

# Как отпугнуть клиента и получить прибыль: фильтрация в системе массового обслуживания

София Сурова<sup>1</sup>, Кирилл Фурманов<sup>2</sup>

<sup>1</sup> Национальный исследовательский университет “Высшая школа экономики”,  
г. Москва, Россия

<sup>2</sup> Центральный экономико-математический институт РАН  
г. Москва, Россия

## Информация о статье

Поступила в редакцию:  
02.03.2024

Принята  
к опубликованию:  
23.04.2024

УДК 303.4

JEL C60

## Аннотация

Рассматривается система массового обслуживания с ограничением на длину очереди и разнородным потоком входящих заявок (клиентов), которые делятся на два типа. Заявки первого типа терпеливы: они встанут в очередь, если это позволяет ёмкость системы и дождутся обслуживания. Заявки второго типа нетерпеливы: они отказываются вставать в очередь, если время ожидания обслуживания оказывается слишком велико. Показывается, что если заявки второго типа приносят меньшую прибыль, чем заявки первого типа, то увеличение прибыли может достигаться за счёт замедления обслуживания, отпугивающего нетерпеливых клиентов.

## How to Scare a Customer Away and Get Profit: Filtration in a Queueing System

Sofia V. Surova, Kirill K. Furmanov

## Abstract

We consider a simple exponential queueing system with a finite capacity and heterogeneous customers. Type I customers are patient, they join the queue if the system capacity allows it. Type II customers are impatient, so that they refuse to join the queue if the waiting time is too large (so-called wait-based balking). The system makes profit servicing customers, and the profit depends on the type of the customer. We show that if type II customers bring less profit than type I customers then the administration may, in certain cases,

## Ключевые слова:

фильтрация, система  
массового обслуживания,  
нетерпеливые заявки

---

## Keywords:

filtration, queueing systems,  
balking

DOI: <https://doi.org/10.24866/2311-2271/2024-1/1113>

*increase profit by slowing down the service (decreasing the service rate). It makes the system unattractive for relatively unprofitable type II customers who stop forming the queue and thus leave the place for more profitable customers who otherwise would not have joined the queue due to capacity restriction.*

### **Введение**

В ряде экономических моделей особое внимание уделяется фильтрации — действиям экономических агентов, косвенным образом препятствующих неблагоприятному отбору. Так, работодатель может поставить барьер для трудоустройства в виде требования образования, которое может не пригодиться на работе непосредственно, но благодаря этому требованию отсеиваются малопродуктивные кандидаты, для которых издержки получения образования относительно высоки [1–3]. Как правило, фильтрация моделируется в рамках теории игр: работодатель делает ходы так, чтобы выигрышной стратегией желательных кандидатов было поведение, отличающее их от нежелательных.

Настоящая статья представляет попытку описания аналогичного процесса в рамках теории массового обслуживания. Мы рассматриваем систему, в которую поступают заявки (клиенты) двух типов, причём заявки первого типа терпеливы и встают в очередь, если это позволяет ёмкость системы, а заявки второго типа нетерпеливы: они отказываются вставать в очередь, если предполагаемое время ожидания в очереди слишком велико. Мы показываем, что, если заявки второго типа приносят меньшую прибыль, чем заявки первого типа, администрация системы в некоторых случаях может увеличивать прибыль за счёт замедления обслуживания и искусственного создания очереди. При этом нетерпеливые и относительно невыгодные клиенты будут “фильтроваться” — отказываться присоединяться к очереди, освобождая место для относительно выгодных и терпеливых.

### **Описание системы**

Рассмотрим механизм фильтрации в простейшей системе массового обслуживания (СМО) — системе M/M/1/2 согласно нотации Кендалла. Это система с простейшим потоком входящих заявок, показательно распределённым временем обслуживания, одним каналом обслуживания и ёмкостью в две заявки (т.е. очередь вмещает не более одной заявки, в это время вторая заявка обслуживается). Система представлена схематически на рис. 1.

