, which always lie in interval [0; 1]. The

kernels $p_s(\cdot)$ monotonically decrease at argument increase.

On l-th step the value \bar{x}_t^l calculated in accordance with formula (4) gets to the individual interval which covering some number of the corresponding value of a discrete variable. So the possible values of all discrete variables are calculated. Calculated (on the previous working step) on a formula (4) interval of a variation $2\Delta \bar{x}_t^l$ of a continuous auxiliary variable allocates numbers of possible values of the discrete variable. The intervals of unit length covering these numbers form a new interval of change of an auxiliary variable. For any component \bar{x}_t the initial values of these variables is equal: $\bar{x}_t^0 = (r_t + 1)/2$, $\Delta \bar{x}_t^0 = r_t/2$.

When approaching to minimum the region of trial movement reduced, and thus, there is more exact tracking of position of extremum. The criterion of stop of search process is the condition of reduction (at some *l*) of size of region of variation of variables to the given value:

$$\max \left\{ \frac{\overline{\Delta x}_{t}^{l}}{\overline{\Delta x}_{t}^{0}}, t = \overline{1, h} \right\} \leq \varepsilon_{2}. \tag{5}$$

The corrected value \overline{x}_t^l in formula (4) directly is not used but $\overline{\Delta x}_t^l$ gives a variation interval $2\overline{\Delta x}_t^l$ for auxiliary coordinate \overline{x}_t (at the subsequent formation of sampling points) and is used in the condition break of search process (5).

4. Numeric example

Let's consider test function with 16 minima which constructed at the expense of operations «min» to 16 degree one-extreme potential functions:

$$f(y_1, y_2) = \min\{z_i(y_1, y_2), i = \overline{1,16}\};$$

$$z_1(y_1, y_2) = 2|y_1 - 9|^2 + 2|y_2 - 9|^2; z_2(y_1, y_2) = 4|y_1 - 9|^{1.5} + 4|y_2 + 1|^{1.8} + 7;$$

$$z_3(y_1, y_2) = 4|y_1 - 6|^{0.8} + 4|y_2 - 5|^{1.6} + 4; z_4(y_1, y_2) = 3|y_1|^{1.1} + 3|y_2 - 2|^{1.8} + 16;$$

$$z_5(y_1, y_2) = 6|y_1 + 4| + 6|y_2 - 7| + 5; z_6(y_1, y_2) = 4|y_1 + 8|^{1.5} + 4|y_2 - 13|^{1.6} + 10;$$

$$z_7(y_1, y_2) = 2|y_1 - 3|^{1.5} + 2|y_2 - 11|^{1.5} + 9; z_8(y_1, y_2) = 4|y_1 - 11|^{0.8} + 4|y_2 - 2|^{0.9} + 8.5;$$

$$z_9(y_1, y_2) = 4|y_1 + 8|^{0.8} + 4|y_2 + 1|^{0.8} + 14; z_{10}(y_1, y_2) = 3|y_1 - 13|^{1.8} + 3|y_2 - 12|^{1.6} + 13;$$

$$z_{11}(y_1, y_2) = 3|y_1 + 13|^{1.3} + 3|y_2 + 4|^{1.3} + 12; z_{12}(y_1, y_2) = 5|y_1 - 6|^{0.8} + 5|y_2 + 1|^{0.6} + 15;$$

$$z_{13}(y_1, y_2) = 5|y_1 + 13|^{1.6} + 5|y_2 - 9|^{1.9} + 8; z_{14}(y_1, y_2) = 6|y_1 - 9|^{0.6} + 6|y_2 + 8|^{0.6} + 18;$$

$$z_{15}(y_1, y_2) = 5|y_1 - 3|^{1.1} + 5|y_2 + 4|^{1.3} + 6; z_{16}(y_1, y_2) = 5|y_1 - 3|^{1.6} + 5|y_2 + 13|^{1.6} + 10.5.$$

The admissible region has the form of square which defined by two constraints-inequalities:

$$\varphi_1(y_1) \equiv |y_1| - 15 \le 0$$
; $\varphi_2(y_2) \equiv |y_2| - 15 \le 0$.

The global minimum thus corresponds to a point $y^* = (9, 9)$ and has a $f(y^*) = 0$.

We enter the additional constraints, one of which cuts the specified minimum:

$$\varphi_3(y) \equiv y_1 + y_2 - 12 \le 0$$
; $\varphi_4(y) \equiv -y_1 - y_2 - 10 \le 0$.

The conditional global extremum is in point $y^* = (6, 5)$, $f(y^*) = 4$.

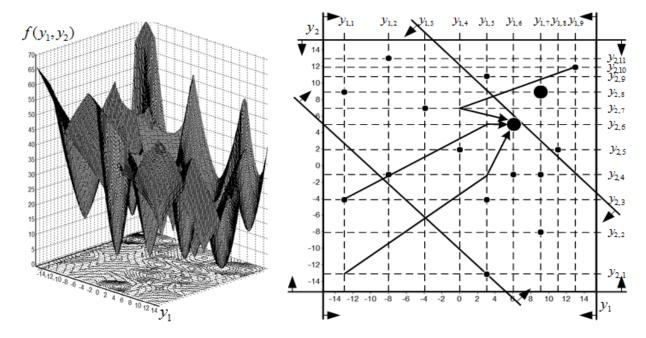


Fig. 2. Multiextremal function f(y): at the left is a perspective view of function; on the right is the position of minima, admissible region of search and movement trajectory in a neighborhood of global minimum at three starting points (-14, -14), (-14, -4), (13, 12)

On Fig. 2 it is shown the perspective view of minimized function, the positions of minima of given function and the constraints which select admissible region. Both figures are constructed in the assumption that the discrete variables are continuous, and then the possible values of the discrete variables are specified. In fact we have 9 cross sections of the presented function at the specified discrete values of the first variable y_1 and the 11 cross sections for the second variable y_2 . Each point in Fig. 2 on the right corresponds to the position of the minimum (small points is local minima, large points is global minimum) of a given objective function, the intersection of dotted lines is admissible points y_μ with one of combinations of possible values of its coordinates.

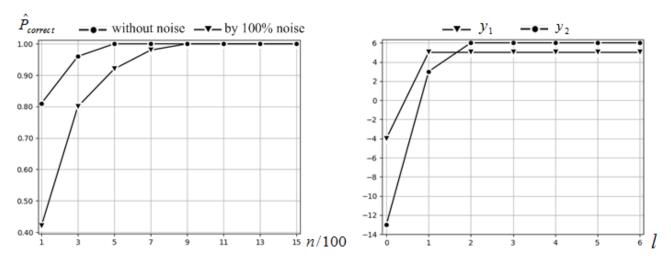


Fig. 3. The left is the dependence of estimate of probability of hitting to neighborhood of true decision from a number of sampling point; the right is the change of variables on iterations at 100% noise

Let's show operability of the offered algorithm. The complex characteristic of process of search of global minimum is the estimate $\widehat{P}_{\text{correct}}$ of probability of hitting of the obtained decision y^* in the neighborhood of true decision y^{**} : $|y^*_{t,N^*} - y^{**}_{t,N^{**}}| = 0$, $t = \overline{1,h}$, i.e. calculated values of discrete variables precisely coincide with the true values in the point of minimum.

Let's consider dependence of received of true decision from a number of sampling point n. For this purpose the M realization of process of search of minimum carried out. The estimate of probability $\widehat{P}_{\text{correct}}$ will be equal to the relative frequency of hitting of the obtained decision y_j^* , $j = \overline{1, M}$ in neighborhood of true decision.

The parameters of the algorithm of minimization: $(y^0, \overline{\overline{x}}^0) = (-13, -4; 1, 3)$, $\overline{\Delta x}^0 = (8, 5; 8, 5)$, $\gamma_q = 1$ (is equal 2 for case by 100% noise), q = 2, kernel on the minimized function parabolic with degree of selectivity s = 300 (is equal 1000 for the case by 100% noise). Research parameters: number of realization M = 101.

On Fig. 3 (the left) the dependence P_{correct} from a number of sampling point is presented. On Fig. 3 (the right) given a typical implementation of reorganization from step to step of discrete variables when search global minimum by a number of sampling points n = 500 (for the case without noise are required 3–5 iterations, for the case with 100% noise are required 7–11 iterations).

Conclusion

The developed method of selective averaging of decision variables and respectively base numerical algorithm of a global optimization has a rather simple structure. Due to the selective averaging of decision variables, the base algorithm has high noise stability. Use in an argument of a kernel of base algorithm of non-dimension (relative) values an essentially increase rate of convergence of algorithm, and reduces number of adjustable parameters.

The dimension of a problem of optimization at the chosen approach doesn't change. The complexity of calculations in comparison with a problem with continuous-discrete variables increases slightly. All main properties of an algorithm are also preserved.

The algorithm is based on execution the same type operations, which can be executed in parallel on multiprocessing computing systems. This feature is important in the presence of a large number of decision variables and different constraints.

REFERENCES

- 1. Schmit, L.A. & Farshi, B. (1974) Some approximation concepts for structural synthesis. AIAA J. 12(5). pp. 692–699.
- 2. Gupta, O.K. & Ravindran, A. (1983) Nonlinear integer programming and discrete optimization. *Journal of Mechanisms, Transmissions, and Automation in Design.* 105(2), pp. 160–164.
- 3. Olsen, G. & Vanderplaats, G.N. (1989) A method for nonlinear optimization with discrete variables. *AIAA J.* 27(11). pp. 1584–1589.
- 4. Rajeev, S. & Krishnamoorthy, C.S. (1992) Discrete optimization of structures using genetic algorithms. *Journal of Structural Engineering*, 118(5), pp. 1233–1250. DOI: 10.1061/(ASCE)0733-9445(1992)118:5(1233)
- 5. Cohn, M.Z. & Dinovitzer, A.S. (1994) Application of structural optimization. *Journal of Structural Engineering*. 120(2). pp. 617-650. DOI: 10.1061/(ASCE)0733-9445(1994)120:2(617)
- 6. Ruban, A.I. (1994) The nonparametric search global optimization method. *Cybernetics and Higher Educational Institution*. vol. 28 (Intellectual information technology). Tomsk: Tomsk Polytechnic University. pp. 107–114.
- 7. Ruban, A.I. (1995) The nonparametric search optimization method. Russian Physics Journal. 38(9). pp. 65-73.
- 8. Salajegheh, E. (1996) Discrete variable optimization of plate structures using dual methods. *Computers & Structures*. 58(6). pp. 1131–1138.
- 9. Gutkowski, W. (1997) Discrete structural optimization: design problems and exact solution methods. Discrete structural optimization. Vienna: Springer. pp. 1–53.
- 10. Ruban, A.I. (1997) Global extremum of continuous functions. Computer Science and Control Systems. 2. pp. 3-11.
- 11. Beckers, M. (2000) Dual methods for discrete structural optimization problems. *International Journal for Numerical Methods in Engineering*. 48. pp. 1761–1784. DOI: 10.1002/1097-0207(20000830)48:12<1761::AID-NME963>3.0.CO;2-R
- 12. Pezeshk, S., Camp, C. & Chen, D. (2000) Design of nonlinear framed structures using genetic optimization. *Journal of Structural Engineering*. 126(3). pp. 382–388. DOI: 10.1061/(ASCE)0733-9445(2000)126:3(382)
- 13. Salajegheh, J. (2001) Continuous and discrete optimization of space structures using approximation concepts. Ph.D. Thesis. University of Kerman. Iran. (In Persian).
- 14. Ruban, A.I. (2004) *Global'naya optimizatsiya metodom usredneniya koordinat* [Global optimization by a method of averaging of coordinates]. Krasnoyarsk: State Technical University.
- 15. Camp, C.V., Bichon, B.J. & Stovall, S.P. (2005) Design of steel frames using ant colony optimization. *Journal of Structural Engineering*. 131(3). pp. 369–379. DOI: 10.1061/(ASCE)0733-9445(2005)131:3(369)
- 16. Kaveh, A. & Talatahari, S. (2008) A hybrid particle swarm and ant colony optimization for design of truss structures. *Asian Journal of Civil Engineering*. 9(4). pp. 329–348.
- 17. Floudas, C.A. & Pardalos, P.M. (2009) Encyclopedia of Optimization. 2th ed. Boston, MA: Springer.
- 18. Rao, S.S. (2009) Engineering optimization: theory and practice. 4th ed. John Wiley & Sons.
- 19. Spillers, W.R. & MacBain, K.M. (2009) Structural Optimization. Boston, MA: Springer.
- 20. Kaveh, A. & Talatahari, S. (2010) An improved ant colony optimization for the design of planar steel frames. *Engineering Structures*. 32(3). pp. 864–873. DOI: 10.1016/j.engstruct.2009.12.012
- 21. Kripakaran, P., Brian, H. & Abhinav, G. (2010) A genetic algorithm for design of moment-resisting steel frames. *Structural Multidisciplinary Optimization*. 32(3). pp. 559–574. DOI: 10.1007/s00158-011-0654-7
- 22. Gandomi, A.H., Yang, X.S. & Alavi, A.H. (2011) Mixed variable structural optimization using firefly algorithm. *Computers & Structures*. 89(23). pp. 2325–2336. DOI: 10.1016/j.compstruc.2011.08.002
- 23. Luh, G.C. & Lin, C.Y. (2011) Optimal design of truss-structures using particle swarm optimization. *Computers & Structures*. 89(23–24). pp. 2221–2232. DOI: 10.1016/j.compstruc.2011.08.013
- 24. Rouban, A.I. (2013) Global optimization method based on the selective averaging coordinates with restrictions. Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika – Tomsk State University Journal of Control and Computer Science. 22. pp. 114–123. DOI: 10.17212/1814-1196-2017-3-126-141
- 25. Baghlani, A. & Makiabadi, M. & Sarcheshmehpour, M. (2014) Discrete optimum design of truss structures by an improved firefly algorithm. *Advances in Structural Engineering*. 17(10). pp. 1517–1530. DOI: 10.1260/1369-4332.17.10.1517
- 26. Carbas, S. (2016) Design optimization of steel frames using an enhanced firefly algorithm. *Engineering Optimization*. pp. 1–19. DOI: 10.1080/0305215X.2016.1145217
- 27. Rouban, A.I. & Mikhalev, A.S. (2017) Global optimization with selective averaging of mixed variables: continuous and discrete with the ordered possible values. *Nauchnyy vestnik Novosibirskogo gosudarstvennogo tekhnicheskogno universiteta Scientific Bulletin of NSTU*. 3(68). pp. 126–141. DOI: 10.17212/1814-1196-2017-3-126-141
- 28. Korn, G.A. & Korn, T. (1973) Mathematical Handbook. New York, San Francisco, Toronto, London, Sydney: [s.n.].

Received: July 8, 2019

Rouban A.I., Mikhalev A.A. (2020) THE GLOBAL OPTIMIZATION METHOD WITH SELECTIVE AVERAGING OF THE DISCRETE DECISION VARIABLES *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50, pp. 47–55

DOI: 10.17223/19988605/50/6

Рубан А.И., Михалев А.С. МЕТОД ГЛОБАЛЬНОЙ ОПТИМИЗАЦИИ С СЕЛЕКТИВНЫМ УСРЕДНЕНИЕМ ДИСКРЕТНЫХ ИСКОМЫХ ПЕРЕМЕННЫХ. Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2020. № 50, С. 47–55

Предложен функционал покомпонентного селективного усреднения искомых дискретных переменных. В положительное убывающее ядро функционала введен положительный коэффициент селективности. При его увеличении усреднение в пределе обеспечивает получение оптимального значения искомых дискретных переменных. На основе оценки функционала селективного усреднения синтезирован базовый алгоритм глобальной оптимизации на множестве дискретных переменных с упорядоченными возможными значениями при наличии ограничений неравенств. На тестовом примере продемонстрированы высокие скорость сходимости и помехоустойчивость базового алгоритма. Статистическое исследование алгоритма показало, что оценка вероятности получения истинного решения достигает единицы.

Ключевые слова: глобальная оптимизация; дискретные переменные; селективное усреднение искомых переменных; много-экстремальная функция; ограничения типа неравенств.

ROUBAN Anatoly Ivanovich (Doktor of Technical Sciences, Professor of Computer Science Department of Institute of Space and Information Technologies, Siberian Federal University, Krasnoyarsk, Russian Federation).
E-mail: ai-rouban@mail.ru

MIKHALEV Anton Sergeevich (Senior Lector of Computer Science Department of Institute of Space and Information Technologies, Siberian Federal University, Krasnoyarsk, Russian Federation).

E-mail: asmikhalev@yandex.ru

2020 Управление, вычислительная техника и информатика

№ 50

УДК 519.218.72

DOI: 10.17223/19988605/50/7

Г.Ш. Цициашвили

СТАЦИОНАРНОЕ РАСПРЕДЕЛЕНИЕ В СИСТЕМЕ $M \mid M \mid 1 \mid \infty$ С ОСТАНАВЛИВАЮЩЕЙСЯ ИНТЕНСИВНОСТЬЮ ВХОДНОГО ПОТОКА

Работа выполнена при частичной поддержке РФФИ, проект № 17-07-00177.

Вычисляется стационарное распределение системы массового обслуживания, в которой интенсивность пуассоновского входного потока является марковским процессом, останавливающимся в некоторых состояниях. Исследуется зависимость стационраного распределения в системе от начального состояния входного потока. Показывается, что подобная модель случайной среды идентична модели разорения игрока, а ее исследование сводится к решению дискретного аналога задачи Дирихле и использованию известных формул для стационарных распределений числа заявок в системе обслуживания с постоянной интенсивностью входного потока.

Ключевые слова: система массового обслуживания; пуассоновский входной поток с останавливающейся интенсивностью; задача Дирихле; игра на разорение игрока.

Моделям массового обслуживания в случайной среде посвящено большое количество работ, затрагивающих как теоретические, так и прикладные аспекты (см., напр.: [1–7]). В этих работах основное внимание уделяется вычислению и анализу стационарных распределений, не зависящих от начального состояния. Однако появляются и другие модели, в которых случайный процесс, определяющий состояние среды, развивается на отрезке с отражающими и поглощающими концами (см., напр.: [8]). Стационарное распределение таких процессов уже зависит от начального состояния. В настоящей работе рассматривается наиболее простой вариант подобной модели случайной среды с помещенной в нее системой массового обслуживания. Вычисляется стационарное распределение системы массового обслуживания в такой случайной среде и исследуется ее зависимость от начального состояния. Показывается, что подобная модель случайной среды идентична модели разорения игрока, а ее исследование сводится к решению дискретного аналога задачи Дирихле [9. Гл. II, § 2–7] и использованию известных формул для стационарных распределений числа заявок в системе обслуживания с постоянной интенсивностью входного потока.

1. Постановка задачи

Рассмотрим одноканальную систему массового обслуживания $M \mid M \mid 1 \mid \infty$ с бесконечным числом мест ожидания и интенсивностью обслуживания μ . Случайный поцесс $\lambda(t)$, характеризующий интенсивность входного потока, задается марковской цепью z_k , k=0,1,..., с множеством состояний $\{0,...,N\}$. Марковская цепь z_k , k=0,1,..., описывает игру на разорение игрока [10. Гл. $[1, \S 9]]$ и определяет кусочно-постоянный случайный процесс $z(t)=z_k$, $k\leq t< k+1$, k=0,1,... Случайный процесс $\lambda(t)=\Lambda(z(t))>0$, $t\geq 0$, с множеством состояний $\{\Lambda(0),...,\Lambda(N)\}$ является марковским, $\Lambda(i)\neq\Lambda(j)$, $i\neq j$. Система массового обслуживания $M\mid M\mid 1\mid \infty$ с так определенной интенсивностью входного потока $\lambda(t)$ будет обозначаться M.

Рассмотрим марковский процесс (z(t),k(t)), $t \ge 0$, характеризующий случайную интенсивность входного потока и случайное число заявок в системе **M** в момент времени t. Далее предполагаем, что выполняются неравенства

$$\Lambda(0) < \mu, \dots, \Lambda(N) < \mu. \tag{1}$$

Нашей задачей является вычисление стационарного распределения второй компоненты k(t) этого процесса.

Очевидно, что стационарное распределение марковской цепи z_k , k=0,1,..., зависит от начального состояния z_0 . На первый взгляд, эта зависимость затрудняет решение поставленной задачи. Однако аналогия марковской цепи z_k , k=0,1,..., с процессом, описывающим игру на разорение игрока, позволяет, наоборот, существенно упростить вычисление стационарного распределения процесса в системе \mathbf{M} . Решение данной задачи может быть сведено к решению дискретного аналога уравнения Дирихле и к использованию известных формул для стацинарного распределения процесса обслуживания в системе с постоянной интенсивностью входного потока.

2. Стационарное распределение марковской цепи z_k

Рассмотрим марковскую цепь z_k , k=0,1,..., с множеством состояний $\{0,1,...,N\}$, с ненулевыми элементами матрицы $\Theta = \parallel \theta_{i,j} \parallel_{i,j=0}^{N}$ переходных вероятностей

$$\theta_{i,i+1} = p, \ \theta_{i,i-1} = q, \ i = 1,...,N-1, \ \theta_{0,0} = \theta_{N,N} = 1, \ 0 (2)$$

Всюду далее вероятность $P_{z_0}(A)$ обозначает вероятность события A при условии, что марковская цепь z_k , k=0,1,..., принимает начальное значение z_0 . Обозначим

$$\pi_{z_0}(k) = \sum_{z=1}^{N-1} P_{z_0}(z_k = z), \ \psi_{z_0}(k) = P_{z_0}(z_k = 0) + P_{z_0}(z_k = N),$$

очевидно, что справедливо соотношение

$$\pi_{z_0}(k) + \Psi_{z_0}(k) = 1. \tag{3}$$

Формулы (2), (3) выполняются при любом $z_0 = 0,...,N$.

Из соотношений (2) следуют неравенства

$$\begin{split} &P_{z_0}(z_0=0) \le P_{z_0}(z_1=0) \le ..., \\ &P_{z_0}(z_0=N) \le P_{z_0}(z_1=N) \le ..., \end{split}$$

и значит существуют пределы

$$\lim_{k\to\infty} P_{z_0}(z_k=0) = \Psi_{z_0}(0), \ \lim_{k\to\infty} P_{z_0}(z_k=N) = \Psi_{z_0}(N),$$

удовлетворяющие соотношению

$$\lim_{k \to \infty} \Psi_{z_0}(k) = \Psi_{z_0}(0) + \Psi_{z_0}(N) \le 1.$$

Очевидно, что при $\gamma = p^N + q^N < 1$ выполняется соотношение

$$\pi_{z_0}(kN) - \pi_{z_0}((k+1)N) \ge \pi_{z_0}(kN)\gamma$$
,

и значит

$$\pi_{z_0}((k+1)N) \le \pi_{z_0}(kN)\gamma, \ k=1,2,...$$

Отсюда получаем предельное соотношение

$$\lim_{k \to \infty} \pi_{z_0}(k) = 0. \tag{4}$$

Из формул (3), (4) следует предельное соотношение

$$\lim_{k \to \infty} \Psi_{z_0}(k) = \Psi_{z_0}(0) + \Psi_{z_0}(N) = 1.$$
 (5)

Из определения марковской цепи z_k , k=0,1,..., следует, что $\Psi_{z_0}(0)$ — это вероятность того, что, выходя из состояния z_0 , марковская цепь z_k , k=0,1,..., когда-нибудь попадет в состояние 0. Поэтому $\Psi_{z_0}(0)$ можно интерпретировать как вероятность разорения игрока с начальной суммой z_0 с вероятностью выигрыша p, вероятностью проигрыша q и суммой выигрыша в игре N.

Тогда функция $\Psi_{z_0}(0)$ удовлетворяет системе линейных алгебраических уравнений

$$\Psi_{z_0}(0) = p\Psi_{z_0-1}(0) + q\Psi_{z_0+1}(0), \ z_0 = 1, ..., N-1, \ \Psi_0(0) = 1, \ \Psi_N(0) = 0. \tag{6}$$

Используя известные формулы для вероятности разорения игрока [10. Гл. I, § 9], можно выписать следующие соотношения:

$$\Psi_{z_0}(0) = \frac{(p/q)^{z_0} - (p/q)^N}{1 - (p/q)^N}, \ p \neq q; \ \Psi_{z_0}(0) = 1 - \frac{z_0}{N}, \ p = q; \ z_0 = 0,...,N,$$
 (7)

В свою очередь, из формулы (5) следует, что вероятность выигрыша

$$\Psi_{z_0}(N) = 1 - \Psi_{z_0}(0). \tag{8}$$

Распределение $\Psi_{z_0}(0),...,\Psi_{z_0}(N)$ является стационарным распределением марковской цепи z_k , k=0,1,..., с начальным состоянием z_0 . Это распределение сосредоточено в точках 0,N, которые являются поглощающими для данной цепи. Отсюда следует, что цепь z_k , k=0,1,..., при любом начальном условии $z_0=0,...,N$ имеет стационарное распределение. Однако эта цепь не является эргодической, так как ее стационарное распределение зависит от начального условия. Этими же свойствами обладает кусочно-постоянный процесс z(t), определяемый по марковской цепи z_k , k=0,1,..., и задающий интенсивность входного потока $\lambda(t)=\Lambda(z(t))$ системы массового обслуживания \mathbf{M} .

3. Стационарное распределение числа заявок в системе массового обслуживания М

При начальном условии $z(0)=z_0$ стационарное распределение процесса $\lambda(t)=\Lambda(z(t))$ удовлетворяет соотношениям

$$\lim_{t \to \infty} P_{z_0}(\lambda(z(t)) = \Lambda(0)) = \Psi_{z_0}(0), \ \lim_{t \to \infty} P_{z_0}(\lambda(z(t)) = \Lambda(N)) = \Psi_{z_0}(N).$$

Иными словами, случайный процесс $\lambda(t)$ имеет на множестве состояний $\{0,...,N\}$ двухточечное стационарное распределение $P(\lambda(t) \equiv \Lambda(0)) = \Psi_{z_0}(0), \ P(\lambda(t) \equiv \Lambda(N)) = \Psi_{z_0}(N),$ зависящее от начального состояния z_0 . Причем процесс $\lambda(t)$, попадая в состояние $\Lambda(0)$ (попадая в состояние $\Lambda(N)$), останавливается в нем и далее не меняется.

Перейдем теперь к вычислению предельного распределения числа заявок в системе **M** при условии (1). Известно, что стационарное распределение p(k) числа заявок в системе $M \mid M \mid 1 \mid \infty$ при постоянной интенсивности входного потока λ и интенсивности обслуживания μ удовлетворяет соотношениям

$$p(k) = (1 - \rho)\rho^k, \ k = 0, 1, ..., \ \rho = \frac{\lambda}{\mu} < 1.$$
 (9)

Обозначим $\rho(0) = \frac{\Lambda(0)}{\mu} < 1$, $\rho(N) = \frac{\Lambda(N)}{\mu} < 1$, тогда стационарное распределение $P_{z_0}(k)$, k = 0,1,...,

числа заявок в системе массового обслуживания М удовлетворяет равенству

$$P_{z_0}(k) = \Psi_{z_0}(0)(1 - \rho(0))\rho^k(0) + \Psi_{z_0}(N)(1 - \rho(N))\rho^k(N), \ k = 0, 1, \dots$$
 (10)

Если предположить, что A(i), i=0,...,N, — начальное распределение случайной величины z_0 , то тогда вследствие формулы (10) стационарное распределение $\Pi(k)$, k=0,1,..., числа заявок в системе обслуживания **M** имеет вид:

$$\Pi(k) = \sum_{i=0}^{N} A(i) \left[\Psi_i(0) (1 - \rho(0)) \rho^k(0) + \Psi_i(N) (1 - \rho(N)) \rho^k(N) \right], \quad k = 0, 1, \dots$$
(11)

Заметим, что если A(1) = ... = A(N-1) = 0, то для существования стационарного распределения достаточно потребовать $\Lambda(0) < \mu$, $\Lambda(N) < \mu$.

Заключение

Перечислим теперь возможные обобщения полученных результатов. От простейшей одноканальной системы $M \mid M \mid 1 \mid \infty$, лежащей в основе модели \mathbf{M} , можно перейти к многоканальным системам, системам с отказами, к открытым сетям массового обслуживания. Марковская цепь z_k , k=0,1,..., определяющая интенсивность входного потока, может быть заменена случайным процессом с непрерывным временем, например диффузионным процессом с поглощениями или частичными поглощениями и отражениями в концах некоторого отрезка. Эта марковская цепь может определять не только интенсивность входного потока, но и интенсивность обслуживания. Тогда вероятность попадания марковской цепи z_k , k=0,1,..., удовлетворяет дискретному аналогу уравнения Дирихле и может быть решена известными методами теории вероятностей и математической физики. Можно также по стационарному распределению процесса k(t), $t \ge 0$, числа заявок в системе массового обслуживания \mathbf{M} (оцениваемому, например, по наблюдениям) определить начальное состояние z_0 процесса z(t), $t \ge 0$, характеризующего интенсивность входного потока в начальный момент времени.

ЛИТЕРАТУРА

- 1. Вишневский В.М., Дудин А.Н., Клименок В.И. Стохастические системы с коррелированными потоками. Теория и применение в телекоммуникационных сетях. М.: Техносфера, 2018. 564 с.
- 2. Вишневский В.М., Семёнова О.В. Системы поллинга: теория и применение в широкополосных беспроводных сетях. М.: Техносфера, 2007. 470 с.
- 3. Vishnevskiy V.M., Evfrosinin D.V., Krishnamurti A. Principles of Construction of Mobile and Stationary Tethered High-Altitude Unmanned Telecommunication Platforms of Long-Term Operation // Communications in Computer and Information Science. 2018. V. 919. P. 561–569.
- 4. Klimenok V.I. Two-Server Queueing System with Unreliable Servers and Markovian Arrival Process // Communications in Computer and Information Science. 2017. V. 800. P. 63–74.
- 5. Klimenok V.I., Dudin A.N., Vishnevskiy V.M. A Retrial Queueing System with Alternating Inter-retrial Time Distribution // Communications in Computer and Information Science. 2018. V. 919. P. 302–315.
- 6. Коротаев И.А., Спивак Л.Р. Системы массового обслуживания в полумарковской случайной среде // Автоматика и телемеханика. 1992. Вып. 7. С. 86–92.
- 7. Жерновой Ю.В. Система массового обслуживания $M \mid M \mid n \mid r$, функционирующая в синхронной случайной среде // Информационные процессы. 2009. Т. 9, вып. 4. С. 352–363.
- 8. Бондрова О.В., Головко Н.И., Жук Т.А. Вывод уравнений типа Колмогорова–Чепмена с интегральным оператором // Дальневосточный математический журнал. 2017. Т. 17, вып. 2. С. 135–146.
- 9. Дынкин Е.Б., Юшкевич А.А. Теоремы и задачи о процессах Маркова. М.: Наука, 1967. 232 с.
- 10. Ширяев А.Н. Вероятность. М.: Наука, 1989. 432 с.

Поступила в редакцию 2 октября 2019 г.

Tsitsiashvili G.Sh. (2019) STATIONARY DISTRIBUTION IN THE SYSTEM WITH THE STAYING INTENSITY OF THE INPUT FLOW. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 56–60

DOI: 10.17223/19988605/50/7

Consider a single server queuing system with an infinite number of waiting places and service intensity representing by some random process. A Markov chain with many states defines the process characterizing the intensity of the input flow. The Markov chain describes a game of player ruin and defines a piecewise constant random process with many states which is Markov process. A queuing system with so defined input flow rate will be denoted **M**.

Consider the Markov process characterizing the random intensity of the input flow and the random number of customers in the system at a time. Then assume that the following inequalities are satisfied: $\Lambda(0) < \mu, ..., \Lambda(N) < \mu, \quad \Lambda(i) \neq \Lambda(j), \quad i \neq j$. Our task is to calculate the stationary distribution of the process, describing a number of customers in the system **M**.

It is obvious that the stationary distribution of the Markov chain depends on the initial state At first glance, this dependence makes it difficult to solve the problem. However, the analogy of the Markov chain process, describing the game to ruin the player, on the contrary, simplifies the computation of the stationary distribution of the process in the system M. The solution of this problem can be reducing to the solution of a discrete analogue of the Dirichlet equations. Then it is possible to use well-known formulas for stationary distribution service process in the system with a constant intensity of the input flow.

Consider a Markov chain with a set of states with nonzero elements of the transition probability matrix

$$\theta_{i,i+1} = p, \ \theta_{i,i-1} = q, \ i = 1, ..., N-1, \ \theta_{0,0} = \theta_{N,N} = 1, \ 0$$

Everywhere further, probability $P_{z_0}(A)$ denotes the probability of an event A provided that the Markov chain z_k , k = 0,1,..., takes the initial value z_0 . Denote

$$\pi_{z_0}(k) = \sum_{z=1}^{N-1} P_{z_0}(z_k = z), \ \psi_{z_0}(k) = P_{z_0}(z_k = 0) + P_{z_0}(z_k = N).$$

It is obvious that the following relation is valid

$$\pi_{z_0}(k) + \psi_{z_0}(k) = 1.$$
 (2)

Formulas (1), (2) are true at any $z_0 = 0,...,N$. Denote $\rho(0) = \frac{\Lambda(0)}{\mu} < 1$, $\rho(N) = \frac{\Lambda(N)}{\mu} < 1$, then stationary distribution $P_{z_0}(k)$,

k = 0,1,..., of a number of customers in the queuing system **M** satisfies the equality

$$P_{z_0}(k) = \Psi_{z_0}(0)(1-\rho(0))\rho^k(0) + \Psi_{z_0}(N)(1-\rho(N))\rho^k(N), \ k = 0,1,...$$

Keywords: queuing system; Poisson input flow with stopping intensity; Dirichlet problem; the game to ruin the player.

TSITSIASHVILI Gurami Shalvovich (Doctor of Physical and Mathematical Sciences, Professor, Institute for Applied Mathematics, Far Eastern Branch of RAS, Vladivostok, Russian Federation).

E-mail: guram@iam.dvo.ru

REFERENCES

- 1. Vishnevsky, V.M., Dudin, A.N. & Klimenok, V.I. (2018) *Stokhasticheskie sistemy s korrelirovannymi potokami. Teoriya i primenenie v telekommunikatsionnykh setyakh* [Stochastic systems with correlated flows. Theory and application in telecommunication networks]. Moscow: Tekhnosfera.
- 2. Vishnevsky, V.M. & Semenova, O.V. (2007) *Sistemy pollinga: teoriya i primenenie v shirokopolosnykh besprovodnykh setyakh* [Polling systems: theory and application in broadband wireless networks]. Moscow: Tekhnosfera.
- 3. Vishnevsky, V.M., Evfrosinin, D.V. & Krishnamurti, A. (2018) Principles of Construction of Mobile and Stationary Tethered High-Altitude Unmanned Telecommunication Platforms of Long-Term Operation. *Communications in Computer and Information Science*. 919. pp. 561–569. DOI: 10.1007/978-3-319-99447-5_48
- 4. Klimenok, V.I. (2017) Two-Server Queueing System with Unreliable Servers and Markovian Arrival Process. *Communications in Computer and Information Science*. 800. pp. 63–74. DOI: 10.1007/978-3-319-68069-9_4
- Klimenok, V.I., Dudin, A.N. & Vishnevsky, V.M. (2018) A Retrial Queueing System with Alternating Inter-retrial Time Distribution. Communications in Computer and Information Science. 919. pp. 302–315. DOI: 10.1007/978-3-319-99447-5
- 6. Korotaev, I.A. & Spivak, L.R. (1992) Queuing Systems in the semi-Markov random environment. *Automatics Remote Control*. 7. pp. 86–92.
- Zhernovoy, Yu.V. (2009) Queuing system functioning in synchronous random environment. *Information Processes*. 9(4). pp. 352–363.
- 8. Bodrova, O.V., Golovko, N.I. & Zhuk, T.A. (2017) Derivation of Kolmogorov-Chapman type equations with integral operator. *Far Eastern Mathematical Journal*. 17(2). pp. 135–146.
- 9. Dynkin, E.B. & Yushkevich, A.A. (1967) *Teoremy i zadachi o protsessakh Markova* [Theorems and problems on Markov processes]. Moscow: Nauka.
- 10. Shiryaev, A.N. (1989) Veroyatnost' [Probability]. Moscow: Nauka.

2020 Управление, вычислительная техника и информатика

№ 50

УДК 519.24

DOI: 10.17223/19988605/50/8

V.M. Chubich, E.V. Filippova

ACTIVE PARAMETRICAL IDENTIFICATION OF STOCHASTIC LINEAR CONTINUOUS-DISCRETE SYSTEMS BASED ON THE EXPERIMENT DESIGN IN THE PRESENCE OF ABNORMAL OBSERVATIONS

The reported study was funded by RFBR according to the research project No. 18-31-00283.

The procedure of active parametrical identification of stochastic linear continuous-discrete systems including robust estimation of parameters and optimal design of input signals is offered. A general case of entering unknown parameters into the equations of state and observation, initial conditions and covariance matrices of system noise and measurements is considered. The efficiency of this procedure is demonstrated by the example of a direct current motor control system.

Keywords: continuous-discrete system; anomalous observations; optimal input signal; robust estimation.

Development of information technologies for identification of complex dynamic systems of stochastic nature is an important area of research and has attracted considerable interest. Application of the optimal experiment theory methods in parametrical identification improves the quality of the results by taking into account more fully the properties of the dynamic object and data collection procedures [1–7]. Thus, given the structure of the mathematical model, the procedure of active parametrical identification involves the following steps:

- Calculation of parameter estimates based on measurement data corresponding to a test input signal.
- Synthesis based on the obtained estimates of the optimal input signal (experiment design).
- Recalculation of estimates of unknown parameters according to the measured data corresponding to the synthesized signal.

