

УДК 004.01:006.72 (470.22)

## МЕТОДЫ РЕГЕНЕРАТИВНОГО МОДЕЛИРОВАНИЯ ДЛЯ АНАЛИЗА МНОГОСЕРВЕРНЫХ СИСТЕМ ОБСЛУЖИВАНИЯ

И. В. Пешкова<sup>1</sup>, А. С. Румянцев<sup>2</sup>

<sup>1</sup> *Петрозаводский государственный университет, Россия*

<sup>2</sup> *Институт прикладных математических исследований КарНЦ РАН,  
ФИЦ «Карельский научный центр РАН», Петрозаводск, Россия*

В статье представлены методы регенеративного моделирования многосерверных систем обслуживания в применении к доверительному оцениванию их характеристик. Кроме хорошо известных методов построения классической и слабой регенерации представлены менее известные и в значительной мере новые методы искусственной регенерации и регенеративных огибающих в применении к современным моделям многосерверных систем с управлением энергоэффективностью и вычислительного кластера.

**Ключевые слова:** регенеративное моделирование; доверительное оценивание; многосерверные системы; метод расщепления; метод регенеративных огибающих; высокопроизводительный вычислительный кластер.

**I. V. Peshkova, A. S. Rumyantsev. REGENERATIVE SIMULATION METHODS FOR MULTISERVER SYSTEMS**

We present regeneration techniques for performance estimation of multiserver system characteristics. Providing a necessary background of classical and weak regeneration notions, we focus on less well-known and mainly new construction of artificial classical regeneration, as well as recent method of regenerative envelopes for modeling modern energy efficient multiserver systems and high-performance computing cluster.

**Keywords:** regenerative simulation; confidence estimation; multiserver systems; splitting; high-performance computing cluster.

### ВВЕДЕНИЕ

Современные вычислительные системы, как правило, имеют многопроцессорную или *многоядерную* архитектуру. Это обусловлено как физическими особенностями современного технологического процесса производства микросхем [34], так и потребностью в создании вычислителей большой мощности, пригодных для решения вычислительноемких задач

(таких, например, как моделирование климата). Развитие каналов связи и коммуникационных сетей позволяет осуществлять перенос части задач в так называемые *облачные* системы, в которых вычисления производятся удаленно, на оборудовании поставщика услуг, на гибко выделяемой части ресурсов из общего пула серверов (как правило, расположенных в нескольких центрах обработки данных). Пе-

редача данных до конечного потребителя в современных сетях также предполагает согласованную работу нескольких устройств, а кроме того, работу с несколькими каналами, в том числе для решения задач снижения задержек и оптимизации пропускной способности.

Исследование характеристик качества обслуживания, надежности, энергоэффективности таких систем можно проводить методами теории массового обслуживания. Модели таких систем в рамках теории массового обслуживания относятся к классу многосерверных. В то же время особенности вычислительных систем и сетей, такие как неконсервативность дисциплины обслуживания (возможность простоя оборудования при наличии ожидающих клиентов) [30] и сложная структура входного потока клиентов [14], затрудняют получение аналитических результатов. Таким образом, необходимо применение методов стохастического моделирования и надежного оценивания характеристик систем.

Регенеративный метод является одним из самых мощных инструментов моделирования неклассических систем обслуживания [2]. Этот метод позволяет проводить доверительное оценивание характеристик системы даже в случае, когда классическое оценивание неприменимо (например, в случае наличия зависимости в исследуемом процессе). Метод основан на выделении так называемых циклов регенерации на траектории исследуемого процесса, между которыми отсутствует либо имеется лишь слабая зависимость, что позволяет применять аналог центральной предельной теоремы для построения доверительного интервала характеристики процесса [9]. Основная цель данной работы состоит в следующем: исследовать применимость методов регенеративного моделирования для анализа новых моделей многосерверных систем обслуживания. Важной задачей работы является развитие регенеративного метода для анализа больших систем, таких как программно определяемые сети и вычислительные кластеры. Указанные объекты обладают большим объемом вспомогательной информации (например, информации об ожидании заявки на узлах сети при прохождении маршрута), которая позволяет идентифицировать регенерацию конструктивно. В свою очередь, конструктивное построение моментов регенерации позволит проводить доверительное оценивание практически важных метрик производительности больших систем.

*Классическая регенерация* является достаточно хорошо изученным типом регенерации

[9, 19, 31, 35], в которой в качестве моментов регенерации используются, например, моменты опустошения системы (когда нет заявок ни в очереди системы, ни на обслуживании), а циклы регенерации являются независимыми. Критерий стационарности односерверной системы обслуживания гарантирует существование бесконечной последовательности таких моментов [9]. В то же время наличие моментов классической регенерации в модели многосерверной системы не гарантировано. *Слабая регенерация* означает попадание исследуемого процесса в определенное множество состояний, а также допускает зависимость на соседних циклах регенерации [9, 36], тем самым расширяя возможности ее использования для моделирования более сложных, в том числе многосерверных систем. Для построения моментов слабой регенерации могут использоваться конструкция обновляющих событий Боровкова [15, 16, 20] или метод расщепления цепей Маркова, возвратных по Харрису [18, 33]. При этом указанные методы позволяют исследовать многосерверные модели в достаточно общем случае (например, при произвольных распределениях интервалов между приходами клиентов и времени их обслуживания). Указанные типы регенерации пригодны для анализа классических систем, однако при анализе современных моделей (таких, как модели высокопроизводительных и распределенных вычислительных систем) возникают трудности, требующие доработки или адаптации методов.

Для успешного применения существующих методов регенеративного оценивания необходимо гарантировать конечность средней длины цикла регенерации. В ряде случаев такое условие обеспечить затруднительно (например, если достаточные условия стационарности системы не гарантируют существования классической регенерации). Одним из перспективных методов разрешения данного затруднения является метод искусственной регенерации с помощью экспоненциального расщепления плотности [8]. Другим новым подходом является метод регенеративных огибающих. Суть метода сводится к замене исходной системы на основе так называемого метода каплинга [6, 7] и свойства монотонности исследуемых процессов парой систем [25] (минорантной и мажорантной), которые регенерируют в классическом смысле и используются для построения доверительного интервала характеристики исходной системы. В данной статье применение указанных методов продемонстрировано на моделях современных вычислительных систем.

Изложение построено по усложнению методов, адекватному усложнению моделей. Более точно структура работы следующая. В первом разделе представлен метод построения классической регенерации, для иллюстрации метода исследуется модель односерверной системы. Далее приводятся два метода построения моментов слабой регенерации: метод обновляющих событий и метод расщепления цепи Маркова, возвратной по Харрису. Применение методов проиллюстрировано на классической многосерверной модели. В следующем разделе рассматривается метод искусственной регенерации с использованием экспоненциального расщепления плотности. Иллюстрация метода проведена на основе исследования энергоэффективности многосерверной системы с управлением активностью серверов (в указанной системе классические моменты регенерации отсутствуют). Далее представлен метод регенеративных огибающих, работа которого продемонстрирована на модели высокопроизводительного вычислительного кластера. Наконец, представлены практические аспекты построения доверительных интервалов для характеристик регенерирующих процессов.

## КЛАССИЧЕСКАЯ РЕГЕНЕРАЦИЯ

**Определение 1.** Случайный процесс  $Z = \{Z(t)\}_{t \geq 0} \in E$  с непрерывными справа траекториями называется *регенерирующим в классическом смысле* [35], если найдется такой процесс восстановления  $\beta = \{\beta_n = \alpha_0 + \dots + \alpha_n, n \geq 0\}$ , что для  $k \geq 0$  процесс

$$\theta_{\beta_k}(Z, \beta) := \{\{Z(\beta_k + t)\}_{t \geq 0}, \{\beta_i - \beta_k\}_{i > k}\} \quad (1)$$

обладает следующими свойствами:

- (1) распределение  $\theta_{\beta_k}(Z, \beta)$  не зависит от  $k$ ;
- (2) процесс  $\theta_{\beta_k}(Z, \beta)$  не зависит от предистории  $\{\{Z(t)\}_{t < \beta_k}, \beta_0, \dots, \beta_k\}$ .

Процесс  $\beta$  называется *процессом восстановления, вложенным в процесс  $Z$* , а моменты  $\beta_k$  называются *моментами классической регенерации*. Заметим, что последовательность моментов регенераций не обязательно уникальна. Значения  $G_k$  процесса между соседними моментами регенерации (м.р.) называют *циклами регенерации*, т. е.

$$G_k = \{Z(t + \beta_k)\}_{0 \leq t < \alpha_{k+1}}, \quad k \geq 0,$$

при этом элемент  $\{Z(t)\}_{t < \alpha_0}$  называют *задержкой* (если  $\alpha_0 = \beta_0 = 0$ , то говорят о процессе без задержки). Будем считать, что длина задержки конечна с вероятностью 1,  $\alpha_0 < \infty$ .

*Длины циклов регенерации*  $\{\alpha_k, k \geq 1\}$  классически регенерирующих процессов являются независимыми и одинаково распределенными (н.о.р.) случайными величинами (с.в.), и циклы процесса являются н.о.р. Регенерирующие процессы с конечной средней длиной цикла регенерации  $E\alpha_1 < \infty$  называются *положительно возвратными*.

