

Количественные оценки информационной чувствительности алгоритмов

В.А. Головешкин, В.Н. Петрушин, М.В. Ульянов

Аннотация. В статье рассматриваются количественные оценки информационной чувствительности алгоритмов по функции трудоемкости и особенности их применения. Предложена новая симметричная по плотности вероятностей количественная оценка информационной чувствительности. Приведены экспериментальные данные по относительным частотам значений трудоемкости для алгоритма поиска подстроки в строке, их аппроксимации функцией бета-распределения и результаты сравнительного анализа предложенной и существующих оценок.

Ключевые слова: алгоритмы, оценки алгоритмов, информационная чувствительность, количественные оценки информационной чувствительности.

Введение

При анализе ресурсной эффективности алгоритмов актуальным является получение результатов, позволяющих прогнозировать ресурсные затраты, требуемые алгоритмом, при решении различных задач в данной проблемной области. Одной из наиболее важных ресурсных функций алгоритма является функция трудоемкости, отражающая требования алгоритма к временному ресурсу механизма реализации. Функция трудоемкости в худшем и лучшем случаях определяется, соответственно, как максимальное и минимальное число базовых операций принятой модели вычислений, заданных алгоритмом на входах фиксированной длины [1].

При исследовании временной эффективности программных реализаций алгоритма функция трудоемкости ассоциирована со временем получения результата на определенных исходных данных. В этом аспекте результаты ресурсного анализа алгоритма должны обеспечивать возможность прогнозирования его трудоемкости как для разных длин входов, так

и для различных входов фиксированной длины. Общепринятым является подход, связанный с введением функции трудоемкости как функции длины входа алгоритма, причем прогнозирование по изменению длин входов осуществляется на основе трудоемкости в среднем, а прогнозирование на фиксированной длине входа — по трудоемкостям в лучшем и худшем случаях [1].

При решении задачи прогнозирования на фиксированной длине входа количественно-зависимые алгоритмы (класс N) [1] представляют собой наиболее благоприятный класс. Но, к сожалению, подавляющее большинство алгоритмов относятся к количественно-параметрическому классу (класс NPR) [1] и обладают ненулевым размахом варьирования трудоемкости при фиксированной длине входа. Для целого ряда алгоритмов на актуальных длинах входа значение этого размаха достаточно велико. Прогнозирование по трудоемкости в худшем случае (гарантированная оценка сверху) дает для этого класса почти всегда сильно завышенные результаты, а прогнозирование по трудоемкости в среднем не учитывает инфор-

мацию о размахе варьирования. При этом очевидно, что качество прогноза во многом определяется тем, насколько велико влияние различных входов фиксированной длины на трудоемкость, т. е. насколько велика информационная чувствительность исследуемого алгоритма [2] в рамках выбранной количественной оценки.

Помимо необходимости обеспечения качества прогнозирования в некоторых случаях, при разработке алгоритмического обеспечения программ для систем реального времени, возникает задача обеспечения их стабильности по времени. Под стабильностью по времени программной реализации алгоритма содержательно понимается слабая зависимость времени выполнения от исходных данных при фиксированной длине входа [1]. Для решения задачи рационального выбора алгоритма с учетом этого требования также необходима детальная информация о влиянии на его трудоемкость различных входов фиксированной длины. При этом в качестве оценки стабильности по времени программной реализации алгоритма может быть использована и количественная оценка его информационной чувствительности.

С точки зрения проблем прогнозирования временной эффективности программных реализаций алгоритмов и оценки их стабильности по времени представляет интерес задача введения и сравнительного анализа различных количественных оценок информационной чувствительности. Настоящая статья кратко излагает результаты наших предыдущих исследований в этой области [2, 3], которые мы подвергаем конструктивной критике и на основе которой вводится новая оценка информационной чувствительности и проводится анализ оценок с рекомендациями по их применению.

1. Терминология и обозначения

Опираясь, в основном, на работу [1], будем использовать далее следующую терминологию и обозначения, связанные с анализом ресурсной эффективности алгоритмов:

D — вход алгоритма A : конечное множество слов фиксированной длины в бинарном алфавите, задающее конкретную решаемую задачу;

$\lambda(D)$ — длина входа алгоритма: $D \rightarrow N$, целочисленная функция, в общем случае определяемая как мощность множества D : $\lambda(D) = |D| = n$;

$f_A(D)$ — трудоемкость алгоритма A на входе D , целочисленная функция, значение которой есть число базовых операций (в принятой модели вычислений), заданных алгоритмом A на входе D ;

D_n — множество входов алгоритма A , имеющих длину n : $D_n = \{D \mid \lambda(D) = n\}$;

$f_A^{\wedge}(n)$ — трудоемкость алгоритма в худшем случае на всех допустимых входах длины n , т. е. максимум $f_A(D)$ на множестве D_n ;

$f_A^{\vee}(n)$ — трудоемкость алгоритма в лучшем случае на всех допустимых входах длины n , т. е. минимум $f_A(D)$ на множестве D_n , при этом для всех классов алгоритмов всегда выполнено:

$$f_A^{\vee}(n) \leq f_A(D \in D_n) \leq f_A^{\wedge}(n). \quad (1)$$

2. Понятие информационной чувствительности

Впервые понятие *информационной чувствительности* алгоритма по трудоемкости введено М.В. Ульяновым и В.А. Головешкиным в [2]. Понятие «информационная чувствительность» отражает тот факт, что алгоритм задает различное число базовых операций принятой модели вычислений $f_A(D)$ на разных входах D , имеющих фиксированную длину n . Ключевым для содержательной интерпретации этого термина является выбор количественной оценки (меры), обладающей свойством сопоставимости, т. е. дающей возможность решения задач сравнения алгоритмов и их рационального выбора.

Мы констатируем информационную чувствительность алгоритма (по трудоемкости), если наблюдаем разброс значений трудоемкости на различных входах фиксированной длины. В общем случае причиной такой вариации является влияние значений элементов конкретного входа и/или их порядка, равно как и других содержательных особенностей входов при их фиксированной длине. Теоретические границы такого разброса на входах длины n задаются

значениями $f_A^{\wedge}(n)$ и $f_A^{\vee}(n)$, и в этих границах наблюдаемые значения трудоемкости могут быть аппроксимированы некоторым законом распределения или же на основе экспериментальных данных могут быть вычислены статистические точечные оценки распределения.