Traditionally, the estimation of unknown parameters is carried out based on the classical Kalman filter, which makes it possible to find estimates of the state vector and corresponding covariance matrices under the assumption of the normality of the noise distribution of the system and measurements. When solving practical problems (for example, problems of communication, navigation, control and radar) there are cases when the usual mechanism of formation of observation data is broken and the appearance of anomalous observations that do not contain information about the object under study. In this case the specified filter can lead to biased estimates or even diverge. At the moment, numerous robust modifications of the Kalman filter, resistant to the appearance of outliers, have been developed. In this regard, it is advisable to consider robust estimation methods that provide good quality results.

This work is devoted to the development of mathematical and software procedures of active identification of stochastic continuous-discrete systems based on robust parameter estimation. The efficiency of the developed procedure is demonstrated by the example of one model structure.

1. Problem statement

Consider the following controlled, observed, identifiable dynamic system model in state space:

$$\frac{d}{dt}x(t) = F(t)x(t) + \Psi(t)u(t) + \Gamma(t)w(t), t \in [t_0, t_N],$$
(1)

$$y(t_{k+1}) = H(t_{k+1})x(t_{k+1}) + v(t_{k+1}), k = 0,1,...,N-1,$$
 (2)

where x(t) is the state *n*-vector; u(t) is deterministic control (input) *r*-vector; w(t) is the process noise *p*-vector; $y(t_{k+1})$ is the measurement (output) *m*-vector; $v(t_{k+1})$ is the measurement error *m*-vector.

Let us suppose that

- the random vectors w(t) and $v(t_{k+1})$ form a white Gaussian noise, for which

$$E[w(t)] = 0, \quad E[w(t)w^{T}(\tau)] = Q(t)\delta(t-\tau),$$

$$E[v(t_{k+1})] = 0, \quad E[v(t_{k+1})v^{T}(t_{i+1})] = R(t_{k+1})\delta_{ki},$$

$$E[v(t_{k+1})w^{T}(\tau)] = 0$$

(here E[] is operator of mathematical expectation, $\delta(t-\tau)$ is delta function, δ_{ki} is the Kronecker symbol);

– initial state $x(t_0)$ has a normal distribution with parameters

$$E[x(t_0)] = \overline{x}(t_0), \quad E\{[x(t_0) - \overline{x}(t_0)][x(t_0) - \overline{x}(t_0)]^T\} = P(t_0)$$

and is uncorrelated with w(t) and $v(t_{k+1})$ for all values of k;

- output data may contain outliers;
- unknown parameters are summarized in the s-vector θ , including the elements of matrices F(t), $\Psi(t)$, $\Gamma(t)$, $H(t_{k+1})$, Q(t), $R(t_{k+1})$, $P(t_0)$ and vector $\overline{x}(t_0)$ in various combinations.

For the mathematical model (1), (2), taking into account the a priori assumptions, it is necessary to develop procedures for the active parametrical identification of stochastic continuous-discrete systems based on robust parameter estimation and conduct a numerical study of the effectiveness of its application.

2. Methods of research

Let us consider the main theoretical aspects of the active identification procedure.

Parameter estimation. Unknown parameters estimations of the mathematical model (1), (2) are carried out according to observational data Ξ by using some criterion of identification χ . The collection of numerical data occurs during identification experiments which are carried out under some discrete design ξ_{ν} :

$$\xi_{v} = \begin{cases} u^{1}(t), u^{2}(t), \dots, u^{q}(t) \\ \frac{k_{1}}{v}, \frac{k_{2}}{v}, \dots, \frac{k_{q}}{v} \end{cases}, u^{i}(t) \in \Omega_{u}, i = 1, \dots, q.$$

Here v is the total number of launches of the system, q is the number of points of the design, k_i is the number of experiments corresponding to the signal $u^i(t)$, Ω_u is the set of design (determined by restrictions on the conditions of the experiment).

Let us denote through $Y_{ij}^{T} = \left[\left[y^{ij}(t_1) \right]^{T}, ..., \left[y^{ij}(t_N) \right]^{T} \right]$ realization of the output signal with number j ($j = 1, ..., k_i$) corresponding to the input signal $u^i(t)$. Then

$$\Xi = \{ (u^i(t), Y_{ij}), j = 1, 2, ..., k_i, i = 1, 2, ..., q \}, \sum_{i=1}^{q} k_i = v.$$

Due to the fact that the measurement data contain anomalous observations, we will calculate quasi-likelihood estimates [8], solving the following optimization problem:

$$\hat{\theta} = \arg\min_{\theta \in \Omega_{\theta}} \left[\chi(\theta; \Xi) \right] = \arg\min_{\theta \in \Omega_{\theta}} \left[-\ln L(\theta; \Xi) \right]. \tag{3}$$

Here

$$\chi(\theta;\Xi) = \frac{Nm\nu}{2} \ln 2\pi + \frac{1}{2} \sum_{i=1}^{q} \sum_{j=1}^{k_i} \sum_{k=0}^{N-1} \left[\varepsilon^{ij}(t_{k+1}) \right]^T \left[B^i(t_{k+1}) \right]^{-1} \left[\varepsilon^{ij}(t_{k+1}) \right] + \frac{1}{2} \sum_{i=1}^{q} k_i \sum_{k=0}^{N-1} \ln \det B^i(t_{k+1}),$$
(4)

where $\varepsilon^{ij}(t_{k+1})$ and $B^i(t_{k+1})$ determined by the recurrent equations of a robust filter.

The calculation of the conditional minimum (3) will be carried out by the method of sequential quadratic programming [9, 10], implemented in the optimization Toolbox MATLAB package and assuming the calculation of the gradient.

Differentiating the equality (4) by θ_{α} ($\alpha = 1,...,s$) taking into account expression

$$\frac{\partial \ln \det B(t_{k+1})}{\partial \theta_{\alpha}} = Sp \left[B^{-1}(t_{k+1}) \frac{\partial B(t_{k+1})}{\partial \theta_{\alpha}} \right];$$

$$\frac{\partial B^{-1}(t_{k+1})}{\partial \theta_{\alpha}} = -B^{-1}(t_{k+1}) \frac{\partial B(t_{k+1})}{\partial \theta_{\alpha}} B^{-1}(t_{k+1})$$

and symmetry of the matrix $B(t_{k+1})$, we obtain

$$\begin{split} \frac{\partial \chi \left(\boldsymbol{\theta};\Xi\right)}{\partial \boldsymbol{\theta}_{\alpha}} &= \sum_{i=1}^{q} \sum_{j=1}^{k_{i}} \sum_{k=0}^{N-1} \left\{ \left[\frac{\partial \boldsymbol{\epsilon}^{ij} \left(t_{k+1}\right)}{\partial \boldsymbol{\theta}_{\alpha}} \right]^{T} \left[\boldsymbol{B}^{i} \left(t_{k+1}\right) \right]^{-1} \left[\boldsymbol{\epsilon}^{ij} \left(t_{k+1}\right) \right] - \\ &- \frac{1}{2} \left[\boldsymbol{\epsilon}^{ij} \left(t_{k+1}\right) \right]^{T} \left[\boldsymbol{B}^{i} \left(t_{k+1}\right) \right]^{-1} \frac{\partial \boldsymbol{B}^{i} \left(t_{k+1}\right)}{\partial \boldsymbol{\theta}_{\alpha}} \left[\boldsymbol{B}^{i} \left(t_{k+1}\right) \right]^{-1} \left[\boldsymbol{\epsilon}^{ij} \left(t_{k+1}\right) \right] \right\} + \\ &+ \frac{1}{2} \sum_{i=1}^{q} k_{i} \sum_{k=0}^{N-1} Sp \left\{ \left[\boldsymbol{B}^{i} \left(t_{k+1}\right) \right]^{-1} \frac{\partial \boldsymbol{B}^{i} \left(t_{k+1}\right)}{\partial \boldsymbol{\theta}_{\alpha}} \right\}. \end{split}$$

Derivatives of $\frac{\partial \varepsilon^{ij}(t_{k+1})}{\partial \theta_{\alpha}}$ and $\frac{\partial B^{i}(t_{k+1})}{\partial \theta_{\alpha}}$ determined by the equations arising from the corresponding rela-

tions of the robust filter.

In [11] the authors conducted a comparative analysis of the efficiency of some modern robust filters for non-stationary linear continuous-discrete systems. While best and quite comparable the results showed correntropy filters Izanloo–Fakoorian–Yazdi–Simon [12] and Chen–Liu–Zhao-Principe [13]. From the point of view of the organization of calculations and, as a consequence, the software implementation of the first of these filters is much easier. In this regard, it seems appropriate to use the Izanloo–Fakoorian–Yazdi–Simon filter when estimating the parameters of models of stochastic linear continuous-discrete systems in the presence of anomalous observations. The corresponding recurrence relations for a single system startup are shown below.

Izanloo-Fakoorian-Yazdi-Simon filter. Initialization:

$$\hat{x}\left(t_0\mid t_0\right) = \overline{x}\left(t_0\right),\ P\left(t_0\mid t_0\right) = P\left(t_0\right);\ \sigma = \sigma_0.$$

To run in a loop on $k = \overline{0, N-1}$:

$$\begin{split} \frac{d}{dt} \hat{x} \Big(t \, | \, t_k \, \Big) &= F \Big(t \Big) \hat{x} \Big(t \, | \, t_k \, \Big) + \Psi \Big(t \Big) u \Big(t \Big), \, t \in \left[t_k, t_{k+1} \right]; \\ \frac{d}{dt} P \Big(t \, | \, t_k \, \Big) &= F \Big(t \Big) P \Big(t \, | \, t_k \, \Big) + P \Big(t \, | \, t_k \, \Big) F^T \Big(t \Big) + \Gamma \Big(t \Big) Q \Big(t \Big) \Gamma^T \Big(t \Big), t \in \left[t_k, t_{k+1} \right]; \\ \varepsilon(t_{k+1}) &= y(t_{k+1}) - H(t_{k+1}) \hat{x}(t_{k+1} \, | \, t_k \,); \quad L \Big(t_{k+1} \Big) = \exp \left(- \frac{\varepsilon^T \Big(t_{k+1} \Big) R^{-1} \Big(t_{k+1} \Big) \varepsilon \Big(t_{k+1} \Big)}{2\sigma^2} \right); \end{split}$$

$$B(t_{k+1}) = H(t_{k+1})P(t_{k+1} | t_k)L(t_{k+1})H^T(t_{k+1}) + R(t_{k+1}); K(t_{k+1}) = P(t_{k+1} | t_k)L(t_{k+1})H^T(t_{k+1})B^{-1}(t_{k+1});$$

$$\hat{x}(t_{k+1} | t_{k+1}) = \hat{x}(t_{k+1} | t_k) + K(t_{k+1})\varepsilon(t_{k+1}); P(t_{k+1} | t_{k+1}) = \left[I - K(t_{k+1})H(t_{k+1})\right]P(t_{k+1} | t_k).$$

End of loop.

Algorithms for calculating the maximum likelihood criterion and its gradient based on robust filtering for linear continuous-discrete models are presented in [14].

Experiment design. Let us consider the features of input signal design for models of continuous-discrete systems (1), (2). Continuous normalized design in this case can be specified as

$$\xi = \begin{cases} u^{1}(t), u^{2}(t), \dots, u^{q}(t) \\ p_{1}, p_{2}, \dots, p_{q} \end{cases}, p_{i} \ge 0, \sum_{i=1}^{q} p_{i} = 1, u^{i}(t) \in \Omega_{u}, i = 1, 2, \dots, q.$$
 (5)

Unlike discrete design ξ_{ν} , weights p_i in continuous design ξ can take any values, including irrational number.

Information matrix $M(\xi)$ for design (5) is determined by the relation

$$M(\xi) = \sum_{i=1}^{q} p_i M(u^i(t), \hat{\theta}),$$

in which the information matrices of single-point design depend on the unknown parameters to be estimated and are calculated in accordance with [15].

We find the optimal experiment design for some convex functional X information matrix $M(\xi)$ by solving the following extremal problem

$$\xi^* = \arg\min_{\xi \in \Omega_{\varepsilon}} X \Big[M \big(\xi \big) \Big]. \tag{6}$$

The construction of optimal design can be associated with the representation of the components of the input signals in the form of linear combinations of basic functions (as such, you can use orthogonal polynomials Legendre, Chebyshev, Walsh function, etc.) and then search for the coefficients of such linear combinations. Another approach is related to the assumption that the input signals are piecewise-constant functions preserving their values on the interval between adjacent measurements. In [16] have demonstrated the effectiveness and applicability of the piecewise-constant approximation of the input signal, which makes it possible to calculate the derivatives of the information matrix Fisher from the components of the input signal by recurrent analytical formulas and, consequently, to apply gradient procedures for the synthesis of optimal signals. This means that

$$u^{i}(t) = \left[u^{i}(t_{0}), u^{i}(t_{1}), ..., u^{i}(t_{N-1})\right]^{T} = U_{i}.$$

Based on the method of sequential quadratic programming, we present a combined procedure for constructing D- or A-optimal continuous design for pre-computed estimates of the parameters $\hat{\theta}$, which using direct and dual approaches [17–19] for solving the extremal problem (6).

1. Set the initial nondegenerate design

$$\xi_0 = \begin{cases} U_1^0, U_2^0, ..., U_q^0, \\ p_1^0, p_2^0, ..., p_q^0 \end{cases}, \quad U_i^0 \in \Omega_u, \quad p_i^0 = \frac{1}{q}, \quad i = 1, 2, ..., q.$$

Calculate the information matrix of a single-point designs $M\left(U_i^0,\hat{\theta}\right)$ for i=1,...,q and put k=0.

2. Counting the design weight $p_1^k, ..., p_q^k$ are fixed, find the design

$$\overline{\xi}_{k+1} = \arg\min_{U_1^k, \dots, U_n^k \in \Omega_n} X [M(\xi_k)].$$

Calculate the information matrix of a single-point designs $M(U_i^{k+1}, \hat{\theta})$, i = 1,...,q.

3. Having fixed the points of the design spectrum $\overline{\xi}_{k+1}$, we find the design

$$\xi_{k+1} = \arg\min_{p_1^k,...,p_q^k} X[M(\overline{\xi}_{k+1})], \ p_i^k \ge 0, \sum_{i=1}^q p_i^k = 1, \ i = 1,...,q.$$

4. If an inequality holds for a small positive δ_1

$$\sum_{i=1}^{q} \left[\left\| U_i^{k+1} - U_i^{k} \right\|^2 + \left(p_i^{k+1} - p_i^{k} \right)^2 \right] \leq \delta_1,$$

then let's put $\xi_0 = \xi_{k+1}$ (the execution of the direct procedure is finished), k = 0 and go to step 5. Otherwise, take k = k + 1 and go to step 2.

- 5. Calculate the information matrix $M(\xi_k)$.
- 6. Find the local maximum

$$U^k = \arg\max_{U \in \Omega_u} \mu(U, \, \xi_k).$$

If the condition $\left|\mu(U^k,\xi_k)-\eta\right| \leq \delta_2$ for a small positive δ_2 is met, then the process is finished. If $\mu(U^k,\xi_k) > \eta$, let's move on the step 7, otherwise, to seek a new local maximum.

7. Find τ_k

$$\tau_k = \arg\min_{0 \le \tau \le 1} X \left[M \left(\xi_{k+1}^{\tau} \right) \right].$$

Here $\xi_{k+1}^{\tau} = (1-\tau)\xi_k + \tau\xi(U^k)$, where $\xi(U^k)$ is a one-point design posted at the point U^k .

8. Make a design

$$\xi_{k+1} = (1 - \tau_k) \xi_k + \tau_k \xi(U^k),$$

let's clean it up by following [17], put k = k+1 and go to step 5.

The correspondence of the parameters $X[M(\xi)]$, $\mu(U, \xi)$, η of the combined procedure to the specified criteria is presented in table 1.

Parameters of the combined procedure

Table 1

Criterion	Parameters			
	$X[M(\xi)]$	$\mu(U, \xi)$	η	
D	$-\ln \det M(\xi)$	$\mathrm{Sp}\Big[M^{-1}(\xi)M(U)\Big]$	S	
A	$\operatorname{Sp}\!\left[M^{-1}(\xi)\right]$	$\operatorname{Sp}\left[M^{-2}(\xi)M(U)\right]$	$\operatorname{Sp}\!\left[M^{-1}(\xi)\right]$	

The analytical expressions of the required for the combined procedure derivatives and its calculation algorithms are presented in [15].

Practical application of the synthesized optimal design is difficult, because weights are arbitrary real numbers, enclosed in the interval from zero to one. In the case of a given number ν of possible system starts, it is necessary to "round" the continuous design to discrete [18]. As a result, we obtain a discrete design

$$\xi_{v} = \begin{cases} U_{1}^{*}, U_{2}^{*}, \dots, U_{q}^{*} \\ \frac{k_{1}^{*}}{v}, \frac{k_{2}^{*}}{v}, \dots, \frac{k_{q}^{*}}{v} \end{cases},$$

we perform identification experiments and recalculate estimates of unknown parameters.

3. Simulation

We assume that all a priori assumptions made in the problem definition are fulfilled. Following [20], consider a position control system consisting of an antenna and a direct current (DC) motor. Let the first com-

ponent of the state vector be responsible for the angular position of the antenna, the second - for its angular velocity. The input signal is the voltage at the input of the DC amplifier controlling the motor. The angular position is measured using a potentiometer. Then the models of state and observation can be determined by the relations:

$$\frac{d}{dt}x(t) = \begin{bmatrix} 0 & 1 \\ 0 & -\theta_1 \end{bmatrix} x(t) + \begin{bmatrix} 0 \\ \theta_2 \end{bmatrix} u(t) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} w(t), \ t \in [0,30];$$
$$y(t_{k+1}) = \begin{bmatrix} 1 & 0 \end{bmatrix} x(t_{k+1}) + v(t_{k+1}), \quad k = 0,..., N-1.$$

Here θ_1 , θ_2 are unknown parameters and $\Omega_{\theta} = \{1 < \theta_1 < 10, 0 < \theta_2 < 1\}$.

Set
$$N = 30$$
, $u(t) = 12$, $Q(t) = Q = 0.01$, $R(t_{k+1}) = R = 0.1$, $\overline{x}(t_0) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $P(t_0) = \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix}$. As-

sume that the measurements are made uniformly every single second and the value of the parameter Izanloo–Fakoorian–Yazdi–Simon filter is σ =10. Chose area of allowable values of input signals

$$\Omega_u = \{ 2 \le u(t) \le 30, t \in [0;30] \}.$$

We simulate samples with anomalous observations using the MATLAB software environment by setting the pollution coefficient of the sample $\lambda = 0.1$ and the noise dispersion of anomalous observations $R_A = 1000R$, assuming that the true values of the parameters $\theta_1^* = 4.600$, $\theta_2^* = 0.787$.

To reduce the dependence of the estimation results on the experimental data, we perform five independent starts of the system and average the obtained estimates of the unknown parameters. The quality of parametric identification will be judged by the value of the relative estimation error δ_{θ} , calculated by formula:

$$\delta_{\boldsymbol{\theta}} = \sqrt{\frac{\left(\boldsymbol{\theta}_{1}^{*} - \hat{\boldsymbol{\theta}}_{1}\right)^{2} + \left(\boldsymbol{\theta}_{2}^{*} - \hat{\boldsymbol{\theta}}_{2}\right)^{2}}{\left(\boldsymbol{\theta}_{1}^{*}\right)^{2} + \left(\boldsymbol{\theta}_{2}^{*}\right)^{2}}} \;,$$

where $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)$ is vector of unknown parameters estimates.

Numerical results of calculations are presented in table 2 (the optimal design is one-point).

 $$\operatorname{Table}\ 2$$ Results of robust procedure of active parametrical identification

	System start	Estimates and estimation errors		
Input signal and values of the D-optimality criterion	number	$\hat{ heta}_1$	$\hat{\theta}_2$	δ_{θ}
u(t)	1	5,012	0,547	0,102
12	2	3,411	0,689	0,255
10 -	3	3,632	0,302	0,232
8 -	4	4,890	0,705	0,064
6 -	5	3,075	0,488	0,332
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	Average value for startups	4,004	0,546	0,137
u*(t)	1	4,277	0,705	0,071
30	2	4,332	0,791	0,057
25	3	5,016	0,589	0,098
20	4	4,170	0,691	0,094
15	5	3,994	0,722	0,130
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	Average value for startups	4,357	0,699	0,055

Analysis of table 2 contents shows that design of D-optimal input signal using the combined procedure make it possible to improve the quality of estimation by 8,2%.

Thus, the authors consider that the applying of the active parametrical identification procedures based on robust estimation and optimal design of input signals is helpful and advisable in the presence of outliers in the measurement data.

Conclusion

Robust procedure of active parametrical identification for models of stochastic linear continuous-discrete systems including robust estimation of parameters based on the Izanloo–Fakoorian–Yazdi–Simon filter and optimal design of input signal is developed. The case of the entry of unknown parameters into the state and observations equations, the initial conditions and the covariance noise matrices of the system and measurements is considered. The efficiency of the developed robust procedure of active parametrical identification is demonstrated on the example of a DC motor control system for the random nature of the outliers.

REFERENCES

- 1. Ljung, L. (1987) System identification: Theory for the User. Englewood Cliffs, NJ: Prentice-Hall.
- 2. Walter, E. & Pronzato, L. (1997) Identification of parametric models from experimental data. Berlin: Springer-Verlag.
- 3. Jauberthie, C., Bournonville, F., Coton, P. & Rendell, F. (2006) Optimal input design for aircraft parameter estimation. *Aerospace Science and Technology*. 10. pp. 331–337. DOI:10.1016/j.ast.2005.08.002
- 4. Brasil, N., Hemerly, E. & Goes, L. (2009) Aircraft parameter estimation experiment design considering measurement colored residuals. *Journal of Aircraft*. 46(6). pp. 1857–1865. DOI:10.2514/1.34133
- Bernaerts, K., Servaes, R.D., Kooyman, S., Versyck, K.J. & Van Impe, J. (2002) Optimal temperature input design for estimation of the Square Root model parameters: parameter accuracy and model validity restrictions. *International Journal of Food Micro-biology*. 73. pp. 145–157. DOI: 10.1016/S0168-1605(01)00645-6
- 6. Galvanin, F., Marchesini, R., Barolo, M., Bezzo, F. & Fidaleo, M. (2016) Optimal design of experiments or parameter identification in electrodialysis models. *Chemical Engineering Research and Design*. 105. pp. 107–119. DOI: 10.1016/j.cherd.2015.10.048
- 7. Chubich, V.M. & Filippova, E.V. (2017) Information technology of active identification of stochastic dynamic systems using parameterization of the input signal. *Evraziyskoe Nauchnoe Ob"edinenie Eurasian Scientific Association*. 11(33). pp. 63–66.
- 8. Mudrov, V.I. & Kushko, V.L. (1983) *Metody obrabotki izmereniy: Kvazipravdopodobnye otsenki* [Methods of measurement processing. Quasi-likelihood estimates]. Moscow: Radio i svyaz'.
- 9. Antoniou, A. & Lu, W-S. (2007) Practical optimization: algorithms and engineering applications. New York: Springer. DOI:10.1007/978-0-387-71107-2
- 10. Izmailov, A.F. & Solodov, M.V. (2008) Chislennye metody optimizatsii [Numerical optimization methods]. Moscow: Fizmat.
- 11. Chubich, V.M. & Filippova, E.V. (2018) Research of the efficiency of some robust filters for non-stationary linear continuous-discrete systems. *Sovremennye naukoemkie tekhnologii Modern High Technologies*. 12(1). pp. 153–161.
- 12. Izanloo, R., Fakoorian, S.A., Yazdi, H.S. & Simon, D. (2016) Kalman filtering based on the maximum correntropy criterion in the presence of non-Gaussian noise. *Annual Conference on Information Science and Systems (CISS), Princeton, USA.* pp. 500–505. DOI: 10.1109/CISS.2016.7460553
- 13. Chen, B., Liu, X., Zhao, H. & Principe, J. (2017) Maximum correntropy Kalman filter. *Automatica*. 76. pp. 70–77. DOI: 10.1016/j.automatica.2016.10.004
- 14. Filippova, E.V. (2018) Algorithms of parametrical identification of stochastic linear continuous-discrete systems based on robust filtering. *Evraziyskoe Nauchnoe Ob"edinenie Eurasian Scientific Association*. 11(45). pp. 86–90.
- 15. Chubich, V.M. & Filippova, E.V. (2017) Aktivnaya identifikatsiya stokhasticheskikh dinamicheskikh sistem. Planirovanie eksperimenta dlya modeley diskretnykh sistem [Active identification of stochastic dynamic systems. Design of experiments for models of discrete systems]. Novosibirsk: NSTU.
- 16. Denisov, V.I., Chubich, V.M. & Filippova, E.V. (2018) The choice of the parameterization of the input signal in the problem of experiment design for model stochastic systems. *Izvestiya Tul'skogo gosudarstvennogo universiteta*. *Tekhnicheskie nauki Proc. of Tula State University*. *Technical Science*. 2. pp. 387–397.
- 17. Fedorov, V.V. (1971) Teoriya optimal'nogo eksperimenta [Theory of optimal experiment]. Moscow: Nauka.
- 18. Ermakov, S.M. & Zhiglyavsky, A.A. (1987) *Matematicheskaya teoriya optimal'nogo eksperimenta* [Mathematical theory of optimal experiment]. Moscow: Nauka.

- 19. Pronzato, L. & Pazman, A. (2013) Design of Experiments in Nonlinear Models. Asymptotic Normality, Optimality Criteria and Small-Sample Properties. New York: Springer. DOI: 10.1007/978-1-4614-6363-4
- 20. Kwakernaak, H. & Sivan, R. (1972) Linear Optimal Control Systems. New York: John Wiley Sons. DOI: 10.1115/1.3426828

Received: April 5, 2019

Chubich V.M., Filippova E.V. (2020) ACTIVE PARAMETRICAL IDENTIFICATION OF STOCHASTIC LINEAR CONTINU-OUS-DISCRETE SYSTEMS BASED ON THE EXPERIMENT DESIGN IN THE PRESENCE OF ABNORMAL OBSERVA-TIONS. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 61–68

DOI: 10.17223/19988605/50/8

Чубич В.М., Филиппова Е.В. АКТИВНАЯ ПАРАМЕТРИЧЕСКАЯ ИДЕНТИФИКАЦИЯ СТОХАСТИЧЕСКИХ ЛИНЕЙ-НЫХ НЕРЕРЫВНО-ДИСКРЕТНЫХ СИСТЕМ НА ОСНОВЕ ПЛАНИРОВАНИЯ ЭКСПЕРИМЕНТА ПРИ НАЛИЧИИ АНОМАЛЬНЫХ НАБЛЮДЕНИЙ. Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2019. № 50. С. 61–68

Предложена процедура активной параметрической идентификации стохастических линейных непрерывно-дискретных систем, включающая робастное оценивание параметров и оптимальное планирование входных сигналов. Рассматривается общий случай вхождения неизвестных параметров в уравнения состояния и наблюдения, начальные условия и ковариационные матрицы шумов системы и измерений. Эффективность данной процедуры продемонстрирована на примере системы управления электродвигателем постоянного тока.

Ключевые слова: непрерывно-дискретная система; аномальные наблюдения; оптимальный входной сигнал; робастное оценивание.

CHUBICH Vladimir Mikhailovich (Doctor of Technical sciences, Professor, Head of the department of theoretical and applied informatics, Novosibirsk State Technical University, Russian Federation).

E-mail: chubich@ami.nstu.ru

FILIPPOVA Elena Vladimirovna (Candidate of Technical sciences, Associate professor, Associate professor of theoretical and applied informatics, Novosibirsk State Technical University, Russian Federation).
E-mail: e.filippova@corp.nstu.ru

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

ИНФОРМАТИКА И ПРОГРАММИРОВАНИЕ

УДК 519.873:519.718.7 DOI: 10.17223/19988605/50/9

Л.А. Золоторевич

АППАРАТНАЯ ЗАЩИТА ЦИФРОВЫХ УСТРОЙСТВ

Рассматривается задача защиты проектов цифровых устройств на структурном уровне от вредоносного искажения и нарушения авторских прав. Предлагается алгоритм управляемого кодирования комбинационных структур на основе применения методов и средств тестового диагностирования. Алгоритм не требует моделирования неисправностей устройства в явном виде, что сокращает объем вычислительных процедур при кодировании схемы. Приводятся особенности проектирования современных СнК. Акцентируется внимание на необходимости создания и развития общего подхода к рассмотрению задач контроля и верификации проектов (таксономии отклонений). Таксономия отклонений включает анализ ошибок, возникающих непосредственно в процессе проектирования, и преднамеренных искажений на этапах проектирования и изготовления. Ключевые слова: цифровая экономика; СБИС; защита авторских прав; кодирование логических схем.

Акцент на цифровую экономику, приоритетная разработка цифровых технологий требуют постоянного совершенствования теории и практики проектирования интегральных схем, систем на кристалле (СнК) как технической базы создания электронных систем различного назначения.

Развитие технологии СБИС, СнК определяющим образом зависит от развития методов и качества применяемых средств автоматизированного проектирования (САПР) [1], в особенности методов и средств контроля, верификации, построения тестов контроля функциональных блоков и систем. Сложность решения указанных задач постоянно возрастает из-за возрастания сложности проектируемых объектов, отсутствия общего подхода к рассмотрению ошибок, вносимых в проект при проектировании, неисправностей реальных объектов, корреляции разного типа ошибок проектирования и неисправностей структурных реализаций. Все проблемы, связанные с разработкой методов и созданием средств верификации проектов и построения тестов контроля объектов в разных классах неисправностей, систем функционального контроля, являются достаточно сложными, но естественными, возникающими непреднамеренно, и должны решаться в режиме благоприятствующего проектирования. Однако в последние годы возникла потребность в дополнительном контроле проектов на предмет несанкционированного внедрения с целью их искажения с разными основополагающими целями. Подобные действия являются преднамеренными и тщательно скрываемыми, что препятствует прямому применению существующих методов тестирования и функционального контроля СБИС.

В связи с этим стала очевидной необходимость защиты проектов на основе создания общего подхода к контролю СБИС, СнК, таксономии нарушений и отклонений, с моделями которых приходится работать при проектировании и организации контроля на всех этапах жизненного цикла цифровой системы с учетом злонамеренных внедрений в цикл проектирования и производства интегральных схем.

Как развитие теории контролепригодного проектирования (Design-for-Testability – DfT) в работе [2] предлагается подход к проектированию Design-for-Trust – DfTr, который дополнительно включает средства для контроля и предотвращения аппаратных атак при проектировании и изготовлении СБИС.

1. Источники угроз в области производства аппаратного обеспечения

В связи с быстрыми темпами роста объемов производства цифровых устройств в настоящее время особую остроту приобретает проблема нарушения авторских прав [2]. Рост степени интеграции интегральных схем и вместе с этим высокая стоимость эксплуатации кремниевых производств расширяет аутсорсинг, который стал важной тенденцией в производстве интегральных схем.

Ущерб от пиратства и других угроз в области производства аппаратного обеспечения составляет около 4 млрд долларов в год, что примерно в 10 раз превышает ущерб от пиратства в области ПО [3]. Кроме пиратства появляются новые виды угроз [4]: внедрение в проект дополнительных вредоносных несанкционированных операций с различной основополагающей целью, изменяющих функциональное наполнение системы; внедрение механизмов деградации схемных решений с целью нарушения системы синхронизации, приводящих к нарушению временной согласованности путей распространения сигналов и в конечном итоге к сбою системы; включение средств для получения конфиденциальной информации (к примеру, получение криптографических ключей) через порты контроля и др.

В работе [5] проанализированы различные модели процесса злонамеренного искажения проекта, описывающие условия, при которых подобное искажение может внедриться в цифровую систему. В числе возможных источников искажений рассматриваются поставщики базовых функциональных блоков интеллектуальной собственности (IP's), которые приобретаются разработчиками СнК, собственно разработчики СнК, а также кремниевые фабрики – изготовители СнК.

В связи с тем, что искажения в проекте могут происходить на разных этапах проектирования, – на RTL-уровне, на уровне структурного описания схем (уровень netlist), в топологическом проекте, существует потребность в разработке методов обнаружения искажений на разных уровнях абстракции. Одной из известных методик защиты исходных кодов программ от обратного проектирования является функциональная обфускация. К сожалению, эффект от применения методов обфускации в случае языка VHDL ограничен, поскольку результаты их применения не приводят к изменению конечного результата синтеза, так как структурные реализации устройств до и после обфускации выглядят одинаково [3].

2. Обфускация и логическое кодирование цифрового устройства на структурном уровне

Одним из методов блокирования попыток внешнего вмешательства в проект цифровой системы на структурном уровне является логическое кодирование структурной реализации, которое обеспечивает доступ к объекту только авторизованным пользователям [7]. Метод предполагает сокрытие функциональности проекта и использование ключа, применение которого выводит систему в область правильного функционирования. Кроме логического шифрования комбинационной схемы известен метод внедрения новых внутренних состояний в граф перехода для последовательностных устройств, эффективность практического применения которого, к сожалению, не установлена [8].

Первый метод основан на включении в логическую сеть дополнительных вентилей, управляемых внешними логическими ключами, т.е. на применении обфускации структуры объекта. В такой постановке если злоумышленник не владеет ключом, то ему недоступна внутренняя реализация объекта. Поэтому задача структурной обфускации и логического кодирования заключается в том, чтобы затруднить или сделать невозможным получение правильного ключа.

Чтобы защитить комбинационную схему с помощью k-разрядного ключа, предлагается простая процедура, которая требует включения в схему k дополнительных вентилей [7]. Выбор линии для включения вентиля, тип вентиля существенно влияют на эффективность кодирования. На рис. 1, a приведен фрагмент логической схемы, а на рис. 1, b проиллюстрирована основная идея логического кодирования. Выход элемента C_1 отключен от нагрузки (элементы D_1 и D_2) и подключен к одному из входов дополнительного «ключевого» элемента типа XOR CC_1 , на второй вход которого поступает

внешний входной сигнал K_1 однобитового ключа. Схема будет работать в требуемом режиме только в том случае, если сигнал на входе K_1 будет равен 0. В противном случае на выходе элемента XOR CC_1 будет формироваться сигнал, инверсный правильному.

Вместо элемента CC_1 типа XOR может быть установлен элемент XNOR. В этом случае однобитовый правильный ключ, поступающий на вход K_1 , равен 1. Заметим, что применение неправильного ключа равносильно появлению неисправности константного типа const 0 (const 1) на выходе элемента C_1 в зависимости от входного набора и истинного значения сигнала на C_1 , равного 1 (0). Этот факт является важным, так как позволяет формализовать задачу обфускации на основе применения методов и средств тестового контроля цифровых устройств.

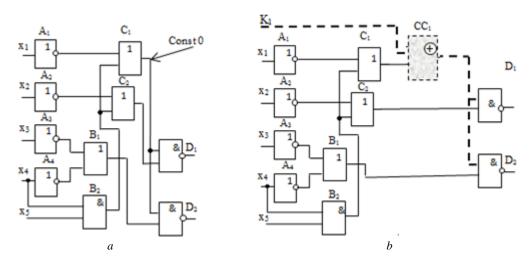


Рис. 1. Фрагмент логической сети: a – исходная комбинационная схема; b – схема с однобитовым ключом Fig. 1. Fragment of a logical network: a – the original combinational circuit; b – scheme with a one-bit key

При воздействии входного набора X=(00000) и неправильного ключа $K_1=1$ (см. рис. 1) на выходах схемы D_1 , D_2 формируются сигналы (11), в то время как при правильном ключе $K_1=0-(00)$. Так же поведет себя схема (см. рис. 1) при неисправности const 0 на выходе элемента C_1 . То есть входной набор X=(00000) является тестом контроля данной неисправности и в то же время при отсутствии неисправности искажает выходное состояние схемы при подаче неправильного ключа.

Таким образом, для сокрытия функциональности схемы необходимо добавить в некоторые линии схемы дополнительные элементы и определить правильный код, искажение которого выводит схему из области правильного функционирования. Заметим, что при воздействии входного набора X = (01110) и неправильного ключа $K_1 = 1$ (см. рис. 1) на выходах схемы D_1 , D_2 появятся сигналы (11) как и при правильном ключе, так как входной набор X = (01110) не является тестом контроля неисправности const 0 на выходе элемента C_1 .

Основная задача, которая должна быть решена при практической реализации данной общей идеи, заключаются в том, чтобы определить оптимальное множество внутренних линий схемы и количество ключевых элементов для создания максимальных трудностей для злоумышленника по поиску правильного ключа.

Положим, что цифровое устройство состоит из n первичных входов, m первичных выходов и k бит ключа шифрования. При воздействии входного вектора $X \in 2^n$ на выходах устройства формируется соответствующий правильный выходной вектор $Z \in 2^m$. Пусть $K \in 2^k$ – правильные значения ключевых сигналов (правильный ключ). Возможны два сценария функционирования устройства при разных значениях переменных шифрования:

- 1) при использовании действительного секретного ключа К функция производит правильные выходы для всех тестовых шаблонов ввода;
- 2) при использовании неправильных значений секретных ключей функция генерирует неправильные выходы соответственно:

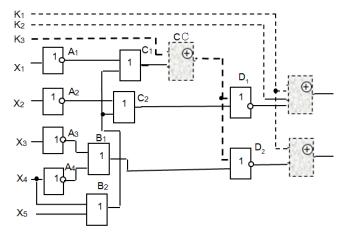
$$F(x,k) = \begin{cases} Z \forall X \in 2^{n}, Z \in 2^{m}, \\ Z \forall X \in 2^{n}, Z \in 2^{m}, Z \neq Z, \end{cases}$$

здесь Z(Z) – правильный (неправильный) выходной вектор.