Одно из важных свойств регенерирующих процессов состоит в том, что предельное распределение регенерирующего процесса существует при достаточно слабых условиях.

**Теорема 1.** [9] Пусть длина первого цикла  $\alpha_1$  положительно возвратного регенерирующего процесса  $Z$  – нерешетчатая, т. е. не существует такого числа  $d > 0$ , что выполняется равенство  $P\left(\sum_{j=1}^{\infty} \{\alpha_1 = a + jd\}\right) = 1$ , где  $a \in \mathbf{R}$ . Тогда существует предельное распределение  $P_e$  процесса  $Z$  при  $t \rightarrow \infty$ , и для ограниченной измеримой функции  $f(Z) : E \rightarrow \mathbb{R}^+$  выполнено соотношение

$$E_e f(Z) = \frac{1}{E\alpha_1} E_0 \int_0^{\alpha_1} f(Z(s)) ds, \quad (2)$$

где  $E_0$  – математическое ожидание процесса без задержки (естественно считать  $f(Z)$  некоторой характеристикой системы).

Кроме того, если  $f$  непрерывна почти всюду, то выполняется соотношение

$$E_0 f(Z(t)) \rightarrow E_e f(Z), \quad t \rightarrow \infty.$$

Классическая регенерация процесса с дискретным временем  $Z = \{Z_n\}_{n \geq 1}$  определяется аналогично, а процесс восстановления  $\beta$  определен на множестве натуральных чисел  $\mathbf{N}$ . Для регенерирующего в классическом смысле процесса с дискретным временем циклы  $G_k = \{Z_i\}_{\beta_k \leq i < \beta_{k+1}}$  независимы и одинаково распределены для всех  $k \geq 1$ . Соотношение (2) в теореме 1 для процесса с дискретным временем [3] принимает вид

$$E_e(f(Z)) = \frac{E_0\left(\sum_{i=0}^{\beta_1-1} f(Z_i)\right)}{E\alpha_1}. \quad (3)$$

## Система $GI/G/1$

В качестве моментов классической регенерации при исследовании систем часто используются моменты опустошения системы (если оно возможно). Примером системы, в которой критерий стационарности гарантирует наступление опустошения, является односерверная система  $GI/G/1$ . Напомним, в такой системе в моменты  $0 = t_0 < t_1, \dots$ , образующие процесс восстановления, поступают клиенты,

при этом клиенту  $i$  требуется  $S_i$  единиц времени для обслуживания на единственном сервере. Клиент  $i$  ожидает в общей неограниченной очереди время  $D_i \geq 0$  и по окончании обслуживания покидает систему. Времена между приходами  $\{T_i := t_i - t_{i-1}\}$  и времена обслуживания  $\{S_i\}$ ,  $i \geq 1$  являются н.о.р.с.в. Время ожидания каждого клиента можно определить с помощью хорошо известной стохастической рекуррентной последовательности, рекурсии Линдли:

$$D_{i+1} = (D_i + S_i - T_i)^+, \quad i \geq 1, \quad (4)$$

где  $D_1 = 0$  есть время ожидания первого клиента (процесс без задержки), а  $(\cdot)^+ = \max(0, \cdot)$ . Отметим, что рекурсия (4) определяет оставшуюся работу (нагрузку)  $D_i$  в системе в момент времени  $t_i^-$  перед приходом клиента  $i$ . Известно [9], что стационарное время ожидания существует тогда и только тогда, когда

$$E(T - S) > 0 \text{ или } \rho := \frac{\lambda}{\mu} < 1, \quad (5)$$

где  $ET = 1/\lambda$ .  $ES = 1/\mu$  (нижних индексов лишены типичные время между приходами и время обслуживания). В качестве моментов регенерации можно определить моменты прихода клиентов в пустую систему [9]:

$$\beta_k = \min\{i > \beta_{k-1} : D_i = 0\}.$$

При этом  $\beta_0 = 0$  и можно показать, что с.в.  $\alpha_1$  – непериодическая, поскольку (5) влечет

$$P(\alpha_1 = 1) = P(T - S > 0) > 0.$$

Отметим, что цикл регенерации в такой системе состоит из *цикла занятости сервера* (в течение которого каждый приходящий клиент вынужден ожидать) и последующего *интервала простоя сервера*. При этом процесс  $\{D_i\}_{i \geq 0}$  с последовательностью  $\beta$  является положительно возвратным регенерирующим процессом в классическом смысле тогда и только тогда, когда выполнен критерий (5) [4, 9].

### СЛАБАЯ РЕГЕНЕРАЦИЯ

С. Асмуссен и Х. Торрисон [9, 36] обобщили понятие классической регенерации до так называемой *регенерации в широком смысле*, при которой допускается зависимость между циклами регенерации. Более точно, в определении 1 свойство (2) заменяется на следующее [35]:

(2') процесс  $\theta_{\beta_k}(Z, \beta)$  не зависит от  $\{\beta_0, \dots, \beta_k\}$ .

На практике широко используется частный случай регенерации в широком смысле, *слабая (однозависимая) регенерация*, при которой допускается зависимость лишь между соседними циклами регенерации. Отметим, что при слабой регенерации процесс  $\beta = \{\beta_n, n \geq 0\}$  по-прежнему является процессом восстановления. Как и в случае классической регенерации, существует предельное распределение  $P_e$  процесса  $Z$  [3, 9], и теорема 1 справедлива.

### Метод обновляющих событий

Для построения моментов слабой регенерации можно использовать метод обновляющих событий А. А. Боровкова [15, 16, 20]. Метод основан на определении так называемых периодов обновления и разработан для стохастических рекуррентных последовательностей, широко применяемых в анализе систем обслуживания (примером такой последовательности является рекурсия Линдли (4)).

Пусть последовательность  $\{Z_n\}_{n \geq 1}$  определяется с помощью наперед заданной *управляющей последовательности*  $\{X_n\}_{n \geq 1}$  следующим одношаговым рекуррентным соотношением:

$$Z_{n+1} = g_1(Z_n, X_n), \quad n \geq 1,$$

где н.о.р.с.в.  $X_n$  принимают значения на некотором пространстве состояний  $\mathfrak{X}$ ,  $g_1 : E \times \mathfrak{X} \rightarrow E$  – измеримая функция. Элементы  $X_n \in \mathfrak{X}$  обычно используются для описания входного потока и длительностей обслуживания клиентов системы, например,  $X_n = (T_n, S_n)$  для системы  $GI/G/1$ . Разворачивая рекурсию от шага  $n$  на  $L$  шагов вперед, получим выражение для вычисления значения  $Z_{n+L}$ :

$$Z_{n+L} = g_L(z, X_n, \dots, X_{n+L-1}),$$

где  $g_L : E \times \mathfrak{X}^L \rightarrow E$  и  $\mathfrak{X}^L = \mathfrak{X} \times \dots \times \mathfrak{X}$ . Если при этом найдутся множество  $C \in E$  и измеримые подмножества  $B_L \in \mathfrak{X}^L$  такие, что для любых значений  $z_1, z_2 \in C$  и  $(x_1, \dots, x_L) \in B_L$  выполняется соотношение

$$g_L(z_1, x_1, \dots, x_L) = g_L(z_2, x_1, \dots, x_L), \quad (6)$$

то говорят о *независимости от прошлого* процесса  $Z$  через  $L$  шагов после попадания процесса во множество  $C$  на шаге  $n$  и элементов  $\{X_n, \dots, X_{n+L-1}\}$  во множество  $B_L$ . При этом момент  $n + L$  наступления события

$$\Omega_n = \{Z_n \in C, (X_n, \dots, X_{n+L-1}) \in B_L\}$$

является моментом регенерации, т. е. моменты регенерации можно определить следующим образом:

$$\beta_k = \min_n \{n + L > \beta_{k-1} : I\{\Omega_n\} = 1\}, \quad k > 1,$$

где  $I\{\Omega\}$  – индикатор случайного события  $\Omega$ . Заметим, что если множество  $C$  состоит только из одной точки, то равенство (6) выполняется автоматически для любых  $L$  и  $B_L$ .

Отметим, что цепочка *обновляющих событий*, приводящая к слабой регенерации, предполагает алгебраическую независимость процесса от предыстории. При этом поведение слабо регенерирующего процесса после момента регенерации зависит от  $L$  значений процесса на предыдущем цикле. При таком определении циклы регенерации являются *однозависимыми* (т. е. зависимы только два соседних цикла) [11, 17], однако сдвинутый процесс  $\{Z_{n+\beta_k}\}_{n \geq 0}$  не зависит от моментов регенерации  $\{\beta_0, \dots, \beta_k\}$ . Следовательно, процесс  $Z$  является слабо регенерирующим. Событие  $\Omega_n$  называется *обновляющим событием* [1]. Если событие  $\Omega_n$  происходит, то интервал  $[n, n+L)$  является *периодом обновления*, момент  $n+L$  является моментом регенерации, а распределение процесса в момент  $n+L$   $\mu(\cdot) := P(Z_{n+L} \in \cdot)$  называется *мерой регенерации*.