Таким образом, в качестве основной идеи построения количественных оценок информационной чувствительности рассматривается описание трудоемкости как дискретной ограниченной случайной величины. В целях решения задач прогнозирования и рационального выбора такие количественные оценки должны, очевидно, вводиться как функции длины входа алгоритма. В рамках данной статьи мы более детально описываем построение оценок для фиксированной длины входа.

3. Оценка информационной чувствительности по энтропии

Классической характеристикой хаотичности некоторой системы является энтропия — понятие, впервые введенное Клаузиусом в термодинамике в 1865 г. для определения меры необратимого рассеивания энергии. В теорию информации энтропия введена К. Шенноном [4] как мера случайности, мера неопределенности какого-либо испытания, имеющего разные исходы. Для дискретной случайной величины X , принимающей n значений с вероятностями $p_i, i = \overline{1, n}$, энтропия $H(X)$ вычисляется по формуле:

$$H(X) = - \sum_{i=1}^n p_i \log_2 p_i \quad (2)$$

и характеризует степень разнообразия состояний системы. Энтропия обращается в нуль лишь в том случае, если состояние системы полностью определено, и максимальна, когда все состояния равновероятны, что легко доказывается методом неопределенных множителей Лагранжа. Значение $H(X)$ характеризует среднюю неопределенность выбора одного состояния из ансамбля и зависит, в силу (2), только от вероятностей состояний.

Для непрерывных случайных величин неопределенность значений состояния системы связана с плотностью распределения вероятностей

этих значений — дифференциальным законом распределения $p(x)$. В связи с этим непрерывный аналог формулы (2) получил название относительной дифференциальной энтропии или просто дифференциальной энтропии [5]:

$$h(X) = - \int_{-\infty}^{\infty} p(x) \cdot \log_2 p(x) dx. \quad (3)$$

Значение $h(X)$ можно трактовать как среднюю неопределенность случайной величины X с произвольным законом распределения по сравнению со средней неопределенностью случайной величины с размахом варьирования, равным единице при равномерном распределении. Дифференциальная энтропия так же не зависит от конкретных значений случайной величины X . Если X является случайной величиной с ограниченной вариацией, то максимальная дифференциальная энтропия соответствует равномерной плотности распределения вероятностей [5].

Использование энтропийного подхода позволило получить ряд значимых результатов как в теории информации, так и целом ряде других областей, например, при анализе динамических систем [6, 7].

Однако значение энтропии зависит только от значений вероятности и не чувствительно к изменению положения моды и математического ожидания распределения. Таким образом, функции плотности, имеющие качественно различный вид, могут иметь одинаковую дифференциальную энтропию. Для иллюстрации этого факта в качестве примера рассмотрим случайную величину X , заданную бета-распределением. Функция плотности бета-распределения имеет вид [8]

$$p(x) = \begin{cases} k \cdot (x)^{\alpha-1} (1-x)^{\beta-1}, & x \in [0,1] \\ 0, & x \notin [0,1] \end{cases} \quad (4)$$

где нормирующая константа k определяется из известного представления бета-функции Эйлера через гамма-функцию

$$B(\alpha, \beta) = \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)},$$

$$\Gamma(s) = \int_0^{\infty} t^{s-1} e^{-t} dt.$$

Таким образом, $k = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)}$.

Качественное поведение функции плотности (4) достаточно разнообразно — от U -образного, уходящего в бесконечность на границах при $\alpha < 1, \beta < 1$ до унимодального при $\alpha > 1, \beta > 1$. При $\alpha = 1, \beta = 1$ бета-распределение совпадает со стандартным непрерывным равномерным распределением.

Утверждение 1. Значение дифференциальной энтропии $h(X)$ случайной величины, заданной бета распределением, имеет вид

$$h(\alpha, \beta) = -\frac{1}{\ln 2} \left(\ln \Gamma(\alpha + \beta) - \ln \Gamma(\alpha) - \ln \Gamma(\beta) + (\alpha - 1)\psi(\alpha) + (\beta - 1)\psi(\beta) - (\alpha + \beta - 2)\psi(\alpha + \beta) \right).$$

Доказательство. Значение дифференциальной энтропии $h(X)$ зависит, очевидно, от параметров α, β функции плотности распределения вероятностей (4), поэтому задача состоит в получении функциональной зависимости $h(\alpha, \beta)$. Подставляя (4) в (2), имеем:

$$h(\alpha, \beta) = -\int_0^1 k \cdot (x)^{\alpha-1} (1-x)^{\beta-1} \cdot \log_2(k \cdot (x)^{\alpha-1} (1-x)^{\beta-1}) dx.$$

Переходя к натуральному логарифму, получаем

$$h(\alpha, \beta) = -c \int_0^1 k \cdot (x)^{\alpha-1} (1-x)^{\beta-1} \cdot \ln(k \cdot (x)^{\alpha-1} (1-x)^{\beta-1}) dx,$$

$$c = \frac{1}{\ln 2}.$$

Данное выражение представимо в виде трех слагаемых

$$h(\alpha, \beta) = h_1(\alpha, \beta) + h_2(\alpha, \beta) + h_3(\alpha, \beta),$$

$$h_1(\alpha, \beta) = -c \left[\int_0^1 kx^{\alpha-1} (1-x)^{\beta-1} \ln k dx \right],$$

$$h_2(\alpha, \beta) = -c \left[\int_0^1 kx^{\alpha-1} (1-x)^{\beta-1} (\alpha-1) \ln x dx \right],$$

$$h_3(\alpha, \beta) = -c \left[\int_0^1 kx^{\alpha-1} (1-x)^{\beta-1} (\beta-1) \ln(1-x) dx \right].$$

Вычисляя $h_1(\alpha, \beta)$, с учетом значений констант k, c , имеем

$$h_1(\alpha, \beta) = -\frac{1}{\ln 2} [\ln \Gamma(\alpha + \beta) - \ln \Gamma(\alpha) - \ln \Gamma(\beta)].$$

Значение $h_2(\alpha, \beta)$ представимо через частную производную бета-функции

$$h_2(\alpha, \beta) = -c \cdot k \cdot (\alpha - 1) \left[\int_0^1 x^{\alpha-1} (1-x)^{\beta-1} \ln x dx \right] = -c \cdot k \cdot (\alpha - 1) \frac{\partial B(\alpha, \beta)}{\partial \alpha}.$$