Каждая из поступающих заявок независимо от прочих принадлежит к одному из двух типов. Заявки первого типа при поступлении входят в систему, если очередь не переполнена. Каждая такая заявка приносит единицу прибыли (допустима также случайная величина прибыли — тогда её математическое ожидание равно единице). Заявки второго типа отличаются нетерпеливостью типа balking [4–5]: они входят в систему, если это позволяет ёмкость — т.е. очередь не переполнена — и, если математическое ожидание времени нахождения в очереди не превышает допустимый для этих заявок предел  $\theta$ . Каждая за-



данием  $\rho$ . Таким образом, нетерпеливые заявки будут присоединяться к очереди в случае  $\rho \leq \theta$  и уйдут без обслуживания при  $\rho > \theta$ .

*Случай 1:  $\rho \leq \theta$ .* Все заявки входят в систему, если позволяет ёмкость. Это классическая система M/M/1/2 с хорошо известным распределением числа заявок (см. например [6, 7]):

$$p_0 = p_1 = p_2 = \frac{1}{3}, \rho = 1;$$

$$p_j = \frac{(1 - \rho)\rho^j}{1 - \rho^3}, j = 0, 1, 2, \rho \neq 1.$$

Здесь  $p_j$  — вероятность пребывания системы в состоянии  $j$  в стационарном режиме, совпадает с долей времени в долгосрочном периоде, в течение которого в системе находится  $j$  заявок.

В дальнейшем удобнее будет использовать единое выражение для вероятностей, доступное, в нашем случае, благодаря точно известной ёмкости системы в две заявки:

$$p_j = \frac{\rho^j}{1 + \rho + \rho^2}, j = 0, 1, 2. \tag{1}$$

*Случай 2:  $\rho > \theta$ .* Заявки второго типа не соглашаются стоять в очереди и входят в систему, только когда она свободна. Это также вариант процесса размножения и гибели, но в отличие от классической системы M/M/1/2 здесь интенсивности “рождений” зависят от текущего состояния системы — графы интенсивностей переходов для обоих случаев представлены на рис. 2. Благодаря нормировке  $\lambda = 1$ , интенсивность обслуживания (“гибели” заявки) обратна приведённой интенсивности входящего потока  $\rho$ .



Рис. 2. Графы интенсивностей переходов между состояниями СМО в случаях (а)  $\rho \leq \theta$  и (б)  $\rho > \theta$

Пользуясь выражениями для стационарного распределения процессов размножения и гибели, получаем:

$$p_1 = p_0\rho, p_2 = p_1s\rho = p_0s\rho^2.$$

Найдём  $p_0$  из условия  $p_0 + p_1 + p_2 = 1$ :

$$p_0(1 + \rho + s\rho^2) = 1,$$

$$p_0 = \frac{1}{1 + \rho + s\rho^2} \quad (2)$$

Отсюда получаем стационарные вероятности остальных состояний:

$$p_1 = \frac{\rho}{1 + \rho + s\rho^2}, \quad (3)$$

$$p_2 = \frac{s\rho^2}{1 + \rho + s\rho^2}. \quad (4)$$

### Ожидаемая прибыль

Прибыль приносят только те заявки, которые входят в систему — не теряются из-за ограничения ёмкости или нетерпения. Как следствие, заявки первого типа приносят прибыль тогда, когда по прибытии застают систему без очереди (0 или 1 заявка в системе). По свойству PASTA простейшего потока (Poisson Arrivals See Time Averages — см., например [7]) это происходит с вероятностью  $p_0 + p_1$ , так что среднее число не потерянных заявок первого типа за единицу времени составляет  $\lambda s(p_0 + p_1)$  и совпадает с ожидаемой прибылью за единицу времени от заявок первого типа. При  $\lambda = 1$  получаем выражение для этой прибыли:

$$\pi_1 = s(p_0 + p_1).$$

Индекс “1” здесь соответствует типу заявок.