### Система $GI/G/m$

В отличие от односерверной системы критерий стационарности многосерверной системы не гарантирует ее опустошения и, следовательно, существования моментов классической регенерации. Действительно, рассмотрим  $m$ -серверную систему обслуживания  $\Sigma$  типа  $GI/G/m$ , адаптировав обозначения потока клиентов из анализа системы  $GI/G/1$ . Рассмотрим вектор Кифера – Вольфовица  $W_n = (W_{n,1}, \dots, W_{n,m})$ , содержащий упорядоченную по возрастанию нагрузку (незавершенную работу) на серверах в момент  $t_n^-$ , включая ожидающих в очереди клиентов (обслуживание которых ведется в порядке поступления). Одношаговое рекуррентное соотношение для вектора нагрузки, предложенное в работе [21], является многомерным аналогом рекурсии Линдли (4):

$$W_{n+1} = R(W_n + lS_n - \mathbf{1}T_n)^+, \quad (7)$$

где  $l = (1, 0, \dots, 0)$ ;  $\mathbf{1} = (1, \dots, 1)$  и оператор  $R(\cdot)$  упорядочивает компоненты вектора в порядке возрастания. Заметим, что значение  $D_n := W_{n,1}$  является временем ожидания клиента  $n$ , при этом моменты прихода клиентов являются последовательными моментами дискретного времени системы.

Критерий стационарности процесса нагрузки обобщает (5) следующим образом:

$$E(mT - S) > 0 \text{ или } \rho < m. \quad (8)$$

Условие (8) гарантирует лишь условие  $P(mT - S > 0) > 0$ , т. е. в «худшем» случае [12]

$$\min\{i \geq 0 : \mathbf{s} < i\mathbf{t}\} = m, \quad (9)$$

где

$$\begin{aligned} \mathbf{s} &= \sup\{x : P(S \geq x) = 1\}, \\ \mathbf{t} &= \inf\{x : P(T \leq x) = 1\}. \end{aligned}$$

Отметим, что если  $\mathbf{t} = \infty$  либо  $\mathbf{s} = 0$ , то  $P(T - S > 0) > 0$  и система может полностью опустошаться, см. [4, 32, 38]. Условие (9) означает, что в системе после шага  $m-1$  всегда находится не менее  $m-1$  заявки. Действительно, рассмотрим систему  $\Sigma^*$  с *детерминированным входным потоком*  $T_i^* = \mathbf{t}, S_i^* = \mathbf{s}$ . Из условия (9) следует, что до шага  $m-1$  в системе не произойдет ни одного ухода, т. е.

$$W_{m-1}^* = \mathbf{v} := (0, \mathbf{s} - (m-1)\mathbf{t}, \dots, \mathbf{s} - \mathbf{t}).$$

Перед приходом клиента  $m$  гарантированно произойдет уход клиента, т. е.

$$W_m^* = R(\mathbf{v} + l\mathbf{s} - \mathbf{1}\mathbf{t})^+ = \mathbf{v}. \quad (10)$$

При этом можно показать, опираясь на свойство монотонности процесса  $W_n$  относительно управляющих последовательностей [37], что поскольку  $T_n \geq \mathbf{t}, S_n \leq \mathbf{s}$ , то  $W_n \geq W_n^*$  (по вероятности). Таким образом, при любых управляющих последовательностях нагрузка в системе  $\Sigma$  всегда *не менее*  $\mathbf{v}$ .

Соотношение (10) подсказывает метод построения обновляющего события: необходимо, чтобы цепь  $W_n$  оказалась «достаточно близко» к значению  $\mathbf{v}$  и при этом управляющие последовательности  $T_n$  и  $S_n$  оказались «достаточно близко» к значениям  $\mathbf{t}$  и  $\mathbf{s}$  соответственно [12, 22, 36]. Для того чтобы приблизить значение  $W_n$  к  $\mathbf{v}$ , необходимо прежде всего показать, что  $W_n$  не растет неограниченно по вероятности, если выполнено (8). Действительно, предположим обратное: для любого  $x > 0$  и любого  $\varepsilon_0 > 0$  найдется  $n_0 = n_0(\varepsilon_0, x) < \infty$  такое, что  $P(D_n \leq x) < \varepsilon_0$  для всех  $n \geq n_0$ . Рассмотрим приращение суммарной работы за один приход:

$$\Delta_n = \sum_{i=1}^m (W_{n+1,i} - W_{n,i}) \leq S_n.$$

Заметим, что если  $D_n = W_{n,1} > \mathbf{t}$ , то все серверы системы работают *без простоев*, и из (7) следует, что  $W_{n+1,i} > 0$  (поскольку уменьшение нагрузки на каждом сервере составит  $T_n \leq \mathbf{t}$ ), поэтому

$$\begin{aligned} E(\Delta_n) &= E(\Delta_n | D_n \leq \mathbf{t})P(D_n \leq \mathbf{t}) + \\ &+ E(\Delta_n | D_n > \mathbf{t})P(D_n > \mathbf{t}) \leq \\ &\leq ES_n\varepsilon_0 + E(S_n - mT_n)(1 - \varepsilon_0). \end{aligned} \quad (11)$$

За счет произвольности выбора  $\varepsilon_0$  и использования условия (8) правая часть может быть сделана строго отрицательной для всех  $n > n_0(\varepsilon_0, \mathbf{t})$ . Это означает, что

$$\mathbb{E} \sum_{i=1}^m W_{n,i} \leq \mathbb{E} \sum_{i=1}^m W_{n_0,i} \leq n_0 \mathbb{E} S < \infty,$$

что противоречит предположению. Таким образом, время ожидания в системе не растет неограниченно по вероятности, т. е. найдется неслучайная последовательность  $\{z_i, i \geq 1\}$ , постоянные  $D_0 < \infty, \varepsilon_1 > 0$  такие, что

$$\inf_i \mathbb{P}(D_{z_i} \leq D_0) \geq \varepsilon_1.$$

Наконец, используя стохастическую ограниченность процесса  $\{W_{n,m} - W_{n,1}\}_{n \geq 1}$  [9, 23], можно показать, что (вообще говоря, для некоторой подпоследовательности  $z_i$ , при необходимости увеличив  $D_0$ , уменьшив  $\varepsilon_1$  и переобозначив элементы последовательности)

$$\mathbb{P}(W_{z_i} \leq \mathbf{1}D_0) \geq \varepsilon_1. \quad (12)$$

Таким образом, процесс  $\{W_n\}_{n \geq 1}$  бесконечно часто с положительной вероятностью  $\geq \varepsilon_1$  попадает в компакт  $\{x \in E : x \leq \mathbf{1}D_0\}$  в неслучайные моменты  $\{z_i\}$ .

Рассмотрим один из таких моментов  $z_1$ . Будем приближать вектор нагрузки, заключенный в компакт  $\mathbf{1}D_0$ , к вектору  $\mathbf{v}$ . Для этого построим циклы разгрузки, состоящие из  $m$  последовательных наступлений следующего события положительной вероятности:

$$H_i = \{T_i \in (\mathbf{t} - \delta, \mathbf{t}], S_i \in [\mathbf{s}, \mathbf{s} + \delta)\},$$

где  $\delta$  выбирается так, чтобы было выполнено неравенство

$$0 < \delta < \frac{m\mathbf{t} - \mathbf{s}}{m + 1}. \quad (13)$$

При этом, даже если в «худшем» случае  $D_0 > (m - 1)(\mathbf{t} - \delta)$ , то уже к моменту  $z_1 + m$  из (7) следует

$$W_{z_1+m} \leq \mathbf{1}(D_0 + \mathbf{s} - m\mathbf{t} + (m + 1)\delta) < \mathbf{1}(D_0 - \delta_0),$$

где величина  $\delta_0 := -(\mathbf{s} - m\mathbf{t} + (m + 1)\delta) > 0$  равна минимальной разгрузке каждого сервера. Заметим, однако, что разгрузка всех серверов происходит лишь через  $m$  шагов, поскольку из (9) следует  $(m - 1)t - s < 0$ . Продолжая реализовывать событие  $H_i$ , через

$$k_1 = m \left\lceil \frac{D_0 - (m - 1)(\mathbf{t} - \delta)}{\delta_0} \right\rceil$$

шагов получим  $W_{z_1+k_1} < D_0^1$  для некоторого  $D_0^1 < \mathbf{1}(m - 1)(\mathbf{t} - \delta)$ . Здесь  $\lceil a \rceil$  означает наименьшее целое больше  $a$ . Выполнив еще  $m - 1$  шаг, заметим, что первая компонента вектора  $W_{z_1+k_1+m-1}$  обращается в нуль, т. е. один из серверов освободится:

$$W_{z_1+k_1+m-1} \leq (0, D_0^1 + \mathbf{s} - (m - 1)\mathbf{t} + m\delta, \dots).$$

На следующем шаге заявка, поступившая без ожидания, займет в векторе последнее место (поскольку из (13)  $D_0^1 - \delta_0 < \mathbf{s} - \mathbf{t} + 2\delta$ ):

$$W_{z_1+k_1+m} \leq (D_0^1 - \delta_0, \dots, D_0^1 - \delta_0, \mathbf{s} - \mathbf{t} + 2\delta).$$