Для определения степени защищенности устройства при его кодировании принимается расстояние Хэмминга (HD) — число, используемое для обозначения меры различия между двумя двоичными последовательностями. HD позволяет количественно определить степень отличия правильной реакции устройства от ошибочной. Если HD (Z, Z) = 0, то это означает, что реакция закодированной схемы не зависит от ключа блокировки. При HD (Z, Z) = m, Z дополняет Z, что упрощает злоумышленнику поиск правильного ключа. Для того чтобы затруднить восстановление правильного ключа, необходимо обеспечить наименьшую корреляцию между правильными и неправильными выходными векторами, что достигается при HD (Z, Z) = m/z, когда на каждом входном воздействии около 50% выходных сигналов в случае применения неправильного ключа принимает логические значения, инверсные правильным.

3. Применение методов и средств тестового диагностирования для защиты цифровых устройств от вредоносных искажений

При включении очередного вентиля при кодировании логических устройств необходимо проводить анализ на появление эффекта маскирования неисправностей, который способен блокировать эффект кодирования. В работе [7] при кодировании логических устройств ключевые вентили помещались в схему случайным образом. При таком подходе применение неправильного ключевого бита не гарантирует появления неправильного выходного сигнала (рис. 2) и не может требуемым образом затруднить злоумышленнику доступ к структуре устройства. Во-первых, возможен эффект маскирования неисправностей, что показано на рис. 2. Схема, зашифрованная тремя битами ключа K_1 , K_2 , K_3 на рис. 2, на входном наборе 00000 как при подаче правильного ключа 000, так и при неправильном ключе 111 вырабатывает одинаковую выходную реакцию 00. Это происходит по причине маскирования неисправностей const 0, которые одновременно возникают на выходах элементов C_1 , D_1 и D_2 . Во-вторых, для некоторых линий отсутствует возможность активизации пути от данной линии к выходам устройства.



Puc. 2. Влияние маскирования неисправностей на результаты кодирования Fig. 2. Impact of masking faults on coding results

На рис. З приведена структура цифрового устройства , реализующего систему булевых функций $D_1 = \overline{x_1}x_3x_4x_5 \vee \overline{x_2}x_3x_4x_5$; $F_1 = \overline{x_1}\overline{x_3} \vee \overline{x_2}\overline{x_3} \vee \overline{x_1}\overline{x_4} \vee \overline{x_2}\overline{x_4} \vee x_6 \vee x_7$. Как было сказано выше, кодирование схемы путем случайного подбора мест вставки в структуру ключевых вентилей оказывается недостаточно эффективным. К примеру, добавление вентиля XOR на выходе элемента B_3 не принесет

ожидаемого эффекта, так как для неисправности const 0 на выходе B_3 не существует проверяющего теста, и применение неправильного ключа, равного 1, не изменит реакции схемы при подаче любой входной последовательности. Поэтому при кодировании структуры устройства необходимо отслеживать эффективность каждого шага. При решении основной задачи — затруднить злоумышленнику доступ к структурной реализации устройства — необходимо обеспечить оптимизацию объема необходимого дополнительного оборудования, учесть влияние задержек дополнительно включенных элементов на функционирование устройства.

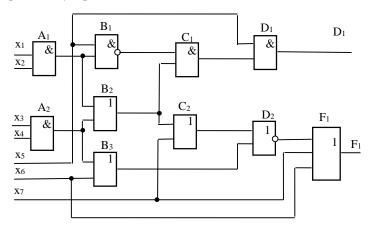


Рис. 3. Структурная реализация цифрового устройства Fig. 3. Structural implementation of a digital device

В работе [9] предложен подход к определению множества линий структуры для кодирования, основанный на моделировании схемы с внесенной i-й неисправностью и вычислении признака $P_i = X_i \times Y_i$, характеризующего линию с точки зрения эффективности ее выбора при кодировании схемы. Здесь X_i — количество входных наборов, которые покрывают анализируемую неисправность, Y_i — количество выходных переменных, которые искажаются при появлении данной неисправности. По результатам анализа полученных признаков определяется множество внутренних линий схемы для кодирования.

Очевидно, что данный подход требует моделирования схемы $M=2s\times 2^n$ раз, где s – общее количество линий схемы (переменных полного состояния схемы), n – количество входных переменных схемы. Для схемы на рис. $3 M = 128\times 34 = 4352$. Для реальных схем подобный подход практически неприемлем по причине высоких вычислительных затрат. С целью оптимизации вычислительных процедур предлагается эвристическое решение – сократить количество моделируемых входных наборов до 100 [9] (в этом случае M = 200k)

Сведем задачу кодирования к поиску неисправностей константного типа кодируемой структуры, обнаруживаемых на большем количестве выходных линий на максимальном количестве входных векторов.

В отличие от решения, принятого в работе [9], рассмотрим более эффективный подход, основанный на применении метода сквозного вычисления неисправностей, покрываемых рассматриваемым входным вектором (метод конкурентно-дедуктивного моделирования) вместо моделирования каждой неисправной модификации схемы на определенном множестве случайных входных наборов с целью оценки степени влияния неисправностей на выходы схемы [10]. Метод конкурентно-дедуктивного моделирования неисправностей основан на моделировании исправной схемы и позволяет за один проход моделирования определять все неисправности константного типа, обнаруживаемые на моделируемом входном наборе. За счет того, что моделируется только исправная схема, эффективность решения существенно повышается по сравнению с моделированием одиночной неисправности на множестве входных векторов.

Вначале вычисляются неисправности, обнаруживаемые на моделируемом ограниченном множестве случайных входных наборов. Затем по результатам анализа определяются те неисправности,

которые обнаруживаются наибольшим числом наборов и указывают преимущественные линии схемы для вставки ключевых вентилей. В то же время численное ограничение количества моделируемых входных воздействий [12] ограничивает возможность поиска наиболее эффективного решения.

Здесь предлагается другой подход, основанный на построении теста в классе неисправностей константного типа [10] и его применении на первом этапе кодирования. В рамках данного подхода вместо использования заранее определенного числа случайных входных воздействий (как, например, 100 в работе [11]) применятся тестовая последовательность входных векторов, которая обеспечивает близкое к полному покрытию неисправностей константного типа кодируемой структуры.

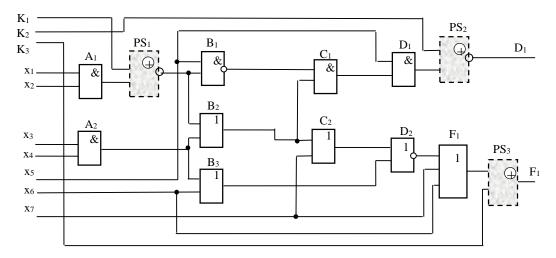
В табл. 1 приведены результаты построения теста для схемы, приведенной на рис. 3, и соответствующие разностные неисправные функции. Первый столбец таблицы содержит входные наборы теста, последующие – идентификаторы неисправностей константного типа всех линий схемы и единичные значения разностных неисправных функций, реализуемых на соответствующем выходе схемы. Здесь X_1^0 – неисправность типа const 0 на входе X_1 , а A_1^1 – неисправность типа const 1 на выходе элемента A_1 . Верхний индекс при единичном значении разностной неисправной функции указывает, на каком выходе схемы реализуется данная функция. В данном случае, значение 1^1 относится к функции, реализуемой на первом выходе схемы, т.е. на выходе элемента D_1 (рис. 3).

Таблица разностных неисправных функций для схемы на рис. 3

Неисправности																	
	X_1^0	X_1^1	X_2^0	X_2^1	X_3^0	X_3^1	X_4^0	X_4^1	X_5^0	X_{5}^{1}	X_6^0	X_6^1	X_7^0	X_7^1	A_1^0	A_1^1	A_2^0
Тест-векторы																	
1010100				12				1^{1}								1 ²	
1011100				11	11		11		1 ¹			12		12		1 ¹	1^{1}
1101100	12		12									12		12	12		
0101111		12				12										1 ²	
0111010										11	12						
0011101					1^{1}		1^{1}		1^{1}				12			1 ¹	1^{1}
Неисправности																	
Неисправности	A_2^1	B ₁ ⁰	B ₁ ¹	$\mathbf{B_2}^0$	B_2^1	B ₃ ⁰	B ₃ ¹	C ₁ ⁰	C_1^1	C_2^0	C_2^1	$\mathbf{D_1}^0$	$\mathbf{D_1}^1$	$\mathbf{D}_2{}^0$	\mathbf{D}_2^1	F ₁ ⁰	F ₁ ¹
Неисправности Тест-векторы	A_2^1	B ₁ ⁰	B ₁ ¹	$B_2{}^0$	B_2^1	B ₃ ⁰	B ₃ ¹	C ₁ ⁰	C_1^1	C_2^0	C_2^1	$\mathbf{D_1}^0$	$\mathbf{D_1}^1$	$\mathbf{D}_2{}^0$	\mathbf{D}_2^1		F ₁ ¹
	A ₂ ¹	B ₁ ⁰	B ₁ ¹	B ₂ ⁰	B ₂ ¹	B ₃ ⁰	B_3^1 12	C ₁ ⁰	C ₁ ¹	C ₂ ⁰	C_{2}^{1} 1^{2}	$\mathbf{D_1}^0$	D ₁ ¹	D_2^0	\mathbf{D}_2^1		F ₁ ¹
Тест-векторы		B ₁ ⁰	B ₁ ¹	B ₂ ⁰		B ₃ ⁰		C ₁ ⁰		C2 ⁰		D ₁ ⁰			D ₂ ¹	F ₁ ⁰	F ₁ ¹
Тест-векторы 1010100			B ₁ ¹			B ₃ ⁰				C ₂ ⁰					-	F ₁ ⁰	
Тест-векторы 1010100 1011100				1 ¹		B ₃ ⁰			11				11		1 ²	F ₁ ⁰	12
Тест-векторы 1010100 1011100 1101100	11			1 ¹	11	B ₃ ⁰	12		11		12		1 ¹	12	1 ²	F ₁ ⁰	12

Первая строка таблицы содержит идентификаторы неисправностей, последующие — единичные значения разностных неисправных функций. Верхний индекс в обозначении разностной неисправной функции (1^2) указывает, что функция относится ко второму выходу схемы, т.е. F_1 . Если неисправность обнаруживается не на одном, а, к примеру, на трех выходах, то верхний индекс может иметь вид $1^{2, 3, 5}$. Из табл. 1 видно, что размещение ключевого вентиля XOR на выходе элемента B_3 не имеет смысла, так как теста контроля неисправности const 0 на выходе элемента B_3 не найдено по причине его отсутствия. Наиболее целесообразно выбрать вначале для последующего кодирования выходы элементов A_1 , D_1 , F_1 , так как столбцы, соответствующие неисправностям A_1^1 , D_1^1 , F_1^0 данных элементов, содержат большее число единичных значений разностных неисправных функций. Это свидетельствует о том, что большее число входных векторов в случае применения неправильного ключа приведет к искажению реакции схемы.

На рис. 4 приведена схема с внесенными ключевыми элементами PS_1 , PS_2 , PS_3 и ключевыми входами K_1 , K_2 , K_3 . В схеме ключевой элемент PS_3 имеет тип XOR, так как неисправность const 0 на выходе элемента F_1 обнаруживается большим числом входных сигналов по сравнению с неисправностью const 1. Ключевые элементы PS_1 и PS_2 имеют тип XNOR, так как соответствуют столбцам с неисправностями типа const 1.



Puc. 4. Схема с вентилями PS₁, PS₂ и PS₃ для логического шифрования Fig. 4. Circuit with PS₁, PS₂, and PS₃ gates for logical encryption

После добавления ключевых элементов в структуру необходимо проанализировать полученные результаты кодирования моделированием полученной частично закодированной структуры на наборах теста на всем булевом интервале множества ключевых входов и сравнением в каждом случае выходных реакций схемы с результатами моделирования исходной схемы. Как было указано выше, для максимального затруднения доступа к получению структуры схемы необходимо обеспечить кодовое расстояние Хэмминга между выходными состояниями схемы в условиях применения правильных и ошибочных ключевых кодов, близкое к 0,5 [7].

4. Управляемое кодирование цифровых устройств на структурном уровне

Очевидно, что результат кодирования проявляется на выходах схемы в зависимости от числа неправильных битов кода [9]. Если ключевой вентиль управляется одним битом ключевого кода, то вероятность того, что данный вентиль будет приведен в действие, P=0,5. Это означает, что только половина ключевых вентилей повлияет на результат функционирования схемы при применении неправильного ключа. Для того чтобы увеличить вероятность P и усилить влияние неправильного бита кодового слова на результат функционирования схемы, применим управляющие вентили, с помощью которых можно объединить биты кодового слова в группы, использовав при этом их выходы в качестве входов ключевых вентилей. В таком случае будет реализовано групповое воздействие нескольких битов кодового слова на активизацию ключевого вентиля. Если хотя бы один из ключевых входов, включенных в группу, принимает неправильное значение, ключевой вентиль окажется активированным. Для этого с каждым ключевым вентилем используется управляющий вентиль. При этом, если применяется двухвходовый управляющий вентиль, то вероятность активизации ключевого вентиля возрастает с 0,5 до 0,75, в случае трехвходового вентиля вероятность составляет 0,88, а при пятивходовом -0,97 (только один ключевой вектор из 32 векторов данной группы является правильным).

На рис. 5, a приведена структура схемы с тремя выходами, а в табл. 2 – тестовая последовательность и соответствующие разностные неисправные функции. На рис. 5, b приведен пример двухуровневого кодирования. В соответствии с результатами табл. 2 в качестве линий для первоочередного включения ключевых вентилей для кодирования выбраны выходы элементов A2 (вентиль PS_1 типа XNOR) и A3 (вентиль PS_2 типа XOR). Тип ключевого вентиля XNOR на выходе элемента A2 выбирается в соответствии с неисправностью const 1, которая покрывается четырьмя из семи входными векторами и обнаруживается на двух из трех выходов. Выбор неисправности A_3^0 обусловлен тем, что по сравнению с неисправностью C_2^1 неисправность A_3^0 очувствляет (приводит к изменению) два выхода.

Дополнительно в схему включены управляющие двухвходовые вентили KK_1 и KK_2 , которые усилили влияние на функционирование схемы каждого бита ключевого входа.

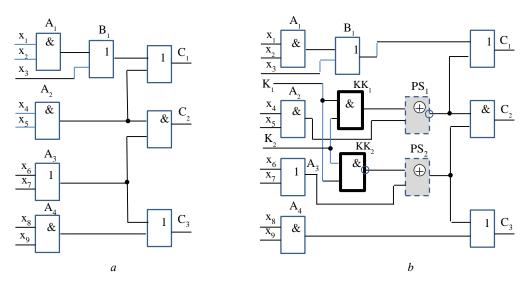


Рис. 5. Пример схемы с двухуровневым кодированием:

a – логическая структура с тремя выходами; b – двухуровневое кодирование схемы

Fig. 5. An example of a scheme with two-level coding: a – logical structure with three outputs; b – two-level coding scheme

Таблица разностных неисправных функций для схемы на рис. 5

Неисправности																	
	X_1^0	X_1^1	X_2^0	X_2^1	X_3^0	X_3^1	X_4^0	X_4^1	X_{5}^{0}	X_{5}^{1}	X_6^0	X_6^1	X_7^0	X_7^1	X_{8}^{0}	X_{8}^{1}	X_9^0
Тест-векторы																	
100100101				1^{1}		11				11,2			13				
110010001	11		11									13		13		13	
001111001							12		12		$1^{2,3}$						
000010011						11		11							13		1^{3}
000110010							1^{1}		1^{1}			$1^{2,3}$		$1^{2,3}$			
101000110					11								13				
010001100		11				11											
Неисправности																	
	X_9^1	A_1^0	A_1^1	A_2^0	A_2^1	A_3^0	A_3^1	A_4^0	A_4^1	$\mathbf{B_1}^0$	$\mathbf{B_1}^1$	C_1^0	C_1^1	C_2^0	C_2^1	C_3^0	C_3^1
Тест-векторы																	
100100101			11		11,2	13					11		11		12	13	
110010001		11					13		13	1^{1}		11			12		1^{3}
001111001				12		$1^{2,3}$						1^{1}		12		1^{3}	
000010011			1^{1}		11			1^{3}			1^{1}		1^{1}		12	1^{3}	
000110010	13			1^{1}			$1^{2,3}$		13			1^{1}			1^{2}		13
101000110					12	13				1^{1}		11			12	13	
010001100			11		11,2	13					1^{1}		11		12	13	

Ниже приводятся основные этапы алгоритма управляемого логического кодирования комбинационных структур при использовании двухвходовых управляющих вентилей.

Исходные данные: описание кодируемой структуры схемы. **Результаты:** описание закодированной структуры схемы; правильный ключ.

- 1) построить тест контроля структуры в классе неисправностей константного типа методом случайного поиска на основе применения метода конкурентно-дедуктивного моделирования неисправностей;
- 2) упорядочить множество FN обнаруживаемых на наборах теста неисправностей по убыванию числа покрывающих входных наборов и активизированных выходов схемы;
 - 3) J := 1;
- 4) из множества FN выбрать j-ю неисправность; в соответствии с типом неисправности включить в структуру схемы ключевой элемент (типа XOR, если неисправность const 0, и элемент XNOR,

если неисправность const 1); включить управляющий вентиль с ключевым входом k_i ; на второй вход управляющего вентиля подключить случайным образом дополнительный ключевой вход;

- 5) моделировать полученную структуру на всех наборах теста при всех возможных комбинациях значений ключа;
- 6) анализировать кодовое расстояние Хэмминга между реакциями исходной схемы и частично закодированной при неправильных битах ключа;
 - 7) если результат анализа кодирования неудовлетворителен, то J := J + 1; перейти к п. 4;
 - 8) выход.

Заключение

В работе акцентирована необходимость развития таксономии отклонений, возникающих по разным причинам в проектах СБИС типа СнК на разных этапах проектирования и изготовления.

Предложен алгоритм управляемого кодирования описаний цифровых устройств комбинационного типа на структурном уровне на основе применения средств тестового контроля. Предложенный алгоритм по сравнению с известными в литературе алгоритмами требует меньших вычислительных затрат и времени и проявляет устойчивость к восстановлению правильного ключа на основе «атаки SAT» [11]. Это обусловлено тем, что ключевые входы не связаны напрямую с ключевыми вентилями, а ключевые вентили активизируются не одним ключевым входом. Применение метода сквозного вычисления множества покрываемых неисправностей на основе моделирования исправной схемы существенно сокращает объем вычислительных процедур.

ЛИТЕРАТУРА

- 1. Zolotorevich L.A. Project verification and construction of superchip tests at the RTL level // Automation and Remote Control. 2013. V. 74, is. 1. P. 113–122.
- 2. Rajendran J., Sam M., Sinanoglu O., Karri R. Security analysis of integrated circuit camouflaging // ACM SIGSAC conference on Computer & communications security. Germany, Berlin. 04–08 November 2013. P. 709–720.
- 3. Сергейчик В.В., Иванюк А.А. Методы лексической обфускации VHDL-описаний // Information Technologies and Systems 2013 (ITS 2013): Proc. of The Int. Conference. BSUIR. Minsk, 2013. C. 198–199.
- 4. Shakya B., Salmani T.H., Forte D., Bhunia S., Tehranipoor M. Benchmarking of hardware Trojans and maliciously affected circuits // J. Hardw. Syst. Secur. (HaSS). 2017. V. 1 (1). P. 85–102.
- 5. Xiao K., Forte D., Jin Y., Karri R., Bhunia S., Tehranipoor M. Hardware Trojans: Lessons learned after one decade of research // ACM transactions on design automation of electronic system. 2016. V. 22, No. 1. P. 1–23.
- 6. Dupuis S., Rouzeyre B., Flottes M.-L., Natale G.D., Ba P.-S. New Testing Procedure for Finding Insertion Sites of Stealthy Hardware Trojans // DATE: Design, Automation and Test in Europe. Grenoble, 2015. P. 776–781.
- 7. Roy J.A., Koushanfar F., Markov I.L. EPIC: Ending Piracy of Integrated Circuits // IEEE Computer. 2010. V. 43, No. 10. P. 30–38.
- 8. Chakraborty R.S., Bhunia S. Security against Hardware Trojan through a Novel Application of Design Obfuscation // IEEE/ACM Int. Conference on Computer-Aided Design. 2009. P. 113–116.
- 9. Karousos N., Pexaras K., Karybali I.G., Kalligeros E. Weighted Logic Locking: a New Approach for IC Piracy Protection // IEEE 23rd Int. Symposium on On-Line Testing and Robust System Design (IOLTS). 2017. P. 221–226.
- 10. Золоторевич Л.А. Исследование методов и средств верификации проектов и генерации тестов МЭС // Проблемы разработки перспективных микроэлектронных систем (МЭС-2006) : сб. науч. тр. всерос. науч.-техн. конф. / под общ. ред. А.Л. Стемпковского. М. : ИППМ РАН, 2006. С. 163–168.
- 11. Yasin M., Rajendran J., Sinanoglu O., Karri R. On Improving the Security of Logic Locking // IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems. 2016. V. 35, No. 9. P. 1411–1424.

Поступила в редакцию 29 мая 2019 г.

Zolotorevich L.A. (2020) HARDWARE PROTECTION OF DIGITAL DEVICES. *Vestnik Tomskogo gosudarstvennogo universiteta*. *Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 69–78

DOI: 10.17223/19988605/50/9

In the last decade, the problem of protection and additional control of VLSI projects with the aim of detecting the consequences of unauthorized third-party interference in a project with different fundamental goals became urgent. Such actions are deliberate and carefully hidden, which prevents the direct application of existing methods for testing and functional control of VLSI.

The task of protecting projects of digital devices at a structural level from malicious misrepresentation and copyright infringement is considered. The algorithm of controlled coding of combinational structures, based on the use of methods and tools for test diagnostics, is proposed. The algorithm does not require the simulation of device malfunctions in an explicit form, which reduces the number of computational procedures for encoding the circuit. The features of SoC design are considered. Attention is focused on the need to create and develop a unified approach to reviewing the tasks of monitoring and verification of projects (taxonomy of deviations). Taxonomy deviations include the analysis of errors that occur directly in the design process, and deliberate distortion during the design and manufacturing stages.

Keywords: digital economy; VLSI design; copyright protection; coding logic circuits.

ZOLOTOREVICH Ludmila Andreevna (Candidate of Technical sciences, docent, Belarusian State University of Informatics and Radioelectronics, Minsk, Republic of Belarus).

E-mail: zolotorevichla@bsuir.by

REFERENCES

- 1. Zolotorevich, L.A. (2013) Project verification and construction of superchip tests at the RTL level. *Automation and Remote Control*. 74(1). pp. 113–122. DOI: 10.1134/S0005117913010104
- 2. Rajendran, J., Sam, M., Sinanoglu, O. & Karri, R. (2013) Security analysis of integrated circuit camouflaging. *ACM SIGSAC conference on Computer & communications security*. Germany. Berlin. pp. 709–720.
- Sergeychik, V.V. & Ivanyuk A.A. (2013) Metody leksicheskoy obfuskatsii VHDL-opisaniy [Methods of lexical obfuscation of VHDL-Descriptions]. Information Technologies and Systems (ITS 2013): Proc. of The Int. Conference. BSUIR. Minsk. pp. 198–199
- 4. Shakya, B., Salmani, T.H., Forte, D., Bhunia, S. & Tehranipoor, M. (2017) Benchmarking of hardware Trojans and maliciously affected circuits. *Journal of Hardware and System Securuty (HaSS)*. 1(1), pp. 85–102. DOI: 10.1007/s41635-017-0001-6
- 5. Xiao, K, Forte, D, Jin, Y, Karri, R., Bhunia, S. & Tehranipoor, M. (2016) Hardware Trojans: Lessons learned after one decade of research. *ACM transactions on design automation of electronic system.* 22(1), pp.1–23. DOI: 10.1145/2906147
- 6. Dupuis, S., Rouzeyre, B., Flottes, M.-L., Natale, G.D. & Ba, P.-S. (2015) New Testing Procedure for Finding Insertion Sites of Stealthy Hardware Trojans. *DATE: Design, Automation and Test in Europe*. France. Grenoble. March 9–13, 2015. pp. 776–781.
- 7. Roy, J.A., Koushanfar, F. & Markov, I.L. (2010) EPIC: Ending Piracy of Integrated Circuits. IEEE Computer. 43(10). pp. 30–38.
- 8. Chakraborty, R.S. & Bhunia, S. (2009) Security against Hardware Trojan through a Novel Application of Design Obfuscation. *IEEE/ACM International Conference on Computer-Aided Design*. pp. 113116.
- 9. Karousos, N., Pexaras, K., Karybali, I.G. & Kalligeros, E. (2017) Weighted Logic Locking: A New Approach for IC Piracy Protection. *IEEE 23rd Int. Symposium on On-Line Testing and Robust System Design (IOLTS)*. pp. 221–226.
- 10. Zolotorevich, L.A. (2006) Issledovanie metodov i sredstv verifikatsii proektov i generatsii testov MES [Research of methods and means of project verification and generation of MES tests]. In: Stempkovsky, A.L. (ed.) Problemy razrabotki perspektivnykh mikroelektronnykh sistem (MES-2006) [Problems of Microelectronic System Development (MIC-2006)]. Moscow: RAS. pp. 163–168.
- 11. Yasin, M., Rajendran, J., Sinanoglu, O. & Karri, R. (2016) On Improving the Security of Logic Locking. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 35(9). pp. 1411–1424. DOI: 10.1109/TCAD.2015.2511144

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

УДК 004.94

DOI: 10.17223/19988605/50/10

О.С. Исаева, Н.В. Кулясов, С.В. Исаев

МЕТОД СТРУКТУРНО-ГРАФИЧЕСКОГО АНАЛИЗА И ВЕРИФИКАЦИИ ИНТЕЛЛЕКТУАЛЬНОЙ ИМИТАЦИОННОЙ МОДЕЛИ

Исследование выполнено при финансовой поддержке РФФИ и Правительства Красноярского края в рамках научного проекта № 18-47-242007.

Представлен метод анализа интеллектуальной модели имитации функционирования бортовой аппаратуры космического аппарата. Модель состоит из графической структуры, отражающей состав элементов бортовой аппаратуры, и базы знаний, описывающей методы ее работы. Выполнена формализация модели и предложены критерии анализа и верификации ее структуры и свойств. Разработаны визуальные компоненты интерактивной инфографики, выполняющие интерпретацию формального описания модели в интерактивные графические образы и формирующие перечень ошибок в функциональных зависимостях базы знаний.

Ключевые слова: имитационное моделирование; бортовая аппаратура космического аппарата; верификация; валидация; базы знаний; инфографика.

Конструирование современных космических аппаратов — это наукоемкий и дорогостоящий процесс, связанный с множеством разнообразных и нередко трудно формализуемых факторов, оказывающих влияние на результаты разработки. Компьютерное моделирование является основным научно обоснованным методом исследования характеристик сложных систем, используемым для принятия решений в различных сферах инженерной деятельности [1]. Существующие и проектируемые системы можно эффективно исследовать с помощью имитационных моделей, выступающих в качестве инструмента экспериментатора взамен проведения дорогостоящих и трудоемких исследований реальных объектов [2]. Выбор метода моделирования и необходимая детализация моделей существенно зависят от этапа жизненного цикла сложной технической системы. Ранние этапы проектирования систем связаны с отсутствием достоверных данных о методах их функционирования. Применяемые в этом случае модели носят описательный характер и преследуют цель наиболее полно представить в компактной форме опыт и знания экспертов предметной области об объекте исследования. Имитационные системы, ориентированные на использование экспертных знаний при решении технологических и функциональных задач, называют интеллектуальными системами имитации [3].

Наиболее распространенный класс интеллектуальных систем, ориентированных на знания как способ тиражирования опыта высококвалифицированных специалистов, — это экспертные системы [4]. Применение экспертных систем способно обеспечить механизмы построения информационной памяти предприятия [5]. Экспертные системы позволяют описывать знания о динамическом поведении анализируемых объектов и используются на этапах построения концепции космической миссии, проектирования оборудования, проверки его работоспособности, а также в процессе эксплуатации для контроля и диагностики отказов [6].

В традиционных исследованиях контроль качества имитационных моделей основан на статистическом анализе и оценке ошибок моделирования, однако при отсутствии данных функциональных испытаний традиционный подход не может быть применен [7]. В этом случае используются методы, основанные на привлечении качественного опыта экспертов предметной области [8]. В работах [7–11] рассматриваются подходы к оценке качества имитационных моделей. В [7] предложен метод валидации, основанный на применении эмпирических данных и знаний, полученных из смежных областей. В [9] решается задача проверки непротиворечивости знаний на основе анализа матрицы инцидентности графа, построенной для обобщенных отношений следования между целевыми установ-

ками базы знаний. В [10] предложен интеллектуальный метод проверки модели, основанный на анализе сходства между временными рядами моделирования из компьютеризированной модели и наблюдаемыми временными рядами из реальной системы. В [11] верификация моделей обеспечивает контроль согласованности на уровне проверки ссылочной целостности между модулями, поиска циклических зависимостей и др.

В настоящей работе решается задача создания методологии структурно-графического анализа имитационной модели, которая может использоваться при моделировании функционирования бортовой аппаратуры космических систем [12–13]. В работе рассмотрены задачи оценки качества имитационных моделей, их структур, решение которых позволяет выполнять графическую интерактивную визуализацию таких свойств модели, как полнота, адекватность, непротиворечивость. Предлагаемые методы позволят повысить качество моделей сложных систем.

1. Формальное описание модели

Интеллектуальная модель имитации функционирования бортовой аппаратуры космического аппарата представляет собой набор множеств, описывающих состояние бортовых систем и их функционирование в каждый момент времени. Модель $S = \langle G, F, T \rangle$ [14], где G — структурнопараметрическое описание (множество элементов структуры), F — функциональное описание (множество методов функционирования), T — моменты времени наблюдения.

Коммутационные интерфейсы $I_i = \{I_i^1, ..., I_i^n\}$, n — количество точек входов и выходов B_i . I_i^n имеет характеристики: тип интерфейса $Tp(I_i^n)$, направленность передачи $Rt(I_i^n)$ и признак состояния $Onf(I_i^n)$. На основе заданных характеристик коммутационных интерфейсов выполняется типизация информационных зависимостей между блоками. Типизация позволяет при проведении моделирования применять единые правила к представлению и обработке данных, передаваемых по однотипным интерфейсам, например задавать вероятности потери сигнала и ошибки передачи.

Множество коммутационных соединений $C_{ij} = \{C_{ij}^{\ l}, ..., C_{ij}^{\ nm}\} \subseteq C$ определяет пути взаимодействия пары моделей B_i и B_j (i, j = [1, ..., l]). $C_{ij}^{\ nm} = \langle I_i^n, I_j^m, \tau_{ij} \rangle$, где I_i^n – интерфейс B_i , I_j^m – интерфейс B_j , τ_{ij} – время прохождения сигнала между интерфейсами. Коммутационное соединение однозначно определяет два элемента модели, обозначим это как $C_{ij} \in B_i \times B_j$, если $\exists I_i^n \in C_{ij}, I_i^n \in C_{ij}, I_i^n \in B_i, I_i^n \in B_j$.

Функциональное описание модели представлено множеством $F = \{R(P,T)\}$, где R — множество правил функционирования. $R = \{A(I_i^n, X, K, T) \to Z(I_i^m, Y, K, T, D(K))\}$, где A — условия, при которых правило должно быть выполнено, Z — действия, вызываемые при выполнении правил. Условия и действия представляют собой выражения над параметрами, заданными в структурно-параметрическом описании, либо функции, осуществляющие изменение состояния модели. Условие A может инициироваться интерфейсом I_i^n элемента модели B_i , на который поступили данные X или команды K, или таймером, определяющим время T наступления события. В этом случае будем говорить, что $I_i^n \in A$, $X \in A$, $K \in A$, $T \in A$. Действие Z изменяет состояние модели: выполняет передачу данных на интерфейс I_i^m элемента модели B_i , изменяет множество выходных параметров Y, команд K, таймеров выполнения T или вызывает выполнение функций, определяемых командой K, таких как переключение на работу по основному или резервному каналу и пр. Для каждого элемента $B_i \in B$ всегда можно однозначно определить его параметры $P_i \in P$, интерфейсы $I_i \in I$ и правила $R_i = A_i \to Z_i$, т.е. можно говорить, что $\forall B_i \exists \{P_i\}, \{I_i\}, \{R_i\}$.

Мощностью множества будем называть количество различных элементов этого множества. Например, количество элементов модели для $B_i \in B$ обозначим как |B|. Цепочкой правил $R'(B_i) = \{R_i^1, ..., R_i^q\} \subseteq R$ назовем взаимосвязанное подмножество правил, выполняемых в процессе логического вывода в блоке модели B_i . Обозначим $L(I_i^n, I_j^m) = R'(B_i) \cup R'(B_j)$ путь от B_i к B_j через интерфейсы $I_i^n \in B_i$ и $I_j^m \in B_j$, определяемый множеством правил $R'(B_i) \cup R'(B_j)$. I_i^n — начало пути, I_j^m — окончание. Длиной пути $L(I_i^n, I_j^m)$ назовем мощность множества правил $R'(B_i) \cup R'(B_j)$, обозначим $|L(I_i^n, I_j^m)|$.

Суть построения модели сводится к созданию упрощенной структуры, свойства и поведение которой соответствуют системе-оригиналу. В зависимости от степени агрегирования структуры модели и характера обобщения методов функционирования для одного и того же оригинала можно получить несколько различных реализаций моделей.

2. Пример построения модели

Бортовые системы решают широкий круг задач, в их числе: обеспечение обмена информацией с наземным комплексом управления, измерение текущих навигационных параметров движения космического аппарата на орбите, сбор, хранение, обработка и передача телеметрической информации, управление работой систем космического аппарата и др. Для построения модели применяются графические инструменты и редактор баз знаний. Они позволяют формировать структуру модели, задавать конфигурацию функциональных блоков и коммутационных связей и описывать логику ее работы в виде правил.

Рассмотрим пример построения имитационной модели. Пусть имеется два бортовых устройства, осуществляющих информационное взаимодействие — имитатор бортового комплекса управления (БКУ) и имитатор командно-измерительной системы (КИС), т.е. $B = \{B_1, B_2\} = \{ «Имитатор БКУ», «Имитатор КИС» \}$. Предположим, одно из устройств формирует и передает данные другому, а то, в свою очередь, обрабатывает данные и возвращает их первому устройству. Пусть в нашей модели передаются аналоговые сигналы и цифровые пакеты данных.

Множество коммутационных интерфейсов первого устройства $I_1 = \{I_1^1, I_1^2, I_1^3\}$, $Tp(I_1^1) = \ll Pen \gg (1)$ (релейный интерфейс), $Tp(I_1^2) = \ll RS-232 \gg (1)$ (recommended standard 232), $Tp(I_1^3) = \ll RS-422 \gg (1)$ (recommended standard 422). $Rt(I_1^1) = \ll Ucx \gg (1)$, $Rt(I_1^2) = (1)$, $Rt(I_1^2) =$

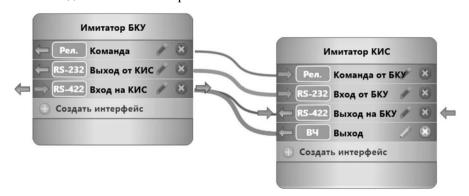


Рис. 1. Графическое представление модели Fig. 1. Graphical presentation of a model

Определим методы функционирования имитационной модели. В функциональном представлении R создадим правило $A_1 \to Z_1$, которое описывает следующие действия: устройство B_2 получает от B_1 данные на интерфейс I_2^2 и часть принятого пакета длиной в два байта возвращает в B_1 по интер-

фейсу I_2 ³. Подобным образом задается взаимодействие между блоками модели и по другим коммутационным интерфейсам. В базе знаний правила описываются в виде конструкций: «Если A то B», текстовый вид правила показан на рис. 2.



Рис. 2. Редактор базы знаний

Fig. 2. Knowledge base editor

Имитационная модель может содержать десятки функциональных блоков, имитирующих бортовые системы, и описывать методы их взаимодействие между собой и с наземным комплексом управления. После построения модели требуется провести анализ ее соответствия целям и задачам моделирования.

3. Метод структурно-графического анализа модели

Методы оценки модели в общем случае сводятся к получению информации о том, насколько хорошо модель описывает реальные процессы, происходящие в исходном объекте, и насколько качественно она будет имитировать развитие данных процессов [15]. Критерии оценки включают в себя такие аспекты существования модели, как адекватность, непротиворечивость, полнота и др. Для имитационных моделей, в основе которых лежит база знаний, такие критерии носят эмпирический характер.