Заявки второго типа обязательно приносят прибыль, если поступают в пустую систему (что происходит с вероятностью  $p_0$ ), и не приносят, если система заполнена (вероятность  $p_2$ ). Если канал обслуживания занят, но очереди нет, то заявка второго типа создаст очередь и принесёт прибыль при  $\rho \leq \theta$  и покинет систему при  $\rho > \theta$ . Учитывая, что каждая такая заявка приносит прибыль  $r$ , получаем выражение для ожидаемой прибыли от заявок второго типа за единицу времени:

$$\pi_2 = \begin{cases} r(1-s)(p_0 + p_1), & \theta \geq \rho; \\ r(1-s)p_0, & \theta < \rho. \end{cases}$$

Чтобы избежать ошибок при истолковании формул для ожидаемой прибыли, надо помнить, что вероятности  $p_0$ ,  $p_1$ ,  $p_2$  и сами зависят от соотношения “предела терпения”  $\theta$  и среднего времени обслуживания  $\rho$ .

*Случай 1:  $\rho \leq \theta$ .* Распределение числа заявок задаётся формулой (1), так что выражения для ожидаемой прибыли приобретают вид:

$$\pi_1 = s(p_0 + p_1) = \frac{s(1 + \rho)}{1 + \rho + \rho^2},$$

$$\pi_2 = r(1-s)(p_0 + p_1) = \frac{r(1-s)(1 + \rho)}{1 + \rho + \rho^2}.$$

Случай 2:  $\rho > \theta$ . Распределение числа заявок задаётся формулами (2)–(4). Выражения для ожидаемой прибыли:

$$\pi_1 = s(p_0 + p_1) = \frac{s(1 + \rho)}{1 + \rho + s\rho^2},$$

$$\pi_2 = r(1 - s)p_0 = \frac{r(1 - s)}{1 + \rho + s\rho^2}.$$

Общая прибыль за единицу времени естественно получается сложением прибылей от заявок первого и второго типа. Будем использовать для неё два обозначения:  $\pi_P$  для случая  $\rho \leq \theta$  (индекс  $P$  от слова “patience” — терпение, так как в этом случае все заявки терпеливо ждут в очереди) и  $\pi_{IP}$  для случая  $\rho > \theta$  (ImPatience — “нетерпение”). В обоих случаях суммарная прибыль выражается с помощью ранее введенных формул следующим образом:

$$\pi_P = \frac{(1 + \rho)(s + r - rs)}{1 + \rho + \rho^2}, \quad (5)$$

$$\pi_{IP} = \frac{s(1 + \rho) + r(1 - s)}{1 + \rho + s\rho^2}. \quad (6)$$

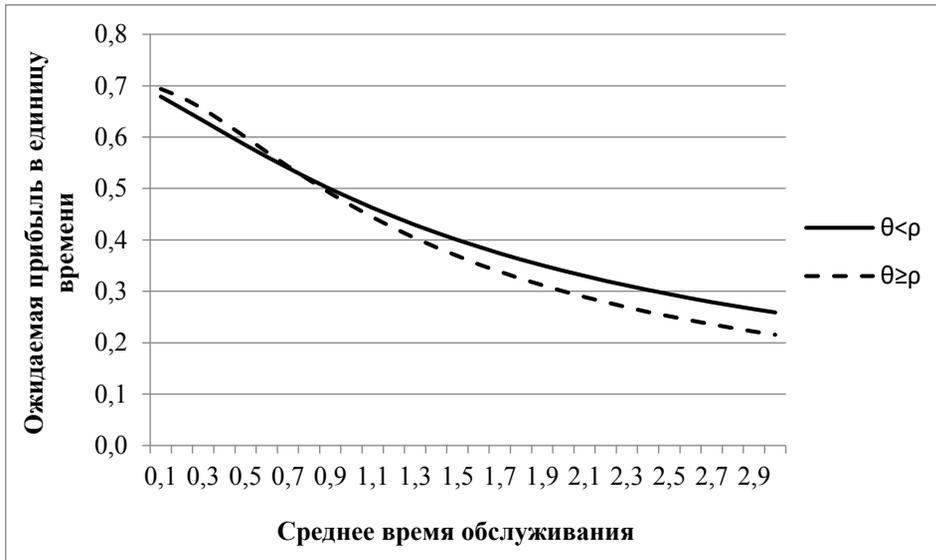
### Как отпугнуть клиента: прибыль и время обслуживания