При этом нагрузка на остальных серверах уменьшилась не менее чем на величину  $\delta_0$ . Заметим также, что из (13) следует, что  $W_{z_1+k_1+im-1,1} = 0$  для любого целого  $i \geq 1$ . Продолжая  $m$ -шаговые циклы разгрузки, не более чем через  $k_2 = m \lceil \frac{\mathbf{t} - \delta}{\delta_0} \rceil$  шагов добьемся, что две последовательные заявки поступят без ожидания, при этом нагрузка на остальных серверах может только уменьшиться. Продолжая процедуру, можно заметить, что число шагов  $k_0$  с момента прихода заявки  $z_1$  (момента первого попадания в компакт  $W_{z_i} \leq \mathbf{1}D_0$ ), необходимое для получения серии из  $m - 1$  последовательных заявок, поступающих на обслуживание без ожидания, определяется соотношением

$$k_0 \leq m \left\lceil \frac{D_0}{\delta_0} \right\rceil. \quad (14)$$

Кроме того, (14) влечет  $D_0 - k_0\delta_0 < 0$ , поэтому в момент времени  $t_{z_1+k_0}^-$  выполнено

$$W_{z_1+k_0} \leq (0, \mathbf{s} - (m - 1)\mathbf{t} + m\delta, \dots, \mathbf{s} - \mathbf{t} + 2\delta).$$

Таким образом, в момент дискретного времени  $z_1 + k_0$  в системе одновременно обслуживаются не более  $m - 1$  последних пришедших заявок, которые поступили на обслуживание без задержки. Это означает, что будущее процесса нагрузки  $\{W_{z_1+k_0+i}\}_{i \geq 0}$  не зависит от его прошлого до момента  $z_1 + k_0 - m$ . Отметим, что данная процедура построения моментов регенерации уточняет конструкции в [12, 36]. Осталось указать, что обновление происходит с вероятностью, не меньшей чем

$$\varepsilon_2 := \mathbb{P}(H_i)^{k_0} > 0. \quad (15)$$

Соответственно, множество  $C$  для процесса  $\{W_n\}$  можно определить так:

$$C = \{x \in E : x \leq \mathbf{1}D_0\},$$

длительность периода обновления  $L \leq k_0$ , и множество  $B_L$  для управляющих последовательностей связано с событием  $H_i$  таким образом:

$$B_L = \{T_i \in (\mathbf{t} - \delta, \mathbf{t}], S_i \in [\mathbf{s}, \mathbf{s} + \delta), i = 1, \dots, L\}.$$

При этом из (12) и (15) следует, что событие  $\Omega_n$  (регенерация  $W_n$ ) происходит на интервале  $[z_i; z_i + k_0]$  с вероятностью не меньше  $\varepsilon_1 \varepsilon_2 > 0$ . Отметим, что можно доказать конечность средней длины цикла регенерации,  $E\alpha_1 < \infty$  [13], см. также [4, 23].

### ВОЗВРАТНОСТЬ ПО ХАРРИСУ

Регенерация в классическом смысле, как правило, связана с попаданием цепи Маркова (процесса) в определенное состояние, например, нулевое (опустошение системы), при этом циклы регенерации являются независимыми. Как показано ранее, такое состояние не всегда существует. Это вызвало необходимость расширить понятие классической регенерации. Слабая регенерация допускает попадание цепи и управляющей последовательности в определенное множество состояний, при этом ослабляя требование независимости циклов регенерации. Подобными свойствами возвращения в некоторое множество с положительной вероятностью обладают цепи Маркова, возвратные по Харрису.

Рассмотрим однородную цепь Маркова  $\{Z_n\}_{n \geq 1}$  с начальным состоянием  $Z_1 = x \in E$  и  $r$ -шаговым переходным ядром

$$\begin{aligned} K^r(x, \cdot) &= P(Z_{r+1} \in \cdot | Z_1 = x) = \\ &= P_x(Z_r \in \cdot), \quad r \geq 1. \end{aligned}$$

Заметим, что  $E$  может быть не счетным. Обозначим через  $\tau(A)$  момент первого возвращения в множество  $A$ :

$$\tau(A) = \inf\{n \geq 1 : Z_n \in A\},$$

который является моментом остановки относительно сигма-алгебры, порожденной траекторией цепи до шага  $n$ , поскольку

$$\{\tau(A) \leq n\} = \bigcup_{k=1}^n \{Z_k \in A\}.$$

**Определение 2.** ([18, 33]) Цепь Маркова  $\{Z_n\}_{n \geq 1}$  называется *возвратной по Харрису*, если существует нетривиальная мера  $\varphi$  ( $\sigma$ -конечная мера на  $E$ , оснащенном борелевской сигма-алгеброй  $\mathfrak{E}$ ), такая, что из положительности меры  $\varphi(A) > 0$  для множества  $A \in \mathfrak{E}$  следует, что

$$P_x(\tau(A) < \infty) = 1$$

для любого начального состояния  $x \in E$ , т. е. каждое достижимое множество  $A$  имеет конечную меру  $\varphi$ .

Можно показать [28], что для такой цепи существует единственная инвариантная мера  $\pi$  и, в случае конечности последней, цепь называется *положительно возвратной по Харрису*. Например, если  $\varphi(A) = \sum_{i=1}^{\alpha_1} I\{X_i \in A\}$  и  $E\alpha_1 = \varphi(E) < \infty$ , то

$$\pi(A) = \frac{\varphi(A)}{\varphi(E)} = \frac{\varphi(A)}{E\alpha_1}.$$

В работе [33] (в условиях сепарабельности  $E$ ) дается эквивалентное определение возвратности по Харрису в терминах так называемых *малых множеств*:

**Определение 3.** Цепь  $\{Z_n\}_{n \geq 1}$  называется возвратной по Харрису, если

1) для некоторого *фиксированного* множества  $V \in \mathfrak{E}$ , *меры регенерации*  $\mu$  (вероятностной меры, такой, что  $\mu(V) = 1$ ) и любого  $x \in E$  выполнено

$$P_x(\tau(V) < \infty) = 1;$$

2) для любого  $x \in V$ ,  $B \in \mathfrak{E}$  и некоторых  $\varepsilon \in (0, 1)$ ,  $r \geq 1$  выполнено *условие миноризации*:

$$K^r(x, B) \geq \varepsilon \mu(B). \quad (16)$$

Взаимосвязь данных определений рассмотрим в следующем разделе.

### Метод расщепления

Поясним связь между однозависимой регенерацией и возвратностью по Харрису. Для этого можно воспользоваться *процедурой расщепления* (рандомизации). В моменты попадания цепи  $Z_n$  во множество  $V$  определим 0-1 н.о.р.с.в.  $\xi_n$  такие, что  $P(\xi_n = 1) = \varepsilon$ . Назовем событие  $\{\xi_n = 1\}$  успехом. В случае успеха строим  $Z_{n+r}$  в соответствии с распределением  $\mu(\cdot)$ , в противном случае строим  $Z_{n+r}$  по условному распределению  $Q(Z_n, \cdot)$ , определяемому равенством:

$$Q(x, \cdot) = \frac{K^r(x, \cdot) - \varepsilon \mu(\cdot)}{1 - \varepsilon}.$$

Таким образом, переходное ядро имеет вид

$$K^r(x, \cdot) = \varepsilon \mu(\cdot) + (1 - \varepsilon) Q(x, \cdot). \quad (17)$$

В силу (16) в момент попадания цепи во множество  $V$  происходит *условное расщепление* распределения вероятности перехода за  $r$  шагов.

Поскольку цепь с вер. 1 попадает во множество  $V$  неограниченное число раз, то можно определить последовательность моментов регенерации  $\{\beta_n\}$  следующим образом:

$$\beta_{n+1} = \inf\{i > \beta_n : Z_i \in V, \xi_i = 1\} + r, \quad n \geq 1,$$

причем в силу положительности вероятности успеха  $\varepsilon$  длина цикла регенерации  $\alpha_n$  конечна [10, 18, 28]. При этом значения процесса  $\{Z_{n+i}, 1 \leq i \leq r-1\}$  определяются граничными значениями  $Z_n, Z_{n+r}$ , распределенными в соответствии с  $\mu$ , вследствие чего возникает зависимость между соседними циклами регенерации. В общем случае реализовать такую конструкцию достаточно трудно (см. [9]). В то же время в некоторых случаях возможно конструирование таких моментов, аналогичное построению обновляющих событий (наступление успешного расщепления в возвратной по Харрису цепи эквивалентно наступлению обновляющего события в слабо регенерирующем процессе). Более того, в случае  $r = 1$  цепь регенерирует в классическом смысле [18, 33].

Вернемся к рассмотрению системы  $GI/G/m$ , где в качестве цепи Маркова рассматривается вектор нагрузки  $\{W_n\}_{n \geq 1}$ . Положим в (17)  $\varepsilon = \varepsilon_1 \varepsilon_2$  и определим меру  $\varphi(\cdot)$  как среднее число попаданий во множество  $V = \{x \in E : \mathbf{v} \leq x \leq (0, \mathbf{s} - (m-1)\mathbf{t} + m\varepsilon, \dots, \mathbf{s} - \mathbf{t} + 2\varepsilon)\}$  на цикле регенерации [9, 10, 17, 28] при условии, что стартовым распределением  $W_1$  является  $\mu(\cdot)$  (мера регенерации), т. е.

$$\begin{aligned} \varphi(\cdot) &:= E_\mu \left[ \sum_{i=1}^{\alpha_1} I\{W_i \in \cdot\} \right] = \\ &= \int E_x \left[ \sum_{i=1}^{\alpha_1} I\{W_i \in \cdot\} \right] d\mu(x). \end{aligned}$$

Слабая регенерация произойдет через  $r = m-1$  шагов цепи, и в качестве меры регенерации  $\mu$  можно выбрать следующую меру:

$$\mu(\cdot) = K^{m-1}(0, \cdot).$$

### Искусственная регенерация

Напомним, что при построении регенерации возвратной по Харрису цепи Маркова методом расщепления в случае  $r = 1$  происходит потеря зависимости между соседними циклами (поскольку выбор из распределения  $\mu$  происходит на следующем шаге после попадания в компакт  $V$ ), т. е. цепь регенерирует в классическом смысле. Одним из конструктивных методов использования этого свойства является искусственное построение моментов регенерации за счет экспоненциального расщепления плотности [8]. Этот метод актуален для моделирования систем, в которых характеристики имеют распределения с тяжелыми хвостами.