Тогда, используя представление бета-функции через гамма-функцию, имеем

$$h_2(\alpha, \beta) = -\frac{1}{\ln 2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} (\alpha - 1) \left[\Gamma(\beta) \frac{\Gamma(\alpha)\Gamma(\alpha + \beta) - \Gamma(\alpha)\Gamma'(\alpha + \beta)}{\Gamma^2(\alpha + \beta)} \right] = -\frac{1}{\ln 2} (\alpha - 1) \left[\frac{\Gamma'(\alpha)}{\Gamma(\alpha)} - \frac{\Gamma'(\alpha + \beta)}{\Gamma(\alpha + \beta)} \right].$$

Воспользуемся далее определением пси-функции Эйлера $\psi(x)$ для представления выражения в квадратных скобках

$$\psi(x) = \frac{d}{dx} (\ln \Gamma(x))$$

и получим окончательное выражение:

$$h_2(\alpha, \beta) = -\frac{1}{\ln 2} (\alpha - 1) [\psi(\alpha) - \psi(\alpha + \beta)].$$

Рассуждая аналогично, получим выражение для функции $h_3(\alpha, \beta)$:

$$h_3(\alpha, \beta) = -c \cdot k \cdot (\beta - 1) \left[\int_0^1 x^{\alpha-1} (1-x)^{\beta-1} \ln(1-x) dx \right] = -c \cdot k \cdot (\beta - 1) \frac{\partial B(\alpha, \beta)}{\partial \beta} = -\frac{1}{\ln 2} (\beta - 1) [\psi(\beta) - \psi(\alpha + \beta)].$$

Суммируя полученные результаты, окончательно получаем искомое выражение для дифференциальной энтропии бета-распределения.

Конец доказательства.

Функция $h(\alpha, \beta)$ задает поверхность в трехмерном пространстве. Следовательно,

уравнение $h(\alpha, \beta) = \text{const}$ имеет в области допустимых значений функции бесконечное множество решений, являясь, по сути, уравнением эквипотенциали, во всех точках которой, т. е. при различных значениях α, β , значения дифференциальной энтропии совпадают.

Эти свойства не обеспечивают достаточной информативности энтропийной оценки информационной чувствительности алгоритмов. Отметим также, что максимум энтропии достигается на равномерном распределении [5], в то время как максимум дисперсии ограниченной случайной величины достигается при равновероятной реализации наибольшего и наименьшего значений [2]. Очевидно, что на основе экспериментальной гистограммы относительных частот трудоемкости можно вычислить энтропийную оценку информационной чувствительности по формуле (2). Однако надо соблюдать определенную осторожность, связанную с выбором числа сегментов гистограммы, который влияет на получаемую оценку.

Энтропийная оценка информационной чувствительности позволяет качественно оценить информационную чувствительность алгоритма, не локализуя при этом положение сегмента значений с наибольшими вероятностями.

4. Статистическая оценка информационной чувствительности

Статистическая оценка информационной чувствительности предложена М.В. Ульяновым и В.А. Головешкиным в [2] на основе следующих рассуждений. На множестве входов фиксированной длины трудоемкость алгоритма рассматривается как дискретная ограниченная случайная величина. Классической точечной мерой рассеяния случайной величины является σ — среднеквадратическое отклонение, которое оценивается по данным выборки через стандартное отклонение s [9]. Корректное определение количественной оценки информационной чувствительности должно учитывать также и длину сегмента варьирования. Это связано с тем, что при одинаковом значении дисперсии достоверно более чувствительным должен быть алгоритм с большей длиной сегмента

возможных значений трудоемкости. Для учета длины этого сегмента используется такое понятие математической статистики, как размах варьирования [9]. Очевидно, что теоретический размах варьирования трудоемкости алгоритма является функцией длины входа, тогда, вводя обозначение $R(n)$ с учетом (1), получаем:

$$R(n) = f_A^{\wedge}(n) - f_A^{\vee}(n).$$

На этой основе в [2] вводится понятие нормированного (относительного) размаха варьирования функции трудоемкости для входов длины n — $R_N(n)$ как отношение половины теоретического размаха варьирования к его середине

$$R_N(n) = \frac{f_A^{\wedge}(n) - f_A^{\vee}(n)}{f_A^{\wedge}(n) + f_A^{\vee}(n)}, \quad (5)$$

при этом значение $R_N(n) = 0$ соответствует ситуации, когда $f_A^{\wedge}(n) - f_A^{\vee}(n) = 0$, т. е. идентифицирует принадлежность алгоритма к классу N .

Одной из общепринятых точечных характеристик выборки является коэффициент вариации V , определяемый как отношение стандартного отклонения к среднему значению [9]. Для выборки, полученной на основе экспериментального исследования трудоемкости алгоритма, коэффициент вариации V также зависит от размерности и имеет вид:

$$V(n) = s_{f_A}(n) / \overline{f_A}(n), \quad 0 \leq V(n) \leq 1, \quad (6)$$

где $s_{f_A}(n)$ — стандартное отклонение трудоемкости, как дискретной ограниченной случайной величины, при фиксированной длине входа n , а $\overline{f_A}(n)$ — выборочное среднее, рассчитываемые по данным выборки. На основе этих рассуждений в [2] вводится *статистическая* количественная оценка (мера) информационной чувствительности алгоритма по трудоемкости, с обозначением $\delta_{IS}(n)$, в виде

$$\delta_{IS}(n) = V(n) \cdot R_N(n), \quad 0 \leq \delta_{IS}(n) \leq 1. \quad (7)$$

Поскольку оценка $\delta_{IS}(n)$ использует только статистические точечные оценки трудоемкости как случайной величины, ее применение возможно в случае отсутствия знаний о законе распределения значений трудоемкости или какой-

либо его аппроксимации. Таким образом, значения оценки $\delta_{IS}(n)$ могут быть получены на основе экспериментальных исследований алгоритма, по результатам которых вычисляется значение $V(n)$, и его теоретического анализа, необходимого для вычисления нормированного размаха варьирования по формуле (5).

Альтернативная оценка — оценка нормированного размаха варьирования по экспериментальным данным приводит к вычислению $R_N(n)$ по формуле

$$R_N(n) = \frac{f_{A \max}(n) - f_{A \min}(n)}{f_{A \max}(n) + f_{A \min}(n)},$$

т.е. на основе вычисления разности максимального и минимального значений трудоемкости в выборке. Этот подход обладает рядом особенностей, которые мы хотим отметить. Не вызывает сомнений неравенство $R_B \leq R_T$: размах варьирования выборки не может превышать размаха варьирования генеральной совокупности. Возникает вопрос о поведении стандартного отклонения *нормированной на размах варьирования* случайной величины при наращивании объема выборочной совокупности. Заметим, что при этом стартовый объем выборки должен быть репрезентативным, т.е. должен позволять почти достоверно установить вид распределения и его параметры. При дальнейшем увеличении объема репрезентативной выборки возможны следующие случаи.