Из формул (5) и (6) следует, что величины  $\pi_P$  и  $\pi_{IP}$  убывают с ростом  $\rho$  (для наглядности соответствующие графики приведены на рис. 3, использованные значения прочих параметров:  $r = 0,4$ ,  $p = 0,5$ ). Таким образом, ожидаемая прибыль отрицательно зависит от времени обслуживания на участках  $\rho \in (0; \theta]$  и  $\rho \in (\theta; +\infty)$ . Это естественно: увеличение времени обслуживания приводит к постоянной загруженности системы и росту доли потерянных заявок. Интересно другое: как видно из рис. 3, начиная со средней длительности обслуживания  $\rho \approx 0,8$  ожидаемая прибыль в “нетерпеливом” случае  $\rho > \theta$  оказывается выше. Таким образом, если предел ожидаемого времени в очереди для заявок второго типа превышает 0,8, а среднее время обслуживания немного меньше этого предела, то администрация системы может увеличить прибыль, замедлив обслуживание.

Действительно, хотя ожидаемая прибыль убывает на множествах  $\rho \in (0; \theta]$  и  $\rho \in (\theta; +\infty)$ , в точке  $\rho = \theta$  происходит скачок, который может как понизить, так и повысить прибыль (система переключается с прерывистой линии, изображённой на рис. 3, на сплошную). Этот скачок изображён на рис. 4, где для примера взято предельное время в очереди  $\theta = 2$  (прочие параметры — как и для рис. 3).

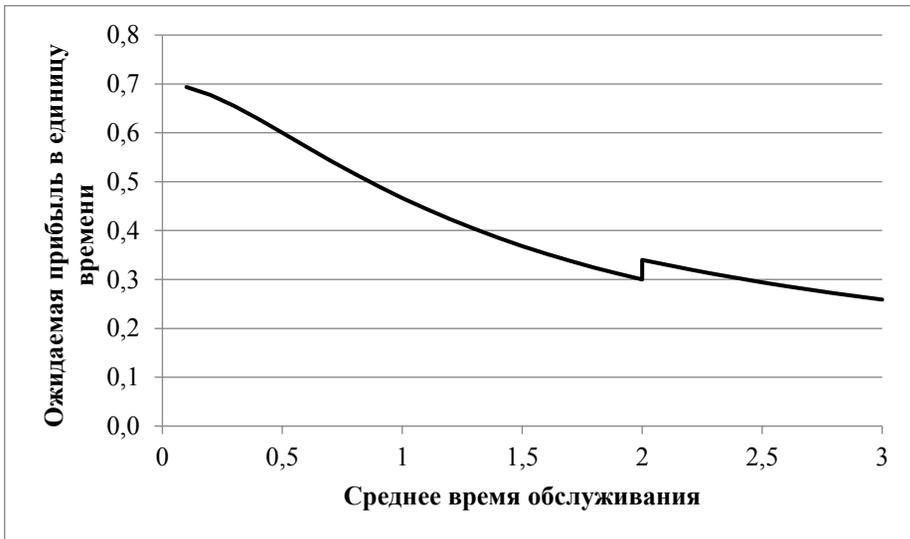
Когда среднее время обслуживания становится настолько велико, что заявки второго типа отказываются стоять в очереди, очередь освобождается для более прибыльных заявок первого типа, что может увеличить прибыль, если только заявки первого типа будут поступать до-

статочно интенсивно, чтобы компенсировать потерю нетерпеливых заявок.



Источник: расчёты авторов по формулам (5), (6).

Рис. 3. Ожидаемая прибыль за единицу времени в случаях  $\theta \geq \rho$  (прерывистая линия) и  $\theta < \rho$  (сплошная линия) в зависимости от среднего времени обслуживания.



Источник: расчёты авторов по формулам (5), (6).