Экспоненциальное расщепление плотности заменяет исходную случайную величину  $S$  (с

плотностью  $f$ ) на комбинацию трех случайных величин: экспоненциально распределенной, индикатора и остатка. Экспоненциальное расщепление возможно, если существуют такие константы  $\theta > 0$ ,  $x_0$  и  $0 < p < 1$ , что выполнено неравенство миноризации

$$f(x) \geq pf_0(x), \quad x \geq x_0, \quad (18)$$

где  $f_0$  – плотность усеченного слева экспоненциального распределения:

$$f_0(x) := \begin{cases} 0, & x \leq x_0; \\ \theta e^{-\theta(x-x_0)}, & x > x_0. \end{cases} \quad (19)$$

Введем плотность  $f_1(x) = \frac{f(x) - pf_0(x)}{1-p}$ . Тогда

$$f_1(x) := \begin{cases} \frac{f(x)}{1-p}, & x \leq x_0; \\ \frac{f(x) - p\theta e^{-\theta(x-x_0)}}{1-p}, & x > x_0. \end{cases} \quad (20)$$

Легко увидеть, что с.в.  $S$  может быть представлена в виде следующей суммы:

$$S = I\{\xi = 1\}S_0 + I\{\xi = 0\}S_1, \quad (21)$$

где независимые с.в.  $S_0$  и  $S_1$  имеют плотности распределения  $f_0$  и  $f_1$  соответственно. С.в.  $\xi$  имеет распределение Бернулли с вероятностью успеха  $P(\xi = 1) = p$  и называется индикатором расщепления.

Поясним вкратце метод построения искусственной регенерации, предложенный в работе [8]. Рассмотрим многокомпонентный процесс  $\Theta = \{X_1(t), T_1(t), \dots, X_M(t), T_M(t)\}_{t \geq 0}$ , имеющий дискретные компоненты  $X_i \in E$  (например, число заявок в очереди системы) и непрерывные компоненты  $T_i \geq 0$ , линейно убывающие до обнуления (например, остаточное время обслуживания). В момент  $t^*$  обнуления  $i$ -й компоненты происходит изменение дискретных компонент процесса в соответствии с некоторой вероятностной мерой

$$P_i(x, x') = P\{X_1(t^*+) = x'_1, \dots, X_M(t^*+) = x'_M \\ | X_1(t^*-) = x_1, \dots, X_M(t^*-) = x_M\},$$

при этом новое значение  $T_i$  имеет некоторую плотность  $f_i(x, x')$ . Непрерывные компоненты в момент  $t^*$  не претерпевают изменения. Напротив, в момент, не соответствующий обнулению непрерывной компоненты, дискретные компоненты не изменяются. Пусть существует такой вектор  $x^* \in E^M$ , что если дискретные компоненты цепи приняли значение  $x^*$ , то каждая непрерывная компонента допускает миноризацию плотности (18) (возможно, со своими параметрами  $p_i(x, x')$ ,  $x_{0,ix'}$ ,  $\theta_i(x, x')$ ).



Пусть в момент времени  $t^*$  произошло попадание некоторой дискретной компоненты  $X_i(t^*+)$  в состояние  $x_i^*$ . Тогда произведем экспоненциальное расщепление соответствующей непрерывной компоненты  $T_i(t^*+)$ . Для этого добавим в процесс  $\Theta$   $M$ -мерный дискретный процесс  $\{(Z_1(t), \dots, Z_M(t))\}_{t \geq 0}$ , где  $Z_i$  есть фаза экспоненциально расщепленной  $i$ -й непрерывной компоненты. В момент расщепления положим  $Z_i(t) = 1$ , если соответствующий индикатор расщепления равен 1, и такие компоненты назовем проходящими предэкспоненциальную фазу. Тогда в момент времени  $t^* + x_{0,ixx'}$  положим  $Z_i(t) = 0$ , а соответствующая компонента попадет в экспоненциальную фазу (и соответствующая непрерывная компонента может в моменты событий переразыгрываться). Положим  $Z_i(t) = 2$ , если в момент  $t^*$  индикатор расщепления равен 0 либо  $X_i(t^*+) \neq x_i^*$ . Тогда событие искусственной регенерации наступает, когда все непрерывные компоненты попадают в экспоненциальную фазу, т. е. в такой момент  $t^*$ , что

$$\{X_1(t^*) = x_1^*, \dots, X_M(t^*) = x_M^*, \\ Z_1(t^*) = Z_M(t^*) = 0\}.$$

### Система $GI/G/m$ с выключением приборов

Поясним технику построения моментов искусственной регенерации на следующем примере. Пусть в многосерверной системе обслуживания в момент освобождения сервера, при условии пустой очереди, сервер уходит в состояние пониженного энергопотребления, в котором обслуживание заявок невозможно (так называемый спящий режим) на случайное время  $C_i$ , имеющее плотность  $c(\cdot)$ , и напомним, что клиенты приходят в систему через моменты времени  $T_i$  (имеющие плотность  $a(\cdot)$ ), обслуживаются время  $S_i$  (с плотностью  $b(\cdot)$ ). Введем следующий  $m + 1$ -мерный процесс:  $X_0(t)$  есть число заявок в очереди в момент времени  $t$ , а  $X_1(t), \dots, X_m(t)$  есть режимы работы сервера  $1, \dots, m$ :  $X_1(t) = 1$ , если сервер занят обслуживанием клиента, и  $X_1(t) = 0$ , если он находится в режиме пониженного энергопотребления. Непрерывные компоненты определим следующим образом:  $T_0(t)$  есть время до прихода следующего клиента,  $T_1(t), \dots, T_m(t)$  есть времена до окончания текущей активности (бездействия) сервера. Заметим, что клиенты никогда не попадают на обслуживание в момент прихода, но вынуждены дожидаться окончания текущей активности (бездействия) сервера. Пусть вектор  $l = (1, 0, \dots, 0)$ . Тогда в момент прихода

заявки происходит следующее изменение дискретной компоненты: из состояния  $x$  в  $x + l$  (попадание заявки в очередь) с вероятностью  $P_0(x, x + l) = 1$ . В момент окончания текущей активности на сервере  $i \geq 1$  происходит либо уход сервера в режим пониженного энергопотребления (если очередь пуста), либо уход заявки из очереди на обслуживание, т. е.  $P_i(x, x') = 1$ , если  $x_0 = 0$ ,  $x'_i = 0$ , либо  $x_0 > 0$ ,  $x'_0 = x_0 - 1$ ,  $x'_i = 1$ , где  $x_j = x'_j, j \neq i$ . При этом непрерывная компонента  $T_0$  в момент прихода выбирается из распределения  $a(x)$ , а компонента  $T_i$  в момент окончания очередного периода активности (бездействия) сервера  $i$  выбирается из распределения  $b(x)$ , если  $x_0 > 0$ , и  $c(x)$  иначе. Зафиксируем состояния дискретного многомерного процесса, например, следующим образом:  $x^* = (k, 1, \dots, 1)$ . Произведем экспоненциальное расщепление плотностей  $T_0, \dots, T_m$  в моменты попадания дискретных компонент в зафиксированное состояние. Пусть  $Z_0, \dots, Z_m$  есть соответствующие фазы каждой непрерывной компоненты, определенные ранее. Если все фазы экспоненциальные, а дискретные компоненты попали в состояние  $x^*$ , то наступил момент регенерации.

Поясним при этом технику построения циклов регенерации методом дискретно-событийного моделирования. Применение данного метода обусловлено свойствами модели, поскольку изменение состояния системы происходит только в моменты прихода и ухода клиентов, а также окончания периодов бездействия, в то время как между данными событиями непрерывные компоненты линейно убывают.