1. В случае равномерного распределения стандартное отклонение и выборочное среднее для нормированной случайной величины не будут претерпевать существенных изменений.

2. Если наблюдаемое распределение U -образно или на одной из границ генеральной совокупности плотность вероятностей бесконечна, то с ростом объема выборки весьма вероятно нарастание численной оценки стандартного отклонения.

3. Чаще всего (а при исследовании чувствительности алгоритмов так и есть) плотность распределения при приближении к границам генеральной совокупности стремится к нулю, наращивание выборки приводит к появлению маловероятных значений случайной величины и

росту размаха варьирования, что в свою очередь вызывает уменьшение стандартного отклонения.

Поскольку статистические точечные оценки сами по себе являются случайными величинами [9], то оценка $\delta_{IS}(n)$ будет варьироваться в экспериментах, что является одним из ее недостатков. Формально этот недостаток устраним — достаточно при этом знать закон распределения случайной функции $\delta_{IS}(n)$ или найти ее достаточно надежную оценку. В силу конечности возможных значений функции трудоемкости такое распределение имеет конечные значения дисперсии и математического ожидания, т.е. интервальная оценка информационной чувствительности алгоритма по трудоемкости в принципе возможна. В случаях 1 и 3 интервальная оценка является достаточной.

Другим, и содержательно более важным, недостатком статистической количественной меры является то, что она не позволяет указать интервал значений трудоемкости при заданной вероятности, а, следовательно — не может быть адаптирована к требованиям разработчиков алгоритмического обеспечения, в частности, при прогнозировании стабильности по времени.

5. Квантильная оценка информационной чувствительности

Идея рассмотрения трудоемкости алгоритма при фиксированной длине входа как ограниченной случайной величины, аппроксимируемой некоторой известной функцией плотности распределения вероятностей и вычисления γ -квантиля этого закона распределения, привела к рассмотрению квантильной оценки информационной чувствительности алгоритмов [3]. Основная идея состоит в определении *длины сегмента* нормированных значений трудоемкости, по которому интеграл от функции плотности равен заданной вероятности (надежности) γ , отражающей требования разработчиков по информационной чувствительности, а по сути — особенности проблемной области применения алгоритма.

Содержательно такая количественная оценка с обозначением $\delta_{IQ}(\gamma)$ [3] есть доля теоретиче-

ского сегмента варьирования трудоемкости, в которой, с заданной вероятностью γ , будут наблюдаться значения трудоемкости алгоритма на произвольных входах фиксированной длины. В общем случае задача определения сегмента, по которому интегрируется заданная вероятность при известной функции плотности, имеет бесконечное множество решений [9]. Традиционное решение этой задачи — определение длины сегмента, задающего вероятность γ , симметрично относительно медианы распределения. При известной функции распределения вероятностей эта задача решается на основе вычисления $1/2 \pm \gamma/2$ -квантилей этого распределения. Такая оценка информационной чувствительности является характеристикой алгоритма для входов фиксированной длины. Очевидно, что в целях практического применения и теоретического исследования алгоритма необходимо рассматривать $\delta_{IQ}(\gamma)$ не только как функцию вероятности γ , но и как функцию размерности n , т.е. $\delta_{IQ} = \delta_{IQ}(\gamma, n)$.

Пусть функция $f(n, x), x \in [0, 1]$ есть непрерывная функция плотности распределения, аппроксимирующая нормированные значения трудоемкости, параметризованная аргументом n — длиной входа алгоритма. Для любой непрерывной не обращающейся в нуль на интервалах функции плотности распределения вероятностей $f(n, x)$ интегральная функция распределения вероятностей $F(n, x)$ является монотонно возрастающей, в силу чего имеет обратную функцию. Обозначим через $F^{-1}(n, x)$ функцию, обратную к $F(n, x)$, тогда согласно [10]

$$\delta_{IQ}(\gamma, n) = F^{-1}\left(n, \frac{1}{2} + \frac{\gamma}{2}\right) - F^{-1}\left(n, \frac{1}{2} - \frac{\gamma}{2}\right). \quad (8)$$

Интуитивно понятно, что с уменьшением разброса значений трудоемкости относительно математического ожидания, приводящего к уменьшению дисперсии, информационная чувствительность алгоритма по трудоемкости будет падать. Заметим, что использование медианы, а не математического ожидания в определении квантильной оценки информационной чувствительности связано с тем, что медиана и математическое ожидание необязательно совпа-

дают для несимметричных функций плотности распределения с ограниченной вариацией. Использование в этом случае значения математического ожидания как центра сегмента интегрирования функции плотности может привести для значений γ , близких к единице, к выходу границ этого сегмента за границы определения функции плотности с ограниченной вариацией.

Использование квантильной оценки позволяет разработчикам, варьируя значение вероятности γ , гибко оценивать информационную чувствительность, учитывая требования технического задания на разработку программ. Заметим, что значение $\gamma = 1$ приводит к значению $\delta_{IQ}(1, n) = 1 \forall n$, т.е. информационной чувствительности по полному нормированному размаху варьирования. Отметим также, что эта оценка не содержит информацию о положении γ -квантиля закона распределения на нормированном сегменте. Этот недостаток легко устраним, но мы сделаем это при построении следующей оценки.

6. Симметричная по плотности вероятностей оценка информационной чувствительности

На первый взгляд квантильная оценка информационной чувствительности позволяет довольно корректно решить задачу рационального выбора алгоритмов по критерию стабильности во времени с учетом требований разработчиков. Однако более детальное рассмотрение правила ее вычисления позволяет указать существенный недостаток этой оценки в случае, если аппроксимирующая функция плотности является унимодальной и асимметричной. Проблема состоит в том, что на границах сегмента

$$\left[F^{-1}\left(n, \frac{1}{2} - \frac{\gamma}{2}\right), F^{-1}\left(n, \frac{1}{2} + \frac{\gamma}{2}\right) \right],$$

значения которых определяются по формуле (8), функция плотности имеет различные значения (Рис. 1). Следовательно, при малых приращениях Δx мы отбрасываем различные вероятности, и квантильная мера в этом смысле не является симметричной — происходит вы-

брасывание из сегмента, определяемого по формуле (8), более вероятных величин в пользу менее вероятных.