Рис. 4. Скачок прибыли при увеличении времени обслуживания

### Обобщение для произвольной ёмкости системы

Введённое ранее фиксированное ограничение на длину очереди (не более одной заявки) удобно для рассмотрения примера фильтрации, но не обязательно. В настоящем разделе приведены более общие

результаты — основные формулы для системы М/М/1/К в нотации Кендалла с разнородными заявками.

Такая система может вместить  $K$  заявок, из которых одна будет обслуживаться, а остальные ждать в очереди. Пусть  $n$  — наибольшее число заявок в системе, при котором нетерпеливые заявки присоединяются к очереди. Из ранее сделанных предпосылок следует, что  $n = \min\left(K - 1, \left\lfloor \frac{\theta}{\rho} \right\rfloor\right)$ , где  $\left\lfloor \frac{\theta}{\rho} \right\rfloor$  — целая часть от деления.

Динамика системы описывается процессом размножения и гибели с состояниями  $0, \dots, K$ , где интенсивности “гибели” (обслуживания) постоянны и равны  $\rho^{-1}$ , а интенсивности “рождений” (поступления заявок) зависят от текущего состояния. Пока в системе не более  $n$  заявок, новые клиенты присоединяются с интенсивностью 1. Когда в системе более  $n$  заявок, присоединяются к очереди только терпеливые клиенты I типа, так что интенсивность притока новых заявок равна  $s$ . Как следствие, стационарные вероятности состояний связаны друг с другом формулами

$$\begin{aligned} p_j &= p_{j-1}\rho = p_0\rho^j, j \leq n + 1, \\ p_j &= p_{j-1}\rho s = p_n(\rho s)^{j-n-1} = p_0\rho^j s^{j-n-1}, n + 1 < j \leq K. \end{aligned}$$

В сумме вероятности всех состояний должны равняться единице, что позволяет выразить вероятность простоя системы  $p_0$ :

$$p_0 = \left( 1 + \sum_{j=1}^{n+1} \rho^j + \sum_{j=n+2}^K \rho^j s^{j-n-1} \right)^{-1}.$$

Следовательно, вероятность произвольного состояния  $j$  можно найти по следующей формуле:

$$p_j = \begin{cases} \frac{\rho^j}{1 + \sum_{j=1}^{n+1} \rho^j + \sum_{j=n+2}^K \rho^j s^{j-n-1}}, & j \leq n + 1, \\ \frac{\rho^j s^{j-n-1}}{1 + \sum_{j=1}^{n+1} \rho^j + \sum_{j=n+2}^K \rho^j s^{j-n-1}}, & n + 1 < j \leq K. \end{cases} \quad (7)$$

Заметим, что в ранее рассмотренном случае  $K = 2$  формула (7) сводится к формуле (1) при  $n = 1$  (т.е.  $\rho \leq \theta$ ) и к формулам (2)–(4) при  $n = 0$  ( $\rho > \theta$ ).

Суммы в формуле (7) имеют разное выражение в зависимости от того, являются ли слагаемые постоянными или образуют геометрическую прогрессию:

$$1 + \sum_{j=1}^{n+1} \rho^j = \begin{cases} n + 2, & \rho = 1, \\ \frac{1 - \rho^{n+2}}{1 - \rho}, & \rho \neq 1. \end{cases}$$

$$\sum_{j=n+2}^K \rho^j s^{j-n-1} = \begin{cases} (K-n-1)\rho^{n+1}, & \rho s = 1, \\ \frac{\rho^{n+2}(1-(\rho s)^K)}{1-\rho s}, & \rho s \neq 1. \end{cases}$$