Авторами данной работы предложен метод оценки соответствия имитационной модели и базы знаний предметной области в объеме, отвечающем целям моделирования. Метод имеет человекомашинную реализацию. Программное обеспечение автоматически выполняет контроль структуры модели и синтаксиса базы знаний, автоматически строит графические представления, описывающие основные элементы модели и их связи в виде интерактивных образов инфографики, формирует списки ошибок, выявленные при анализе, и текстовые рекомендации для их устранения. Метод включает анализ таких характеристик, как валидация модели и данных, верификация и полнота модели.

Валидация модели — соответствие между поведением модели и исследуемого реального объекта. Для интеллектуальных систем имитации валидация заключается в оценке логической непротиворечивости знаний. Для валидации модели требуется произвести следующие действия:

- 1. Выполнить поиск правил, которые при одинаковых условиях содержат разные действия при применении правил. Если такие правила найдены, то база знаний противоречива.
- 2. Выбрать правила, для которых в условии содержатся команды (из базы команд). Если существуют такие команды, для которых есть несколько правил, действия которых задают переходы к разным элементам модели, то база знаний содержит ошибки.
- 3. Для выбранных в п. 2 правил построить зависимые цепочки, просмотрев все переходы между правилами. Если существуют такие команды, для которых построены различные последовательности действий, то такие цепочки являются ошибочными.
- 4. Сформировать перечень найденных ошибок с указанием элементов модели, команд и правил базы знаний.

Валидация данных заключается в анализе реакций модели на изменения входных параметров. Валидация проводится с помощью имитационных экспериментов при штатных значениях переменных X и значениях на границе их области определения. Автоматическая валидация предполагает, что для всех данных, передаваемых по интерфейсам модели, есть методы их получения и обработки. Для валидации данных требуется выполнить следующие действия:

- 1. Выполнить поиск элементов модели, для которых в базе знаний не содержится правил для получения или передачи данных по коммутационным соединениям, в зависимости от типа интерфейса (вх., исх.).
- 2. Найти все коммутационные соединения элементов модели. Выбрать не связанные правилами соединения.
- 3. Выполнить поиск команд из базы команд, для которых нет правил их приема / обработки / передачи.
- 4. Сформировать перечень найденных ошибок с указанием элементов модели, коммутационных соединений или команд.

Верификация модели предполагает проверку на соответствие модели замыслу исследователя на уровне структуры и методов. Для верификации модели метод структурно-графического анализа выполняет следующие действия:

- 1. Выбрать соединения элементов модели, у которых типы коммутационных интерфейсов не совпадают.
- 2. Выбрать соединения элементов модели, у которых одинаковые направления передачи данных (вх., исх.).
 - 3. Выполнить поиск элементов модели, для которых не заданы правила в базе знаний.
- 4. Выполнить поиск коммутационных соединений, которые не заданы в правилах (в части «условия» или «действия»).
- 5. Сформировать перечень найденных ошибок с указанием элементов модели, коммутационных соединений или интерфейсов.

Пункты 1, 2 позволяют выявлять ошибки структуры, 3, 4 — наличие методов функционирования для всех элементов модели.

Полнота модели — проверка всех возможных вариантов развития моделируемых процессов. Для интеллектуальных систем имитации полнота недостижима, метод структурно-графического анализа определяет характеристики модели, на основе которых конструктор может судить о допустимой степени детализации модели при ее декомпозиции. В данном случае метод может построить рекомендации, на основании которых конструктор сам принимает решение об изменении модели. Для формирования рекомендаций выполняются следующие действия:

- 1. Выполнить расчет функциональной нагрузки на элементы модели, коммутационные интерфейсы и соединения элементов. Выбрать элементы с наибольшей нагрузкой. Сформировать рекомендации о резервировании устройств и линий связи.
- 2. Построить все пути передачи данных между элементами модели с указанием интерфейсов, по которым происходит взаимодействие. Выбрать «несвязные» пути, «замыкания» и «тупики». Сформировать рекомендации о проверке корректности коммутации устройств.
- 3. Выбрать альтернативные пути передачи данных. При их отсутствии сформировать рекомендации о резервировании путей передачи данных.
- 4. Сформировать перечень рекомендаций с указанием элементов модели, коммутационных соединений или интерфейсов.

4. Реализация метода анализа имитационной модели

Реализация метода структурно-графического анализа выполнена на основе библиотек интерактивной инфографики d3.js и sigma.js [16]. Для верификации модели строится таблица ошибок, которая содержит перечень элементов модели, интерфейсов и описание найденных при анализе ошибок. Фрагмент таблицы приведен на рис. 3.

Таблица (ошибок			
Передающ	ий блок	Принимающи	й блок	Описание ошибки
ми кис	На ПРД (RS-422)	Передатчик	Выход (ВЧ)	Соединены интерфейсы различного типа
МИ КИС	От БАТС (RS-232)	Передатчик	Команда (RS-232)	Соединенные интерфейсы не синхронизированны по направлению передачи

Рис. 3. Фрагмент таблицы ошибок структуры модели Fig. 3. Model's structure error table

Для анализа функциональных связей, заданных в базе знаний применяется круговая диаграмма зависимостей (dependency wheel). Секции диаграммы обозначают элементы модели, например имитаторы бортовых устройств, а лучи между ними отображают взаимодействия, заданные в базе знаний. Диаграмма позволяет интерактивно выбирать зависимости отдельных элементов модели, определять ошибки базы знаний, выявлять недостающие или избыточные данные и структуры, для которых не заданы правила в базе знаний, обеспечивая контроль полноты функционального представления. Пример диаграммы зависимостей представлен на рис. 4.

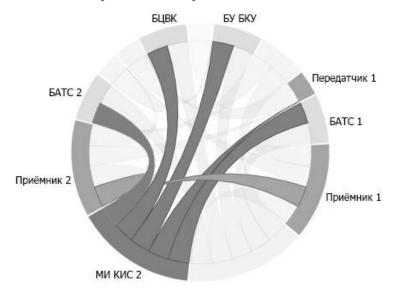
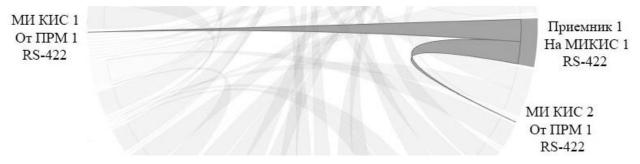


Рис. 4. Круговая диаграмма зависимостей между элементами модели Fig. 4. Dependency circle chart of interaction between model's elements

Для исключения ошибок структуры, при коммутации элементов модели применяется графическое представление, отображающее связи между интерфейсами, показанное на рис. 5 (фрагмент диаграммы). На основе исследования диаграмм делается вывод о соответствии структуры модели и базы знаний.



Puc. 5. Отображение связей между интерфейсами имитационной модели Fig. 5. Links between interfaces of a simulation model

На основе проведенного анализа программное обеспечение формирует перечень ошибок в функциональных зависимостях. Фрагмент таблицы ошибок показан на рис. 6.

Таблица ошибок			
Передающий блок	Интерфейс	Принимающий блок	Интерфейс
Описание ошибки передающего блока		Описание ошибки принимающего блока	
МИ КИС 2	От БАТС 1	БАТС	Выход
Для передающего интерфейса описано пр	равило приёма.	Принимающий интерфейс, в блоке указан пе Для принимающего интерфейса нет правил.	редающим.
МИ КИС 1	На БЦВК	БЦВК	Вход
Для передающего интерфейса нет правил	п.	Для принимающего интерфейса нет правил.	
МИ КИС 1	На БУ БКУ 2	БУБКУ	Вход 2
		Для принимающего интерфейса описано пра	вил передачи.
МИ КИС 1	От БАТС 2	БATC 2	Выход
		Для принимающего интерфейса описано пра	вил передачи.

Рис. 6. Фрагмент таблицы ошибок функциональных зависимостей Fig. 6. Functional dependency error table

Таблица ошибок содержит наименования элементов модели, команд, перечень правил базы знаний и текстовое описание найденных при анализе проблем.

Исследование коммуникационной нагрузки модели выполняется на специальном графе, узлы которого представляют элементы модели, дуги — пути их взаимодействия с другими подмоделями. Пример графа показан на рис. 7. Размер узлов и ширина дуг отображают степень нагрузки, вычисляемую как мощность множества правил для элементов модели. Высокая загруженность элементов модели может послужить поводом к ее пересмотру или дополнительному резервированию оборудования и коммутации.

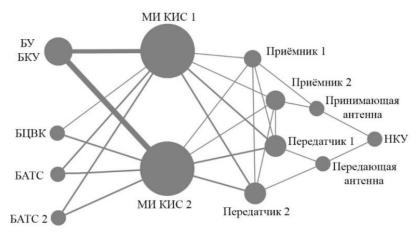


Рис. 7. Граф нагрузки на элементы модели Fig. 7. Model's workload graph

Взаимодействие между правилами базы знаний, описывающими обработку данных на отдельных интерфейсах модели, отображается на графе покрытия (рис. 8). Изменение масштаба отображения графа позволяет рассматривать модель в целом, детализировать отдельные элементы модели или все коммутационные соединения, описанные в базе знаний, с указанием начальных и конечных интерфейсов. Интерактивные графические элементы позволяют визуализировать показатели покрытия структуры модели правилами базы знаний. Выделяются элементы модели, в составе которых все интерфейсы описаны логическими правилами либо отдельные узлы не имеют правил. Это свидетельствует о наличии коммутационных соединений, по которым не предусмотрена передача данных, что является ошибкой проектирования модели. Фрагмент графа, представленный на рис. 8, демонстрирует связи между интерфейсами блоков и позволяет детализировать покрытие правилами по каждому из интерфейсов.

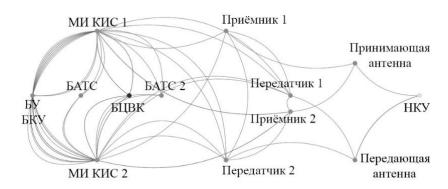


Рис. 8. Граф покрытия правилами Fig. 8. Coverage graph

На основе построенных графов программное обеспечение формирует таблицу предупреждений о превышении средней нагрузки на элементы модели. Фрагмент таблицы показан на рис. 9.

езервирование устройства.
езервирование устройства.
езервирование интерфейса.
езервирование интерфейса.
езервирование соединения.
езервирование соединения.

Рис. 9. Фрагмент таблицы предупреждений превышения средней нагрузки Fig. 9. Overload warning table

В таблице приводятся перечень элементов, коммутационных интерфейсов и соединений, для которых нагрузка превышает среднюю, рассчитанную по всем элементам модели, и рекомендации о резервировании устройств или линий связи.

Графическая визуализация модели, а также формируемые перечни ошибок и предупреждений позволяют конструктору бортовой аппаратуры выполнять анализ характеристик базы знаний.

Заключение

В работе предложен человеко-машинный метод анализа структуры и свойств имитационной модели функционирования бортовой аппаратуры космического аппарата. Программное обеспечение строит наглядные диаграммы, отражающие свойства и взаимосвязи элементов модели, формирует перечни ошибок и рекомендаций, на основании которых конструктор может принимать решение о внесении изменений в модель. Реализованный метод позволяет выявлять зависимости отдельных элементов модели, ошибки базы знаний, недостающие или избыточные данные и структуры, для которых не заданы правила в базе знаний, обеспечивая контроль полноты функционального представления. Помимо автоматического формирования списка ошибок модели, метод может использоваться для ручной проверки соответствия моделей техническим описаниям, заданным в конструкторской документации.

ЛИТЕРАТУРА

1. Аксенов К.А., Гончарова Н.В. Моделирование и принятие решений в организационно-технических системах. Екатеринбург: Изд-во Урал. ун-та, 2015. 104 с.

- 2. Советов Б.Я., Яковлев С.А. Моделирование систем. М.: Высш. шк., 2009. 343 с.
- 3. Литвинов В.В., Марьянович Т.П. Методы построения имитационных систем. Киев : Наукова думка, 1991. 120 с.
- 4. Остроух А.В. Интеллектуальные системы. Красноярск: Науч.-инновационный центр, 2015. 110 с.
- 5. Tan C.F., Wahidin L.S., Khalil S.N., Tamaldin N. The application of expert system: a review of research fnd applications // ARPN Journal of Engineering and Applied Sciences. 2016. No. 11 (4). P. 2448–2453.
- 6. Eickhoff J. Simulating Spacecraft System. Springer, 2009. 376 p.
- 7. Лычкина Н.Н. Имитационное моделирование экономических процессов / под ред. В.В. Година. М.: Академия IT, 2005. 165 с.
- 8. Min F., Yang M., Wang Z. Knowledge-based method for the validation of complex simulation models // Simulation Modelling Practice and Theory. 2010. No. 18 (5). P. 500–515.
- 9. Василенко Д.Е., Обидин Д.Н., Бердник П.Г. Разработка процедуры контроля непротиворечивости знаний для открытой экспертной системы реального времени // Системи обробки інформації. 2016. № 9 (146). С. 90–93.
- 10. Yuchen Zhou, Ke Fang, Ming Yang, Ping Ma. An intelligent model validation method based on ECOC SVM // Proc. of the 10th Int. Conference on Computer Modeling and Simulation. 2018. P. 67–71.
- 11. Zanon O. The SimTG simulation modeling framework a domain specific language for space simulation // Proc. of the 2011 Symposium on Theory of Modeling & Simulation: DEVS Integrative M&S Symposium. 2011. P. 16–23.
- 12. Ноженкова Л.Ф., Исаева О.С., Евсюков А.А. Инструменты компьютерного моделирования функционирования бортовой аппаратуры космических систем // Тр. СПИИРАН. 2018. № 56. С. 144–168. DOI: 10.15622/sp.56.7.
- 13. Исаева О.С. Разработка методики автоматизации испытаний на основе имитационной модели функционирования бортовой аппаратуры космического аппарата // Вестник компьютерных и информационных технологий. 2018. № 10 (172). С. 30–38.
- 14. Ноженкова Л.Ф., Исаева О.С., Грузенко Е. А. Метод системного моделирования бортовой аппаратуры космического аппарата // Вычислительные технологии. 2015. № 20 (3). С. 33–44.
- 15. Антонов А.В. Системный анализ: учебник для вузов. 3-е изд. М.: Высш. шк., 2008. 454 с.
- 16. Bostock M. Data-Driven Documents. URL: https://d3js.org/ (accessed: 28.04.2019).

Поступила в редакцию 16 мая 2019 г.

Isaeva O.S., Kulaysov N.V., Isaev S.V. (2020) METHOD OF STRUCTURAL AND GRAPHICAL ANALYSIS AND VERIFICATION OF INTELLECTUAL SIMULATION MODEL. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 79–88

DOI: 10.17223/19988605/50/10

The paper presents a method for analyzing an intelligent model of simulating the functioning of the spacecraft's onboard equipment. The model consists of a graphic structure that presents the elements of onboard equipment and a knowledgebase that describes the methods of its operation. The authors formalized the model and proposed criteria for the analysis and verification of its structure and properties. We have developed visual components of interactive infographics that performing interpretation of the formal description of the model in interactive graphic images for analysis and detect errors of the knowledgebase. In addition to automatic control, interactive graphic tools can be used to manually check the completeness and consistency of knowledge, as well as the compliance of models with the technical descriptions given in the design documentation.

Intelligent simulation model consists of a graphical structure duplicating the composition of the elements of the onboard equipment and a knowledge base describing the methods of its operation. Model $S = \langle G, F, T \rangle$, where G is a structural-parametric description (a set of structure elements), F is a functional description (a set of functioning methods), T is the time of observation. A structural-parametric description of $G = \langle B, I, C, D, P \rangle$, where B is a set of elements representing properties or functions of physical devices, I is a set of commutation interfaces, C is typed information dependencies describing connections between elements from B, D is a set of data structures, P is a set of parameters. Functional description $F = \{R(P, T)\}$, where R is the set of rules of the knowledge base. $R = \{A \rightarrow Z\}$, where A is the conditions for rule's fulfillment and Z is the actions required for changing the model's state.

The method allows to interactively to build graphs of dependencies of different elements of the model, detect errors of the knowledge base, reveal lacking or excessive data and structures that do not have rules set in the knowledge base and provide control of completeness of functional presentation. For example, the completeness property of a model: $\forall B_i \in B | R(B_i)| \neq 0$, where $|R(B_i)|$ is the number of rules for element of the model B_i , is visualized on the workload graph. The size of nodes and arcs of the graph show the intensity of load calculated as the number of rules for the model's elements. High workload of the elements may be a reason for revision of the model or for additional reservation of equipment and commutation.

 $Keywords: simulation\ modeling;\ spacecraft\ onboard\ equipment;\ verification;\ validation;\ knowledge\ base;\ inforgafika.$

ISAEVA Olga Sergeevna (Candidate of Technical Sciences, Senior Researcher, Institute of Computational Modelling of the Siberian Branch of the Russian Academy of Sciences, Krasnoyarsk, Russian Federation).

E-mail: isaeva@icm.krasn.ru

KULYASOV Nikita Vladimirovich (Engineer, Institute of Computational Modelling of the Siberian Branch of the Russian Academy of Sciences, Krasnoyarsk, Russian Federation).

E-mail: razor@icm.krasn.ru

ISAEV Sergey Vladislavovich (Candidate of Technical Sciences, Assistant Professor, Deputy Director, Institute of Computational Modelling of the Siberian Branch of the Russian Academy of Sciences, Krasnoyarsk, Russian Federation).

E-mail: si@icm.krasn.ru

REFERENCES

- 1. Aksenov, K.A. & Goncharova, N.V. (2015) *Modelirovanie i prinyatie resheniy v organizatsionno-tekhnicheskikh sistemakh* [Modeling and decision making in organizational and technical systems]. Ekaterinburg: Ural State University.
- 2. Sovetov, B.Ya. & Yakovlev, S.A. (2009) Modelirovanie sistem [System modeling]. Moscow: Vysshaya shkola.
- 3. Litvinov, B.V. & Maryanovich, T.P. (1991) *Metody postroeniya imitatsionnykh sistem* [Methods for constructing simulation systems]. Kiev: Naukova dumka.
- 4. Ostroukh, A.V. (2015) Intellektual'nye sistemy [Intellectual systems]. Krasnoyarsk: Nauch.-innovatsionnyy tsentr.
- 5. Tan, C.F., Wahidin, L.S., Khalil, S.N. & Tamaldin, N. (2016) The application of expert system: a review of research fnd applications. *ARPN Journal of Engineering and Applied Sciences*. 11(4). pp. 2448–2453.
- 6. Eickhoff, J. (2009) Simulating Spacecraft System. Springer.
- Lychkina, N.N. (2005) Imitatsionnoe modelirovanie ekonomicheskikh protsessov [Simulation modeling of economic processes].
 Moscow: IT Academy.
- 8. Min, F., Yang, M. & Wang, Z. (2010) Knowledge-based method for the validation of complex simulation models. *Simulation Modelling Practice and Theory*. 18(5) pp. 500–515. DOI: 10.1016/j.simpat.2009.12.006
- 9. Vasilenko, D.E., Obidin, D.N. & Berdnik, P.G. (2016) Establishing procedures for the contradiction control for open knowledge of expert system real time. *Sistemi obrobki informatsii Information Processing Systems*. 9(146). pp. 90–93.
- 10. Zhou, Y, Fang K., Yang, M. & Ma, P. (2018) An intelligent model validation method based on ECOC SVM. *Proc. of the 10th Int. Conference on Computer Modeling and Simulation*. Sydney, Australia. pp. 67–71.
- 11. Zanon, O. (2011) The SimTG simulation modeling framework a domain specific language for space simulation. *Proc. of the 2011 Symposium on Theory of Modeling & Simulation: DEVS Integrative M&S Symposium*. pp. 16–23.
- 12. Nozhenkova, L.F., Isaeva, O.S. & Evsyukov, A.A. (2018) Tools of computer modeling of the space systems' onboard equipment function. *Tr. SPIIRAN SPIIRAS Proceedings*. 56. pp. 144–168. DOI: 10.15622/sp.56.7.
- 13. Isaeva, O.S. (2018) Development of a method for automation of testing based on the simulation model of functioning of the onboard equipment of the spacecraft. *Herald of Computer and Information Technologies*. 10(172). pp. 30–38.
- 14. Nozhenkova, L.F., Isaeva, O.S. & Gruzenko, E.A. (2015). The method for system modelling of the spacecraft on-board equipment. *Vychislitel'nye tekhnologii Computational Technologies*. 20(3). pp. 33–44.
- 15. Antonov, A.V. (2008) Sistemnyy analiz [System Analysis]. Moscow: Vysshaya shkola.
- 16. Bostock, M. (n.d.) Data-Driven Documents. [Online] Available from: https://d3js.org/. (Accessed: 28th April 2019).

2020 Управление, вычислительная техника и информатика

№ 50

УДК 681.324

DOI: 10.17223/19988605/50/11

П.Х. Карим, П.А. Михеев, В.В. Поддубный, С.П. Сущенко

ЧИСЛЕННЫЕ ИССЛЕДОВАНИЯ ПРОПУСКНОЙ СПОСОБНОСТИ ТРАНСПОРТНОГО ПРОТОКОЛА С МЕХАНИЗМОМ ПРЯМОЙ КОРРЕКЦИИ ОШИБОК В МЕЖСЕГМЕНТНОМ ПРОСТРАНСТВЕ

Предложена математическая модель транспортного протокола с механизмом прямой коррекции ошибок. Показано, что применение механизма ведет к росту пропускной способности канала при определенных значениях многомерного признакового пространства протокольных параметров, характеристик тракта передачи данных и механизма прямой коррекции ошибок. Представлены результаты численного анализа зависимости прироста пропускной способности транспортного протокола с прямой коррекцией ошибок от параметров корректирующего механизма, достоверности передачи протокольных блоков данных и длительности круговой задержки. Ключевые слова: транспортное соединение; пропускная способность; цепь Маркова; прямая коррекция ошибок; размер окна; длительность тайм-аута; круговая задержка.

Пропускная способность транспортного соединения является крайне важной характеристикой компьютерных сетей. Данный показатель определяет качество сетевых сервисов для абонентов и определяется значениями протокольных параметров (размер окна, длительность тайм-аута), характеристиками тракта передачи данных (длительность круговой задержки, достоверность передачи данных в различных направлениях транспортного соединения) [1]. В настоящее время получают применение технологии прямой коррекции ошибок [2–7] в виде дополнительного сервиса в транспортных протоколах наряду с методом решающей обратной связи для снижения объема повторно передаваемого трафика. Исследование транспортного соединения и анализ его потенциальных возможностей выполнялся в [2-14], но аналитические результаты получены только для однозвенного тракта передачи данных [8-10] либо при существенных ограничениях на параметры протокола [11-14]. Примером одной из модификаций транспортного протокола является подключаемый к UDP протоколу дополнительный механизм под названием QUIC (Quick UDP Internet Connections) [7]. Это новое дополнение протокола пока еще не является стандартом, и в настоящее время эксперименты с ним и исследование его эффективности продолжаются. Протокол QUIC развернут как дополнительный сервис в кампании Google. Следует отметить, что протокол QUIC имеет ряд недостатков, к числу которых относится использование ограниченного множества значений параметров механизма прямой коррекции ошибок. Исследование протокола QUIC выполнялось только в натурных экспериментах (тестирование на оборудовании) [7]. Как правило, исследования преимуществ метода прямой коррекции ошибок проводятся на качественном уровне и для некоторых частных случаев численно. Известные исследования не определяют области признакового пространства параметров протокола и транспортного соединения, в которых применение метода опережающего исправления ошибок дает положительные результаты. Скрытные возможности транспортного протокола с применением метода прямой коррекции ошибок не изучены полностью. Отсутствуют аналитические зависимости комплексного влияния протокольных параметров, характеристик тракта передачи данных и параметров метода коррекции на быстродействие транспортного соединения. Не исследовано влияние соотношений между длительностью круговой задержки и протокольными параметрами на пропускную способность тракта передачи данных, управляемого транспортным протоколом.

В работе предложена математическая модель процесса передачи данных с прямой коррекцией ошибок в фазе информационного переноса в виде цепи Маркова с дискретным временем. Проведен числительный анализ пропускной способности транспортного канала с применением механизма прямой коррекции ошибок, показано преимущество протокола с прямой коррекцией ошибок по сравнению с классическим протоколом с решающей обратной связью для определенных областей признакового пространства протокольных параметров, характеристик транспортного соединения и параметров механизма кодирования.

1. Модель транспортного соединения

Рассмотрим процесс переноса данных между абонентами транспортного протокола, основанного на алгоритме с решающей обратной связью. Примером семейства таких надежных протоколов является доминирующий в современных компьютерных сетях протокол ТСР [1]. Полагаем, что взаимодействующие абоненты имеют неограниченный поток данных для передачи, а обмен выполняется сегментами данных транспортного протокола одинаковой длины. Считаем, что участки переприема вдоль тракта передачи данных имеют одинаковое быстродействие в обоих направлениях, а длительность цикла передачи сегмента в отдельном звене составляет t. В общем случае длина пути от источника до адресата, переносящего информационный поток, и длина обратного пути, по которому передаются подтверждения на принятые сегменты, могут быть различными. Полагаем, что длина тракта передачи данных, выраженная в количестве участков переприема, в прямом направлении равна $D_n \ge 1$. Обратный тракт, по которому доставляются подтверждения отправителю о корректности приема последовательности блоков сегментов, имеет длину $D_o \ge 1$. Заданы вероятности искажения сегмента в каналах связи для прямого $R_{\Pi}(d)$, $d=\overline{1,D_n}$ и обратного $R_{\sigma}(d)$, $d=\overline{1,D_n}$ направлений передачи каждого участка переприема. Тогда достоверность передачи сегментов вдоль тракта от источника до адресата и обратно составит $F_n = \prod_{d=1}^{D_n} (1 - R_{\Pi}(d))$ и $F_o = \prod_{d=1}^{D_o} (1 - R_{O}(d))$ соответственно. Считаем, что потерь сегментов из-за отсутствия буферной памяти в узлах тракта не происходит. Полагаем, что передача данных отправителем реализуется блоками, содержащими B сегментов, из которых $1 \le A \le B$ являются информационными, а B - A – дополнительными (избыточными) той же длины. Полагаем, что все сегменты имеют контрольные суммы, позволяющие обнаружить ошибки в каждом из них. Потеря (искажение) до B-A произвольных сегментов в блоке позволяет восстановить все сегменты блока (например, передачей дублей при $A = 1, B \ge 2$, оправкой избыточного сегмента с поразрядной четностью всех информационных сегментов по технологии RAID-массивов [15] при A > 1, B = A + 1 и др.). Управление потоком данных реализуется механизмом скользящего окна [1] с протокольным параметром ширины окна $\omega \ge 1$, выраженным в количестве блоков. Полагаем, что подтверждения о корректности полученных адресатом блоков сегментов переносятся в каждом сегменте встречного потока. При невозможности прямого восстановления переданных сегментов блока (искажение более B-A сегментов в блоке) весь блок передается повторно.

Тогда процесс информационного переноса в виртуальном соединении, управляемом транспортным протоколом, может быть описан марковским процессом с дискретным временем (с длительностью такта t) в силу того, что время между получениями подтверждений имеет геометрическое распределение с параметром F_o . Данная модель является обобщением формализаций процесса передачи данных, предложенных в [11–14], на случай транспортного соединения произвольной длины и механизма прямой коррекции ошибок. Область возможных состояний цепи Маркова определяется длительностью тайм-аута ожидания подтверждения S, выраженной в количестве циклов продолжительности t. Размер тайм-аута связан с длиной тракта, шириной окна и размером блока неравенствами $S \ge \omega B + 1$, $S \ge D_n + D_o + B - 1$. Очевидно, что сумму длин прямого и обратного трактов можно интерпретировать как круговую задержку одиночного сегмента $D = D_n + D_o$ в детерминированном

транспортном соединении, выраженную в длительностях t. Круговая задержка для блока сегментов составит D+B-1. Состояниям цепи Маркова $i=\overline{0,\omega B}$ соответствует размер очереди переданных, но не подтвержденных сегментов в источнике потока, а состояниям $i = \overline{\omega B + 1, S - 1}$ — время, в течение которого отправитель не активен и ожидает получения подтверждения о корректности приема переданной последовательности из ω блоков сегментов. Из нулевого состояния в D+B-2 источник продвигается с каждым тактом t с вероятностью детерминированного события. В состояниях $i \ge D + B - 2$ после истечения очередного дискретного цикла t к отправителю начинают прибывать подтверждения и, в зависимости от результатов доставки блоков сегментов с учетом технологии прямой коррекции ошибок, отправитель передает новые блоки сегментов (при положительном подтверждении) либо повторно – искаженные (не поддающиеся прямому восстановлению). Завершение щикла пребывания в состоянии D+B-2 соответствует времени доведения первого блока сегментов до адресата и получения на него подтверждения. Дальнейший рост номера состояния происходит с вероятностью искажения подтверждения $1 - F_o$ в обратном тракте. Получение подтверждения в состояниях $i \ge D + B - 2$ в предположении отсутствия расщепления точек возврата, обусловленных конвейерным эффектом, вызывает переход в D-1 состояние при $\omega \ge K+2$ только в случае успешной передачи доставленных адресату блоков, в противном случае следует переход в 0 состояние. Здесь $K = \left\lfloor \frac{D-2}{B} \right\rfloor$, где $\lfloor \dots \rfloor$ означает «целая часть» дроби.

В силу того что в состояниях $i \geq \omega B$ источник приостанавливает отправку блоков сегментов, получение подтверждений при $\omega \geq K+2$ в состояниях $i=\overline{(\omega+k)B-1,(\omega+k+1)B-2},\ k=\overline{1,K}$ приводит к переходу в состояния $D-kB-1,k=\overline{1,K}$, только при успешной доставке данных (иначе — в 0 состояние). В состояниях $i=\overline{(\omega+K+1)B-1,S-2}$ выполняется переход в нулевое состояние, поскольку размер очереди переданных, но не подтвержденных информационных сегментов при этом обнуляется. В состоянии S-1 истекает тайм-аут ожидания подтверждения от получателя о корректности принятых блоков сегментов и происходит безусловный переход в нулевое состояние.

2. Операционные характеристики транспортного протокола с механизмом прямой коррекции ошибок

Переходные вероятности π_{ij} из исходного состояния i в результирующее j цепи Маркова, описывающей процесс передачи информационного потока с технологией прямой коррекцией ошибок в режиме группового отказа для $\omega \ge K+2$, $S \ge D+B(\omega+1)-2$, имеют вид:

$$\pi_{ij} = \begin{cases} 1, i = \overline{0, D + B - 3}, j = i + 1; \\ 1 - F_o, i = \overline{D + B - 2, S - 2}, j = i + 1,; \\ F_o \psi^k, i = \overline{D + Bk - 2, D + (k + 1)B - 3}, k = \overline{1, G}, j = D - 1; \\ F_o (1 - \psi^k), i = \overline{D + Bk - 2, D + (k + 1)B - 3}, k = \overline{1, G}, j = 0; \\ F_o \psi^G, i = \overline{D + B(G + 1) - 2}, B(\omega + 1) - 2, j = D - 1; \\ F_o (1 - \psi^G), i = \overline{D + B(G + 1) - 2}, B(\omega + 1) - 2, j = 0; \\ F_o \psi^{G+k}, i = \overline{B(\omega + k) - 1}, B(\omega + k + 1) - 2, k = \overline{1, K}, j = D - Bk - 1; \\ F_o (1 - \psi^{G+k}), i = \overline{B(\omega + k) - 1}, B(\omega + k + 1) - 2, k = \overline{1, K}, j = 0; \\ F_o, i = \overline{B(\omega + K + 1) - 1}, S - 2, j = 0; \\ 1, j = 0, i = S - 1. \end{cases}$$

Здесь $G = \left\lfloor \frac{B(\omega+1)-2-(D+B-2)+1}{B} \right\rfloor = \omega - \left\lfloor \frac{D-1}{B} \right\rfloor$ — расстояние между моментами начала прекращения активности отправителя $B(\omega+1)-2$ (завершения оправки ω блоков) и начала поступления ему квитанций D+B-2, выраженное в размерах $B, \ \psi = \sum\limits_{i=A}^{B} C \int\limits_{B}^{i} F_{n}^{i} (1-F_{n})^{B-i}$. Решая систему уравнений равновесия находим вероятности состояний цепи Маркова и далее получаем показатель пропускной способности транспортного соединения с применением механизма прямой коррекции ошибок:

$$Z(D, \omega, S, A, B, F_n, F_o) = \frac{P_0 A \psi (1 - \overline{F}_0^B) (1 - \overline{F}_0^B \psi)}{B \mathcal{I} F_o (1 - \psi)} \left\{ \left[1 - \overline{(F}_0^B \psi)^{\omega} \right] \left(\frac{1 - \psi}{1 - \overline{F}_0^B \psi} \right) - \overline{F}_0^{S - D - B + 2} (1 - \psi^{\omega}) \right\}, \quad (1)$$

где

$$\begin{split} P_0 &= \overline{D} F_o (1 - \overline{F}_0^B \psi) \bigg/ \bigg\{ \mathcal{A} F_o (D - BK - 1) (1 - \overline{F}_0^B \psi) + \bigg(1 - \overline{F}_0^B \psi \bigg)^2 \bigg[1 + F_o (B - 1) - \overline{F}_0^{S - D - B + 2} \bigg] + \\ &+ B F_o \bigg[K (1 - \psi) (1 - \overline{F}_0^B \psi) + \psi^{G + 1} [1 - (\overline{F}_0^B \psi)^K] (\overline{F}_0^{B \circ D - D + 1} (1 - \overline{F}_0^B \psi) - (1 - \psi) \overline{F}_0^{G + 1}) \bigg] \bigg\}, \\ \bar{D} &= (1 - \overline{F}_0^B \psi) [\psi^{G + 1} \overline{F}_0^{B \circ D - D + 1} (\overline{F}_0^B \psi)^K] + (1 - \psi) [1 - (\overline{F}_0^B \psi)^{K + G + 1}], \ \overline{F} = 1 - F_o. \end{split}$$

Очевидно, что поиск в многомерном пространстве признаков (D, ω , S, A, B, F_n , F_o) областей, обеспечивающих превосходство механизма прямой коррекции ошибок перед классической протокольной процедурой с решающей обратной связью по критерию пропускной способности, является сложной задачей. Решением данной проблемы может быть редуцирование размерности признакового пространства. Эффективными вариантами редуцирования размерности признакового пространства являются случаи абсолютно надежного обратного тракта передачи данных (F_o = 1), неограниченной ширины окна ($\omega \rightarrow \infty$), а следовательно, и длительности тайм-аута ($S \rightarrow \infty$), а также случай однородного тракта передачи данных в прямом и обратном направлениях (F_n = F_o = F).

3. Численный анализ пропускной способности транспортного соединения

Рассмотрим случай абсолютно надежного обратного тракта $F_o = 1$. Тогда пропускная способность транспортного соединения принимает вид:

$$Z(D, \omega, S, A, B, F_n, 1) = \frac{A\psi}{B(D-1)(1-\psi) + B^2}.$$

Чтобы определить области, где механизм с применением технологии прямой коррекции ошибок дает преимущество, сравним показатель пропускной способности классического транспортного протокола и протокола с коррекцией ошибок. В классическом варианте пропускная способность имеет вид [11]:

$$Z_{\text{KJI}}(D, \omega A, S, F_n, F_o = 1) = \frac{F_n}{(D-1)(1-F_n)+1}$$
.

Определим выигрыш от применения технологии прямой коррекции ошибок перед классическим транспортным протоколом при $F_o = 1$:

$$\Delta Z = Z(D, \omega, S, A, B, F_n, F_o) - Z_{\text{KJI}}(D, A\omega, S, F_n, F_o) = \frac{A.\psi}{B(D-1)(1-\psi) + B^2} - \frac{F_n}{(D-1)(1-F_n) + 1}.$$

Численные исследования выигрыша при наборе параметров механизма прямой коррекции ошибок A=1, B=A+1 показывают (рис. 1, a), что область положительных значений выигрыша с ростом длительности круговой задержки D от 51 до 66 расширяется от $F \in (0,05;0,93)$ до

 $F \in (0,04;0,95)$. Кроме того, из рис. 1, 2 следует, что экстремальное значение выигрыша растет с увеличением D и достигается практически в одной и той же точке F = 0,88. На рис. 1, b проиллюстрированы сходные зависимости для набора параметров A = 2, B = A + 1.

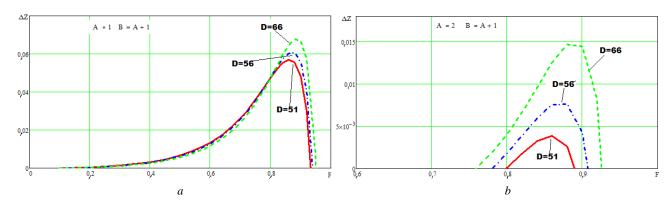
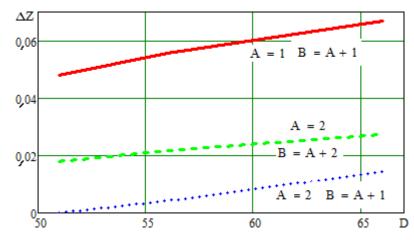


Рис. 1. Зависимость выигрыша ΔZ от достоверности передачи данных в прямом тракте $F=F_n$ при $F_o=1$ Fig. 1. Dependence of winning ΔZ on reliability data transfer in the forward path $F=F_n$ at $F_o=1$

Из результатов, приведенных на рис. 2, нетрудно видеть, что режим дублирования данных B = 2A имеет преимущество перед параметрами A = 2, B = A + 1 механизма прямой коррекции ошибок, при этом наибольший выигрыш достигается для параметрического набора A = 1, B = A + 1.