Преодоление этого недостатка приводит к введению новой, уже симметричной по плотности вероятностей количественной оценке информационной чувствительности, которую мы будем далее обозначать через $\delta_{IP}(\gamma)$. Основной принцип построения этой оценки для унимодальной и асимметричной функции плотности иллюстрируется на Рис.2.

Таким образом, мы хотим определить такой сегмент нормированной функции плотности распределения вероятностей $f(x)$, по которому интегрируется заданная вероятность γ при условии совпадения значений функции $f(x)$ на границах сегмента. Длина этого сегмента и есть значение симметричной по вероятности количественной оценки информационной чувствительности $\delta_{IP}(\gamma)$.

Опишем построение такой оценки формально для случая унимодальной кусочно-монотонной непрерывной функции плотности распределения вероятностей $f(x)$, нормированной в сегмент $[0, 1]$. Определим x^* как $x^* = \arg \max_{x \in [0,1]} f(x)$, при этом функция $f(x)$ моно-

тонно возрастает на сегменте $[0, x^*]$ до максимального значения $f(x^*)$ и монотонно убывает на сегменте $[x^*, 1]$ (Рис.2). В силу высказанных предположений значение x^* единственно и, следовательно, $f(x)$ представима в виде:

$$f(x) = \begin{cases} f_+(x), & 0 \leq x \leq x^* \\ f_-(x), & x^* \leq x \leq 1 \end{cases}$$

где обозначения компонент функции $f(x)$ выбраны по знаку их производных. В силу монотонности функций $f_+(x)$ и $f_-(x)$ на области их определения для них существуют обратные функции: $f_+^{-1}(x)$ и $f_-^{-1}(x)$. Для рассматриваемой функции $f(x)$ через $F(x)$ обозначим интегральную функцию распределения вероятностей

$$F(x) = \int_0^x f(t) dt.$$

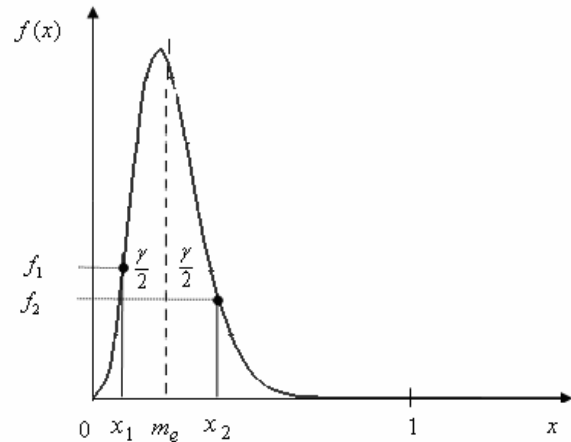


Рис.1. Квантильная количественная мера информационной чувствительности

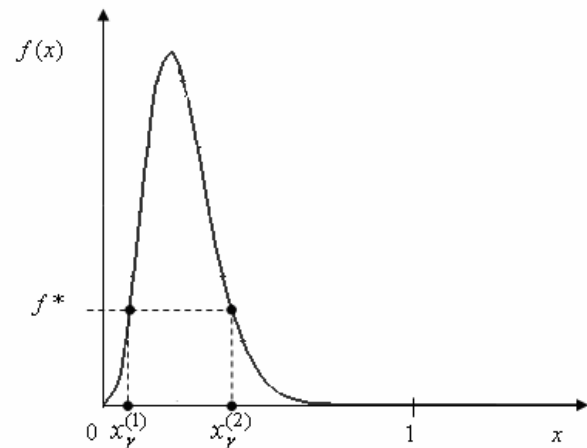


Рис.2. Симметричная по плотности вероятностей количественная оценка информационной чувствительности

Используя данные обозначения, введем в рассмотрение функцию $I(s)$, определяемую выражением

$$I(s) = F(f_-^{-1}(s)) - F(f_+^{-1}(s)), \quad 0 \leq s \leq f(x^*).$$

Таким образом, значение $I(s)$ — это вероятность попадания случайной величины, имеющей плотность распределения вероятностей $f(x)$, в сегмент, на границах которого $f(x)$ имеет значение s . Существование такого сегмента следует из того, что функция $I(s)$ монотонно убывает в силу предположений о $f(x)$ и ограничений на значения s , причем $I(f(x^*)) = 0, I(0) = 1$. Исполь-

зую функцию $I(s)$, определим симметричную по плотности вероятностей оценку информационной чувствительности при фиксированной длине входа в виде:

$$\delta_{IP}(\gamma) = x_{\gamma}^{(2)} - x_{\gamma}^{(1)}, \quad (9)$$

$$x_{\gamma}^{(1)} = f_+^{-1}(f_{\gamma}), \quad x_{\gamma}^{(2)} = f_-^{-1}(f_{\gamma}), \quad f_{\gamma} : I(f_{\gamma}) = \gamma. \quad (10)$$

Если воспользоваться обозначением $I^{-1}(p)$ для функции, обратной к $I(s)$, то (9) представимо в явном функциональном виде:

$$\delta_{IP}(\gamma) = f_-^{-1}(I^{-1}(\gamma)) - f_+^{-1}(I^{-1}(\gamma)),$$

Вычисление этой оценки связано с вычислением по заданной вероятности γ значения $f_{\gamma} = I^{-1}(\gamma)$ и значений границ сегмента $f_+^{-1}(f_{\gamma}), f_-^{-1}(f_{\gamma})$, для которых $f_+(x_{\gamma}^{(1)}) = f_-(x_{\gamma}^{(2)})$. Как правило, для большинства функций плотности вычисление границ сегмента в соответствии с (10) приводит к необходимости итерационного решения задачи определения нуля функции (Рис. 3).

Мы выбираем два начальных значения f_1 и f_2 (Рис. 3), определяющих границы двух сегментов, для которых $I(f_1) < \gamma, I(f_2) > \gamma$, и, используя любой численный метод нахождения нуля функции, решаем уравнение $I(f) - \gamma = 0$. Решение этого уравнения $f_{\gamma} = I^{-1}(\gamma)$ и определяет необходимый нам сегмент $[x_{\gamma}^{(1)}, x_{\gamma}^{(2)}]$, границы которого вычисляются по формуле (10), а длина этого сегмента и есть искомое значение $\delta_{IP}(\gamma)$.