Предположим, что изначально цепь стартует в момент регенерации. Выберем ближайшее по времени событие  $t^*$  и изменим состояние системы, пусть оно связано с обнулением компоненты  $T_i(t^*-) = 0, i \geq 0$ . Если при этом  $X_i(t^*+) = x^*$ , то можно провести экспоненциальное расщепление плотности *нового* времени  $T_i(t^*+)$ . Для этого разыграем бернуллиевскую случайную величину  $\xi_i$  с соответствующей вероятностью успеха. Если  $\xi_i = 1$ , начнем отсчет времени предэкспоненциальной фазы, положив  $Z_i = 1$  и временно положив  $T_i(t^*+) = x_{0,ixx'}$ . Далее, как только истечет предэкспоненциальная фаза, величину  $T_i$  можно будет каждый раз разыгрывать заново из классического экспоненциального распределения с параметром  $\theta_i(x, x')$ , полагая  $Z_i = 0$ . Если же  $\xi_i = 0$ , то разыграем  $T_i(t^*+)$  из распределения  $f_{1,i}(x, x')$  и полагая  $Z_i = 2$ . Дальнейшая эволюция системы связана с переходом к следующему ближайшему событию. Наконец за

метим, что если в момент  $t^*$  дискретная компонента не попала в  $x^*$ , то новое время  $T_i(t^*+)$  выбирается из исходного распределения ( $a, b, c$  соответственно), и полагаем  $Z_i = 2$ . Более подробное пояснение алгоритма представлено в работе [8]. Отметим, что на построенных таким образом циклах регенерации можно применять технику регенеративного оценивания, рассмотренную ниже, например, для оценивания некоторой ценовой функции (энергопотребления, среднего числа заявок и т. п.).

## МЕТОД РЕГЕНЕРАТИВНЫХ ОГИБАЮЩИХ

В этом разделе мы рассмотрим метод регенеративных огибающих для построения моментов регенерации высокопроизводительного вычислительного кластера. Суть метода заключается в следующем. Предположим, что описывающий исходную систему случайный процесс не является регенерирующим, либо число моментов слишком мало для эффективного оценивания. Тогда для исходной системы  $\Sigma$  строятся две новые системы: мажорантная система  $\bar{\Sigma}$  и минорантная система  $\underline{\Sigma}$ , входные потоки в которые совпадают с входным потоком исходной системы  $\Sigma$ . В мажорантной (минорантной) системе времена обслуживания стохастически больше (меньше), чем в исходной. Обе новые системы регенерируют в классическом смысле. Тогда оценки характеристик (например, средней загрузки, средней длины очереди, среднего числа заявок), получаемых в этих системах, являются верхней (оценка, получаемая в мажорантной системе) и нижней (оценка, получаемая в минорантной системе) границами для соответствующих характеристик исходной системы. Соответствующие величины в новых системах  $\bar{\Sigma}$  и  $\underline{\Sigma}$  будем также обозначать верхней и нижней чертой соответственно.

Рассмотрим модель высокопроизводительного вычислительного кластера  $\Sigma$  с  $m$  серверами, работающими параллельно. Обозначим  $T_n = t_{n+1} - t_n$  времена между приходами клиентов (с интенсивностью  $\lambda$ ), где  $t_n$  – момент прихода  $n$ -го клиента, которому требуется для обслуживания случайное время  $S_n$  (с интенсивностью  $\mu$ ) на  $N_n$  серверах *одновременно* с распределением  $\{p_k := P(N_i = k)\}$ . При этом  $n$ -й клиент занимает  $N_n$  наименее загруженных серверов. Обозначим  $\nu_n(Q_n)$  число заявок в кластере (размер очереди) в момент времени  $t_n^-$ . Условия стационарности для  $\nu, Q$  приведены в работах [24, 30].

Построим для кластера две новые системы  $\bar{\Sigma}$  (мажорантную) и  $\underline{\Sigma}$  (минорантную), с таким же (используя каплинг) входным по-

током  $\{T_i\}$ , как в системе  $\Sigma$ , и соответственно увеличенными (уменьшенными) временами обслуживания  $\{\bar{S}_i\}$  ( $\{\underline{S}_i\}$ ),  $i \geq 1$ . Такие изменения времен обслуживания происходят в моменты, когда мажорантная (минорантная) система попадает в некоторое фиксированное состояние, приводящее к регенерации в классическом смысле. Это позволяет построить доверительные оценки, даже если исходная система не являлась регенерирующей. (Более подробно этот метод представлен в [26].)

Метод базируется на использовании свойства монотонности процессов [25]. Предположим, что времена обслуживания стохастически упорядочены,  $\underline{S}_i \leq S_i \leq \bar{S}_i$ ,  $i \geq 1$ . Тогда

$$\underline{\nu}_i \leq \nu_i \leq \bar{\nu}_i, \underline{Q}_i \leq Q_i \leq \bar{Q}_i, i \geq 1. \quad (22)$$

Для начала определим для мажорантной системы  $\bar{\Sigma}$  моменты регенерации. Рассмотрим процесс  $\bar{Z}_n := \{\bar{\nu}_n, \bar{S}_i(n), i \in \bar{M}_n\}$ ,  $n \geq 1$ , где  $\bar{M}_n = \{i : t_i \leq t_n < \bar{z}_i\}$  – набор номеров заявок, обслуживаемых в  $\bar{\Sigma}$  в момент времени  $t_n$ ,  $\bar{z}_i$  – момент ухода  $i$ -й заявки из системы,  $\bar{S}_i(n)$  – остаточное время обслуживания  $i$ -й заявки в момент  $t_n$ . Фиксированное состояние определим двумя следующими способами.

1 способ: зафиксируем целое  $Q_0$  (размер очереди) и  $N_0$  (число серверов, требуемых для обслуживания первой в очереди заявки), константы  $0 \leq a \leq b < \infty$  и определим моменты

$$\begin{aligned} \bar{\beta}_{n+1} = \inf \{ k > \bar{\beta}_n : \bar{Q}_k = Q_0 > 0, \\ N_{k-Q_0} = N_0, \\ \bar{S}_i(k) \in (a, b), i \in \bar{M}_k, \\ \bar{M}_{\bar{\beta}_n} \cap \bar{M}_k = \emptyset \}, n \geq 0. \end{aligned} \quad (23)$$

2 способ: зафиксируем число занятых серверов  $m_0$  при свободной очереди и определим моменты

$$\begin{aligned} \bar{\beta}_{n+1} = \inf \{ k > \bar{\beta}_n : \bar{Q}_k = 0, \\ \sum_{i \in \bar{M}_k} N_i = m_0 \leq m, \\ \bar{S}_i(k) \in (a, b), i \in \bar{M}_k, \\ \bar{M}_{\bar{\beta}_n} \cap \bar{M}_k = \emptyset \}, n \geq 0. \end{aligned} \quad (24)$$

Условие  $\bar{M}_{\bar{\beta}_n} \cap \bar{M}_k = \emptyset$  означает, что все заявки, обслуживаемые в момент  $\bar{\beta}_n$ , покидают систему до момента  $\bar{\beta}_{n+1}$ .

Аналогично определяются моменты  $\underline{\beta}_n$ ,  $n \geq 0$ , для минорантной системы (константы  $a, b, Q_0, N_0, m_0$  могут быть другими).

Когда система попадает в фиксированное состояние, остаточные времена обслуживания

$\bar{S}_i(k), i \in M_k$ , заменяются на верхнюю границу компакта  $b$  (в мажорантной системе  $\bar{\Sigma}$ ), остаточные времена обслуживания  $S_i(k), i \in M_k$ , заменяются на нижнюю границу компакта  $a$  (в минорантной системе  $\underline{\Sigma}$ ). Построенные моменты времени  $\{\bar{\beta}_n\}, \{\underline{\beta}_n\}$  являются моментами регенерации, и кроме того, эти преобразования сохраняют свойство монотонности исследуемых характеристик.

Процесс  $\bar{Z}_n (Z_n)$  является марковским и в каждый момент времени  $\bar{\beta}_n^+, (\underline{\beta}_n^+)$  переходит в фиксированное состояние  $(\nu_0, b\mathbf{1}), ((\nu_0, a\mathbf{1}))$ . Таким образом, процесс  $\{\bar{Z}_n\}$  регенерирует в классическом смысле в моменты времени  $\{\bar{\beta}_k\}$  и имеет н.о.р. циклы регенерации  $\{\bar{Z}_k, \bar{\beta}_k \leq k < \bar{\beta}_{k+1}\}$  с н.о.р. длинами циклов регенерации  $\bar{\beta}_{k+1} - \bar{\beta}_k, n \geq 1$  (то же самое верно для минорантной системы).

Из (22) следует, что  $\underline{Q}_n \leq Q_n$  и

$$\underline{EQ} \leq EQ \leq \overline{EQ}.$$

Для построения моментов регенерации границы  $a, b$  для остаточных времен обслуживания можно брать разные  $a_i$  и  $b_i$  в зависимости от номера сервера  $i$ . Заметим, что, как и в методе обновляющих событий, для построения моментов слабой регенерации частота моментов регенерации регенеративных огибающих зависит от выбора констант  $a$  и  $b$ .

## РЕГЕНЕРАТИВНОЕ ОЦЕНИВАНИЕ

Предположим, что необходимо оценить некоторый функционал от исследуемого процесса (в дискретном времени)  $f(Z_n)$ . Например,  $Z_n$  – число клиентов в системе, а  $f(Z_n)$  – задержка в системе или ее энергопотребление.

Пусть  $Y_j = \sum_{i=\beta_j}^{\beta_{j+1}-1} f(Z_i)$  – сумма значений характеристики процесса на  $j$ -м цикле, для классически регенерирующего процесса, с.в.  $(Y_j)$  независимы и одинаково распределены. В силу (3) статистическое оценивание характеристики процесса сводится к следующему: при наличии независимых и одинаково распределенных наблюдений  $(Y_j, \alpha_j), j \geq 0$ , необходимо оценить величину

$$r = \frac{E_0 Y_1}{E_0 \alpha_1},$$

где, напомним,  $E_0$  есть математическое ожидание процесса без задержки.