Puc. 2. Зависимость значений максимального выигрыша от длительности круговой задержки при $F_o = 1$ Fig. 2. The dependence of the values of the maximum winnings from the duration of the circular delay at $F_o = 1$

Проанализируем выигрыш пропускной способности канала в условиях неограниченной ширины окна ($\omega \to \infty$) и стохастической однородности прямого и обратного трактов передачи данных ($F_n = F_o = F$). Пропускная способность для протоколов с применением механизма прямой коррекции ошибок и без его использования согласно (1) и [11] определится соответственно следующим образом:

$$Z(D, \infty, \infty, A, B, F, F) = \frac{A\psi(1 - (1 - F)^{B})}{B\{(1 - \psi)F(D - 1) + (1 - (1 - F)^{B}\psi)(1 + F(B - 1))\}},$$

$$Z_{\text{KJI}}(D, \infty, \infty, F, F) = \frac{F^{2}}{1 + F(D - 2)(1 - F)}.$$

Согласно значениям выигрыша, приводимым на рис. 3, нетрудно видеть, что область положительных значений с ростом круговой задержки D незначительно расширяется, а максимальные значения — увеличиваются.

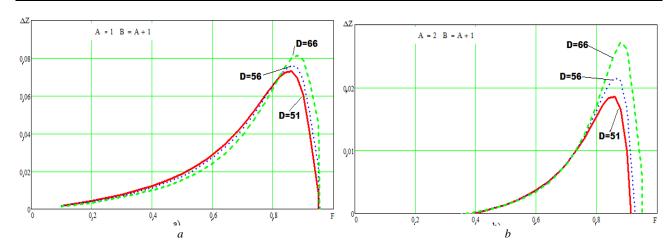


Рис. 3. Зависимость выигрыша ΔZ от достоверности передачи данных в прямом тракте F при $F_n = F_o = F$ Fig. 3. Dependence of winning ΔZ on reliability data transfer in the forward path F at $F_n = F_o = F$

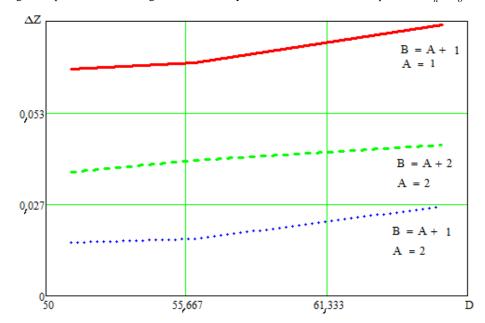


Рис. 4. Зависимость значений максимального выигрыша от длительности круговой задержки при $F_n = F_o = F$ Fig. 4. The dependence of the values of the maximum winnings from the duration of the circular delay when $F_n = F_o = F$

Из зависимостей, приводимых на рис. 4, следует, что лучшими значениями параметров механизма прямой коррекции ошибок являются A = 1, B = A + 1.

Заключение

В работе предложена модель процесса переноса сегментов данных в транспортном соединении, управляемом надежным транспортным протоколом с механизмом прямой коррекции ошибок и подтверждением данных, принятых получателем, в режиме группового повтора. Математическая модель основана на описании очереди переданных, но не подтвержденных сегментов данных цепью Маркова с конечным числом состояний и дискретным временем. Представлен численный анализ пропускной способности транспортного соединения. Численные исследования выполнены для абсолютно надежного обратного тракта и для неограниченного размера окна. При данных условиях показано, что выигрыш пропускной способности увеличивается с ростом круговой задержки. Показано, что для группового режима повтора применение механизма прямой коррекции ошибок целесообразно на транспортных соединениях с большой круговой задержкой.

ЛИТЕРАТУРА

- 1. Fall K., Stevens R. TCP/IP Illustrated. 2nd ed. Addison-Wesley Professional, 2012. V. 1: The Protocols 1017 p. (Addison-Wesley Professional Computing Series)
- 2. Lundqvist H., Karlsson G. TCP with end-to-end FEC // Communications Int. Zurich Seminar. 2004. P. 152-156.
- 3. Barakat Ch., Altman E. Bandwidth tradeoff between TCP and link-level FEC // Computer Networks. 2002. No. 39. P. 133-150.
- 4. Shalin R., Kesavaraja D. Multimedia Data Transmission through TCP/IP using Hash Based FEC with AUTO-XOR Scheme // ICTACT Journal on Communication Technology. 2012. V. 03, is. 03. P. 604–609.
- 5. Flach T., Dukkipati N., Terzis A., Raghavan B., Cheng Yu., Cardwell N., Jain A., Hao S., Katz-Bassett E., Govindan R. Reducing Web Latency: the Virtue of Gentle Aggression // ACM SIGCOMM. 2013. P. 159–170.
- 6. Herrero R. Modeling and comparative analysis of Forward Error Correction in the context of multipath redundancy // Telecommunication Systems. Modelling, Analysis, Designand Management. 2017. V. 65, No. 4. P. 783–794.
- 7. Langley A., Riddoch A., Wilk A., Vicente A., Krasic C., Zhang D., Ang F., Kouranov F., Swett I., Iyengar J., Bailey J., Dorfman J., Roskind J., Kulik J., Westin P., Tenneti R., Shade R., Hamilton R., Vasiliev V., Chang W.-T., Shi Z. The QUIC transport protocol: Design and internet-scale deployment // SIGCOMM"17, August, 2017, Los Angeles, CA, USA. P. 183–196.
- 8. Boguslavsky L.B., Gelenbe E. Analytical models transmission link control procedures for data computer networks with packet // Automation and Remote Control. 1980. No 7. P. 181–192.
- Gelenbe E., Labetoulle J., Pujolle G. Performance Evaluation of the HDLC Protocol // Computer Networks. 1978. V. 2, No. 4/5. P. 409

 –415.
- 10. Кокшенев В.В., Сущенко С.П. Анализ быстродействия асинхронной процедуры управления звеном передачи данных // Вычислительные технологии. 2008. Т. 15, спец. вып. № 5. С. 61–65.
- 11. Kokshenev V.V., Mikheev P.A., Sushchnenko S.P. Comparative Analysis of the Performance of Selective and Group Repeat Transmission Models in a Transport Protocol // Automation and Remote Control. 2017. V. 78, No 2. P. 65–81.
- 12. Кокшенев В.В., Михеев П.А., Сущенко С.П. Анализ селективного режима отказа транспортного протокола в нагруженном тракте передаче данных // Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2013. № 3 (24). С. 78–94.
- 13. Кокшенев В.В., Сущенко С.П. Моделирование сеансов связи цепями Маркова // Теория вероятностей, случайные процессы, математическая статистика и приложения : материалы Междунар. науч. конф., посвященной 80-летию проф. Г.А. Медведева. Минск (23–26 февраля 2015). Минск : РИВШ, 2015. С. 311–316.
- 14. Mikheev P.A., Sushchenko S.P, Tkachev R.V. Estivation of High-Speed Performance of the Transport Protocol with the Mechanism of Forward Error Correction // Communications in Computer and Information Science. 2017. V. 700. P. 259–268.
- 15. Олифер В.Г., Олифер Н.А. Компьютерные сети. Принципы, технологии, протоколы: учебник для вузов. 5-е изд. СПб.: Питер, 2016. 862 с.

Поступила в редакцию 23 марта 2019 г.

Karim P.Kh., Mikheev P.A., Poddubny V.V., Sushchenko S.P. (2020) NUMERICAL STUDIES OF TRANSPORT PROTOCOL THROUGHPUT WITH FORWARD ERROR CORRECTION MECHANISM IN INTERSEGMENT SPACE. *Vestnik Tomskogo gosudarstvennogo universiteta*. *Upravlenie, vychislitelnaja tehnika i informatika* [Tomsk State University Jounal of Control and Computer Science]. 50. pp. 89–96

DOI: 10.17223/19988605/50/11

A mathematical model of the data transfer process with forward error correction at the level of the transport protocol with crucial feedback in the phase of information transfer in the form of a Markov chain with discrete time is proposed. The article analyzes the possibility of increasing the throughput of a reliable transport connection in the presence of interference in the communication channels of the data transmission path through the use of non-laborious methods of noise-resistant coding in a space of grouped protocol data units. Data transmission by the sender is implemented with blocks containing B segments of the transport protocol, where $1 \le A \le B$ are informational, and B - A are redundant. In this case, distortion of B - A arbitrary segments in the block allows the recipient to restore all segments of the block. The Markov chain describes the dynamics of the queue of transmitted, but not confirmed blocks of segments. To identify areas of superiority of the transport protocol with the error correction mechanism over the classical transport protocol, the gain function in the multidimensional feature space of protocol parameters, transport connection characteristics and error correction method parameters is constructed by the performance criterion. To reduce the complexity of the numerical analysis of the increase in throughput of the data transmission path controlled by the transport protocol using the forward error correction mechanism, methods are proposed for reducing the dimension of the parametric space that determines the speed of the transport connection.

Effective options for reducing the dimension of the feature space are cases of absolutely reliable delivery of receipts to the sender of the data stream, unlimited protocol parameters for the window width and timeout duration, as well as the case of a uniform data transmission path in the forward and reverse directions. Under these conditions, it is shown that the gain in throughput increases with increasing the round trip delay of the protocol data units.

Keywords: transport connection; bandwidth; Markov chain; direct error correction; window size; timeout duration; round-trip delay.

KARIM Peshang Hassan (Post-graduate Student, National Research Tomsk State University, Tomsk, Russian Federation). Email: peshangh@yahoo.com

MIKHEEV Pavel Andreevich (Candidate of Technical Sciences, Senior Researcher, A. Alexandrov Scientific and Research Technological Institute. St. Petersburg, Russian Federation).

E-mail: doka.patrick@gmail.com

PODDUBNY Vasily Vasilyevich (Doctor of Technical Sciences, Professor, National Research Tomsk State University, Tomsk, Russian Federation).

Email: vvpoddubny@gmail.com

SUSHCHENKO Sergey Petrovich (Doctor of Technical Sciences, Head of the Department of Applied Informatics, National Research Tomsk State University, Tomsk, Russian Federation).

Email: ssp.inf.tsu@gmail.com

REFERENCES

- 1. Fall, K. & Stevens, R. (2012) TCP/IP Illustrated, vol. 1: The Protocols (2nd Edition). Addison-Wesley Professional Computing Series
- 2. Lundqvist, H. & Karlsson, G. (2004) TCP with end-to-end FEC. *Communications. Int. Zurich Seminar*. Zurich, Switzerland, 2004. pp. 152–156.
- 3. Barakat, Ch. & Altman, E. (2002) Bandwidth tradeoff between TCP and link-level FEC. Computer Networks. 39. pp. 133–150. DOI: 10.1007/3-540-47734-9_10
- 4. Shalin, R. & Kesavaraja, D. (2012) Multimedia Data Transmission through TCP/IP using Hash Based FEC with AUTO-XOR Scheme. *ICTACT Journal on Communication Technology*. 3. pp. 604–609. DOI: 10.21917/ijct.2012.0086. 604
- 5. Flach, T., Dukkipati, N., Terzis, A., Raghavan, B., Cheng, Yu., Cardwell, N., Jain, A., Hao, S., Katz-Bassett, E. & Govindan, R. (2013) Reducing Web Latency: the Virtue of Gentle Aggression. *ACM SIGCOMM*. pp. 159–170. DOI: 10.1145/2486001.2486014
- 6. Herrero, R. (2017) Modeling and comparative analysis of Forward Error Correction in the context of multipath redundancy. *Tele-communication Systems. Modelling, Analysis, Designand Management.* 65(4). pp. 783–794. DOI: 10.1007/s11235-016-0267-y
- 7. Langley, A., Riddoch, A., Wilk, A., Vicente, A., Krasic, C., Zhang, D., Ang, F., Kouranov, F., Swett, I., Iyengar, J., Bailey, J., Dorfman, J., Roskind, J., Kulik, J., Westin, P., Tenneti, R., Shade, R., Hamilton, R., Vasiliev, V., Chang, W.T. & Shi, Z. (2017) The QUIC Transport Protocol: Design and Internet-Scale Deployment. SIGCOMM"17. Los Angeles, CA, USA. pp. 183–196.
- 8. Boguslavsky, L.B. & Gelenbe, E. (1980) Analytical models transmission link control procedures for data computer networks with packet. *Automation and Remote Control*. 7. pp. 181–192.
- 9. Gelenbe, E., Labetoulle, J. & Pujolle, G. (1978) Performance Evaluation of the HDLC Protocol. *Computer Networks*. 2(4/5). pp. 409–415. DOI: 10.1016/0376-5075(78)90019-3
- 10. Kokshenev, V.V. & Sushchenko, S.P. (2008) Analysis of the asynchronous performance management procedures link transmission data. *Vychislitel'nye tekhnologii Computational Technologies*. 15(5). pp. 61–65.
- 11. Kokshenev, V.V., Mikheev, P.A. & Sushchnenko, S.P. (2017) Comparative Analysis of the Performance of Selective and Group Repeat Transmission Models in a Transport Protocol. *Automation and Remote Control.* 78(2). pp. 65–81. DOI: 10.1134/S0005117917020059
- 12. Kokshenev, V.V., Mikheev, P.A. & Sushchenko, S.P. (2013) Transport protocol selective acknowledgements analysis in loaded transmission data path. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie, vychislitel'naya tekhnika i informatika Tomsk State University Journal of Control and Computer Science*. 3(24). pp. 78–94.
- 13. Kokshenev, V.V. & Sushchenko, S.P. (2015) [Modeling sessions with Markov's chains]. *Teoriya veroyatnostey, sluchaynye protsessy, matematicheskaya statistika i prilozheniya* [Theory of probability, random processes, mathematical statistics and applications]. Proc. of the International Conference. Minsk, February 23–26, 2015. pp. 311–316.
- 14. Mikheev, P.A., Sushchenko, S.P. & Tkachev, R.V. (2017) Estimation of High-Speed Performance of the Transport Protocol with the Mechanism of Forward Error Correction. *Communications in Computer and Information Science*. 700. pp. 259–268.
- 15. Olifer, V.G. & Olifer, N.A. (2016) Komp'yuternye seti. Printsipy, tekhnologii, protokoly [Computer Networks. Principles, Technologies, Protocols]. 5th ed. St. Petersburg: Piter.

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

УДК 004.272

DOI: 10.17223/19988605/50/12

А.А. Пазников

РАСПРЕДЕЛЕННАЯ ОЧЕРЕДЬ С ОСЛАБЛЕННОЙ СЕМАНТИКОЙ ВЫПОЛНЕНИЯ ОПЕРАЦИЙ В МОДЕЛИ УДАЛЕННОГО ДОСТУПА К ПАМЯТИ

Публикация выполнена при финансовой поддержке $P\Phi\Phi U$ и СИТМА в рамках научных проектов № 19-07-00784, № 18-57-34001 и Совета по грантам Президента $P\Phi$ (проект СП-4971.2018.5).

Модель удаленного доступа к памяти (RMA) является перспективным средством повышения эффективности и упрощения разработки параллельных программ для распределенных вычислительных систем. Модель реализована в стандарте MPI (Message Passing Interface) и применяется в языках семейства PGAS (Partitioned Global Address Space). Предлагается оригинальный подход для решения актуальной задачи разработки в модели RMA масштабируемых распределенных структур данных. Основная идея заключается в ослаблении (relaxation) семантики выполнения операций. Исследуется эффективность созданной ослабленной распределенной очереди; экспериментально показано, что подход обеспечивает большую эффективность по сравнению со структурами строгой семантики.

Ключевые слова: распределенная очередь; ослабленные структуры данных; масштабируемость; удаленный доступ к памяти; MPI; RMA.

При разработке параллельных программ для вычислительных систем (BC) одной из ключевых является задача синхронизации процессов (потоков), обращающихся к разделяемым (concurrent, thread-safe) структурам данных. Разделяемые структуры данных являются базовым элементом в параллельном программировании, поэтому эффективность синхронизации существенно влияет на время выполнения программ. Такие структуры должны обеспечивать доступ параллельных процессов (потоков) в произвольные моменты времени [1–3].

Синхронизация в ВС реализуется средствами блокировок (locks) и неблокируемых (non-blocking) структур данных. Блокировки обладают интуитивной семантикой и часто не менее эффективны по сравнению с неблокируемыми методами. В то же время программирование без блокировок позволяет избежать тупиковых ситуаций (deadlocks), инверсий приоритетов (priority inversion) и обеспечивает гарантии выполнения. Независимо от реализации большая часть структур данных характеризуется наличием узких мест (bottlenecks) для операций, таких как вставка и удаление элементов (очереди, стеки), удаление максимального элемента (очереди с приоритетом).

Большая часть работ в области разделяемых структур данных ориентированы на ВС с общей памятью, производительность которых может быть недостаточной для решения современных задач. Например, размеры графов социальных сетей достигают нескольких петабайт, а число вершин в графах из теста Graph500 — нескольких триллионов. Проекты в области физики высоких энергий, такие как CMS и Atlas, производят десятки петабайт ежегодно. Планируется, что Большой обзорный телескоп будет каждую ночь генерировать около 20 терабайт. Поскольку ВС с общей памятью имеют технологический предел числа процессорных ядер, для решения таких задач требуется использовать ВС с распределенной памятью (кластерные ВС, ВС с массовым параллелизмом). В процессе программирования таких систем обрабатываемые данные представляются в виде

распределенных структур данных, для которых необходимо обеспечить масштабируемую синхронизацию.

Одной из перспективных моделей параллельного программирования для ВС с распределенной памятью является модель удаленного доступа к памяти (Remote Memory Access, RMA), реализованная в стандарте MPI [4, 5]. В рамках RMA процессы непосредственно обращаются к памяти других процессов вместо отправки и получения сообщений. В отличие от модели разделенного глобального адресного пространства (Partitioned Global Address Space, PGAS), представленной языками UPC, CAF, Chapel, X10, модель RMA тесно интегрирована с библиотеками MPI и может быть использована наравне с моделью передачи сообщений. Программы в модели RMA характеризуются меньшим временем выполнения по сравнению с моделью передачи сообщений и PGAS. Большая часть современных коммуникационных сетей (Infiniband, PERCS, Gemini, Aries, RoCE over Ethernet) обеспечивает поддержку RMA с помощью технологии RDMA [4], реализующей обращение к удаленным сегментам памяти без участия центрального процессора.

Опишем программную модель RMA в MPI. Основными являются неблокируемые функции MPI_Put (запись в память удаленного процесса) и MPI_Get (чтение из удаленной памяти), атомарные MPI_Accumulate, MPI_Get_accumulate, MPI_Compare_and_swap. RMA-вызовы должны находиться внутри областей (эпохи, epochs), в рамках которых выполняется синхронизация. В работе применяется пассивный метод синхронизации (passive target synchronization), реализованный в стандарте MPI [5]. При пассивной синхронизации процесс открывает эпоху реализации удаленного доступа (access epoch) посредством вызова функций MPI_Win_lock/MPI_Win_lockall, после чего он может выполнять RMA-операции для доступа к зарегистрированным сегментам памяти (окна, windows) других процессов. Таким образом, RMA-операции выполняются в одностороннем порядке, без явного вызова функций синхронизации другими процессами.

Основная часть работ в области разделяемых структур данных направлена на создание средств синхронизации для ВС с общей памятью. К ним относятся алгоритмы блокировки потоков [1, 6] (TTS, Backoff, CLH, MCS, Oyama, Flat Combining, RCL и др.). Хотя некоторые методы (Hierarchical Backoff (CLH, MCS), Cohorting и др.) учитывают отдельные иерархические уровни, они неприменимы в ВС с распределенной памятью. Неблокируемые структуры [1-3, 7] также разработаны для многоядерных ВС и неприменимы в распределенных ВС. Перспективным методом повышения масштабируемости разделяемых структур данных является ослабление их семантики (relaxation) [8-14]. Например, в ослабленной очереди с приоритетом извлекается не максимальный элемент, а элемент, близкий к максимальному. В ослабленной очереди (стеке) удаляется не первый (последний) добавленный элемент, а элемент в его окрестности. Ослабленные структуры обеспечивают высокую пропускную способность и приемлемый уровень упорядоченности операций в реальных программах. В работах [8, 9] для построения потокобезопасной очереди с приоритетом предлагается использовать набор последовательных очередей. Аналогичная реализация стека, основанного на временных метках, предлагается в [10]. Также построены аналитические модели ослабления [13, 14], включая квазилинеаризуемость (quasi linearizability), количественное ослабление (quantitative relaxation).

Насколько известно, для распределенных ВС не разработаны эффективные масштабируемые разделяемые структуры данных. Методы, предложенные для распределенных ВС, включают простые спинлоки, блокировки чтения-записи и МСS-блокировки [15]. Работы, посвященные распределенным структурам данных [16–18], неприменимы в модели RMA. В языках PGAS реализованы отдельные примитивы синхронизации и распределенные структуры, но они характеризуются наличием узких мест и высокими накладными расходам. Исходя из вышесказанного, задача разработки эффективных разделяемых структур данных для распределенных ВС является востребованной и нерешенной в настоящее время. В данной статье предлагается метод построения масштабируемых распределенных структур на основе ослабления их семантики, рассмотренный на примере очереди.

1. Распределенная ослабленная очередь

1.1. Ослабление семантики выполнения операций для распределенных очередей

Очередь – коллекция объектов, реализующая дисциплину FIFO («первым вошел – первым вышел»). Основные операции: добавление (insert, enqueue) элемента в последнюю позицию (хвост, tail) и извлечение (remove, dequeue) элемента из первой позиции (голова, head). В классических реализациях распределенных очередей необходимо обеспечивать актуальность данных о расположении (ранг процесса) головного и хвостового элементов. Процесс перед выполнением операций при необходимости обновляет данные о расположении головного (хвостового) элемента. Это приводит к дополнительным накладным расходам и увеличивает время выполнения операций. Другим значимым недостатком является наличие узких мест при одновременном обращении нескольких процессов к процессу, в памяти которого находится головной (хвостовой) элемент.

С целью увеличения масштабируемости распределенной очереди предлагается ослабить ее семантику и допустить извлечение элемента из окрестности первого добавленного элемента. Для этого распределенная структура представляется в виде множества последовательных структур, распределенных между процессами. Каждый процесс асинхронно обращается к удаленным сегментам посредством RMA-вызовов (рис. 1). Данный подход не предполагает выполнения операций для актуализации данных о расположении головы и хвоста очереди и позволяет избежать возникновения узких мест. Кроме того, за счет низкой латентности односторонних коммуникаций и аппаратной поддержки RDMA метод гарантирует снижение времени выполнения операций.

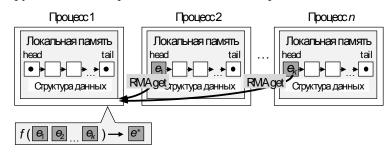


Рис. 1. Операция извлечения элемента в ослабленной распределенной очереди Fig. 1. Execution of item remove operation for relaxed distributed queue

Опишем операции добавления и удаления элементов. Обозначим p — число процессов. Считаем, что часы процессов синхронизированы и каждому элементу e очереди соответствует временная метка t(e) — момент добавления его в очередь. При выполнении операции insert процесс случайным образом выбирает очередь $s \in \{1, 2, ..., p\}$ и помещает в нее элемент. Отметим, что вместо случайного выбора можно использовать другие схемы для локализации обращений к памяти и других оптимизаций.

При выполнении операции удаления remove процесс выбирает k очередей-кандидатов $R \subseteq \{1, 2, ..., p\}$ и посредством RMA-операций получает значения элементов $\{e_1, e_2, ..., e_k\}$, находящихся в голове соответствующих очередей. Далее среди элементов-кандидатов определяется «лучший» элемент с минимальной временной меткой: $e^* = \operatorname{argmin}_{i \in \{1, 2, ..., k\}} t(e_i)$ (рис. 1).

Рассмотрим детально реализацию распределенной ослабленной очереди. При инициализации структуры данных на каждом процессе организуется циклический буфер фиксированного достаточно большого размера (100 000 в данной реализации) и синхронизируются часы процессов [19].

1.2. Операция добавления элементов

Входными данным операции insert добавления элемента являются значение элемента val, число процессов p, окно для выполнения RMA-операций win и массив блокировок locks для защиты данных в очередях на каждом процессе. Блокировки могут быть реализованы с помощью любого спинлока, в данной работе используется простой алгоритм TASLock [1]. Под MPI_Put_Atomic и MPI_Get_atomic

здесь и далее понимаются атомарные операции MPI_Put и MPI_Get, реализованные на основе MPI_Accumulate и MPI_Get_accumulate.

Таблица 1

Алгоритмы выполнения операций для распределенной ослабленной очереди: a — операция insert добавления элемента в очередь, δ — операция remove извлечения элемента из очереди

Входные val — добавляемый элемент locks — массив блокировок для защиты очередей p — число процессов коммуникатора win — окно для выполнения RMA-

```
операций
1
      nqueues = p
2
      do
3
        rank = GETRAND(p)
4
        elem.val = val
5
        elem.ts = GETTIMESTAMP()
        MPI WIN LOCK(rank, win)
6
7
        Lock(rank, locks[rank], win)
8
        MPI_GET_ATOMIC(rank, state, win)
9
        MPI_WIN_FLUSH(win)
10
        if IsFull(state) then
11
          nqueues = nqueues - 1
12
          if nqueues = 0 then
13
             UnLock(rank, locks[rank], win)
14
             MPI_WIN_UNLOCK(rank, win)
15
             return ErrQueueFull
16
          end if
17
        else
18
          MPI_Put(rank, elem, win)
19
          state.tail = (state.tail + 1) \mod size
20
          MPI_PUT_ATOMIC(rank, state.tail, win)
21
        end if
22
        UNLOCK(rank, locks, win)
23
        MPI_WIN_UNLOCK(rank, win)
24
      while IsFull(state)
```

```
Вхолные
             ncand – число очередей-кандидатов
             locks - массив блокировок для защиты очере-
данные:
            р – число процессов коммуникатора
             win – окно для выполнения RMA-операций
     MPI_WIN_LOCK_ALL(win)
2
     ncurr = 0
3
     navail = p
4
     nattempts = 0
5
     while ncurr < ncand do
        rank = GETRAND(p)
6
7
       rc = TRYLOCK(rank, locks[rank], win)
8
       if LockIsAcquired(rc) then
9
          MPI_GET_ATOMIC(rank, states[ncurr], win)
10
          MPI_WIN_FLUSH(win)
11
          if IsEmpty(state) then
12
            UNLOCK(rank, locks[rank], win)
13
            navail = navail - 1
14
            if navail < ncand then
15
              if ncand = 0 then
16
                 MPI_WIN_UNLOCK_ALL(win)
                 return ErrQueueEmpty
17
18
              end if
19
              ncand = navail
            end if
20
21
          else
22
            MPI_Get(rank, elems[ncurr], win)
23
            MPI_WIN_FLUSH(win)
24
            ADDCAND(rank, cands)
25
            ncurr = ncurr + 1
26
            nattempts = 0
27
          end if
2.8
        else if LOCKISBUSY(rc) then
29
          if ncurr > 0 then
30
            nattempts = nattempts + 1
31
            if nattempts = max\_nattempts then
32
              for i = 0 to ncurr do
33
                 UNLOCK(cands[i], locks[cands[i]], win)
34
              end for
35
              ncurr = 0
            end if
36
37
          end if
38
        end if
39
     end while
40
     bestrank = GETBESTRANK(cands, elems)
41
     states[bestrank].tail = (states[bestrank].tail + 1) mod
42
43
     MPI_PUT_ATOMIC(bestrank, states[bestrank].tail, win)
44
     for i = 0 to ncand do
45
       UnLock(cands[i], locks[cands[i]], win)
46
```

MPI_WIN_UNLOCK_ALL(win) **return** *elems*[*bestrank*].*val*

Основные шаги алгоритма (табл. 1, a):

- 1. Проинициализировать число доступных очередей *nqueues* (строка 1).
- 2. Случайным образом выбрать процесс *rank* (строка 3). Установить поля элемента (строки 4, 5). Начать эпоху синхронизации для выбранного процесса (строка 6). Заблокировать очередь процесса *rank* (строка 7) и получить ее состояние (строка 8).
- 3. Если очередь заполнена (строка 10), уменьшить *nqueues*. Если нет доступных очередей (строка 12), разблокировать очередь, завершить эпоху синхронизации, вернуть код ошибки (строки 13–15).
- 4. Если очередь не полна, добавить в нее элемент (строка 18), увеличить указатель *state.tail* на хвост очереди (строка 19) и установить новое значение состояния очереди (строка 20).
 - 5. Разблокировать очередь (строка 22) и завершить эпоху пассивной синхронизации (строка 23).
 - 6. Выполнять шаги 2-5, пока как минимум одна очередь не будет найдена (строка 24).

1.3. Операция удаления элементов

Входными данными для функции remove удаления элемента являются число кандидатов *ncand* для выбора элемента, количество процессов p, окно для выполнения RMA-операций win и массив блокировок locks для защиты данных очередей. Операция включает следующие шаги (табл. 1, b):

- 1. Начать эпоху пассивной синхронизации для всех процессов (строка 1), проинициализировать текущее число найденных очередей *ncurr* (строка 2), число доступных очередей *navail* (строка 3) и число попыток блокировки очереди *nattempts* (строка 4).
- 2. Если текущее число найденных кандидатов *ncurr* равно *ncand*, перейти на шаг 7. Если нет, случайно выбрать очередь *rank* (строка 6). Попытаться заблокировать очередь (строка 7).
 - 3. Если очередь заблокирована, получить ее состояние *state* (строка 9). Если нет, переход на шаг 6.
- 4. Если очередь пуста (строка 11), разблокировать ее (строка 12), уменьшить *navail* (строка 13). Если *navail* < *ncand*, сбросить *ncand* до значения *navail*. Если не осталось доступных кандидатов (строка 15), завершить эпоху синхронизации (строка 16) и вернуть код ошибки (строка 17).
- 5. Если очередь не пуста, получить элемент в голове (строка 22), добавить в список кандидатов (строка 24), увеличить *ncurr* (строка 25), сбросить *nattempts* (строка 26). Перейти на шаг 2.
- 6. Если очередь не захвачена, увеличить *nattempts* (если это не первая очередь-кандидат) (строка 30). Если *nattempts* достигло максимального (строка 31), разблокировать захваченные очереди (строки 32–34), сбросить *ncurr* (строка 35), перейти на шаг 2. Данный шаг необходим для избежания взаимной блокировки, когда два процесса пытаются заблокировать уже захваченные очереди.
- 7. Выбрать кандидата с минимальным значением временной метки (строка 40). Для данной очереди инкрементировать указатель на голову и обновить состояние очереди (строка 41, 42). Разблокировать все очереди-кандидаты (строки 43, 45) и завершить эпоху синхронизации (строка 46).

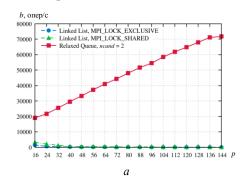
2. Проведение экспериментов

Экспериментальное исследование ослабленной очереди проводилось на вычислительном кластере Јеt Центра параллельных вычислительных технологий Сибирского государственного университета телекоммуникаций и информатики. Кластер укомплектован 18 вычислительными узлами, оборудованными двумя 4-ядерными процессорами Intel Xeon E5420 (суммарное число ядер 144). В качестве МРІ-библиотеки применялась МРІСН 3.2.1.

Разработан синтетический тест, выполняющий $n=100\,000$ операций вставки / удаления (тип операции выбирается случайно). Число процессов p варьировалось от 16 до 144. Созданная распределенная ослабленная очередь Relaxed Queue сравнивалась со связным списком, реализованным в библиотеке MPICH (MPICH linked list) в модели RMA. Использовалось два типа списка: на основе эксклюзивной и разделяемой пассивной синхронизации (MPI_LOCK_EXCLUSIVE и MPI_LOCK_SHARED). Тип синхронизации определяет, допускается ли одновременное обращение нескольких процессов

к памяти удаленного процесса. Также исследовалось влияние числа очередей-кандидатов *ncand* на эффективность очереди. Для этого *ncand* варьировалось от 1 до 4. Кроме того, анализировалась зависимость эффективности очереди от типа пассивной синхронизации в функции вставки. Измерялась пропускная способность b = t / n, где t – время проведения эксперимента.

Пропускная способность разработанной очереди значительно превосходит пропускную способность линейного списка строгой семантики (рис. 2, a). Оптимизация достигается за счет сокращения накладных расходов, возникающих при выполнении доступа к элементам. Недостатки классических распределенных списков — необходимость актуализации данных о расположении головного (хвостового) элементов и возможность образования узких мест при одновременном обращении нескольких процессов к ним. Разработанная ослабленная очередь не требует поддержания согласованного состояния головы (хвоста) очереди, поскольку каждый раз процесс-кандидат и соответствующая очередь выбираются случайно. Такой подход также позволяет распределить нагрузку между процессами и избежать возникновения узких мест. Линейный список строгой семантики на основе разделяемой синхронизации более эффективен, по сравнению с эксклюзивным режимом (рис. 2, b), поскольку во время вставки / удаления несколько процессов одновременно обращаются к одному процессу, в памяти которого находится головной (хвостовой) элемент.



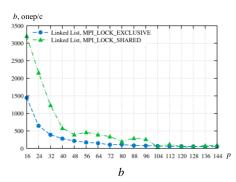
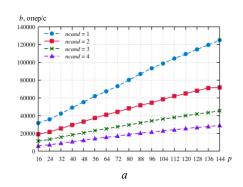


Рис. 2. Сравнение эффективности структур данных: a – пропускная способность; b – пропускная способность линейного списка строгой семантики для разных режимов синхронизации Fig. 2. Comparison of efficiency of data structures: a – throughput; b – throughput of list with strong semantics for different synchronization modes



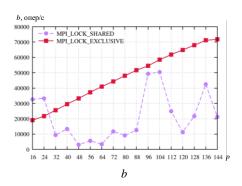


Рис. 3. Анализ эффективности Relaxed Queue: a – пропускная способность в зависимости от числа очередей-кандидатов; δ – пропускная способность ослабленной очереди для разных режимов синхронизации Fig. 3. Analysis of Relaxed Queue efficiency: a – throughput depending on number of candidates; b – throughput of Relaxed Queue for different synchronization modes

Как и ожидалось, пропускная способность уменьшается с увеличением числа кандидатов *ncand* (рис. 3, *a*). Это объясняется дополнительными накладными расходами на выполнение блокировки очередей, получение состояний и значений элементов очередей-кандидатов. На наш взгляд, *ncand* = 2 является достаточным для большинства случаев и обеспечивает приемлемый для практики уровень упорядоченности операций. Данные выводы согласуются с результатами аналогичной структуры для ВС с общей памятью [8]. Тем не менее для повышения близости порядка вставки / удаления элемен-

тов к порядку FIFO можно увеличить *ncand* до 3 и 4. В данной работе не выполняется оценка близости к FIFO, но это планируется сделать в будущем.

В отличие от распределенных списков со строгой семантикой, для ослабленной очереди эксклюзивный режим пассивной синхронизации обеспечивает большую пропускную способность по сравнению с разделяемым режимом (рис. 3, b). Это объясняется тем, что организация разделяемого режима является дорогостоящей операцией. Вместе с тем в ослабленной очереди одновременное обращение к одному процессу (последовательной очереди) является редким событием. Поэтому мы полагаем, что разделяемый режим является избыточным.

Заключение

В данной статье разработаны эвристические алгоритмы реализации ослабленных распределенных очередей в модели RMA. Созданная очередь основана на множестве последовательных очередей, распределенных между процессами. Очередь характеризуется значительно большей пропускной способностью по сравнению с линейными списками строгой семантики в модели RMA. Оптимизация достигается за счет устранения узких мест при выполнении операций. При реализации очереди рекомендуется использовать 2 или 3 очереди-кандидата и эксклюзивный тип пассивной синхронизации (MPI_LOCK_EXCLUSIVE).