Для разработчиков представляет интерес более детальный анализ алгоритма — получение не только собственно длины сегмента, но и его положения, равно как и функциональная зависимость этих величин от длины входа. Таким образом, мы вводим симметричную по плотности оценку информационной чувствительности в виде:

$$\delta_{IP}^*(\gamma, n) = (x_{\gamma}^{(1)}(n), x_{\gamma}^{(2)}(n), x_{\gamma}^{(2)}(n) - x_{\gamma}^{(1)}(n)). \quad (11)$$

Таким образом, количественная оценка $\delta_{IP}^*(\gamma, n)$ задается тремя числами: значениями

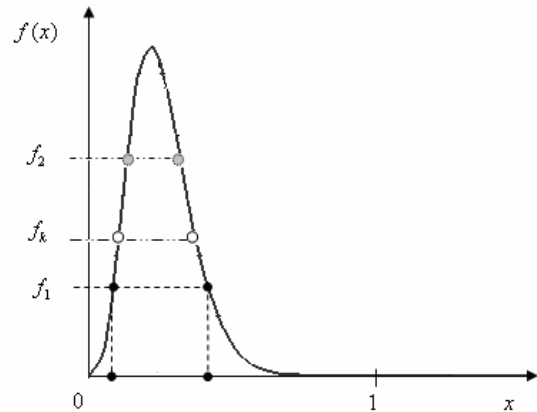


Рис.3. К вычислению симметричной по плотности вероятностей количественной оценки информационной чувствительности

границ нормированного сегмента, по которому интегрируется заданная вероятность γ , и длиной этого сегмента.

7. Методика вычисления оценки

$$\delta_{IP}^*(\gamma, n)$$

На основе вышеизложенного предлагается методика вычисления оценки $\delta_{IP}^*(\gamma, n)$, состоящая из следующих этапов:

- фиксация длины входа n_1 и экспериментальное исследование алгоритма на множестве входов длины n_1 , целью которого является получение значений функции трудоемкости;
- нормирование значений трудоемкости в сегмент $[0, 1]$ на основе теоретических функций трудоемкости в лучшем и худшем случаях или на основе репрезентативной выборки и построение гистограммы относительных частот трудоемкости в нормированном сегменте;
- аппроксимация полученной гистограммы некоторой функцией плотности распределения с ограниченной вариацией;
- задание вероятности γ и вычисление количественной оценки информационной чувствительности на входах длины n_1 по формуле (9) на основе решения уравнения (10);
- вычисление оценки для других длин входа n_i и интерполяция полученных данных в функциональную зависимость от длины входа — $\delta_{IP}^*(\gamma, n)$ по методике, предложенной в [5].

В случае если аппроксимирующая функция плотности вероятностей является асимметричной, авторы считают целесообразным использовать для определения информационной чувствительности именно оценку $\delta_{IP}^*(\gamma, n)$, так как она предусматривает включение в интервал более вероятных значений и исключение менее вероятных. Ниже мы покажем, что получаемый сегмент одновременно является наименьшим из всех возможных.

8. Минимальность оценки $\delta_{IP}(\gamma)$

Решение задачи определения симметричной по плотности вероятностей количественной оценки информационной чувствительности связано с задачей поиска минимума значения $\delta_{IP}(\gamma)$ при условии

$$\int_{x_1}^{x_2} f(t) dt = \gamma, \tag{12}$$

где $f(t)$ — некоторая функция плотности, т. е. с задачей поиска условного экстремума.

Утверждение 2. Предложенная оценка $\delta_{IP}(\gamma)$ (9) доставляет минимум функции $y(x_1, x_2) = x_2 - x_1$ при условии (12).

Доказательство. Найдем минимум функции $y(x_1, x_2) = x_2 - x_1$ методом множителей Лагранжа. Введем в рассмотрение функцию

$$z(x_1, x_2, \lambda) = x_2 - x_1 + \lambda \left(\int_{x_1}^{x_2} f(t) dt - \gamma \right) \tag{13}$$

и найдем ее частные производные:

$$\frac{\partial z}{\partial x_1} = -1 - \lambda f(x_1), \quad \frac{\partial z}{\partial x_2} = 1 + \lambda f(x_2), \quad \frac{\partial z}{\partial \lambda} = \int_{x_1}^{x_2} f(t) dt - \gamma.$$

Приравнивая частные производные к нулю, получаем

$$f(x_1) = f(x_2) = -\frac{1}{\lambda}, \quad \int_{x_1}^{x_2} f(t) dt = \gamma,$$

заметим, что при этом $\lambda < 0$, поскольку значения $f(t) > 0$.

Проверим выполнение достаточных условий минимума функции $z(x_1, x_2, \lambda)$ методом анализа определителя вторых частных производных [11].

Составим определитель из вторых частных производных функции $z(x_1, x_2, \lambda)$ в предположении о дифференцируемости функции $f(x)$:

$$\Delta = \begin{vmatrix} \frac{\partial^2 z}{\partial x_1^2} & \frac{\partial^2 z}{\partial x_1 \partial x_2} & \frac{\partial^2 z}{\partial x_1 \partial \lambda} \\ \frac{\partial^2 z}{\partial x_2 \partial x_1} & \frac{\partial^2 z}{\partial x_2^2} & \frac{\partial^2 z}{\partial x_2 \partial \lambda} \\ \frac{\partial^2 z}{\partial \lambda \partial x_1} & \frac{\partial^2 z}{\partial \lambda \partial x_2} & \frac{\partial^2 z}{\partial \lambda^2} \end{vmatrix} = \begin{vmatrix} -\lambda f'(x_1) & 0 & -f(x_1) \\ 0 & \lambda f'(x_2) & f(x_2) \\ -f(x_1) & f(x_2) & 0 \end{vmatrix}.$$

Вычисляя значение определителя, получаем:

$$\Delta = \lambda f'(x_1) f^2(x_2) - \lambda f'(x_2) f^2(x_1).$$

С учетом того факта, что $f(x_1) = f(x_2)$,

$$\Delta = \lambda f^2(x_1) \cdot (f'(x_1) - f'(x_2)).$$

Для унимодальной функции плотности $f(t)$ значение $f'(x_1) \geq 0$, а $f'(x_2) \leq 0$ в силу (10). Таким образом, если хотя бы для одного из значений аргумента x_1, x_2 значение производной функции $f(t)$ не равно нулю и поскольку если $\lambda < 0$, то и $\Delta < 0$, то, функция $z(x_1, x_2, \lambda)$ имеет минимум, что влечет минимальность $\delta_{IP}(\gamma)$. Если же $f'(x_1) = f'(x_2) = 0$, то возможны варианты: решений будет либо одно, либо бесконечно много.