Предположим, что  $E(Y_1 + \alpha_1)^2 < \infty$ , тогда на основе центральной предельной теоремы можно получить следующий  $(1 - 2\gamma)\%$  до-

верительный интервал для  $r$  [3]:

$$\left[ \bar{r}_n \pm \frac{h_\gamma \sqrt{\text{Var}(n)}}{\sqrt{n \bar{\alpha}_n}} \right], \quad (25)$$

где  $\bar{\alpha}_n (\bar{Y}_n)$  – выборочное среднее для  $\alpha_i (Y_i)$ ,  $n$  – число циклов регенерации,  $\text{Var}(n)$  – оценка  $\text{Var}(Y_1) - 2r \text{cov}(Y_1, \alpha_1) + r^2 \text{Var}(\alpha_1)$ ,  $h_\gamma = \Phi^{-1}((1 - \gamma)/2)$ ,  $\Phi(x)$  – функция Лапласа. Отметим, что данный метод применим также для слабо зависимых циклов регенерации. В этом случае доверительный интервал имеет вид:

$$\left[ \bar{r}_n \pm \frac{h_\gamma \sqrt{\text{Var}(n) + 2(t_1(n) - \bar{r}_n t_2(n))}}{\sqrt{n \bar{\alpha}_n}} \right], \quad (26)$$

где  $t_1(n)$  – выборочная оценка ковариации  $\text{cov}(Y_0, Y_1)$ ,  $t_2(n)$  – выборочная оценка ковариации  $\text{cov}(\alpha_1, Y_0)$ .

Аналогичная конструкция доступна и для цепей Маркова, положительно возвратных по Харрису. Обозначим  $E_\mu(\cdot) = \int E_x(\cdot) \mu(dx)$ ,  $E_x(\cdot) = E(\cdot | Z_0 = x)$ . Пусть  $\{Z_n\}_{n \geq 0}$  положительно возвратна по Харрису, т. е.  $E_\mu \beta_1 < \infty$ . Тогда [9]

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=0}^N f(Z_n) = \frac{E_\mu \sum_{n=0}^{\beta_1-1} f(Z_n)}{E_\mu \beta_1}.$$

Использование методов уменьшения дисперсии оценки, таких как, например, метод общих случайных чисел или метод противоположных случайных чисел [29], позволило бы управлять знаком ковариации в полученном доверительном интервале и, следовательно, сократить его длину. Они применимы при условии монотонности соответствующих функционалов, что требует дополнительного изучения и, с точки зрения авторов, является перспективным направлением дальнейших исследований.

## ЗАКЛЮЧЕНИЕ

В работе представлены методы регенеративного моделирования многосерверных систем обслуживания. Высокая сложность современных многосерверных систем обслуживания является не только вызовом, требующим доработки классических методов, но и возможностями, требующими исследовать новые методы оценивания и типы регенерации, в том числе за счет качественно новой информации о состоянии системы (доступной, например, в программно конфигурируемых сетях [27]). При этом, с учетом дополнительной

информации о типах распределения управляющих последовательностей системы, представляется перспективным применение методов сокращения дисперсии оценки, методов ускоренного моделирования. Отметим также, что для ускорения моделирования возможно применение параллельных и распределенных вычислений. Исследование этих возможностей представляется важной темой для дальнейшей работы.

*Финансовое обеспечение исследований осуществлялось из средств федерального бюджета на выполнение государственного задания КарНЦ РАН (Институт прикладных математических исследований КарНЦ РАН) и при финансовой поддержке РФФИ (18-07-00147, 18-07-00156, 16-07-00622), гранта Президента РФ МК-1641.2017.1.*

## ЛИТЕРАТУРА

1. Боровков А. А. Теория вероятностей. М., 1972.
2. Иглхард Д. Л., Шедлер Д. С. Регенеративное моделирование сетей массового обслуживания. М.: Радио и связь, 1984.
3. Крайн М., Лемуан О. Введение в регенеративный метод анализа моделей. М.: Наука, 1982.
4. Морозов Е. В., Дельгадо Р. Анализ стационарности регенеративных систем обслуживания // Автоматика и телемеханика. 2009. Т. 70. С. 42–58.
5. Морозов Е. В., Румянцев А. С. Модели многосерверных систем для анализа вычислительного кластера // Труды Карельского научного центра РАН. 2011. Вып. 5. С. 75–85.
6. Нуммелин Э. Общие неприводимые цепи Маркова и неотрицательные операторы / Пер. с англ. М.: Мир, 1989.
7. Morozov E. V. Coupling and stochastic monotonicity of queueing processes. Петрозаводск: Изд-во ПетрГУ, 2013. 72 с.
8. Andronov A. Artificial regeneration points for stochastic simulation of complex systems // Simulation Technology: Science and Art. 10th European Simulation Symposium ESS'98, Proceedings. 1998. P. 34–40.
9. Asmussen S. Applied probability and queues. Springer-Verlag, New York, 2003. doi: 10.1007/b97236
10. Athreya K. B., Ney P. A new approach to the limit theory of recurrent Markov Chains // Transactions of the American Mathematical Society. 1978. Vol. 245. P. 493–501. doi: 10.1090/S0002-9947-1978-0511425-0
11. Borovkov A. A. Asymptotic Methods in Queueing Theory. Wiley, New York, 1984.
12. Charlot F., Ghidouche M., Hamami M. Irréductibilité et récurrence au sens de Harris des «Temps d'attente» des files GI/G/q // Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete. 1978. Vol. 43. P. 187–203.
13. Feller W. An introduction to probability theory and its applications. Wiley, New York, 1950.
14. Feitelson D. G. Workload Modeling for Computer Systems Performance Evaluation. Cambridge University Press, New York, 2015. doi: 10.1017/CBO9781139939690
15. Foss S. G. On the ergodicity conditions for stochastically recursive sequences // Queueing Systems. 1992. Vol. 12. P. 287–296. doi: 10.1007/BF01158804
16. Foss S., Kalashnikov V. Regeneration and renovation in queues // Queueing Systems. 1991. No. 8. P. 211–224. doi: 10.1007/BF02412251
17. Foss S., Konstantopoulos T. An overview of some stochastic stability methods // Journal of the Operations Research Society of Japan. 2004. Vol. 47, no. 4. P. 275–303.
18. Glynn P. Wide-sense regeneration for Harris recurrent Markov processes: an open problem // Queueing Systems. 2011. Vol. 68, no. 3–4. P. 305–311. doi: 10.1007/s11134-011-9238-x
19. Glynn P., Iglehart D. Simulation methods for queues: an overview // Queueing Systems. 1988. Vol. 3. P. 221–256. doi: 10.1007/BF01161216
20. Kalashnikov V. V. Regenerative queueing processes and their qualitative and quantitative analysis // Queueing Systems. 1990. Vol. 6. P. 113–136. doi: 10.1007/BF02411469
21. Kiefer J., Wolfowitz J. On the theory of queues with many servers // Transactions of the American Mathematical Society. 1955. P. 1–18. doi: 10.1090/S0002-9947-1955-0066587-3
22. Morozov E. Stochastic boundness of some queueing systems. Preprint No R-95-2022, ISSN 0908-1216, Dept. Math. and Computer Sci., Aalborg Univ., Aalborg, Denmark, 1995.
23. Morozov E. The tightness in the ergodic analysis of regenerative queueing processes // Queueing Systems. 1997. Vol. 27. P. 179–203. doi: 10.1023/A:1019114131583
24. Morozov E., Rumyantsev A. Stability Analysis of a MAP/M/s Cluster Model by Matrix-Analytic Method // Lecture Notes in Computer Science. 2016. Vol. 9951. P. 63–76. doi: 10.1007/978-3-319-46433-6\_5

25. Morozov E., Rumyantsev A., Peshkova I. Monotonicity and stochastic bounds for simultaneous service multiserver systems // 8th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops. Institute of Electrical and Electronics Engineers. 2016. P. 294–297. doi: 10.1109/ICUMT.2016.7765374
26. Morozov E., Rumyantsev A., Nekrasova R., Peshkova I. A Regeneration-Based Estimation of High Performance Multiserver Systems // Communications in Computer and Information Science. 2016. Vol. 608. P. 271–282. doi: 10.1007/978-3-319-39207-3\_24
27. Nguyen T. A. et al. Availability Modeling and Analysis for Software Defined Networks // 2015 IEEE 21st Pacific Rim International Symposium on Dependable Computing (PRDC), Zhangjiajie. 2015. P. 159–168. doi: 10.1109/PRDC.2015.27
28. Nummelin E. A Splitting Technique for Harris Recurrent Markov Chains // Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete. 1978. Vol. 43. P. 309–318. doi: 10.1007/BF00534764
29. Ross S. Simulation. Academic Press, New York, 1997.
30. Rumyantsev A., Morozov E. Stability criterion of a multiserver model with simultaneous service // Annals of Operations Research. 2017. Vol. 252. No. 1. P. 29–39. doi: 10.1007/s10479-015-1917-2
31. Shedler G. Regeneration and networks of queues. Springer-Verlag, 1987. doi: 10.1007/978-1-4612-1050-4
32. Sigman K. Queues as Harris recurrent Markov chains // Queueing Systems. 1988. Vol. 3. P. 179–198. doi: 10.1007/BF01189048
33. Sigman K., Wolff R. W. A review of regenerative processes // SIAM Review. 1993. Vol. 35, no. 2. P. 269–288. doi: 10.1137/1035046
34. Sutter H. The free lunch is over: A fundamental turn toward concurrency in software // Dr. Dobbs's Journal. 2005. P. 1–9.
35. Kalashnikov V. Topics on regenerative processes. CRC Press, Boca Baton, 1994.
36. Thorrisson H. Coupling, stationarity, and regeneration. Springer-Verlag, New York, 2000. doi: 10.1007/978-1-4612-1236-2
37. Whitt W. Comparing counting processes and queues // Advances in Applied Probability. 1981. Vol. 13, no. 1. P. 207–220. doi: 10.2307/1426475
38. Whitt W. Embedded renewal processes in the GI/G/s queue // Journal of Applied Probability. 1972. Vol. 9. P. 650–658. doi: 10.1017/S0021900200035944