Наконец, получим выражение для ожидаемой прибыли в единицу времени. Заявки первого типа приносят прибыль, если система не заполнена (в ней меньше  $K$  заявок). Они поступают с интенсивностью  $s$  заявок в единицу времени, и каждая приносит единичную прибыль, так что ожидаемая прибыль за единицу времени находится следующим образом:

$$\pi_1 = s \sum_{j=0}^{K-1} p_j = s(1 - p_K).$$

Заявки второго типа поступают с интенсивностью  $1 - s$ , приносят  $r$  единиц прибыли и входят в систему, если в ней не более  $n$  заявок, так что

$$\pi_2 = (1 - s)r \sum_{j=0}^n p_j.$$

Выражение для суммарной ожидаемой прибыли за единицу времени:

$$\pi = \pi_1 + \pi_2 = s(1 - p_K) + (1 - s)r \sum_{j=0}^n p_j. \quad (8)$$

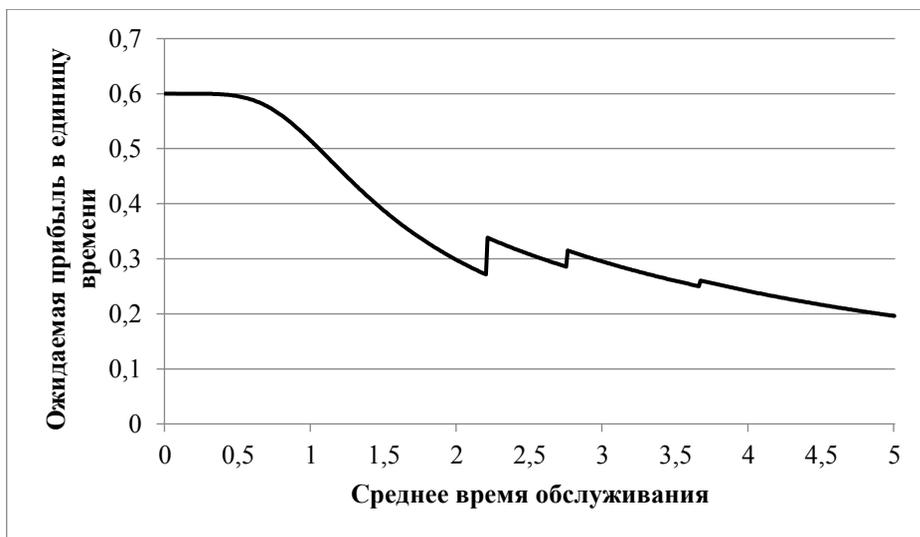


Рис. 5. Скачки прибыли в системе M/M/1/6

На рис. 5 представлена зависимость ожидаемой прибыли  $\pi$  от времени обслуживания  $\rho$  для случая  $K = 6, r = 0.2, s = 0.5, \theta = 11$ . Обратим внимание, что скачков в этом случае несколько. При увеличении времени обслуживания сначала “отпугиваются” те клиенты II типа, которые застают систему почти заполненной (пять заявок из мак-

симальных шести), затем к системе перестают присоединяться те нетерпеливые клиенты, кто застаёт четыре заявки и так далее — в каждом случае прибыль совершает скачок.

### **Заключение**

Как следует из примеров, изображённых на рис. 4 и 5, администрация системы массового обслуживания может увеличить прибыль за счёт фильтрации относительно неприбыльных заявок при увеличении времени обслуживания. Отметим характеристики системы, которые делают фильтрацию возможной.

1. Относительно невыгодные заявки должны быть нетерпеливы.

2. Если прибыль, которую приносят заявки, положительна, то у системы должна быть ограниченная ёмкость — иначе все заявки первого типа будут обслуживаться вне зависимости от того, соглашаются ли заявки второго типа стоять в очереди. Ограничение ёмкости можно заменить нетерпеливостью заявок первого типа, но тогда они должны быть терпеливее невыгодных заявок второго типа.

3. Если прибыль, которую приносят нетерпеливые заявки, отрицательна, то выгодная фильтрация возможна и в системах с неограниченной очередью — “отпугивание” заявок второго типа само по себе увеличивает прибыль, при этом терпеливые заявки первого типа не будут теряться.

Мы рассматривали и другой вариант нетерпения — *reneging*, уход из очереди заявок, не дождавшихся обслуживания [8, 9], в отечественной литературе известность получил частный случай *reneging*, при котором каждая заявка имеет ограничение на время пребывания в очереди [10, 11]. При однородном входящем потоке заявок вывести стационарное распределение и ожидаемую прибыль легко, для разнородных заявок дело усложняется тем, что вероятность ухода заявки из очереди начинает зависеть не только от длины очереди, но и от её состава (числа нетерпеливых заявок). Мы не выводили аналитического решения, но имитационным моделированием обнаружили скачки прибыли, аналогичные приведённым на рис. 4.