Таким образом, возвращаясь к функции (13), можно утверждать, что при нарушении условия $f(x_1) = f(x_2)$ в обоих возможных случаях: $f(x_1) < f(x_2)$ или $f(x_1) > f(x_2)$, - происходит увеличение длины сегмента, по которому осуществляется интегрирование, что доказывает минимальность предложенной оценки $\delta_{IP}(\gamma)$ при условии унимодальности функции плотности распределения вероятностей.

В дополнение к сказанному заметим, что решение задачи поиска минимума оценки информационной чувствительности возможно не только для унимодальной функции плотности распределения. При этом надо отметить, что для различных многомодальных распределений таких решений может быть не одно, а задача становится более общей — задачей поиска наименьшего значения двумерной функции $y(x_1, x_2) = x_2 - x_1$ на ограниченном множестве при выполнении условий (10). В зависимости

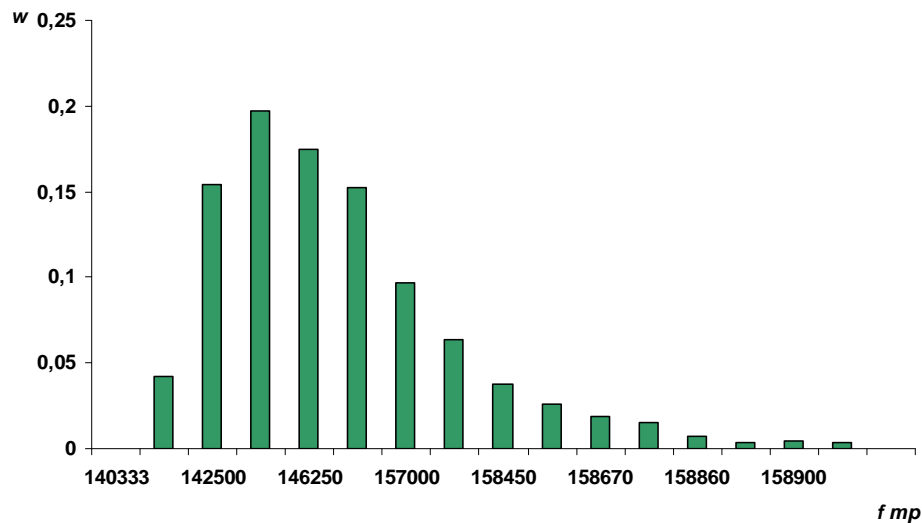


Рис 4. Ненулевая часть гистограммы относительных частот значений функции трудоемкости для алгоритма Кнута-Морриса-Пратта, $n=10000$, $m=20$

от вида функции $f(t)$ возможных решений (x_1, x_2) уравнения (10) может быть одно, несколько или бесконечно много.

9. Модельный пример

Проиллюстрируем различные количественные меры информационной чувствительности на данных экспериментального исследования¹ алгоритма Кнута-Морриса-Пратта для поиска подстроки в строке [12]. Этот алгоритм использует префиксную функцию для преобработки строки поиска. Теоретические функции трудоемкости алгоритма Кнута-Морриса-Пратта в случае однократного вхождения подстроки, необходимые для нормирования значений трудоемкости, получены в [13], при этом трудоемкость в лучшем случае имеет вид:

$$f_A^\vee(n, m) = 14n + 17m - 7,$$

где n — длина строки, m — длина подстроки. Заметим, что функция трудоемкости для этого алгоритма имеет два аргумента.

Трудоемкость в худшем случае задается формулой:

$$f_A^\wedge(n, m) = 24n + 17m - 27.$$

Для значений $n = 10000$, $m = 20$ получаем теоретические минимум и максимум значений функции трудоемкости для случая однократного вхождения, равные:

$$f_A^\vee(10000, 20) = 140333, \quad f_A^\wedge(10000, 20) = 240313.$$

Экспериментальное исследование трудоемкости алгоритма Кнута-Морриса-Пратта было проведено при $n = 10000$ и $m = 20$ по методике, предложенной в [14]. Генерировались случайные входы, для которых случайно выбиралась подстрока с однократным вхождением. Было проведено 20000 экспериментов, в каждом из которых фиксировалось значение трудоемкости — числа базовых операций, заданных алгоритмом. На основе этих данных была построена гистограмма относительных частот W значений трудоемкости (Рис.4). Отметим, что поведение значений трудоемкости как ограниченной случайной величины имеет для этого алгоритма и данных параметров входа ярко выраженную асимметрию, в данном случае — левую.

Нормирование значений трудоемкости в сегмент $[0, 1]$ проведено по теоретически полученным границам. Была выдвинута гипотеза H_0 об аппроксимации гистограммы нормированных относительных частот функцией плотности бета-распределения. Методом моментов [14] были

¹ Экспериментальные данные получены А.С. Алексеенко и М.В. Ульяновым

определены параметры аппроксимирующего бета-распределения: $\alpha = 1,42$; $\beta = 18,47$. После вычисления теоретических относительных частот для функции плотности бета-распределения с данными параметрами, гипотеза H_0 была подтверждена критерием согласия Пирсона на уровне значимости 0,95.

По формуле (8) для вероятности $\gamma = 0,95$ было вычислено значение квантильной количественной оценки информационной чувствительности $\delta_{IQ}(\gamma) = 0,2104$. При этом границы сегмента, соответствующего $1/2 \pm \gamma/2$ -квантилям бета-распределения (Рис.1), оказались равными $x_1 = 0,0049$, $x_2 = 0,2153$. Поскольку аппроксимирующая функция плотности бета-распределения с параметрами $\alpha = 1,42$; $\beta = 18,47$ лево-асимметрична, то значения функции плотности в этих точках не совпадают (Рис.1), и принимают следующие значения: $b(x_1, \alpha, \beta) \approx 0,008$; $b(x_2, \alpha, \beta) = 0,0007$.