Поступила в редакцию 10.04.2018

## REFERENCES

1. Borovkov A. A. Probability Theory. Moscow, 1972 (in Russian).
2. Iglehart D., Shedler G. Regenerative simulation of response times in networks of queues. Moscow: Radio i Svyaz, 1984 (in Russian).
3. Crane M. A., Lemoine A. J. An Introduction to the Regenerative Method for Simulation Analysis. Moscow: Nauka, 1982 (in Russian).
4. Morozov E. V., Delgado R. Stability analysis of regenerative queueing systems. *Avtomatika i telemekhanika*. 2009. Vol. 70. P. 42–58 (in Russian).
5. Morozov E. V., Rumyantsev A. S. Multiserver system models for high performance cluster analysis. *Transactions of Karelian Research Centre of RAS*. 2011. Vol. 5. P. 75–85 (in Russian).
6. Nummelin E. General Irreducible Markov Chains and Non-Negative Operators. Moscow: Mir, 1989 (in Russian).
7. Morozov E. V. Coupling and stochastic monotonicity of queueing processes. Petrozavodsk: PetrSU, 2013.
8. Andronov A. Artificial regeneration points for stochastic simulation of complex systems. *Simulation Technology: Science and Art. 10th European Simulation Symposium ESS'98, Proceedings*. 1998. P. 34–40.
9. Asmussen S. Applied probability and queues. Springer-Verlag, New York, 2003. doi: 10.1007/b97236
10. Athreya K. B., Ney P. A new approach to the limit theory of recurrent Markov Chains. *Transactions of the American Mathematical Society*. 1978. Vol. 245. P. 493–501. doi: 10.1090/S0002-9947-1978-0511425-0
11. Borovkov A. A. Asymptotic Methods in Queueing Theory. Wiley, New York, 1984.
12. Charlot F., Ghidouche M., Hamami M. Irréductibilité et récurrence au sens de Harris des «Temps d'attente» des files GI/G/q. *Zeitschrift*

für Wahrscheinlichkeitstheorie und verwandte Gebiete. 1978. Vol. 43. P. 187–203.

13. Feller W. An introduction to probability theory and its applications. Wiley, New York, 1950.
14. Feitelson D. G. Workload Modeling for Computer Systems Performance Evaluation. Cambridge University Press, New York, 2015. doi: 10.1017/CBO9781139939690
15. Foss S. G. On the ergodicity conditions for stochastically recursive sequences. *Queueing Systems*. 1992. Vol. 12. P. 287–296. doi: 10.1007/BF01158804
16. Foss S., Kalashnikov V. Regeneration and renovation in queues. *Queueing Systems*. 1991. No. 8. P. 211–224. doi: 10.1007/BF02412251
17. Foss S., Konstantopoulos T. An overview of some stochastic stability methods. *Journal of the Operations Research Society of Japan*. 2004. Vol. 47, no. 4. P. 275–303.
18. Glynn P. Wide-sense regeneration for Harris recurrent Markov processes: an open problem. *Queueing Systems*. 2011. Vol. 68, no. 3-4. P. 305–311. doi: 10.1007/s11134-011-9238-x
19. Glynn P., Iglehart D. Simulation methods for queues: an overview. *Queueing Systems*. 1988. Vol. 3, P. 221–256. doi: 10.1007/BF01161216
20. Kalashnikov V. V. Regenerative queueing processes and their qualitative and quantitative analysis. *Queueing Systems*. 1990. Vol. 6. P. 113–136. doi: 10.1007/BF02411469
21. Kiefer J., Wolfowitz J. On the theory of queues with many servers. *Transactions of the American Mathematical Society*. 1955. P. 1–18. doi: 10.1090/S0002-9947-1955-0066587-3
22. Morozov E. Stochastic boundness of some queueing systems. Preprint No R-95-2022, ISSN 0908-1216, Dept. Math. and Computer Sci., Aalborg Univ., Aalborg, Denmark, 1995.
23. Morozov E. The tightness in the ergodic analysis of regenerative queueing processes. *Queueing Systems*. 1997. Vol. 27. P. 179–203. doi: 10.1023/A:1019114131583
24. Morozov E., Rumyantsev A. Stability Analysis of a MAP/M/s Cluster Model by Matrix-Analytic Method. *Lecture Notes in Computer Science*. 2016. Vol. 9951. P. 63–76. doi: 10.1007/978-3-319-46433-6\_5
25. Morozov E., Rumyantsev A., Peshkova I. Monotonocity and stochastic bounds for simultaneous service multiserver systems. *8th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops. Institute of Electrical and Electronics Engineers*. 2016. P. 294–297. doi: 10.1109/ICUMT.2016.7765374
26. Morozov E., Rumyantsev A., Nekrasova R., Peshkova I. A Regeneration-Based Estimation of High Performance Multiserver Systems. *Communications in Computer and Information Science*. 2016. Vol. 608. P. 271–282. doi: 10.1007/978-3-319-39207-3\_24
27. Nguyen T. A. et al. Availability Modeling and Analysis for Software Defined Networks. *2015 IEEE 21st Pacific Rim International Symposium on Dependable Computing (PRDC), Zhangjiajie*. 2015. P. 159–168. doi: 10.1109/PRDC.2015.27
28. Nummelin E. A Splitting Technique for Harris Recurrent Markov Chains. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*. 1978. Vol. 43. P. 309–318. doi: 10.1007/BF00534764
29. Ross S. Simulation. Academic Press, New York, 1997.
30. Rumyantsev A., Morozov E. Stability criterion of a multiserver model with simultaneous service. *Annals of Operations Research*. 2017. Vol. 252, no. 1. P. 29–39. doi: 10.1007/s10479-015-1917-2
31. Shedler G. Regeneration and networks of queues. Springer-Verlag, 1987. doi: 10.1007/978-1-4612-1050-4
32. Sigman K. Queues as Harris recurrent Markov chains. *Queueing Systems*. 1988. Vol. 3. P. 179–198. doi: 10.1007/BF01189048
33. Sigman K., Wolff R. W. A review of regenerative processes. *SIAM Review*. 1993. Vol. 35, no. 2. P. 269–288. doi: 10.1137/1035046
34. Sutter H. The free lunch is over: A fundamental turn toward concurrency in software. *Dr. Dobbs's Journal*. 2005. P. 1–9.
35. Kalashnikov V. Topics on regenerative processes. CRC Press, Boca Baton, 1994.
36. Thorrisson H. Coupling, stationarity, and regeneration. Springer-Verlag, New York, 2000. doi: 10.1007/978-1-4612-1236-2
37. Whitt W. Comparing counting processes and queues. *Advances in Applied Probability*. 1981. Vol. 13, no. 1. P. 207–220. doi: 10.2307/1426475
38. Whitt W. Embedded renewal processes in the GI/G/s queue. *Journal of Applied Probability*. 1972. Vol. 9. P. 650–658. doi: 10.1017/S0021900200035944

Received April 10, 2018

## СВЕДЕНИЯ ОБ АВТОРАХ:

### **Пешкова Ирина Валерьевна**

и.о. зав. кафедрой прикладной математики  
и кибернетики, к. ф.-м. н.  
Петрозаводский государственный университет  
пр. Ленина, 33, Петрозаводск, Республика Карелия,  
Россия, 185910  
эл. почта: iaminova@petsu.ru  
тел.: (8142) 719606

### **Румянцев Александр Сергеевич**

научный сотрудник, к. ф.-м. н.  
Институт прикладных математических  
исследований КарНЦ РАН,  
Федеральный исследовательский центр  
«Карельский научный центр РАН»  
ул. Пушкинская, 11, Петрозаводск,  
Республика Карелия, Россия, 185910  
эл. почта: ar0@krc.karelia.ru  
тел.: (8142) 763370

## CONTRIBUTORS:

### **Peshkova, Irina**

Petrozavodsk State University  
33 Lenin St., 185910 Petrozavodsk, Karelia, Russia  
e-mail: iaminova@petsu.ru  
tel.: (8142) 719606

### **Rumyantsev, Alexander**

Institute of Applied Mathematical Research,  
Karelian Research Centre, Russian Academy of Sciences  
11 Pushkinskaya St., 185910 Petrozavodsk,  
Karelia, Russia  
e-mail: ar0@krc.karelia.ru  
tel.: (8142) 763370