ЛИТЕРАТУРА

- 1. Herlihy M., Shavit N. The art of multiprocessor programming. Morgan Kaufmann, 2012. 537 p.
- 2. Mark M., Shavit N. Concurrent Data Structures. Chapman and Hall / CRC Press, 2004. 32 p.
- 3. Shavit N. Data structures in the multicore age // Communications of the ACM. 2011. V. 54. P. 76–84.
- 4. Liu J., Wu J., Panda D.K. High performance RDMA-based MPI implementation over InfiniBand // International Journal of Parallel Programming, 2004. V. 32. P. 167–198.
- 5. Hoefler T., Dinan J., Thakur R., Barrett B., Balaji P., Gropp W., Underwood K. Remote memory access programming in MPI-3 // ACM Transactions on Parallel Computing. 2015. V. 2, No. 2. P. 9.
- 6. Пазников А.А. Оптимизация делегирования выполнения критических секций на выделенных процессорных ядрах // Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2017. № 38. С. 52–58.
- 7. Аненков А.Д., Пазников А.А. Алгоритмы оптимизации масштабируемого потокобезопасного пула на основе распределяющих деревьев для многоядерных вычислительных систем // Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2017. № 39. С. 73–84.
- 8. Rihani H., Sanders P., Dementiev R. Brief announcement: Multiqueues: Simple relaxed concurrent priority queues // Proc. of the 27th ACM symposium on Parallelism in Algorithms and Architectures. 2015. P. 80–82.
- 9. Табаков А.В., Пазников А.А. Алгоритмы оптимизации потокобезопасных очередей с приоритетом на основе ослабленной семантики выполнения операций // Известия СПбГЭТУ «ЛЭТИ». 2018. № 10. С. 42–49.
- 10. Dodds M., Haas A., Kirsch C.M. A scalable, correct time-stamped stack // ACM SIGPLAN Notices. 2015. V. 50, No. 1. P. 233-246.
- 11. Alistarh D., Kopinsky J., Li J., Shavit N. The Spray List: a scalable relaxed priority queue // ACM SIGPLAN Notices. 2015. V. 50, No. 8. P. 11–20.
- $12.\ Wimmer\ M.\ et\ al.\ The\ lock-free\ k-LSM\ relaxed\ priority\ queue\ //\ ACM\ SIGPLAN\ Notices.\ 2015.\ V.\ 50,\ No.\ 8.\ P.\ 277-278.$
- 13. Henzinger T.A., Kirsch C.M., Payer H., Sezgin A., Sokolova A. Quantitative relaxation of concurrent data structures // ACM SIGPLAN Notices. 2013. V. 48, No. 1. P. 317–328.
- 14. Afek Y., Korland G., Yanovsky E. Quasi-Linearizability: Relaxed Consistency for Improved Concurrency // Int. Conf. on Principles of Distributed Systems. 2010. P. 395–410.
- 15. Schmid P., Besta M., Hoefler T. High-Performance Distributed RMA Locks // Proc. of the 25th ACM Int. Symposium on High-Performance Parallel and Distributed Computing, HPDC 2016, Kyoto, Japan, May 31 June 04, 2016. ACM 2016. P. 19–30.
- 16. Mans B. Portable distributed priority queues with MPI // Concurrency Practice and Experience. 1998. V. 10, No. 3. P. 175–198.
- 17. Brodal G.S., Traff J.L., Zaroliagis C.D. A parallel priority queue with constant time operations // J. of Parallel and Distributed Computing. 1998. Vol. 49, No. 1. P. 4–21.
- 18. Zanny R. Efficiency of distributed priority queues in parallel adaptive integration. MS thesis. Kalamazoo, MI : Western Michigan University, 1999. 148 p.
- 19. Курносов М.Г. MPIPerf: пакет оценки эффективности коммуникационных функций стандарта MPI // Вестник Нижегородского университета им. Н.И. Лобачевского. 2012. № 5 (2). С. 385–391.

Поступила в редакцию 5 апреля 2019 г.

Paznikov A.A. (2020) DISTRIBUTED RELAXED QUEUE IN REMOTE MEMORY ACCESS MODEL. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 97–105

DOI: 10.17223/19988605/50/12

Remote memory access (RMA) technique is a very attractive way for improving efficiency of high-performance computations (HPC) and simplifying parallel program development. Unlike message-passing, in RMA programs processes communicate by directly accessing each other's memories. RMA model is implemented in MPI standard and offers Partitioned Global Address Space (PGAS). Many applications have been shown to benefit from RMA communications, where a process directly accesses the memory of another process instead of sending and receiving messages. To the best of knowledge, there are no efficient scalable concurrent data structures designed in RMA model.

In modern computer systems, multiple processes (threads) execute concurrently and synchronize their activities through shared (concurrent) data structures. Such data structures are therefore key building blocks in parallel programming, and their efficiency is crucial for the overall performance. Concurrent data structures are, however, much more difficult to design than their sequential counterparts, as processes executing concurrently might interleave their shared-memory steps in very complicated ways, with often unexpected outcomes. Coming up with efficient concurrent data structures for distributed environments with deep hierarchy, such as computer clusters and data centers, is challenging. A lot of prior work focused on designing efficient synchronization techniques for shared-memory systems. However, in such systems the shared memory itself may become an impediment for scalability. Moreover, shared-memory systems are not sufficient for processing of large data volumes in current applications. Hence, there is growing demand for efficient concurrent data structures in hierarchical distributed systems (supercomputers, clusters, grids).

Regardless of the design, many data structures are subject to inherent sequential bottlenecks for some operations, such as the delete-min operation for priority queues or insert and remove operations for queues and stacks. A promising way to alleviate the bottleneck problem is relaxing the consistency requirements for such operations. There are evidences that, on most workloads, relaxed data structures outperform data structures with strict semantics and ensure acceptable degrees of operation reordering. However, to the best of our knowledge, nobody looked at relaxed concurrent structures in the distributed environment.

As data structures to study, we consider relaxed queues. Relaxed queues do not guarantee strict FIFO order: remove operation might not remove exactly the first inserted element, but an element close to it. We propose an approach based on multiple sequential data structures distributed among the processes. This approach is well approved for shared-memory systems and outperform data structures with strong ordering. Every process can asynchronously access to the remote segment via RMA calls. Thanks to low latency of one-sided communications and hardware support of RDMA this scheme will guarantee high performance. When a process executes insert operation for the relaxed queue, it sets timestamp value and just picks (randomly or by a specified algorithm) the remote process and inserts the element along with timestamp to its data structure. When it executes a remove operation, the calling process selects a subset of other processes and remotely gets "candidate" elements from the set of processes. Finally, it chooses among candidate elements the "best" one with minimal timestamp and returns it.

We evaluated developed relaxed queue on computer cluster. In experiments we compared developed distributed queue with the linked list, implemented in MPICH library (MPICH linked list) in RMA model. Throughput of developed relaxed distributed queue substantially outperforms MPICH linked lists. Optimization was achieved by reducing overheads for communications while accessing the elements of the structures. The main drawback of common distributed lists that the head pointers become sequential bottlenecks. Unlike them developed distributed queue has multiple access points distributed among the processes and no bottlenecks. In the work we also investigate the influence of candidate elements number and chosen type of synchronization on the efficiency of the queues. As expected, throughput is decreasing with increase of number of candidates. This is explained the additional overheads for locking, getting states and values from candidate queues. We also found for the relaxed queue exclusive mode of synchronization provides better throughput compared with shared mode.

Thus, proposed decentralized asynchronous approach for designing relaxed distributed data structures, as we expected, eliminates bottlenecks, minimizes latency of the operations and provides high scalability of parallel programs.

Keywords: distributed queue; relaxed data structures; remote memory access; MPI; RMA.

PAZNIKOV Alexey Aleksandrovich (Candidate of Technical Sciences, Senior Researcher, Department of Computer Science and Engineering, Saint Petersburg Electrotechnical University "LETI", Saint Petersburg, Russian Federation).

E-mail: apaznikov@gmail.com

REFERENCES

- 1. Herlihy, M. & Shavit, N. (2012) The Art of Multiprocessor Programming. Morgan Kaufmann.
- 2. Mark, M. & Shavit, N. (2004) Concurrent Data Structures. Chapman and Hall/CRC Press.
- 3. Shavit, N. (2011) Data structures in the multicore age. *Communications of the ACM*. 54. pp. 76–84. DOI: 10.1145/1897852.1897873

- 4. Liu, J., Wu, J. & Panda, D.K. (2004) High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming*. 32. pp. 167–198. DOI: 10.1023/B:IJPP.0000029272.69895.c1
- Hoefler, T., Dinan, J., Thakur, R., Barrett, B., Balaji, P., Gropp, W. & Underwood, K. (2015) Remote memory access programming in MPI-3. ACM Transactions on Parallel Computing. 2(2). pp. 9. DOI: 10.1145/2780584. 30
- 6. Paznikov, A.A. (2017) Optimization method of remote core locking. Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie, vychislitel'naya tekhnika i informatika Tomsk State University Journal of Control and Computer Science. 38. pp. 52–58. DOI: 10.17223/19988605/38/8
- 7. Anenkov, A.D. & Paznikov, A.A. (2017) Algorithms of optimization of scalable thread-safe pool based on diffracting trees for multicore computing systems. Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie, vychislitel'naya tekhnika i informatika Tomsk State University Journal of Control and Computer Science. 39. pp. 73–84. DOI: 10.17223/19988605/39/10
- 8. Rihani, H., Sanders, P. & Dementiev, R. (2015) Brief announcement: Multiqueues: Simple relaxed concurrent priority queues. Proceedings of the 27th ACM symposium on Parallelism in Algorithms and Architectures. pp. 80-82. DOI: 10.1109/EIConRus.2019.8657105
- 9. Tabakov, A.V. & Paznikov, A.A. (2018) Algorithms for optimization of relaxed concurrent priority queues in multicore systems. *Izvestia SPbGETU "LETI"*. 10. pp. 42–49.
- Dodds, M., Haas, A. & Kirsch, C.M. (2015) A scalable, correct time-stamped stack. ACM SIGPLAN Notices. 50(1). pp. 233–246.
 DOI: 10.1145/2775051.2676963
- 11. Alistarh, D., Kopinsky, J., Li, J. & Shavit, N. (2015) The Spray List: A scalable relaxed priority queue. *ACM SIGPLAN Notices*. 50(8). pp. 11–20.
- 12. Wimmer, M., Gruber, J., Träff, J.L., & Tsigas, P. (2015) The lock-free k-LSM relaxed priority queue. *ACM SIGPLAN Notices*. 50(8). pp. 277–278. DOI: 10.1145/2858788.2688547
- 13. Henzinger, T.A., Kirsch, C.M., Payer, H., Sezgin, A. & Sokolova, A. (2013) Quantitative relaxation of concurrent data structures. *ACM SIGPLAN Notices*. 48(1). pp. 317–328. DOI: 10.1145/2480359.2429109
- 14. Afek, Y, Korland, G. & Yanovsky, E. (2010) Quasi-Linearizability: Relaxed Consistency for Improved Concurrency. In: Lu, C., Masuzawa, T. & Mosbah, M. (eds) *Principles of Distributed Systems. OPODIS 2010. Lecture Notes in Computer Science*. Vol 6490. Berlin, Heidelberg: Springer.
- 15. Schmid, P., Besta, M. & Hoefler, T. (2016) High-Performance Distributed RMA Locks. *Proc. of the 25th ACM Int. Symposium on High-Performance Parallel and Distributed Computing, HPDC 2016.* Kyoto, Japan, May 31 June 4, 2016. ACM 2016. pp. 19–30.
- 16. Mans, B. (1998) Portable distributed priority queues with MPI. Concurrency Practice and Experience. 10(3). pp. 175-198.
- 17. Brodal, G.S., Traff, J.L. & Zaroliagis, C.D. (1998) A parallel priority queue with constant time operations. *Journal of Parallel and Distributed Computing*. 49(1). pp. 4–21. DOI: 10.1006/jpdc.1998.1425
- 18. Zanny, R. (1999) Efficiency of distributed priority queues in parallel adaptive integration. MS Thesis. Western Michigan University.
- 19. Kurnosov, M.G. (2012) MPIPerf: paket otsenki effektivnosti kommunikatsionnykh funktsiy standarta MPI [MPIPerf: a toolkit for benchmarking MPI libraries]. *Vestnik Nizhego-rodskogo universiteta im. N.I. Lobachevskogo Vestnik of Lobachevsky University of Nizhni Novgorod.* 5(2). pp. 385–391.

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

УДК 004.7:519.872

DOI: 10.17223/19988605/50/13

О.В. Семёнова, З.Т. Буй

ПАКЕТ ПРИКЛАДНЫХ ПРОГРАММ ДЛЯ ИССЛЕДОВАНИЯ СИСТЕМ ПОЛЛИНГА

Работа выполнена при финансовой поддержке фонда РФФИ, проект № 18-5700002.

Приведено описание пакета прикладных программ, предназначенного для расчета характеристик систем массового обслуживания с несколькими очередями и общим сервером, а именно систем поллинга с различными видами порядка опроса очередей: циклическим, адаптивным циклическим, случайным; и различными дисциплинами обслуживания: шлюзовой, исчерпывающей, глобально-шлюзовой, ограниченной, пороговой и случайной. Входной поток может быть простейшим, потоком фазового типа или коррелированным МАР-потоком. Описаны основные структуры и механизм работы пакета прикладных программ. Приведены численные примеры, иллюстрирующие результаты исследования систем поллинга с входным коррелированным потоком.

Ключевые слова: системы поллинга; пакет прикладных программ; имитационное моделирование; МАРпоток.

Системы поллинга представляют собой системы массового обслуживания с несколькими очередями (или несколькими потоками заявок) [1, 2]. Обслуживающий прибор по определенному правилу посещает очереди и обслуживает находящиеся в них заявки. Системы поллинга эффективно используются для оценки производительности, проектирования и оптимизации структуры телекоммуникационных систем и сетей, транспортных систем и систем управления дорожным движением, производственных систем и систем управления запасами [3]. Для исследования математических моделей систем поллинга применяют различные методы: метод производящих функций, метод средних, метод ветвящихся процессов и другие [4] для получения точных формул вычисления характеристик производительности систем, а также различные эвристические методы для систем поллинга специального вида [5–8] и имитационное моделирование в случаях, когда проведение точного анализа системы не представляется возможным либо когда необходимо оценить точность приближенных результатов исследования.

Разработка пакета прикладных программ для расчета характеристик и имитационного моделирования систем поллинга на данный момент становится весьма актуальной задачей, поскольку в литературе удается найти лишь одну работу [9], описывающую пакет программ имитационного моделирования для систем поллинга с несколькими обслуживающими устройствами (серверами), но порядок обслуживания очередей рассматривается лишь циклический. Представляемый же в данной работе пакет программ предполагает наличие только одного сервера в моделях, но рассматривается широкий класс типов опроса очередей в системах поллинга: циклический, адаптивный циклический, случайный. Пакет создан с помощью OMNeT++ Discrete Event Simulator, а также OMNeT++ Simulation Manual Version 4.6. OMNeT++ (Objective Modular Network Testbed на C++) - это модульная, основанная на компонентах библиотека моделирования и C++, в основном для создания сетевых симуляторов.

1. Модуль аналитических расчетов

Разработанный пакет программ делится на два крупных модуля (рис. 1): модуль имитационного моделирования и модуль аналитических расчетов характеристик производительности систем поллинга, реализующий формулы их расчета. Модуль аналитических расчетов характеристик систем поллинга реализован с помощью Matlab 2013 на основе точных результатов анализа [1, 3], полученных для следующих систем поллинга:

- Система циклического опроса с шлюзовой дисциплиной.
- Система циклического опроса с исчерпывающей дисциплиной.
- Система циклического опроса с глобально-шлюзовой дисциплиной.
- Система адаптивного циклического опроса с шлюзовой дисциплиной.
- Система адаптивного циклического опроса с исчерпывающей дисциплиной.

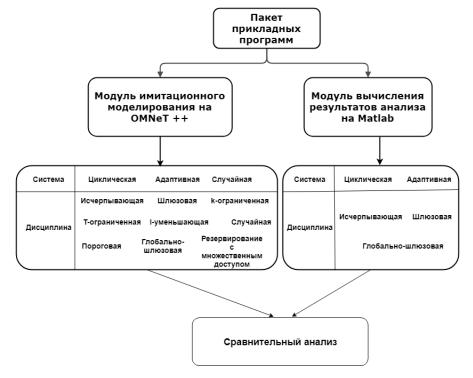


Рис. 1. Схема пакета прикладных программ

Fig. 1. The scheme of the application package

Данный модуль позволяет получить следующие характеристики:

- Вероятность опроса каждой очереди при адаптивном опросе.
- Вероятность того, что сервер пропускает k очередей после завершения обслуживания i-й очереди (переходит в (i+k)-ю очередь), $k=\overline{1,N}$.
 - Первые и вторые моменты длины каждой очереди в момент ее опроса сервером.
 - Вторые моменты длины каждой очереди в момент опроса сервера.
 - Среднее время цикла для каждой очереди.
- Среднее время пребывания заявок в системе адаптивного динамического поллинга с исчерпывающей и шлюзовой дисциплиной.

2. Модуль имитационного моделирования

Для имитационного моделирования, был использован OMNeT++ Discrete Event Simulator, а также OMNeT++ Simulation Manual Version 4.6. OMNeT++ (Objective Modular Network Testbed на C++) – это модульная, компонентно-ориентированная C++ библиотека и фреймворк для дискретно-событийного моделирования, используемая прежде всего для создания сетевых симуляторов. OMNeT++ представляет архитектуру компонентов для моделей. Компоненты (модули) запрограммированы на C++, а затем собраны в более крупные компоненты и модели с использованием языка высокого уровня (NED). OMNeT++ имеет обширную поддержку графического интерфейса, и благодаря

его модульной архитектуре ядро моделирования (и модели) может быть легко внедрено во многие приложения.

В пакете имитационного моделирования реализованы следующие модели систем поллинга: циклический опрос с шлюзовой, исчерпывающей, глобально-шлюзовой, пороговой, *k*- и *Т*-ограниченными дисциплинами, адаптивный циклический опрос с шлюзовой, исчерпывающей и глобально-шлюзовой дисциплиной, упорядоченный адаптивный динамический опрос с шлюзовой и исчерпывающей дисциплиной, циклический опрос с резервированием с множественным доступом, система циклического поллинга с пороговой исчерпывающей дисциплиной и простоями севера, система поллинга с приоритетом, система поллинга со случайным порядком опроса при шлюзовой дисциплине, система поллинга со случайным порядком опроса при исчерпывающей или шлюзовой дисциплине, а также система циклического поллинга со случайной дисциплиной обслуживания очередей. Подробную информацию об этих дисциплинах опроса и обслуживания очередей, а также описание основной модели системы поллинга можно найти в [1, 3].

Структура пакета имитационного моделирования. Программа состоит из трех основных типов файлов (рис. 2):

— Входные файлы: здесь задаются параметры модели системы поллинга: количество очередей, тип (простейший, фазового типа или MAP) и параметры входных потоков, тип и параметры процесса обслуживания заявок в каждой очереди, параметры процессов переключения сервера между очередями (табл. 1). Подробную информацию о MAP-потоках, их обобщении BMAP-потоках и распределении фазового типа можно найти в [10].



Рис. 2. Основная структура программы Fig. 2. The basic structure of the program

Таблица 1

Выход

7. График флуктуации среднего времени ожидания

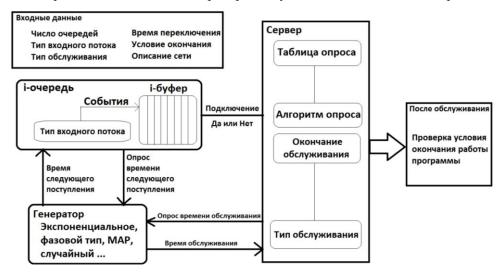
Пример входных и выходных данных

Вход

1. Delay time in Queues i: среднее время пребывание заявок 1. *.numQueues = 2: определение количества очередей. в очереди і 2. *.generator[0].interArrivalTime = exponential(1/30): ин-2. InterVisit time in Queues i: среднее время между последутенсивность поступления заявок для первой очереди. юшими посешениями в очереди і 3. *.generator[1].interArrivalTime = MAP("MAP.txt"): тип 3. Number of cycle in Queues i: общее количество циклов поступления заявок для второй очереди. в этой очереди (посещений сервером этой очереди) 4. *.threshold_gated_queue[*].serviceTime = 4. Cycle time: среднее время цикла PhaseType("Ph.txt"): тип обслуживания заявок для всех 5. Ave delay: среднее время пребывания заявок в системе очередей. 6. Number of packet: количество пакетов данных, прошед-5. **.switching time = "0.001 0.002": интенсивность переших сквозь систему в течение всего времени моделирования ключения между очередями

- Файлы алгоритмов: эти файлы содержат информацию о типе опроса очередей и дисциплинах их обслуживания, генерируют поток входных данных и процесс обслуживания для каждой очереди, а также генерируют случайное время переключения сервера между очередями. Далее в этих файлах производится расчет требуемых характеристик системы.
- Выходные файлы: эти файлы содержат результаты работы программы: среднее время ожидания в очереди системы, среднее время цикла, среднее время между посещениями очереди сервером, а также другие характеристики производительности (см. табл. 1).

Механизм работы пакета программ. На рис. 3 показан механизм работы пакета программ. После определения входных данных и запуска программы для каждой очереди отдельно генерируется момент времени, в который поступит следующая заявка. Сервер по выбранному правилу опроса в момент завершения обслуживания очередной очереди принимает решение, какая следующая очередь подлежит обслуживанию, и начинает к ней переключаться. Заявка, получившая обслуживание, отправляется в Sink, далее по моменту завершения обслуживания происходит перерасчет текущего среднего времени пребывания в системе и проверяется условие остановки моделирования.



Puc. 3. Механизм работы пакета программ Fig. 3. The mechanism of the software package

Задача остановки моделирования. Проблема остановки программного моделирования, когда поведение системы достигает режима, близкого к стационарному, является важной задачей для экономии ресурсов процессора и памяти. Рассмотрим следующий метод, который описан в работе [11].

Пусть $W_i, i \ge 1$ — последовательность значений среднего времени пребывания заявок после каждого обслуживания,

$$\overline{W}(n) = \sum_{i=1}^{n} \frac{W_i}{n}, \ \varepsilon_i = \frac{\Delta w}{\overline{W}(n)}, \Delta W_i - W_{i-1}, i \ge 1,$$

где ε_i - относительная точность доверительного интервала. Пусть вероятность

$$P(|\overline{W}(n)-W|<\Delta w_i)=1-\alpha,$$

где W — неизвестное среднее время пребывания заявок в системе. Критерием остановки моделирования является выполнение неравенства $\varepsilon_i < \varepsilon_{\max}$, где ε_{\max} - требуемая предельная относительная точность результатов при уровне достоверности $100(1-\alpha)\%$ $(0<\varepsilon_{\max}<1)$.

3. Некоторые результаты моделирования

При проведении имитационного моделирования часто требуется больше времени (особенно в случае большой загрузки системы), чем для того, чтобы рассчитать характеристики этих же систем

по точным или приближенным формулам (для тех моделей, для которых получены такие формулы). Это также касается случаев, когда входной поток не является простейшим (МАР- или РН-потоком в нашем случае), поскольку такое усложнение системы ухудшает сходимость результатов моделирования. Однако аналитические результаты получены лишь для относительно небольшого числа систем поллинга, и зачастую доступны только приближенные формулы для расчета характеристик систем поллинга, представляющих интерес для практического применения. В данном разделе для некоторых моделей представим результаты обоих способов расчета их характеристик, а также продемонстрируем результаты имитационного моделирования для систем поллинга, не имеющих на данный момент удобных для численной реализации способов расчета их характеристик (например, для систем поллинга с коррелированными входными потоками).

Рассмотрим систему поллинга с адаптивным циклическим опросом, тремя очередями и исчерпывающим обслуживанием, время обслуживания для всех очередей имеет одинаковое экспоненциальное распределение со средним b=0.01, время переключения между очередями имеет одинаковое экспоненциальное распределение со средним c=0.001, время простоя сервера (в случае, когда все очереди должны быть пропущены согласно адаптивной дисциплине) имеет экспоненциальное распределение со средним $\beta=0.002$, потоки заявок в очереди – простейшие с параметрами $\lambda_1=30$, $\lambda_2=20$, а параметр входного потока в третью очередь изменяется от $\lambda_3=2$ до $\lambda_3=16$ (табл. 2). Очереди обслуживаются согласно шлюзовой дисциплине. Аналитические формулы для расчета характеристик такой системы получены в работе [12].

 $\label{eq:Tadiff} T\, a\, б\, \pi\, u\, \mu\, a \quad 2$ Сравнение аналитических результатов и результатов имитационного моделирования

λ ₃	Аналитические результаты			Имитационное моделирование			Относительная погрешность, %		
	Q1	Q2	Q3	Q1	Q2	Q3	Q1	Q2	Q3
2	0,0239	0,0263	0,0318	0,0236	0,0260	0,0305	13	1,1	4,1
4	0,0247	0,0273	0,0321	0,0242	0,0268	0,0314	2	1,8	2,2
6	0,0256	0,0284	0,0328	0,0254	0,0281	0,0323	0,8	1,1	1,5
8	0,0267	0,0296	0,0337	0,0263	0,0291	0,0333	1,5	1,7	1,2
10	0,0279	0,0311	0,0347	0,0274	0,0305	0,0342	1,8	1,9	1,4
12	0,0292	0,0327	0,0358	0,0288	0,0321	0,0355	1,4	1,8	0,8
14	0,0308	0,0345	0,0371	0,0303	0,0338	0,0366	1,6	2	1,4
16	0,0325	0,0366	0,0386	0,0321	0,0358	0,0379	1,2	2,2	1,9

Рассмотрим далее симметричную систему циклического (не адаптивного) опроса с двумя очередями и шлюзовой дисциплиной. Параметры системы те же, что и в предыдущем случае, а параметры входных потоков идентичны, но изменяются от $\lambda_1 = \lambda_2 = 2$ до $\lambda_1 = \lambda_2 = 40$ (табл. 3). Аналитические формулы для расчета характеристик такой системы получены в [13].

Таблица 3 Сравнение аналитических результатов и результатов имитационного моделирования

$\lambda_2 = \lambda_1$	Имитационное	моделирование	Аналитически	ие результаты	Относительная погрешность, %		
	λ_1	λ_2	λ_1	λ_2	λ_1	λ_2	
2	0,01209	0,01188	0,01198	0,01198	0,918196995	0,834724541	
10	0,01434	0,01439	0,01438	0,01438	0,278164117	0,069541029	
30	0,02879	0,02877	0,02875	0,02875	0,139130435	0,069565217	
40	0,05746	0,05743	0,0575	0,0575	0,069565217	0,12173913	

Высокую практическую значимость имеют модели систем поллинга с коррелированными входными потоками (например, MAP- или BMAP-потоками). Однако известные на данный момент аналитические результаты их исследования [14, 15] ставят дополнительную задачу численной реализации этих результатов, и эта задача авторами работ не решена. Поэтому в пакете прикладных программ мы ограничились лишь имитационным моделированием таких систем. И для примера опишем влияние коэффициента корреляции на характеристики системы, в частности на среднее время пребы-

вания заявки в системе, полученное как сумма средних времен пребывания во всех очередях системы, взвешенных по доле загрузки соответствующих очередей.

Рассмотрим систему 5 очередей с адаптивным циклическим опросом и шлюзовой дисциплиной обслуживания очередей. Входящие потоки в очереди могут быть детерминированными (D), простейшими (M) и МАР-потоками с параметрами $\lambda_1 = \lambda_2 = \lambda_3 = 10$ и $\lambda_4 = \lambda_5 = 20$. Рассмотрим три варианта МАР-потоков в очереди (*МАРі* с коэффициентом корреляции c_i , $i = \overline{1,3}$):

$$\begin{split} \mathit{MAP}_1 \colon D_0 = & \begin{bmatrix} -60 & 20 \\ 0 & -10 \end{bmatrix}, D_1 = \begin{bmatrix} 0 & 40 \\ 0 & 10 \end{bmatrix}, c_1 = 0, \\ \mathit{MAP}_2 \colon D_0 = & \begin{bmatrix} -30 & 0 \\ 0 & -6 \end{bmatrix}, D_1 = & \begin{bmatrix} 20 & 10 \\ 2 & 4 \end{bmatrix}, c_1 = 0,07843, \\ \mathit{MAP}_3 \colon D_0 = & \begin{bmatrix} -18,75 & 0,625 \\ 0,625 & -2,5 \end{bmatrix}, D_1 = & \begin{bmatrix} 18,125 & 0 \\ 0 & 1,875 \end{bmatrix}, c_1 = 0,2704. \end{split}$$

В каждом из случаев предполагаем, что входные потоки заявок в очереди характеризуются одними и теми же матрицами MAP, но для 4-й и 5-й очередей эти матрицы соответствующим образом масштабируются (умножаются на скаляр), чтобы получить требуемые значения λ_4 и λ_5 .

Среднее время обслуживания заявки для всех очередей одинаково и равно 0,01. Рассмотрим различные типы обслуживания: детерминированное (D) время обслуживания заявки в очереди, фазового типа (PH) и марковское (MAP) обслуживание, задаваемое матрицами MAP_1 и MAP_3 , отмасштабированными соответствующим образом для сохранения одинаковых параметров обслуживания. Распределение фазового типа задано вектором α и неприводимой матрицей S следующим образом:

$$\alpha = \begin{bmatrix} 0.5 & 0.5 \end{bmatrix}, S = \begin{bmatrix} -10 & 0 \\ 10 & -20 \end{bmatrix}.$$

На рис. 4 представлена зависимость среднего времени пребывания заявки в системе для различных моделей системы поллинга с 5 очередями от типа выбранной модели и степени коррелированности входных потоков и процесса обслуживания при изменении интенсивности потока заявок в очереди с номерами 4 и 5 от 2 до 34.

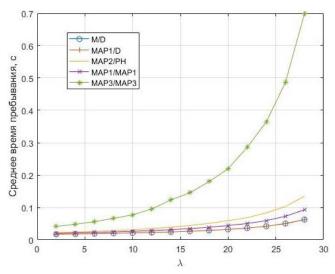


Рис. 4. Зависимость времени пребывания заявок в системе адаптивного поллинга от интенсивности поступления Fig. 4. Dependence of the residence time of applications in the system adaptive polling on the intensity of admission

Заметим, что при отсутствии корреляции среднее время пребывания заявок в системах M/D и MAP_1/D практически совпадает. При введении коррелированного процесса обслуживания среднее время пребывания заявок начинает возрастать и в случае MAP_3 с наибольшим коэффициентом корреляции существенно превышает результаты для остальных рассматриваемых моделей.

Заключение

В данной работе представлен пакет прикладных программ для оценки характеристик широкого класса систем поллинга с различными типами порядка опроса и дисциплинами обслуживания очередей. Пакет программ состоит из модуля расчета характеристик систем поллинга с помощью формул, реализованных с помощью Matlab, а также модуля имитационного моделирования на OMNeT++ для моделей систем поллинга, для которых не удается получить точные формулы расчета характеристик их производительности, в том числе слабоизученных в настоящее время систем поллинга с коррелированными входными потоками. Данный пакет программ зарегистрирован в Реестре программ для ЭВМ [16].

ЛИТЕРАТУРА

- 1. Вишневский В.М., Семенова О.В. Математические методы исследования систем поллинга // Автоматика и телемеханика. 2006. №2. С. 3–56.
- 2. Borst S.C., Boxma O. Polling: past, present, and perspective // TOP. 2018. V. 26, No. 3. P. 335–369.
- 3. Boon M.A.A., van der Mei R.D., Winands E.M.M. Applications of polling systems // Surveys in Operations Research and Management Science. 2011. V. 2011, No. 16. P. 67–82.
- 4. Вишневский В.М., Семёнова О.В. Системы поллинга: теория и применение в широкополосных беспроводных сетях. М.: Техносфера, 2007. 312 с.
- 5. Van der Mei R.D., Winands E. Heavy traffic analysis of polling models by mean value analysis // Performance Evaluation. 2008. V. 65, No. 6-7. P. 400–416.
- 6. Van Vuuren M., Winands E.M.M. Iterative approximation of k-limited polling systems // Queueing Systems. 2007. V. 55, No. 3. P. 161–178.
- 7. Van der Mei R.D., Roubos A. Polling models with multi-phase gated service // Annals of Operations Research. 2012. V. 198, No. 1. P. 25–56.
- 8. Bekker R., Vis P., Dorsman J.L., Van der Mei R.D., Winands E.M.M. The impact of scheduling policies on the waiting-time distributions in polling systems // Queueing Systems. 2015. V. 79, No. 2. P. 145–172.
- 9. Сонькин М.А., Моисеев А.Н., Сонькин Д.М., Буртовая Д.А. Объектная модель приложения для имитационного моделирования циклических систем массового обслуживания // Вестник Томского государственного университета. Управление, вычислительная техника и информатика. 2017. № 40. С. 71–80.
- 10. Вишневский В.М., Дудин А.Н., Клименок В.И. Стохастические системы с коррелированными потоками. Теория и применение в телекоммуникационных сетях. М.: Техносфера, 2018. 564 с.
- 11. Pawlikowski K. Steady state simulation of queueing processes: a survey of problems and solutions // ACM Computing Surveys. 1990. V. 22, No. 2. P. 123–170.
- 12. Semenova O.V., Bui D.T. Method of generating functions for performance characteristic analysis of the polling systems with adaptive polling and gated service // Communications in Computer and Information Science. 2018. V. 912. P. 348–359.
- 13. Yechiali U. Analysis and control of polling systems // Performance Evaluations of Computer and Communication Systems / ed. L. Donatielo, R. Nelson. Springer-Verlag, 1993. P. 630–650.
- 14. Saffer Z., Telek M. Unified analysis of BMAP/G/1 cyclic polling models // Queueing Systems. 2010. V. 64, No. 1. P. 69–102.
- 15. Saffer Z. BMAP/G/1 cyclic polling model with binomial disciplines // Modern Probabilistic Methods for Analysis of Telecommunication Networks. Communications in Computer and Information Science. 2013. V. 356. P. 157–166.
- 16. Вишневский В.М., Семёнова О.В., Буй З.Т. Программный комплекс оценки характеристик систем стохастического поллинга: свидетельство о государственной регистрации программы для ЭВМ № 2019614554 РФ; зарег. 08.04.2019.

Поступила в редакцию 28 апреля 2019 г.

Semenova O.V., Bui D.T. (2020) THE SOFTWARE PACKAGE AND ITS APPLICATION TO STUDY THE POLLING SYSTEMS. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 106–113

DOI: 10.17223/19988605/50/13

The paper describes the application package designed to calculate the characteristics of queuing systems with multiple queues and a common server, namely, polling systems with various types of the polling order: cyclic, adaptive cyclic, random; and various service disciplines: gated, exhaustive, globally-gated, limited, threshold and random. The input flow of customers may be a Poisson, a phase-type, or a correlated MAP. The basic structures and mechanism of the application package are described. Numerical examples illustrate the results of the study of polling systems with correlated input.

Polling systems are the queuing systems with multiple queues (or multiple input flows). The server visits the queues and serves the queued customers accordingly to a certain rule. To study the mathematical models of polling systems, various methods are used:

the method of generating functions, the mean-value method, descendant set approach, etc. to obtain exact formulas for calculating the performance characteristics, as well as various heuristic methods for special types of polling systems and simulation modeling in case when an exact analysis of the system is not possible to provide, or when it is necessary to estimate the accuracy of the approximate investigation results.

The developed software package consists of two large modules: a simulation module based on OMNeT++ Discrete Event Simulator, as well as OMNeT++ Simulation Manual Version 4.6, and a module for analytical calculations of the performance characteristics of the polling systems based on Matlab for some polling systems allowing exact analysis, in particular, polling systems with cyclic and adaptive cyclic polling and gated, exhaustive and globally-gated service of customers.

The following models of the polling systems are implemented in the simulation package: a cyclic polling system with gated, exhaustive, globally-gated, threshold and limited service disciplines, adaptive and ordered adaptive cyclic polling, priority polling system, polling system with random polling and other models.

Keywords: polling systems; software package; simulation.

SEMENOVA Olga Valeryevna (Candidate of Physical and Mathematical Sciences, Senior Researcher, Institute of Control Sciences of RAS, Moscow, Russian Federation).

E-mail: olgasmnv@gmail.com.

BUI Duy Tan (Post-graduate Student, Moscow Institute of Physics and Technology, Dolgoprudny, Moscow, Russian Federation). E-mail: duytan@phystech.edu.

REFERENCES

- 1. Vishnevskii, V.M. & Semenova, O.V. (2006) Mathematical methods to study the polling systems. *Automation and Remote Control*. 67. pp. 173–220. DOI: 10.1134/S0005117906020019
- 2. Borst, S.C. & Boxma, O. (2018) Polling: past, present, and perspective. TOP. 26(3). pp. 335-369. DOI: 10.1007/s11750-018-0490-7
- 3. Boon, M.A.A., van der Mei, R.D. & Winands, E.M.M. (2011) Applications of polling systems. Surveys in Operations Research and Management Science. 16(2). pp. 67–82. DOI: 10.1016/j.sorms.2011.01.001
- Vishnevsky, V. & Semenova O. (2012) Polling Systems: Theory and Applications for Broadband Wireless Networks. LAMBERT Academic Publishing.
- 5. Van der Mei, R.D. & Winands, E. (2008) Heavy traffic analysis of polling models by mean value analysis. *Performance Evaluation*. 65(6-7). pp. 400–416. DOI: 10.1016/j.peva.2007.12.002
- Van Vuuren, M. & Winands, E.M.M. (2007) Iterative approximation of k-limited polling systems. Queueing Systems. 55(3). pp. 161–178. DOI: 10.1007/s11134-007-9010-4
- 7. Van der Mei, R.D. & Roubos, A. (2012) Polling models with multi-phase gated service. *Annals of Operations Research*. 198(1). pp. 25–56. DOI: 10.1007/s10479-011-0921-4
- 8. Bekker, R., Vis, P., Dorsman, J.L., Van der Mei, R.D. & Winands, E.M.M. (2015) The impact of scheduling policies on the waiting-time distributions in polling systems. *Queueing Systems*. 79(2). pp. 145–172. DOI: 10.1007/s11134-014-9416-8
- 9. Sonkin, M.A., Moiseev, A.N., Sonkin, D.M. & Burtovaya, D.A. (2017). Object model of application for simulation of cyclic queueing systems. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie, vychislitel'naya tekhnika i informatika Tomsk State University Journal of Control and Computer Science.* 40. pp. 71–80. DOI: 10.17223/19988605/40/8
- 10. Vishnevsky, V.M., Dudin, A.N. & Klimenok, V.I. (2018) Stochastic Systems with Correlated Arrivals. Theory and Applications in Telecommunication Networks. Moscow: Tekhnosfera.
- 11. Pawlikowski, K. (1990) Steady state simulation of queueing processes: a survey of problems and solutions. *ACM Computing Surveys*. 22(2). pp. 123–170. DOI: 10.1145/78919.78921
- 12. Semenova, O.V. & Bui, D.T. (2018) Method of generating functions for performance characteristic analysis of the polling systems with adaptive polling and gated service. *Communications in Computer and Information Science*. 912. pp. 348–359. DOI: 10.1007/978-3-319-97595-5_27
- 13. Yechiali, U. (1993) Analysis and control of polling systems. In: Donatielo, L. & Nelson, R. (eds) *Performance Evaluations of Computer and Communication Systems*. Springer-Verlag. pp. 630–650.
- 14. Saffer, Z. & Telek, M. (2010) Unified analysis of BMAP/G/1 cyclic polling models. *Queueing Systems*. 64(1). pp. 69–102. DOI: 10.1007/s11134-009-9136-7
- Saffer, Z. (2013) BMAP/G/1 cyclic polling model with binomial disciplines. Modern Probabilistic Methods for Analysis of Telecommunication Networks. Communications in Computer and Information Science. 356. pp. 157–166. DOI: 10.1007/978-3-642-35980-4_18
- 16. Vishnevsky, V.M., Semenova, O.V. & Bui, D.T. (2019) Programmnyy kompleks otsenki kharakteristik sistem stokhasticheskogo pollinga: svidetel'stvo o gosudarstvennoy registratsii programmy dlya EVM № 2019614554 RF [Software system for evaluating the characteristics of stochastic polling systems: Certificate of the State Registration of Computer Programs No. 2019614554 of Russian Federation]. Registered 08.04.2019.