Симметричная по плотности вероятностей количественная оценка информационной чувствительности, вычисленная по предложенной в статье методике, при вероятности $\gamma = 0,95$ равна $\delta_{IP}(\gamma) = 0,1785$, что почти на 20% меньше, чем значение квантильной оценки. При этом границы соответствующего сегмента (Рис.2) $x_1 = 0,00065$, $x_2 = 0,1792$, а значения функции плотности, в соответствии с определением (10), совпадают и равны

$$b(x_1, \alpha, \beta) = b(x_2, \alpha, \beta) = 0,0015.$$

Таким образом, в соответствии с (11) значение $\delta_{IP}(\gamma, n, m)$ при $\gamma = 0,95$, $n = 10000$, $m = 20$ равно

$$\delta_{IP}^*(0,95, 10000, 20) = (x_\gamma^{(1)} = 0,00064, x_\gamma^{(2)} = 0,1792, \delta_{IP}(0,95) = 0,1785).$$

Отметим меньшее, по сравнению с квантильной оценкой, значение $\delta_{IP}(\gamma)$, что обусловливается сокращением длины сегмента, по которому интегрируется заданная вероятность до минимального значения, как это было доказано выше. Это уменьшение повышает точность оценки информационной чувствительности с сохранением ее надежности.

Заключение

На основе вероятностного подхода к описанию трудоемкости на входах фиксированной длины, в статье приведены известные количественные оценки информационной чувствительности алгоритмов и предложена новая симметричная по плотности вероятностей оценка информационной чувствительности алгоритма, показана ее минимальность при заданном значении доверительной вероятности.

Предложенная оценка для асимметричных функций плотности позволяет более точно прогнозировать поведение алгоритма на реальных входах с заданной вероятностью, по сравнению с прогнозированием на основе квантильной меры, за счет корректного решения задачи вычисления границ сегмента на основе аппроксимации распределения значений трудоемкости как ограниченной случайной величины.

Для обсуждаемых в статье количественных оценок информационной чувствительности сформулированы рекомендации по их применению в процессе разработки и выбора алгоритмического обеспечения программ. Показано, что если на данный момент невозможно теоретически определить границы сегмента варьирования значений трудоемкости, то достаточной оценкой информационной чувствительности является оценка на основе репрезентативной выборки.

Экспериментальные данные для алгоритма поиска подстроки в строке согласуются с теоретическими результатами сравнительного анализа симметричной и квантильной оценок.

Введенная симметричная по плотности вероятностей количественная оценка информационной чувствительности может быть использована как более достоверная в случае асимметрии функции плотности, аппроксимирующей значения трудоемкости. Такой выбор рекомендуется для тех алгоритмов, экспериментальные данные трудоемкости которых имеют явно выраженную асимметрию, при этом авторы отмечают, что теоретическое предсказание асимметрии распределения представляет собой сложную и интересную дополнительную задачу.

Рассмотренные количественные оценки могут быть полезны при прогнозировании вре-

менной эффективности компьютерных алгоритмов, при решении задачи рационального выбора алгоритмов по критерию стабильности по времени, в частности для балансировки загрузки кластеров, равно как и при решении других задач прикладной теории алгоритмов.

Литература

1. Ульянов М. В. Ресурсно-эффективные компьютерные алгоритмы. Разработка и анализ. — М.: ФИЗМАТЛИТ, 2008. — 304 с.
2. Ульянов М. В., Головешкин В. А. Информационная чувствительность функции трудоемкости алгоритмов к входным данным // Новые информационные технологии: Сборник трудов VII Всероссийской НТК — М.: МГАПИ, 2004. С. 19–26.
3. Ульянов М. В., Алексеенко А. С. Вероятностный подход к определению количественной меры информационной чувствительности компьютерных алгоритмов // Автоматизация и современные технологии.—2009.— №10.— С.24-32.
4. Шеннон К., Математическая теория связи // в кн.: "Работы по теории информации и кибернетике", — М., 1963. —С. 242-332.
5. Вернер М. Дифференциальная энтропия // Основы кодирования. — ЗАО «РИЦ «Техносфера», 2004. — С. 109—114.
6. Попков Ю. С. Основы теории динамических систем с энтропийным оператором и ее приложения // Автоматика и телемеханика, 2006, № 6, 75–105.
7. Попков Ю. С., Качественный анализ динамических систем с V_q -энтропийным оператором // Автоматика и телемеханика, 2007, № 1, 41–56.
8. Прохоров Ю. В., Розанов Ю. А. Теория вероятностей (Основные понятия. Предельные теоремы. Случайные процессы). — М.: Наука, 1973. — 494 с.
9. Гмурман В. Е. Теория вероятностей и математическая статистика – 9-е изд., стер.- М.: Высш. шк., 2003.- 479 с.
10. Лагутин М. Б. Наглядная математическая статистика — М.: БИНОМ. Лаборатория знаний, 2007. — 472 с.
11. Ильин В. А., Садовничий В. А., Сендов Бл. Х. Математический анализ. — М.: Наука. Главная редакция физико-математической литературы, 1979. — 720 с.
12. Гасфилд Д. Строки, деревья и последовательности в алгоритмах: Информатика и вычислительная биология. — СПб.: Невский диалект, 2003. — 654 с.
13. Алексеенко А. С. Информационная чувствительность алгоритма Кнута-Морриса-Пратта // Задачи системного анализа, управления и обработки информации: межвуз. сб. науч. тр. Вып. 3. — М.: МГУП, 2010. С. 7-10.
14. Петрушин В. Н., Ульянов М. В. Планирование экспериментального исследования трудоемкости алгоритмов на основе бета-распределения // Информационные технологии и вычислительные системы, 2008. № 2. С. 81–91.

Головешкин Василий Адамович. Профессор кафедры «Высшая математика» Московского государственного университета приборостроения и информатики. Окончил Московский государственный университет им. М. В. Ломоносова в 1974 году. Доктор технических наук (2005), профессор (2006). Автор более 70 научных и учебно-методических работ, в том числе 3-х монографий. Область научных интересов: механика деформируемого твердого тела, теория рекурсии, аналитическое исследование алгоритмов. E-mail: nikshevolog@yandex.ru

Петрушин Владимир Николаевич. Доцент кафедры «Прикладная математика и моделирование систем» Московского государственного университета печати. Московский государственный университет им. М. В. Ломоносова в 1974 году. Кандидат физико-математических наук (1988), доцент (1991). Автор более 75 научных работ, в том числе одной монографии. Область научных интересов: теория вероятностей, математическая статистика, теория эксперимента.

Ульянов Михаил Васильевич. Профессор кафедры «Управление разработкой программного обеспечения» Государственного университета - Высшей школы экономики, профессор кафедры «Прикладная математика и моделирование систем» Московского государственного университета печати. Окончил Московский институт электронного машиностроения в 1979 году. Доктор технических наук (2005), профессор (2006). Автор более 70 научных работ, в том числе 4-х монографий. Область научных интересов: анализ и разработка ресурсно-эффективных компьютерных алгоритмов. E-mail: muljanov@mail.ru