На практике определить удачный момент для увеличения времени обслуживания вряд ли возможно, но на предлагаемую модель можно посмотреть с иной стороны: она показывает, почему ускорение обслуживания может иметь эффект ниже желаемого и даже приводить к снижению прибыли.

### *Список источников*

1. Stiglitz J.E. The Theory of “Screening”, Education, and the Distribution of Income // *The American Economic Review*. 1957. Vol. 65. No. 3. P. 283–300.
2. Аистов А.В. О фильтрующей роли образования в России // *Экономический журнал Высшей школы экономики*. 2009. Т. 13. № 3. С. 452–481.
3. Grubb W.N. Further tests of screening on education and observed ability // *Economics of Education Review*. 1993. Vol. 13. No. 2. P. 125–136.

4. Haight F.A. Queueing with balking // *Biometrika*. 1957. Vol. 44. No. 3. P. 360–369.
5. Liu L.Q. Service systems with balking based on queueing time // PhD Thesis. University of North Carolina at Chapel Hill. 2007. — URL: <https://research.tue.nl/en/publications/service-systems-with-balking-based-on-queueing-time>
6. Shortle J.F., Thompson J.M., Gross D. [et al.]. *Fundamentals of Queueing Theory*. — 5<sup>th</sup> ed. — Wiley, 2018.
7. Ross S.M. *Introduction to Probability Models*. — 10<sup>th</sup> ed. — Elsevier, 2010.
8. Pazzal A.I., Radas S. Comparison of customer balking and renegeing behavior to queueing theory predictions: An experimental study // *Computers & Operations Research*. 2008. Vol. 35. No. 8. P. 2537–2548.
9. Choudhury A., Medhi P. Balking and renegeing in multiserver Markovian queueing system // *International Journal of Mathematics in Operational Research*. 2011. Vol. 3. No. 4. P. 377–394.
10. Гнеденко Б.В., Коваленко И.Н. *Введение в теорию массового обслуживания*. — М.: Наука, 1966. — 432 с.
11. Кирпичников А.П. *Методы прикладной теории массового обслуживания*. — М.: ЛЕНАНД, 2018. — 224 с.

#### Сведения об авторах

**Сурова София Валерьевна**, студент факультета компьютерных наук, Национальный исследовательский университет “Высшая школа экономики”. 101000, Россия, г Москва, ул. Мясницкая, 20. E-mail: [svsurova@edu.hse.ru](mailto:svsurova@edu.hse.ru).

*Sofia V. Surova*, student, Faculty of Computer Science, National Research University Higher School of Economics. 20, Myasnitskaya St., Moscow, 101000, Russia. E-mail: [svsurova@edu.hse.ru](mailto:svsurova@edu.hse.ru).

**Фурманов Кирилл Константинович**, кандидат экономических наук, старший научный сотрудник отделения Эконометрики и прикладной статистики, Центральный экономико-математический институт РАН. 117418, Россия, г Москва, Нахимовский проспект 47. E-mail: [kfurmanov@hse.ru](mailto:kfurmanov@hse.ru).

*Kirill K. Furmanov*, PhD in Economics, Senior Researcher, Department of Econometrics and Applied Statistics, Central Economic and Mathematical Institute of the Russian Academy of Sciences. 47, Nakhimovsky pr., Moscow, 117418, Russia. E-mail: [kfurmanov@hse.ru](mailto:kfurmanov@hse.ru).

#### Ссылка для цитирования

Сурова С.В., Фурманов К.К. Как отпугнуть клиента и получить прибыль: фильтрация в системе массового обслуживания // *Известия Дальневосточного федерального университета*. 2024. № 1. С. 80–90. — DOI: <https://doi.org/10.24866/2311-2271/2024-1/1113>.