УДК 004.713, 621.318.57 DOI: 10.17223/19988605/50/14

А.И. Солдатов, А.Ю. Матросова, О.Х. Ким, А.А. Солдатов, М.А. Костина

ПРОГРАММИРУЕМАЯ КОММУТАЦИОННАЯ СРЕДА

Работа выполнена при финансовой поддержке фонда РФФИ, проект № 18-5700002.

Рассмотрены основные тенденции современного развития электронной компонентной базы. Показаны основные проблемы, с которыми сталкиваются производители микросхем. Одним из путей решения указанных проблем может служить разработка микросхем 3D-интеграции с раздельным расположением чипов и межчиповых соединений в пространстве. Предлагаемый подход позволит создать унифицированные базовые технологии и конструкции микросхем 3D-интеграции, обеспечивающие упрощение технологических и алгоритмических проблем.

Ключевые слова: электронная компонентная база; матричный коммутатор; микросхема 3D-интеграции; технология 3D-TSV; технология SiP.

Современный уровень технологии производства интегральных схем ведущих иностранных фирм составляет (28–22) нм с планируемым переходом на (16–10) нм, в то время как современный уровень технологии ведущих отечественных фирм составляет всего 90 нм с планируемым переходом на 45 нм [1–4]. Технологии уровня (22–16–10) нм относят к сфере национальной безопасности. Тренд на повышение уровня технологии сохранился, несмотря на то что еще более 10 лет назад многие иностранные аналитики и эксперты отрасли считали, что 22 нм технология будет последней планарной кремниевой технологией, поскольку «...на таком уровне разрешения разработка и внедрение техпроцессов требуют таких затрат, что вряд ли оправданы» [5].

Другим основным трендом современной микроэлектроники является создание технологий микросхем 3D-интеграции. Для фирм-производителей микросхем быть конкурентоспособным на рынке становится невозможным без овладения технологией микросхем 3D-интеграции [6–13]. Из существующих более 50 программ по созданию микросхем 3D-интеграции абсолютным большинством ведущих иностранных и российских специалистов и экспертов технология 3D-TSV (рис. 1, *a*) признана ключевой [14–18].

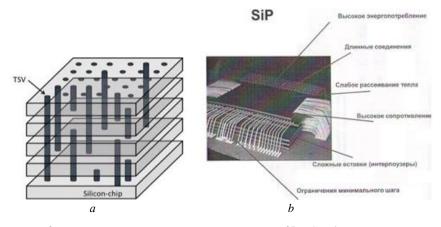


Рис. 1. Микросхемы 3D-интеграции: a — конструкция микросхемы 3D-TSV; b — конструкция микросхемы SiP Fig. 1. 3D integration chips: a — 3D-TSV chip design; b — SiP chip design

На Международной конференции (IEEE International Solid-State Circuits Conference – ISSCC), прошедшей в феврале 2009 г. в Сан-Франциско, компании впервые представили разработки в области 3D-интеграции, выполненные по технологии TSV. Многие европейские компании, специализирующиеся на полупроводниках и являющиеся в настоящий момент мировыми лидерами в производстве 3D-интегральных схем, также считают, что технология TSV является ключевой технологией 3D-интеграции микросхем.

Микросхемы по технологии 3D-TSV получают путем накладывания чипов друг на друга («стекирование») и соединения их между собой проводниками через переходные отверстия (through-silicon-vias – TSV), выполненными в виде медных проводников с сечением 50 мкм. Дальнейшее уменьшение сечения соединений TSV специалистами принято нецелесообразным из-за хрупкости медных проводников [19].

Основные проблемы технологии 3D-TSV связаны с выполнением большого количества соединений TSV и необходимостью обеспечения стабильной работы электрической схемы в многослойной структуре. Увеличение количества стекируемых чипов ведет к увеличению соединений TSV. При 5–6 «стекируемых» чипах количество соединений TSV составляет более 10 000. Встраивание в структуру «стекируемых» чипов технологий уровня (28–22–16–10) нм такого большого количества соединений TSV с сечением 50 мкм приводит к увеличению как физических размеров «стекируемых» чипов, так и физических размеров микросхемы 3D-TSV, в частности к ограничению уровня интеграции.

Второй популярной технологией микросхем 3D-интеграции является технология SiP (рис. 1, b) [20, 21]. Основные недостатки технологии SiP связаны с увеличением длины внешних выводов микросхемы по мере увеличения количества интегрируемых чипов.

Кроме перечисленных недостатков в этих двух технологиях существует еще одна важная проблема большого количества внешних выводов у микросхем. Например, в базовой конструкции технологии SiP все выводы интегрируемых чипов выводятся на внешние выводы микросхемы. В базовой конструкции технологии 3D-TSV часть межчиповых соединений выполняется посредством выполнесоединений TSV, остальные межчиповые соединения выводятся на внешние выводы микросхемы, количество которых также остается высоким.

Развитие микроэлектроники по перечисленным основным трендам ведет к общему резкому усложнению технологических и алгоритмических проблем и в ближайшей перспективе – к системному тупику.

Ускоренному приближению к технологическому тупику способствуют резкий рост энергопотребления по мере увеличения уровня интеграции, снижение отказоустойчивости и снижение процента выхода годных схем, так как «...интенсивность тепла, выделяемого густой "чащей" транзисторов в ходе работы, может достигнуть уровня, когда сами элементы "сварятся"...» [4]. Кроме того, по мере повышения уровня технологии и частоты тактовых импульсов на физические принципы работы электронных схем все большее влияние начинают оказывать «волновые свойства» логических элементов и электрических проводников, располагаемых совместно на «жесткой» плоскостной или пространственной конструктивной среде [22].

«Волновые свойства» элементов схемы и их взаимное влияние способствуют лавинообразному нарастанию новых видов неисправностей, в том числе кратковременных, с нарастанием вероятности их появления. При этом поведение электрической схемы становится теоретически и практически непредсказуемым. Возникающие новые и сложные виды неисправностей приводят к необходимости разработки новых, более сложных моделей неисправностей, а изменения в характере проявления неисправностей и все более усложняющиеся проблемы их локализации приводят к необходимости разработки новых, все более сложных математических, программных и аппаратных методов и средств решения задач проектирования и диагностики. В совокупности все это приводит к увеличению размерности и сложности математических моделей электрических схем и неисправностей и значительному усложнению задач проектирования и диагностики, при этом уже наработанные математические, программные и аппаратные методы и средства становятся мало приемлемыми для практического применения.

1. Постановка задачи

По мнению авторов статьи, основным источником проблем в существующих современных технологиях микроэлектроники является принцип совместного расположения логических элементов и электрических проводников на «жесткой конструктивной среде» 2D или 3D-интеграции.

Существующий принцип создания технологий микроэлектроники приводит к тому, что каждое изделие как объект проектирования и производства сохраняет индивидуальность и целостность. По мере повышения уровня технологии, уровня интеграции, а также частоты тактовых импульсов размерность и сложность каждого изделия нелинейно нарастают, и, соответственно, усложняются все проблемы и задачи проектирования, диагностики и производства.

Индивидуальность и целостность каждого изделия вызывают сложные проблемы в унификации и стандартизации выпускаемой продукции. В настоящее время номенклатура выпускаемых микросхем велика и, по некоторым источникам, составляет 500 и более типов. При этом, учитывая, что каждое новое изделие ЭКБ надо разрабатывать практически заново, выполняя полный цикл проектирования и подготовки производства, все более актуальной становится скорость создания нового продукта, требующая мощных инструментальных средств быстрого бездефектного проектирования ЭКБ, которые в настоящее время в основном базируются на математическом моделировании электрических схем.

В целом постановка задачи сводится к необходимости разработки базовой конструкции микросхемы 3D-интеграции нового типа, в которой был бы нивелирован «шлейф» технологических и алгоритмических проблем, следующий за существующими технологиями микросхем 3D-интеграции, в частности технологии 3D-TSV. В такой постановке задачи сформулированы следующие основные требования к базовой конструкции микросхемы 3D-интеграции нового типа:

- 1. В новой базовой конструкции должно быть осуществлено раздельное размещение конструктивных элементов схемы (интегрируемых чипов и межчиповых соединений), в частности, для нивелирования влияния «волновых» свойств элементов схемы.
- 2. Новая базовая конструкция должна обладать возможностью программирования конструкции, что позволит устранить индивидуальность и целостность каждого изделия, в частности, для получения унифицированной базовой технологии и конструкции.

Ближайшими аналогами, частично отвечающими поставленным требованиям, являются базовая конструкция технологии 2,5D-TSV, которая многими экспертами рассматривается как локомотив технологии 3D-TSV [23], и устройство FPIC (Field Programmable Inter Connection), представленное на рынке фирмой Арtix в 1997 г. [24].

По этой технологии несколько кристаллов объединяются с помощью промежуточной монтажной пластины (Interposer). Такая монтажная пластина, кремниевая или стеклянная с вертикальными сквозными отверстиями и медными контактными столбиками, позволяет присоединять микросхемы с малым шагом на печатные платы с большим шагом проводников. По плотности монтажа кремниевая промежуточная монтажная пластина в 20 раз превосходит платы, предназначенные для монтажа электронных компонентов.

Устройство FPIC представляет собой программируемый кристалл, выполненный в 1 024-контактном корпусе, обеспечивающий до 1 000 внешних соединительных выводов, к которым могут быть подключены цифровые электронные компоненты: микропроцессоры, ПЛИС (программируемые логические интегральные схемы) и т.п. Изменение конфигурации синтезируемых электронных схем осуществляют программированием электрических соединений между выводами электронных компонент. Принцип выполнения электрических соединений в устройстве FPIC в целом совпадает с принципом выполнения электрических соединений в ПЛИС.

Устройство FPIC, так же как ПЛИС, имеет регулярную структуру, представленную матрицей из программируемых ячеек CLB. В ПЛИС программируемая коммутационная среда (ПКС) входит в состав программируемых ячеек CLB. Основное отличие устройства FPIC состоит в том, что в програм-

мируемой ячейке CLB отсутствуют логические элементы, а электрические соединения выполняются на внешних соединительных выводах. В функции макроячейки входит только осуществление транзита электрических сигналов по восьми направлениям в соседние ячейки. Устройство FPIC фактически представляет собой ПКС, выведенное из структуры ПЛИС в отдельное функциональное устройство, предназначенное для выполнения программируемых электрических соединений. Попытка реализовать электрические соединения на внешних соединительных выводах приводит к значительному увеличению количества последовательно соединенных ключей в реализуемых электрических соединениях и соответствующему увеличению времени задержки электрических сигналов.

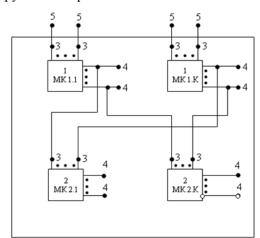
В настоящее время устройство FPIC является единственной практической реализацией универсальной ПКС. Области применения устройства FPIC ограничены его техническими характеристиками. В частности, большие физические размеры кристалла не позволяют использовать его для создания базовой конструкции микросхемы 3D-интеграции с программируемой конструкцией нового типа.

Основные технические и алгоритмические проблемы коммутации современной электроники изложены в [25].

2. Принципиальная схема коммутатора

Электрическая схема разработанной ПКС [26] представлена на рис. 2.

ПКС, имеющая N соединительных выводов, содержит первую группу матричных коммутаторов МК1.1 ... МК1.К и вторую группу матричных коммутаторов МК2.1 ... МК2.К, каждый из которых имеет размерность $n \times n$, соединенных между собой соответствующим образом. Выводы 5 первой группы матричных коммутаторов МК1.1 ... МК1.К образуют группу N внешних соединительных выводов ПКС, к которым можно подключать разные типы чипов: микропроцессоры, ОЗУ, ПЗУ, ПЛИС и другие электронные компоненты.



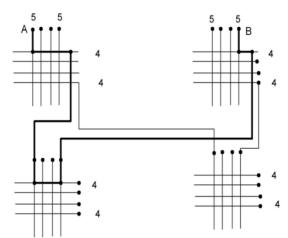


Рис. 2. Электрическая схема ПКС Fig. 2. The block diagram of the PSA

Рис. 3. Принцип выполнения электрических соединений в ПКС Fig. 3. Principle of electrical connections in PCSA

Принцип выполнения электрических соединений на группе внешних соединительных выводов 5 ПКС представлен на рис. 3.

Контакты A и B, объединяемые в одну электрическую цепь путем программирования соответствующих матричных коммутаторов первой группы МК1.1 ... МК1.К, коммутируются на одноименные горизонтальные шины 4, которые через образованное многомерное пространство попадают на соответствующие вертикальные шины одного из матричных коммутаторов второй группы матричных коммутаторов МК2.1 ... МК2.К. Путем программирования этого матричного коммутатора вертикальные шины коммутируются на одну из свободных горизонтальных шин. Таким образом, осуществляется соединение в электрическую цепь контактов A и B. При этом все выполняемые электрические соединения содержат одинаковое количество ключей, равное 4.

3. Оценка сложности ПКС

Суммарное количество ключей в ПКС (рис. 3), содержащем N соединительных выводов, равно $C_1 = (n \cdot n) \cdot k_1 + (n \cdot n) \cdot k_2$, (1)

где: n — количество вертикальных и горизонтальных линий связи одного матричного коммутатора, $(n \cdot n)$ — количество ключей в одном матричном коммутаторе, k_1 — количество матричных коммутаторов первой группы матричных коммутаторов, k_2 — количество матричных коммутаторов второй группы матричных коммутаторов.

С учетом того, что $k_1 = k_2 = n$, выражение (1) примет вид:

$$C_1 = 2 \cdot n^3 = 2 \cdot \sqrt{N^3} \ . \tag{2}$$

Суммарное количество ключей в обычном матричном коммутаторе, содержащем N соединительных выводов, равно

$$C_2 = N^2. (3)$$

Общее сокращение количества ключей в предлагаемой ПКС в сравнении с обычным матричным коммутатором, составит

$$C = \frac{C_2}{C_1} = \frac{N^2}{2 \cdot (\sqrt{N})^3} = \frac{\sqrt{N}}{2}.$$
 (4)

Результаты расчета сокращения количества ключей в предлагаемой ПКС в сравнении с обычным матричным коммутатором, рассчитанным по формуле (4), представлены на рис. 4.

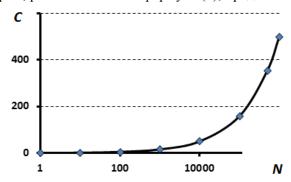


Рис. 4. Сравнение количества ключей в предлагаемой ПКС и матричном коммутаторе Fig. 4. Comparison of the number of switches in PSA and the matrix switch

Как видно из рис. 4, использование предлагаемой ПКС дает возможность существенно уменьшить количество ключей в коммутаторе без снижения его функциональных возможностей. Это позволяет уменьшить размеры кристалла, увеличить надежность и снизить потребляемую мощность. Следовательно, снижаются требования по отводу тепла. Следует заметить, что с увеличением размерности матричного коммутатора увеличивается выигрыш в количестве ключей.

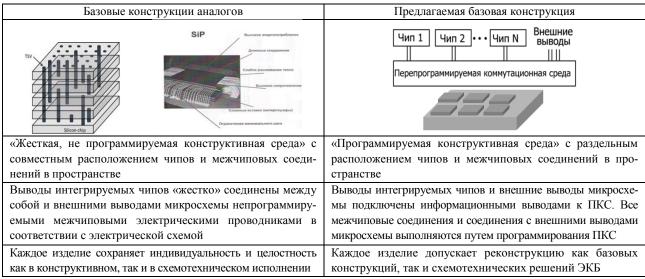
Основные отличия предлагаемой базовой конструкции микросхемы 3D-интеграции нового типа от аналогов представлены в таблице.

Базовая конструкция предлагаемой микросхемы 3D-интеграции нового типа представляет собой универсальную аппаратно-программную платформу, на которую устанавливаются чипы, интегрируемые в микросхему 3D-интеграции. Упрощение технологических проблем обеспечивается следующими преимуществами:

1. Создание унифицированных базовых конструкций и технологий микросхем 3D-интеграции с низкими требованиями к уровню технологии.

На заданном наборе чипов путем программирования ПКС можно получить микросхему 3D-интеграции с набором функциональных назначений, получаемую путем перепрограммирования межчиповых соединений. Чтобы получить микросхему 3D-интеграции с другим набором функциональных назначений, достаточно сменить состав интегрируемых чипов.

Базовые структуры микросхем



2. Многократное увеличение уровня интеграции микросхемы 3D-интеграции при многократном сокращении количества внешних выводов.

В предлагаемой микросхеме 3D-интеграции выводы всех интегрируемых чипов и внешние выводы микросхемы подключены к ПКС, при этом все межчиповые соединения и соединения выводов чипов с внешними выводами микросхемы в соответствии с реализуемой электрической схемой выполняются путем программирования ПКС. На внешние выводы микросхемы выводится только необходимое количество выводов. Общее количество чипов, интегрируемых в микросхему 3D-интеграции, в зависимости от размерности ПКС, рассчитывается по следующей формуле

$$m = \frac{N}{\sum_{i=1}^{n} n_i} - P , \qquad (5)$$

где m — общее количество интегрированных чипов; N — общее количество внешних соединительных выводов ПКС; P — общее количество внешних выводов микросхемы 3D-интеграции в соответствии с реализуемой электрической схемой; $\sum_{i=1}^n n_i$ — суммарное количество выводов чипов, интегрируемых в микросхему 3D-интеграции.

Из формулы следует, что получение требуемого уровня интеграции может быть достигнуто изменением размерности ПКС, например, в диапазоне 256–1 024 и более.

- 3. Многократное сокращение номенклатуры ЭКБ.
- Прогнозируемое сокращение номенклатуры микросхем 3D-интеграции составит 10 раз и более.
- 4. Высокая скорость и низкая стоимость создания нового продукта.

По существу, любой новый продукт можно рассматривать как некоторую модификацию уже существующих продуктов, в которых достаточно только заменить состав интегрируемых чипов: микропроцессоров / микроконтроллеров, ОЗУ, ПЗУ, ПЛИС и других электронных компонентов и радиоэлементов.

Заключение

ПКС может быть как перепрограммируемой, так и однократно программируемой. В связи с этим можно выделить два основных направления создания микросхем 3D-интеграции нового типа:

- 1) разработка микросхем 3D-интеграции с перепрограммируемой архитектурой;
- 2) разработка микросхем 3D-интеграции с однократно программируемой архитектурой.

Предлагаемый подход позволит создать унифицированные базовые технологии и конструкции микросхем 3D-интеграции, обеспечивающие упрощение технологических и алгоритмических про-

блем. На заданном наборе чипов путем программирования ПКС можно получить ЭКБ различного функционального назначения. Многократное увеличение уровня интеграции микросхем возможно при многократном сокращении количества внешних выводов. В базовых конструкциях ЭКБ на основе предлагаемой ПКС выводы всех чипов и внешние выводы ЭКБ подключены к ПКС, при этом все соединения в соответствии с реализуемой электрической схемой выполняются путем программирования ПКС. На внешние выводы ЭКБ выводится только необходимое количество информационных выводов.

ЛИТЕРАТУРА

- 1. Куликова Н.Н. Современное состояние и тенденции развития электронной промышленности в России // Теория и практика общественного развития. 2017. № 12. С. 87–92.
- 2. Микроэлектронная промышленность России: состояние и перспективы развития. URL: https://pandoraopen.ru/2017-06-18/mikroelektronnaya-promyshlennost-rossii-sostoyanie-i-perspektivy-razvitiya/ (дата обращения: 20.10.2019).
- 3. Тенденции и перспективы глобального и российского рынка микроэлектроники. URL: https://www.crn.ru/news/detail.php?ID=119695. (дата обращения: 20.10.2019).
- 4. Перспективы развития микросхем. URL: http://iqrate.com/infotech/perspektivy-razvitiya-mikroshem/ (дата обращения: 20.10.2019).
- 5. Ежов В. Тенденции развития электронных технологий. Ближайшие перспективы. URL: http://www.russianelectronics.ru/leader-r/review/521/doc/40568/ (дата обращения: 20.10.2019).
- 6. Аракелян В.А. Проблемы и перспективы в трехмерном проектировании интегральных схем // SWorld. 2014. Т. 4, № 1. С. 71–78.
- 7. Строгонов А., Цыбин С., Быстрицкий А. Трехмерные интегральные схемы 3D БИС // Компоненты и технологии. 2011. № 1. С. 118–121.
- 8. Патент РФ № 2461911. Многокристальный модуль. Патентообладатель: Минпромторг РФ.
- 9. Патент РФ № 2463684(73). Многокристальный модуль. Патентообладатели: Минпромторг РФ, ЗАО «НПО "НИИТАЛ"».
- 10. Патент РФ № 2335821, Трехмерный электронный модуль. Патентообладатель: ОАО «НПК "ЭЛАРА" им. Г.А. Ильенко» (ОАО «ЭЛАРА»).
- 11. Патент РФ № 2336595 Способ изготовления объемных мини-модулей Патентообладатель: Завадский Александр Ивано-
- 12. Гольцова М. Международная конференция ISSCC 2011. От микросхем больших объемов до имплантируемых устройств // Электроника: наука, технология, бизнес. 2011. № 3. С. 32–45.
- 13. Haran B.S., Kumara A., Adam L. et al. 22 nm Technology Compatible Fully Functional 0.1 μ m² 6T-SRAM Cell // 2008 IEEE International Electron Devices Meeting. 2008. P. 615–619.
- 14. 3D сборка микросхем: «Сделано в России» // SEMICON Russia 2014. URL: https://www.semiconeuropa.org/en/sites/semiconrussia.org/files/docs/3.%20SEMI_2014_GS%20Group_3D%20packaging.pdf (дата обращения: 20.10.2019).
- 15. Liu Z., Tian Q., Li J., Liu X., Zhu W. An Efficient and High Quality Chemical Mechanical Polishing Method for Copper Surface in 3D TSV Integration // IEEE Transactions on Semiconductor Manufacturing_ 2019. V. 32 (3). 8738875. P. 346–351.
- 16. Meng M., Cheng L., Yang K., Sun M., Luo Y. A novel seedless TSV process based on room temperature curing silver nanowires ECAs for MEMS packaging // Micromachines. 2019. V. 10 (6). P. 351.
- 17. Wang Z. Microsystems using three-dimensional integration and TSV technologies: Fundamentals and applications // Microelectronic Engineering. 2019. V. 210. P. 35–64.
- 18. Golishnikov A.A., Kostyukov D.A., Putrya M.G., Shevyakov V.I. Features of silicon deep plasma etching process at 3D-TSV structures producing // Proc. of SPIE The International Society for Optical Engineering/ 2019. 11022. 11022.12.
- 19. Юдинцев В. Трехмерная кремниевая технология. Что, где, когда // Электроника: наука, технология, бизнес. 2011. N 4 (110). С. 70–75.
- 20. Чипы-гибриды парадоксально уплотняют схемы без уплотнения. URL: http://www.membrana.ru/articles/technic/ 2009/09/21/173100.html. (дата обращения: 20.10.2019).
- 21. Miao M., Wang L., Chen T., Sun L., Duan X., Zhang J., Liu H., Jin Y. Modeling and Design of a 3D Interconnect Based Circuit Cell Formed with 3D SiP Techniques Mimicking Brain Neurons for Neuromorphic Computing Applications // Proc. Electronic Components and Technology Conf. 2018. 8429591. P. 490–497.
- 22. U.S. Patent No. 5,414,638 entitled Programmable Interconnect Architecture. Filed May. 2011. by Aptix Corporation.
- 23. U.S. Patent No. 6,272,646 entitled Programmable Logic Device Having an Integrated Phase Lock Loop. Filed February 18. 2011. by Xilinx.
- 24. U.S. Patent No. 6,188,578 entitled Intergrated Circuit Package with Multiple Heat Dissipation Paths and owned by HTC Corporation. Filed February 14. 2011.
- 25. Солдатов А.И., Ким О.Х. Технические и алгоритмические проблемы коммутации современной электроники // Известия высших учебных заведений. Физика. 2010. Т. 53, № 9-3. С. 308–310.

26. Патент РФ № 2 402 061. Пространственная коммутационная среда (варианты). Патентообладатель: ООО «Хонбин». 16.03.2009.

Поступила в редакцию 20 октября 2019 г.

Soldatov A.I., Matrosova A.Yu., Kim O.H., Soldatov A.A., Kostina M.A. (2020) PROGRAMMABLE SWITCHING AREA. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 114–122

DOI: 10.17223/19988605/50/14

Currently, there is a long-term movement to increase the level of integration of the electronic component base. Another major trend in modern microelectronics is the development of 3D integration chip technologies. 3D-TSV technology for the production of 3D integration chips is recognized as key by the absolute majority of leading foreign and Russian experts. The main problems in 3D-TSV technology are related to the implementation of a large number of internal connections and the need to ensure stable operation of the electrical circuit in a multilayer structure. Increasing the number of stackable chips leads to an increase in the number of TSV connections, which leads to increase in physical size "stackable" chips, and physical dimensions of a 3D-TSV chip.

The second most popular technology of 3D chip is SiP technology. The main disadvantages of SiP technology are associated with the increase in the length of the external pins of the chip, as the number of integrated chips increases. In addition to these disadvantages there is another important problem is a large number of pins of chips.

According to the authors, the main source of problems in the existing modern microelectronics technologies is the principle of joint arrangement of electronic components and electrical conductors on a "rigid structure".

The authors propose the concept of building of a multidimensional programmable switching area (PSA), which has a conflict-free, reduced total number of switches compared to a conventional matrix switch, all performed electrical connections contain the same number of switches equal to 4. Significantly reducing the number of switches in PSA without reducing its functionality allows you to reduce the size of the chip, increase reliability and reduce power consumption. As the dimension of the matrix switch increases, the benefit in the number of switches increases.

The total number of switches in the PSA containing N connecting pins is equal to

$$C_1 = 2 \cdot n^3 = 2 \cdot \sqrt{N^3}.$$

The total number of switches in a cross-bar containing N connection pins is equal to

$$C_2 = N^2$$
.

The reduction in the number of switches in the PSA compared to a cross-bar, will be

$$C = \frac{C_2}{C_1} = \frac{N^2}{2 \cdot (\sqrt{N})^3} = \frac{\sqrt{N}}{2}.$$

The use of the proposed PCC allows you to create 3D-chips of high integration, while different types of chips can be connected to the pins of the PSA: microprocessors, RAM, ROM, FPGA and other electronic components, and external the pins will only output the necessary signals that will reduce the number of pins of the chip. All connections in accordance with the implemented electrical circuit are performed by programming the PSA. The proposed approach will allow to create unified basic technologies and designs of 3D integration chips, providing simplification of technological and algorithmic problems.

Keyword: electronic component base; cross-bar switch; 3D-integration chip; 3D-TSV technology; SiP technology.

SOLDATOV Alexey Ivanovich (Doctor of Technical Sciences, Professor, Tomsk state University of control systems and Radioelectronics, Professor, National Research Tomsk Polytechnic University, Tomsk, Russian Federation). E-mail: asoldatof@mail.ru

MATROSOVA Anzhela Yurevna (Doctor of Technical Sciences, Professor, National Research Tomsk State University, Tomsk, Russian Federation).

E-mail: mau11@yandex.ru

KIM Oleg Honbinovich (National Research Tomsk State University, Tomsk, Russian Federation, Tomsk, Russian Federation). E-mail: oh.kim@mail.ru

SOLDATOV Andrey Alexeevich (Candidate of Technical Sciences, Associate Professor, Tomsk state University of control systems and Radioelectronics, Associate Professor, National Research Tomsk Polytechnic University, Tomsk, Russian Federation). E-mail: Soldatov.88@bk.ru

KOSTINA Mariya Alexeevna (Candidate of Technical Sciences, Associate Professor, Tomsk state University of control systems and Radioelectronics, Associate Professor, National Research Tomsk Polytechnic University, Tomsk, Russian Federation). E-mail: mashenkasoldatova@mail.ru

REFERENCES

- 1. Kulikova, N.N. (2017) Modern state and development trends of electronic industry in Russia. *Teoriya i praktika obshchestvennogo razvitiya Theory and Practice of Social Development*. 12. pp. 87–92. DOI: 10.24158/tipor.2017.12.19
- Feedmatic.bot. (2017) Mikroelektronnaya promyshlennost' Rossii: sostoyanie i perspektivy razvitiya [Microelectronic industry of Russia: state and prospects of development]. [Online] Available from: https://pandoraopen.ru/2017-06-18/mikroelektronnaya-promyshlennost-rossii-sostoyanie-i-perspektivy-razvitiya (Accessed: 20th October 2019).
- 3. CRN. (n.d.) *Tendentsii i perspektivy global'nogo i rossiyskogo rynka mikroelektroniki* [Trends and prospects of the global and Russian microelectronics market]. [Online] Available from: https://www.crn.ru/news/detail.php?ID=119695. (Accessed: 20th October 2019).
- 4. IQRATE.com. (n.d.) *Perspektivy razvitiya mikroskhem* [Prospects of development of chips]. [Online] Available from: http://iqrate.com/infotech/perspektivy-razvitiya-mikroshem (Accessed: 20th October 2019).
- 5. Ezhov, V. (2019) *Tendentsii razvitiya elektronnykh tekhnologiy. Blizhayshie perspektivy* [Tendencies of development of electronic technologies. Immediate prospect]. [Online] Available from: http://www.russianelectronics.ru/leader-r/review/521/doc/40568. (Accessed: 20th October 2019).
- 6. Arakelyan, V.A. (2014) Problemy i perspektivy v trekhmernom proektirovanii integral'nykh skhem [Problems and prospects in three-dimensional design of integrated circuits]. *SWorld*. 4(1). pp.71–78.
- 7. Strogonov, A., Tsybin, S. & Bystritsky, A. (2011) Trekhmernye integral'nye skhemy 3D BIS [Three-Dimensional integrated circuits 3D LSI circuit]. *Komponenty i tekhnologii Components and Technologies*. 1. pp. 118–121.
- 8. The Russian Federation. (n.d.) *Mnogokristal'nyy modul'* [Multi-chip module]. RU Patent No. 2461911. Patent Holder: Ministry of Industry and Trade of Russia.
- The Russian Federation. (n.d.) Mnogokristal'nyy modul' [Multi-chip module]. RU Patent No. 2463684(73). Patent Holder: Ministry of Industry and Trade of Russia, ZAO NPO NIITAL.
- 10. The Russian Federation. (n.d.) *Trekhmernyy elektronnyy modul'* [Three-dimensional electronic module]. RU Patent No. 2335821. Patent Holder: OAO NPK ELARA im. G.A. Il'enko.
- 11. The Russian Federation. (n.d.) *Sposob izgotovleniya ob"emnykh mini-moduley* [Method of manufacturing three-dimensional mini-modules]. RU Patent No. 2336595. Patent Holder: Zavadsky Alexander Ivanovich.
- 12. Goltsova, M. (2011) From microchips to large amounts of implantable devices. *Elektronika: nauka, tekhnologiya, biznes Electronics: STB.* 3. pp. 32–45.
- 13. Haran, B.S., Kumara, A., Adam, L. et al. (2008) 22 nm Technology Compatible Fully Functional 0.1 μm² 6T-SRAM Cell. 2008 *IEEE International Electron Devices Meeting*. pp. 615–619. DOI: 10.1109/IEDM.2008.4796769
- 14. Semiconeuropa.org. (2014) 3D sborka mikroskhem: "Sdelano v Rossii" [3D chip: "Made in Russia"]. [Online] Available from: https://www.semiconeuropa.org/en/sites/ semiconrussia.org/files/docs/3.%20SEMI_2014_GS%20Group_3D%20packaging.pdf (Accessed: 20th October 2019)
- Liu, Z., Tian, Q., Li, J., Liu, X. & Zhu, W. (2019) An Efficient and High Quality Chemical Mechanical Polishing Method for Copper Surface in 3D TSV Integration. *IEEE Transactions on Semiconductor Manufacturing*. 32(3). pp. 346–351. DOI: 10.1109/TSM.2019.2923427
- 16. Meng, M., Cheng, L., Yang, K., Sun, M. & Luo, Y. (2019) A novel seedless TSV process based on room temperature curing silver nanowires ECAs for MEMS packaging. *Micromachines*. 10(6). pp. 351. DOI: 10.3390/mi10060351
- 17. Wang, Z. (2019) Microsystems using three-dimensional integration and TSV technologies: Fundamentals and applications. *Microelectronic Engineering*. 210. pp. 35–64. DOI: 10.1016/j.mee.2019.03.009
- 18. Golishnikov, A.A., Kostyukov, D.A., Putrya, M.G. & Shevyakov, V.I. (2019) Features of silicon deep plasma etching process at 3D-TSV structures producing. *Proc. of SPIE. The International Society for Optical Engineering*. 11022. 110221Z. DOI: 10.1117/12.2521969
- 19. Yudintsev, V. (2011) Trekhmernaya kremnievaya tekhnologiya. Chto, gde, kogda [Three-Dimensional silicon technology. What, where, when]. *Elektronika: nauka, tekhnologiya, biznes Electronics: STB.* 4(110). pp. 70–75.
- 20. Membrana.ru. (n.d.) *Chipy-gibridy paradoksal'no uplotnyayut skhemy bez uplotneniya* [Hybrid chips paradoxically seal circuits without sealing]. [Online] Available from: http://www.membrana.ru/articles/technic/2009/09/21/173100.html (Accessed: 20th October 2019).
- 21. Miao, M., Wang, L., Chen, T., Sun, L., Duan, X., Zhang, J., Liu, H. & Jin, Y. (2018) Modeling and Design of a 3D Interconnect Based Circuit Cell Formed with 3D SiP Techniques Mimicking Brain Neurons for Neuromorphic Computing Applications. *Proc. Electronic Components and Technology Conf.* 8429591. pp. 490–497.
- 22. USA. (2011a) U.S. Patent No. 5,414,638 Entitled Programmable Interconnect Architecture. Filed May, 2011, by Aptix Corporation.
- 23. USA. (2011b) U.S. Patent No. 6,272,646 Entitled Programmable Logic Device Having an Integrated Phase Lock Loop. Filed February 18, 2011, by Xilinx.
- 24. USA. (2011c) U.S. Patent No. 6,188,578 Entitled Intergrated Circuit Package with Multiple Heat Dissipation Paths and owned by HTC Corporation. Filed February 14, 2011.
- 25. Soldatov, A.I. & Kim, O.H. (2010) Technical and algorithmic problems of switching of modern electronics. *Izvestiya vysshikh uchebnykh zavedeniy. Fizika*. 53(9/3). pp. 308–310.
- 26. The Russian Federation. (n.d.) *Prostranstvennaya kommutatsionnaya sreda (varianty)* [Spatial switching area (options). RU Patent No. 2 402 061. Registered on March 16, 2009. Patent Holder: Ltd. Honbin.

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

ОРЗОЬРІ

УДК 004.75

DOI: 10.17223/19988605/50/15

Н.В. Петухова, М.П. Фархадов, Д.Л. Качалов

РАЗГРУЗКА И КОНСОЛИДАЦИЯ ВЫЧИСЛИТЕЛЬНЫХ РЕСУРСОВ В СРЕДЕ ТУМАННЫХ И ГРАНИЧНЫХ ВЫЧИСЛЕНИЙ

Исследование выполнено при частичной финансовой поддержке РФФИ в рамках научного проекта № 19-29-06044 (разделы 2, 3, 4) и при поддержке Президиума РАН в рамках научного проекта по Программе 7 (1.30).

Анализируются задачи управления ресурсами, возникающие в среде туманных и граничных вычислений при организации обслуживания запросов, посылаемых устройствами интернета вещей. На основании обзора публикаций выявлены основные темы исследований в этой области и представлены примеры решения задач, относящихся к распределению заявок на обслуживание, к методам стимулирования поставщиков ресурсов к предоставленыю их в общий пул, к проблемам энергосбережения.

Ключевые слова: интернет вещей; туманные вычисления; граничные вычисления; распределение задач; консолидация ресурсов; энергосбережение; стохастические модели; теория игр.

Туманные вычисления (fog computing) представляют собой технологию, которая позволяет использовать для хранения и обработки информации память и вычислительные ресурсы, расположенные ближе к устройствам интернета вещей, где генерируются данные, за счет чего снижется нагрузка на сеть и облачные дата-центры и уменьшается время реакции, потребление энергии и эксплуатационные расходы. Кроме того, для выполнения сторонних вычислений могут быть привлечены свободные ресурсы самих устройств интернета вещей, как предполагает концепция граничных вычислений (edge computing). В результате среда для обработки данных оказывается состоящей из многих устройств, различных по техническим характеристикам и программным возможностям. При этом необходимо учитывать, что устройства, выполняющие задачи, зачастую являются мобильными и предоставляются для решения «чужих» задач без всякого постоянного расписания. Очевидно, что такие структуры значительно отличаются от ставших уже привычными облачных дата-центров и ставят перед разработчиками целый ряд новых задач.

Цель работы – выявить на основании обзора публикаций основные задачи управления ресурсами в среде туманных и граничных вычислений и подходы к их решению, учитывающие особенности таких сетей.

1. Цели и методы решения задач управления ресурсами

Обзоры работ, посвященных архитектуре, технологиям и программно-техническим средствам, используемым в среде туманных и граничных вычислений, регулярно появляются в открытом доступе и показывают, что интерес к этой тематике постоянно растет. Обсуждаются как общие вопросы: концепция, архитектура, технические решения, приложения [1–4], так и специальные, например безопасность [5, 6]. Определены ключевые особенности таких сетей: гетерогенность технических средств, их географическая распределенность, мобильность, высокие требования потребителей к времени реакции, ограниченность ресурсов, энергосбережение.

Чаще всего темами исследований в области управления распределенными ресурсами туманных и граничных сетей, как показывает обзор публикаций, являются задачи выбора ресурса для выполнения возникающих заявок на обслуживание, разгрузка ресурсов, способы привлечения ресурсов в общий пул, проблемы энергосбережения.

Целевые функции и методы решения зависят от конкретной системы. Чаще всего рассматриваются следующие показатели и их комбинации: время отклика, расход энергии, доступность сервиса, экономическая эффективность.

Для обоснования и оценки предлагаемых подходов и алгоритмов управления авторами используется различный математический аппарат: аналитические модели, основанные на теории очередей, линейное и нелинейное программирование, теория графов, теория игр, имитационное моделирование.

2. Распределение заявок на обслуживание и разгрузка ресурсов

Наибольшее количество исследований относится к решению задачи выбора места для переноса выполнения вычислений с устройства, где возникла такая необходимость, на другое устройство или сервер. Авторы работы [7] приводят обзор публикаций по данной тематике, и называют следующие причины, вызывающие перенос задач с одних устройств на другие:

- приложение требует больше вычислений, чем может обеспечить собственное устройство;
- сокращение сетевых задержек на передачу больших объемов данных;
- балансировка нагрузки;
- необходимость длительного хранения накапливаемой информации;
- выгрузка данных с целью обеспечения необходимой конфиденциальности, безопасности, надежности и доступности;
 - выгрузка задач с облачных серверов на вычислительные средства более низких уровней.

Классифицируются алгоритмы переноса задач и сформулированы задачи, требующие дальнейших исследований.

Авторы статьи [8] процесс распределения задач разделили на два этапа. На первом этапе запросы на обслуживание распределяются между fog кластерами, на втором этапе происходит выбор узла кластера, который будет выполнять поступившую заявку. При распределении запросов между кластерами для обслуживания заявки выбирается ближайший кластер, загрузка которого позволяет предоставить нужные ресурсы. Выбор узла внутри кластера осуществляется с учетом нестабильности присутствия узлов в туманной сети. Для этого авторы вводят показатель, названный ими репутацией узла, который количественно выражает степень надежности устройства и представляет собой статистическую вероятность того, что узел останется в кластере на время, требуемое данной задачей. Если статистическая оценка является малодостоверной из-за того, что число выполненных узлом задач недостаточно велико, то показатель репутации узла определяется типом устройства. Авторы разделили устройства сети на три категории: мобильные узлы (смартфоны), статические узлы (настольные компьютеры, серверы) и полумобильные узлы (планшеты, ноутбуки). Показатели репутации этих узлов, по оценкам авторов, находятся в соотношении 1:1,5:1,25. Оптимальная схема распределения ресурсов кластера формулируется авторами как многоцелевая задача, которая минимизирует время обслуживания заявки и максимизирует общую стабильность функционирования кластера.

Для решения поставленной задачи авторы используют улучшенный генетический алгоритм сортировки без доминирования (NSGA-II), предложив свою формулу для вычисления меры близости к соседу, учитывающую разнонаправленность целевых функций. Для оценки эффективности предложенного алгоритма авторы сравнили его с алгоритмом М. Аазат и Е. Huh [9] и с алгоритмом случайного выбора ресурса. Эксперименты показали, что при малых нагрузках все три схемы практически эквивалентны. Однако при увеличении количества задач среднее время обработки задач с использованием предложенного алгоритма было меньше, чем при двух других алгоритмах, а средняя стабильность успешного выполнения задач превосходила этот показатель для других алгоритмов. Авторы объясняют этот выигрыш вводом в рассмотрение показателя репутации устройств.

В работе [10] предложен алгоритм распределения задач с целью минимизации задержки обслуживания приложений интернета вещей и разработана аналитическая модель для оценки предложенного алгоритма. Рассмотрены централизованный и распределенный подходы к маршрутизации задач внутри домена. Выбор узла для передачи задачи осуществляется на основе оценки ожидаемого времени обслуживания запроса, которое включает время ожидания в очереди и время обслуживания. Для упрощения модели авторы рассматривают два класса запросов на обслуживание: «легкие», требующие короткого времени обслуживания, и «тяжелые», сложные для обработки. Такое упрощение недалеко от реальности. Например, сенсоры и датчики могут регулярно отправлять в ближайший узел свои показания, образуя поток «легких» запросов. В то же время запрос на распознавание номера автомобиля, посылаемый дорожной видеокамерой, является «тяжелым» запросом. Для расчета средних длин очередей и ожидаемого времени обслуживания fog узлы рассматриваются авторами как марковские системы обслуживания с неограниченным буфером и пуассоновскими входными потоками легких и тяжелых запросов, поступающих в узел j с интенсивностями λ_i и λ_i соответственно. Времена обработки этих запросов распределены экспоненциально с интенсивностями µ, и µ, соответственно. На основании описания алгоритма составляется система уравнений состояния и выводятся выражения для интенсивностей перехода, что позволяет найти стационарные вероятности состояний и ожидаемые времена обслуживания заявок. Эксперименты показали хорошее совпадение расчетных данных и результатов моделирования.

В работе [11] авторы также применяют аналитическую модель, чтобы количественно оценить выигрыш, получаемый от взаимодействия дата-центров.

Авторы работы [12] ставят перед собой задачу сокращения пространства поиска узловкандидатов для размещения вычислительной нагрузки в гетерогенной среде туманных вычислений и обосновывают свои правила размещения нагрузки. Результатом работы является получение достаточно качественного решения задачи размещения заявок на вычисления в ограниченные временные сроки.

3. Игровые модели для консолидации ресурсов

Yan Sun и Nan Zhang в работе [13] представляют алгоритм краудфандинга для привлечения ресурсов в сеть. Чтобы обеспечить постоянный объем пула ресурсов, авторы разработали механизм стимулирования поставщиков, основанный на теории повторяющихся игр. Согласно разработанному алгоритму, брокер (локальный дата-центр) обеспечивает активным поставшикам (владельцам ресурсов) больший поток задач, в результате чего они имеют возможность увеличить свой доход а. Если задача завершена, то поставщик ресурсов получает дополнительный доход β ($\alpha < \beta$), но несет издержки ϕ , а брокер получает доход θ (θ > 0). Если задача не завершена, то доход брокера равен 0. Чтобы стимулировать владельцев ресурсов выполнять задачи до конца и вовремя, брокер заносит «нерадивого» владельца ресурсов в черный список и не предоставляет ему новых задач, лишая его дохода. Авторы разработали также триггерную стратегию в соответствии с концепцией повторяющейся игры с целью заставить владельцев ресурсов корректно выполнять задачи. Стратегия на стороне брокера: на первой стадии платить более высокое вознаграждение β^* ; если доход брокера всегда равен θ на фазе (t-1), то продолжать платить β^* ; в противном случае не платить вознаграждения, $\beta^* = 0$. Стратегия на стороне поставщика: если доход выше а, принять задачу. Если доход на предыдущей (t-1) фазе всегда равен β^* , поставщик продолжает выполнять задачи на фазе t, в противном случае он задачу не выполняет. Показано, что эта комбинация стратегий представляет собой равновесие Нэша, совершенное на подиграх.

Авторы работы [14] предложили эффективную схему стимулирования пользователей предоставлять свои ресурсы для разгрузки серверов облака на основе теории некооперативных игр. Взаимодействие между серверами облака и устройствами, обладающими свободными ресурсами, формулируется как игра Штакельберга. На первом шаге облако назначает платеж. На втором шаге каждое устройство называет объем предоставляемых ресурсов, ориентируясь на величину платежа. Авторы

доказывают, что для предложенной ими схемы игры Штакельберга между облаком и устройствами существует равновесие Нэша и оно уникально. Был рассмотрен также более общий случай, когда устройство может присоединиться к общему пулу вычислительных ресурсов или покинуть его, и показано, как изменяется при этом равновесие, что очень важно для изучения динамики процесса разгрузки серверов облака. На основе полученных теоретических результатов авторы разработали два алгоритма разгрузки серверов и показали их эффективность с точки зрения уменьшения времени отклика.

Другие модели применения теории игр можно найти в работах [15–17].

4. Энергосбережение

Проблема энергосбережения важна для технических средств туманной сети, а для устройств интернета вещей она является едва ли не первостепенной.

В работе [18] авторы рассматривают задачу максимизации времени жизни облака мобильных устройств за счет перераспределения задач между устройствами. Такие задачи могут возникать при проведении группой людей аварийно-спасательных работ, в обстановке военных действий и чрезвычайных ситуаций, при управлении группой роботов. Рассмотрена сеть из n мобильных устройств, каждое из которых имеет запас энергии E_t^u и доступную вычислительную емкость C_t^u в момент времени t. Каждое устройство u выполняет набор задач T_k^u , k=1,...,j. Каждая задача характеризуется временем жизни TTL_{Tk} и затратами энергии на вычисления C_{Tk} и на передачу данных D_{Tk} . Задачи независимы и могут выполняться одновременно. $E_{Tk}^{u,v}$ обозначает ожидаемые затраты энергии на запуск задачи T_k^u на устройстве v. Эти затраты энергии являются функцией затрат на выполнение вычислений и на передачу данных в другой узел, и задача заключается в том, чтобы продлить время жизни сети путем распределения задач. Для оценки параметров задач авторы разработали специальную платформу, которая позволяет генерировать профили расхода энергии для различных запросов на выполнение вычислений и на передачу данных. Эта информация используется затем как входные данные для системы моделирования. Предложены различные алгоритмы распределения задач и исследованы различные структуры мобильного облака.

В работе [19] ставится задача нахождения компромисса между задержками и электропотреблением в среде туманно-облачных вычислений. Определены функции энергозатрат и задержек и составлена целевая функция задачи разделения нагрузки. Для ее решения применен метод декомпозиции, для чего поставленная задача разделена на три подзадачи, каждая из которых решается одним из известных методов оптимизации. Затем итерационно решается исходная задача.

В работе [20] авторы рассматривают взаимодействие туманной сети и облака с точки зрения удовлетворения требований по времени отклика с учетом выбросов углерода. Дата-центры облака разделены на две категории в зависимости от того, какую энергию они используют: возобновляемую – зеленые дата-центры, имеющие нулевые выбросы углерода, или традиционную – коричневые датацентры с выбросами углерода, пропорциональными потребленной энергии. В облачные дата центры поступают запросы трех типов: обработка данных, хранение информации, программное обеспечение как сервис (PaaS, StaaS, SaaS), отличающиеся средним временем обслуживания. Рассматриваются различные стратегии маршрутизации запросов: «ближайший дата-центр», «ближайший зеленый датацентр», «ближайший зеленый, но со штрафом». Последняя стратегия означает, что выбираемый зеленый центр находится дальше, чем ближайший коричневый, что будет оцениваться штрафом за увеличение длины пути. Трафик fog узлов рассматривается как марковский модулированный пуассоновский поток. Основной результат, полученный авторами, состоит в том, что предложенный ими подход позволяет уменьшить как задержки обслуживания запросов, так и углеродный след облачных структур без значительного ухудшения производительности сети.

В работе [21] авторы рассмотрели взаимодействие между конечными пользователями и облаком, центры которого обладают источниками возобновляемой энергии, резервируемыми энергией из электросетей. Целью являлась минимизация общего потребления энергии, что достигается исполь-

зованием схемы миграции задач на основе учета стоимости миграции при обеспечении необходимого качества обслуживания задач.

В работе [22] авторы решают задачу разгрузки вычислений с мобильных периферийных устройств, где плотное развертывание точек радиодоступа обеспечивает высокую скорость доступа к вычислительным ресурсам туманной сети, но также увеличивает межсотовые помехи. Для совместной оптимизации использования радио- и вычислительных ресурсов сформулирована задача распределения нагрузки с целью минимизации энергопотребления мобильных станций при ограничениях по времени ожидания и по энергопотреблению. Для простейшего случая получен аналитический результат. Для более общего многопользовательского сценария разработаны централизованные и распределенные алгоритмы на основе методов последовательного выпуклого приближения с доказуемой сходимостью к локальным оптимальным решениям.

Анализ публикаций показывает, что вопросам энергосбережения уделяется все больше внимания. Во многих моделях решается как задача обеспечения необходимого времени отклика, так и минимизации расхода энергии.

Заключение

Организация взаимодействия устройств в сетевой среде интернета вещей на фоне постоянного развития технических средств и растущих потребностей подключаемых приложений требует дальнейших разработок в сфере управления и оптимизации функционирования сетевых и вычислительных ресурсов. Задачи управления ресурсами особенно актуальны для туманных и граничных сетей. Обзор публикаций, относящихся к этой области исследований, показывает, что авторы учитывают при разработке моделей как особенности устройств, обеспечивающих сервисы, в первую очередь разнотипность, мобильность, ограниченную энергоемкость, так и особенности потребителей сервисов, такие как высокие требования ко времени реакции и к надежности исполнения заявок. Формулировка целевых функций и выбор методов решения во многом зависят от особенностей конкретных приложений. Чаще всего задача решается как оптимизационная для доминирующего параметра при условии выполнения требований по другим характеристикам.

Первоочередными задачами управления ресурсами считаются в настоящее время поиск более эффективных схем стимулирования для объединения ресурсов, вопросы безопасности, надежности и доступности ресурсов, алгоритмы и модели энергосбережения, учет собственных затрат на реализацию процедур, решающих перечисленные задачи.

Потребность в решении этих задач велика и будет, без сомнения, продолжать расти.

ЛИТЕРАТУРА

- 1. Fernando N., Loke S.W., Rahayu W. Mobile cloud computing: a survey // Future Generation Computer Systems. 2013. V. 29, No. 1. P. 84–106.
- 2. Ahmed A., Ahmed E. A survey on mobile edge computing // Proc. of the 10th Int. Conf. on Intelligent Systems and Control (ISCO), IEEE. 2016. P. 1–8.
- 3. Liu H., Eldarrat F., Alqahtani H., Reznik A., de Foy X., Zhang Y. Mobile Edge Cloud System: Architectures, Challenges, and Approaches // IEEE Systems Journal. 2017. V. 12 (3). P. 2495–2508.
- 4. Gonzalez N., Goya W., Pereira R., Langona K., Silva E. et al. Fog computing: Data analytics and cloud distributed processing on the network edges // Proc. of the 2016 35th Int. Conf. of the Chilean Computer Science Society (SCCC). 2016. P. 1–9.
- 5. Рябоконь В.В., Кузькин А.А., Тутов С.Ю., Махов А.С. Обзор угроз информационной безопасности в концепции граничных вычислений // Вестник Евразийской науки. 2018. № 3. URL: https://esj.today/PDF/79ITVN318.pdf (дата обращения: 12.08.2019).
- 6. Roman R., Lopez J., Mambo M. Mobile edge computing, Fog et al.: a survey and analysis of security threats and challenges // Future Generation Computer Systems. 2018. V. 78, pt. 2. P. 680–698.
- 7. Aazam M., Zeadally S., Harras K.A. Offloading in fog computing for IoT: Review, enabling technologies, and research opportunities // Future Generation Computer Systems. 2018. V. 87, Oct. P. 278–289.
- 8. Sun Y., Lin F, Xu H. Multi-objective Optimization of Resource Scheduling in Fog Computing Using an Improved NSGA-II // Wireless Personal Communications. 2018. V. 102 (2). P. 1369–1385.

- 9. Aazam M., Huh E. Fog computing micro datacenter based dynamic re-source estimation and pricing model for IoT // IEEE 29th Int. Conf. on Advanced Information Networking and Applications. 2015. P. 687–694.
- 10. Yousefpour A., Ishigaki G., Gour R., Jue J.P. On Reducing IoT Service Delay via Fog Offloading // IEEE Internet Things J. 2018. V. 5, is. 2. P. 998–1010.
- 11. Fricker C., Guillemin F., Robert P., Guilherme T. Analysis of an Of-floading Scheme for Data Centers in the Framework of Fog Computing. 2016. 8 p. URL: https://arxiv.org/pdf/1507.05746.pdf (accessed: 12.08.2019).
- 12. Клименко А.Б., Сафроненкова И.Б. Решение задачи распределения вычислительной нагрузки в средах туманных вычислений на базе онтологий // Известия ЮФУ. Технические науки. 2018. № 8. С. 83–94
- 13. Sun Y., Zhang N. A resource-sharing model based on a repeated game in fog computing // Saudi Journal of Biological Sciences, 2017. V. 24 (3). P. 687–694.
- 14. Liua Y., Xua C., Zhanb Y., Liuc Z., Guana J., Zhang H. Incentive mechanism for computation offloading using edge computing: a Stackelberg game approach // Computer Networks. 2017. V. 129, pt. 2. P. 399–409.
- 15. Cao Z, Zhang H., Liu B, Sheng B. A Game-theoretic Framework for Revenue Sharing in Edge-Cloud Computing System. 2018. 10 p. URL: https://arxiv.org/pdf/1711.10102.pdf (accessed: 12.08.2019).
- 16. Zhang H., Liu B., Susanto H., Xue G., Sun T. Incentive mechanism for proximity-based mobile crowd service systems // IEEE INFOCOM: The 35th Annual IEEE Int. Conf. on Computer Communications. 2016. P. 1–9.
- 17. Chen Y., Chen H., Yang S., Gao X., Guo Y., Wu F. Designing Incentive Mechanisms for Mobile Crowdsensing with Intermediaries // Wireless Communications and Mobile Computing. 2019. Article ID 8603526. 20 p. URL: https://doi.org/10.1155/2019/8603526 (accessed: 12.08.2019).
- 18. Mtibaa A., Fahim A., Harras K., Ammar M. Towards resource sharing in mobile device clouds: Power balancing across mobile devices // Proc. of the 2013 2nd ACM SIGCOMM Workshop on Mobile Cloud Computing, MCC 2013. P. 51–56.
- 19. Самуйлов К.Е. К построению модели разделения нагрузки в системе туманных вычислений // Информационные технологии и телекоммуникации. 2017. Т. 5, № 1. С. 8–14
- 20. Borylo P., Lason A., Rzasa J., Szymanski A., Jajszczyk A. Energy-aware fog and cloud interplay supported by wide area software defined networking // Proc. of the 2016 IEEE Inter. Conf. on Communications (ICC). 2016. P. 1–7.
- Fan Q., Ansari N., Sun X. Energy Driven Avatar Migration in Green Cloudlet Networks // IEEE Communications Letters. 2017.
 V. 21, No. 7. P. 1601–1604.
- 22. Sardellitti S., Scutari G., Barbarossa S. Joint optimization of radio and computational resources for multicell mobile-edge computing // IEEE Transactions on Signal and Information Processing over Networks. 2015. V. 1, No. 2. P. 89–103.

Поступила в редакцию 12 августа 2019 г.

Petukhova N.V., Farkhadov M.P., Kachalov D.L. (2020) UNLOADING AND CONSOLIDATION OF COMPUTING RESOURCES IN THE ENVIRONMENT OF FOG AND BOUNDARY COMPUTING. *Vestnik Tomskogo gosudarstvennogo universiteta. Upravlenie vychislitelnaja tehnika i informatika* [Tomsk State University Journal of Control and Computer Science]. 50. pp. 123–129

DOI: 10.17223/19988605/50/15

The article analyzes the management tasks in the fog and edge computing when requests of Internet of Things (IoT) devices are processing. The features of the managed objects are determined: the hardware heterogeneity, the mobility of resources, short response time, resource constraints, energy saving. Based on the review of publications over the past five years, the most urgent management problems in the fog and edge networks have been identified. These tasks include the routing of requests, resources offloading, the control of energy savings. The classification of the reasons for the transfer of tasks from one device to another is given. Examples of researches related to the problem of selecting a resource for service requests are given. Traditional models of resources and algorithms description are supplemented with new characteristics relevant to these networks. The example is an indicator of the instability of the presence of a resource in a fog or edge network. One more example is requests separating into two or more inputs with different characteristics, which is typical for many IoT applications. Power saving is the subject of many studies. Challenges are often a reflection of new realities. For example, the aim can be to distribute requests between the data centers in order to minimize the carbon footprint and to meet the requirements for response time. A specific problem is to attract resources to the common pool. Various algorithms are proposed based on game theory. The future priority tasks are listed in conclusion.

Keywords: internet of things; fog computing; edge computing; offloading; resource consolidation; power saving stochastic models; game theory.

PETUKHOVA Nina Vasilievna (Senior Researcher, V.A. Trapeznikov Institute of Control Sciences of Russian Academy of Sciences, Moscow, Russian Federation).

E-mail: nvpet@ipu.ru

FARKHADOV Mais Pasha (Doctor of Technical Sciences, V.A. Trapeznikov Institute of Control Sciences of Russian Academy of Sciences, Moscow, Russian Federation).

E-mail: mais@ipu.ru

KACHALOV Dmitry Leonidovich (Post-graduate Student, V.A. Trapeznikov Institute of Control Sciences of Russian Academy of Sciences, Moscow, Russian Federation).

E-mail: comdcompdim@mail.ru

REFERENCES

- 1. Fernando, N., Loke, S.W. & Rahayu, W. (2013) Mobile cloud computing: a survey. Future Generation Computer Systems. 29(1). pp. 84–106. DOI: 10.1016/j.future.2012.05.023
- 2. Ahmed, A. & Ahmed, E. (2016) A survey on mobile edge computing. *Proc. of the 10th Int. Conf. on Intelligent Systems and Control (ISCO). IEEE.* pp. 1–8. DOI: 10.1109/ISCO.2016.7727082
- 3. Liu, H., Eldarrat, F., Alqahtani, H., Reznik, A., de Foy, X. & Zhang, Y. (2017) Mobile Edge Cloud System: Architectures, Challenges, and Approaches. *IEEE Systems Journal*. pp. 1–14. DOI: 10.1109/JSYST.2017.2654119
- 4. Gonzalez, N., Goya, W., Pereira, R., Langona, K., Silva, E. et al. (2016) Fog computing: Data analytics and cloud distributed processing on the network edges. *Proc. of the 2016 35th Int. Conf. of the Chilean Computer Science Society (SCCC)*. pp. 1–9. DOI: 10.1109/SCCC.2016.7836028
- 5. Ryabokon, V.V., Kuzkin, A.A., Tutov, S.Yu. & Mahov, A.S. (2018) Review of information security threats in the concept of edge computing. *Vestnik Evraziyskoy nauki The Eurasian Scientific Journal*. 3(10). [Online] Available from: https://esj.today/PDF/79ITVN318.pdf (Accessed: 12th August 2019).
- 6. Roman, R., Lopez, J. & Mambo, M. (2018) Mobile edge computing, Fog et al.: a survey and analysis of security threats and challenges. *Future Generation Computer Systems*. 78(2). pp. 680–698. DOI: 10.1016/j.future.2016.11.009
- 7. Aazam, M., Zeadally, S. & Harras, K.A. (2018) Offloading in fog computing for IoT: Review, enabling technologies, and research opportunities. *Future Generation Computer Systems*. 87. pp. 278–289. DOI: 10.1016/j.future.2018.04.057
- 8. Sun, Y., Lin, F. & Xu, H. (2018) Multi-objective Optimization of Resource Scheduling in Fog Computing Using an Improved NSGA-II. Wireless Personal Communications. 102(2). pp. 1369–1385. DOI: 10.1007/s11277-017-5200-5
- 9. Aazam, M. & Huh, E. (2015) Fog computing micro datacenter based dynamic re-source estimation and pricing model for IoT. *IEEE 29th Int. Conf. on Advanced Information Networking and Applications*. pp. 687–694.
- 10. Yousefpour, A., Ishigaki, G., Gour, R. & Jue, J.P. (2018) On Reducing IoT Service Delay via Fog Offloading. *IEEE Internet Things J.* vol.5. issue 2. pp. 998–1010. DOI: 10.1109/JIOT.2017.2788802
- 11. Fricker, C., Guillemin, F., Robert, P & Guilherme, T. (2016) Analysis of an Of-floading Scheme for Data Centers in the Framework of Fog Computing. [Online] Available from: https://arxiv.org/pdf/1507.05746.pdf (Accessed: 12th August 2019).
- Klimenko, A.B. & Safronenkova, I.B. (2018) Ontology based work-load allocation problem solving in fog computing environment. *Izvestiya SFedU. Tekhnicheskie nauki Izvestiya SFedU. Engineering Sciences*. 8. pp. 83–94. DOI: 10.23683/2311-3103-2018-8-83-94
- 13. Sun, Y. & Zhang, N. (2017) A resource-sharing model based on a repeated game in fog computing. *Saudi Journal of Biological Sciences*. 24(3). pp. 687–694. DOI: 10.1016/j.sjbs.2017.01.043
- 14. Liua, Y., Xua, C., Zhanb, Y., Liuc, Z., Guana, J. & Zhang, H. (2017) Incentive mechanism for computation offloading using edge computing: A Stackelberg game approach. *Computer Networks*. 129(2). pp. 399–409. DOI: 10.1016/j.comnet.2017.03.015
- 15. Cao, Z., Zhang, H., Liu, B. & Sheng, B. (2018) A Game-theoretic Framework for Revenue Sharing in Edge-Cloud Computing System. Pages 10. [Online] Available from: https://arxiv.org/pdf/1711.10102.pdf
- 16. Zhang, H., Liu, B., Susanto, H., Xue, G. & Sun, T. (2016) Incentive mechanism for proximity-based mobile crowd service systems. *In IEEE INFOCOM. The 35th Annual IEEE International Conference on Computer Communications*. pp. 1–9. DOI: 10.1109/INFOCOM.2016.7524549
- 17. Chen, Y., Chen, H., Yang, S., Gao, X., Guo, Y. & Wu, F. (2019) Designing Incentive Mechanisms for Mobile Crowdsensing with Intermediaries. *Wireless Communications and Mobile Computing*. Article ID 8603526. DOI: 10.1155/2019/8603526
- 18. Mtibaa, A., Fahim, A., Harras, K. & Ammar, M. (2013) Towards resource sharing in mobile device clouds: Power balancing across mobile devices. *Proc. of the 2013 2nd ACM SIGCOMM Workshop on Mobile Cloud Computing, MCC-2013.* pp. 51–56. DOI: 10.1145/2534169.2491276
- 19. Samuylov, K. (2017) On the Construction of Workload Allocation Model in Fog Computing System. *Informatsionnye tekhnologii i telekommunikatsii*. 5(1). pp. 8–14.
- Borylo, P., Lason, A., Rzasa, J., Szymanski, A. & Jajszczyk, A. (2016) Energy-aware fog and cloud interplay supported by wide area software defined networking. *Proc. of the 2016 IEEE Int. Conf. on Communications (ICC)*. pp. 1–7. DOI: 10.1109/ICC.2016.7511451
- 21. Fan, Q., Ansari, N. & Sun, X. (2017) Energy Driven Avatar Migration in Green Cloudlet Networks. *IEEE Communications Letters*. 21(7). pp. 1601–1604. DOI: 10.1109/LCOMM.2017.2684812
- 22. Sardellitti, S., Scutari, G. & Barbarossa, S. (2015) Joint optimization of radio and computational resources for multicell mobile-edge computing. *IEEE Transactions on Signal and Information Processing over Networks.* 1(2). pp. 89–103.

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

2020 Управление, вычислительная техника и информатика

№ 50

СВЕДЕНИЯ ОБ АВТОРАХ

БАТРАЕВА Инна Александровна — доцент, кандидат физико-математических наук, заведующая кафедрой технологий программирования на базе филиала ООО «Мирантис ИТ» в г. Саратове факультета компьютерных наук и информационных технологий Саратовского национального исследовательского государственного университета им. Н.Г.Чернышевского. E-mail: BatraevaIA@info.sgu.ru

БУЙ Зуи Тан – аспирант Московского физико-технического института. E-mail: duytan@phystech.edu

ВОРОБЬЕВ Андрей Владимирович – кандидат технических наук, доцент кафедры геоинформационных систем факультета информатики и робототехники Уфимского государственного авиационного технического университета. E-mail: geomagnet@list.ru

ВОРОБЬЕВА Гульнара Равилевна – кандидат технических наук, доцент кафедры вычислительной математики и кибернетики факультета информатики и робототехники Уфимского государственного авиационного технического университета. E-mail: gulnara.vorobeva@gmail.com

ГУБИНА Оксана Викторовна – аспирант Национального исследовательского Томского государственного университета. E-mail: gov7@mail.ru

ДОМБРОВСКИЙ Владимир Валентинович – профессор, доктор технических наук, заведующий кафедрой информационных технологий и бизнес аналитики Института экономики и менеджмента Национального исследовательского Томского государственного университета. E-mail: dombrovs@ef.tsu.ru

ЗОЛОТОРЕВИЧ Людмила Андреевна — кандидат технических наук, доцент Белорусского государственного университета информатики и радиоэлектроники (г. Минск, Беларусь). E-mail: zolotorevichla@bsuir.by

ИСАЕВ Сергей Владиславович – доцент, кандидат технических наук, заместитель директора Института вычислительного моделирования СО РАН (г. Красноярск). E-mail: si@icm.krasn.ru

ИСАЕВА Ольга Сергеевна – кандидат технических наук, старший научный сотрудник отдела прикладной информатики Института вычислительного моделирования СО РАН (г. Красноярск). E-mail: isaeva@icm.krasn.ru

КАРИМ Пешанг Хасан – аспирант кафедры прикладной информатики Института прикладной математики и компьютерных наук Национального исследовательского Томского государственного университета. E-mail: peshangh@yahoo.com

КАЧАЛОВ Дмитрий Леонидович – аспирант Института проблем управления им. В.А. Трапезникова РАН (г. Москва). E-mail: comdcompdim@mail.ru

КИМ Олег Хонбинович – научный сотрудник Национального исследовательского Томского государственного университета. E-mail: oh.kim@mail.ru

КОПАТЬ Дмитрий Ярославович – аспирант Гродненского государственного университета им. Я. Купалы (г. Гродно, Беларусь). E-mail: dk80395@mail.ru

КОСТИНА Мария Алексеевна — кандидат технических наук, доцент Томского государственного университета систем управления и радиоэлектроники, доцент Национального исследовательского Томского политехнического университета. E-mail: mashenkasoldatova@mail.ru

КОШКИН Геннадий Михайлович – профессор, доктор физико-математических наук, профессор Национального исследовательского Томского государственного университета. E-mail: kgm@mail.tsu.ru

КУЛЯСОВ Никита Владимирович – инженер отдела информационно-телекоммуникационных технологий Института вычислительного моделирования СО РАН (г. Красноярск). E-mail: razor@icm.krasn.ru

ЛЕЗГЯН Артем Саркисович – студент Саратовского национального исследовательского государственного университета им. Н.Г. Чернышевского. E-mail: lezgyan@yandex.ru

МАТАЛЬЩКИЙ Михаил Алексеевич – профессор, доктор физико-математических наук, профессор кафедры фундаментальной и прикладной математики Гродненского государственного университета им. Я. Купалы (г. Гродно, Беларусь). E-mail: m.matalytski@gmail.com

МАТРОСОВА Анжела Юрьевна – профессор, доктор технических наук, профессор Национального исследовательского Томского государственного университета. E-mail: mau11@yandex.ru

МИХАЛЕВ Антон Сергеевич — старший преподаватель кафедры информатики Института космических и информационных технологий Сибирского федерального университета (г. Красноярск). E-mail: asmikhalev@yandex.ru

МИХЕЕВ Павел Андреевич – кандидат технических наук, старший научный сотрудник лаборатории расчетных исследований отдела нейтронно-физических исследований Научно-исследовательского технологического института им. А.П. Александрова (г. Санкт-Петербург). E-mail: doka.patrick@gmail.com

НАРЦЕВ Андрей Дмитриевич — студент Саратовского национального исследовательского государственного университета им. Н.Г. Чернышевского. E-mail: narcev.andrey@gmail.com

ПАЗНИКОВ Алексей Александрович – доцент, кандидат технических наук, старший научный сотрудник кафедры вычислительной техники Санкт-Петербургского государственного электротехнического университета «ЛЭТИ» им. В.И. Ульянова (Ленина). E-mail: apaznikov@gmail.com

ПОДДУБНЫЙ Василий Васильевич – профессор, доктор технических наук, профессор кафедры прикладной информатики Института прикладной математики и компьютерных наук Национального исследовательского Томского государственного университета. E-mail: vvpoddubny@gmail.com

РУБАН Анатолий Иванович – профессор, доктор технических наук, заслуженный деятель науки РФ, профессор кафедры информатики Института космических и информационных технологий Сибирского федерального университета (г. Красноярск). E-mail: ai-rouban@mail.ru

СЕМЁНОВА Ольга Валерьевна – кандидат физико-математических наук, старший научный сотрудник Института проблем управления им. В.А. Трапезникова РАН (г. Москва). E-mail: olgasmnv@gmail.com

СОЛДАТОВ Андрей Алексеевич – кандидат технических наук, доцент Томского государственного университета систем управления и радиоэлектроники, доцент Национального исследовательского Томского политехнического университета. E-mail: Soldatov.88@bk.ru

СОЛДАТОВ Алексей Иванович – доктор технических наук, профессор Томского государственного университета систем управления и радиоэлектроники, профессор Национального исследовательского Томского политехнического университета. E-mail: asoldatof@mail.ru

СУЩЕНКО Сергей Петрович – профессор, доктор технических наук, заведующий кафедрой прикладной информатики Института прикладной математики и компьютерных наук Национального исследовательского Томского государственного университета. E-mail: ssp.inf.tsu@gmail.com

ПАШИНСКАЯ Татьяна Юрьевна – кандидат физико-математических наук, доцент кафедры информационных технологий и бизнес аналитики института экономики и менеджмента Национального исследовательского Томского государственного университета. E-mail: tatyana.obedko@mail.ru

ПЕТУХОВА Нина Васильевна — старший научный сотрудник Института проблем управления им. В.А.Трапезникова РАН (г. Москва). E-mail: nvpet@ipu.ru

ФАРХАДОВ Маис Паша – доктор технических наук, главный научный сотрудник, заведующий лабораторией Института проблем управления им. В.А.Трапезникова РАН (г. Москва). E-mail: mais@ipu.ru

ФИЛИППОВА Елена Владимировна – доцент, кандидат технических наук, доцент кафедры теоретической и прикладной информатики Новосибирского государственного технического университета. E-mail: e.filippova@corp.nstu.ru

ЧУБИЧ Владимир Михайлович – профессор, доктор технических наук, заведующий кафедрой теоретической и прикладной информатики Новосибирского государственного технического университета. E-mail: chubich@ami.nstu.ru

ЦИЦИАШВИЛИ Гурами Шалвович – профессор, доктор физико-математических наук, профессор кафедры алгебры, геометрии и анализа Дальневосточного федерального университета, главный научный сотрудник Института прикладной математики ДВО РАН (г. Владивосток). E-mail: guram@iam.dvo.ru

Научный журнал

ВЕСТНИК ТОМСКОГО ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА

УПРАВЛЕНИЕ, ВЫЧИСЛИТЕЛЬНАЯ ТЕХНИКА И ИНФОРМАТИКА

TOMSK STATE UNIVERSITY JOURNAL OF CONTROL AND COMPUTER SCIENCE

2020. № 50

Редактор Е.Г. Шумская Оригинал-макет Е.Г. Шумской Редакторы-переводчики: Г.М. Кошкин; В.Н. Горенинцева Дизайн обложки Л.Д. Кривцовой

Подписано к печати 24.03.2020 г. Формат $60x84^{1}/_{8}$. Гарнитура Times. Усл. печ. л. 15,4. Тираж 250 экз. Заказ № 4280. Цена свободная.

Дата выхода в свет 27.03.2020 г.

Журнал отпечатан на полиграфическом оборудовании Издательского Дома Томского государственного университета 634050, г. Томск, Ленина, 36

Тел. 8(382-2)—52-98-49; 8(382-2)—52-96-75 Сайт: http://publish.tsu.ru; E-mail: rio.tsu@mail.ru