

Вероятностный анализ задачи восстановления зависимостей по интервальным данным: оправдание допускового множества решений

В.Я. КРЕЙНОВИЧ, С.П. ШАРЫЙ

Университет Техаса в Эль-Пасо, Эль-Пасо, США
ФИЦ ИВТ и НГУ, Новосибирск, Россия

29 сентября 2025 г.

В предыдущей статье авторов —



Kreinovich V., Shary S.P. Interval methods for data fitting under uncertainty: a probabilistic treatment // Reliable Computing. – 2016. – Vol. 23. – P. 105–141

— показано, что упрощённый вероятностный подход к описанию интервальной неопределённости в задаче восстановления зависимостей приводит к объединённому множеству решений.

Сейчас мы показываем, что реалистичный теоретико-вероятностный анализ задачи делает необходимым привлечение другого известного понятия — допускового множества решений интервальных уравнений.

Общая формулировка задачи

При обработке неточных данных практики часто используют подходы и методы, основанные на вероятностной модели неопределённости.

Как эти методы соотносятся с методами интервального анализа данных в случае интервального описания неопределённости?



Bounding approaches to system identification / Milanese M., Norton J., Piet-Lahanier H., Walter E. (eds.) – New York: Springer, Plenum Press, 1996.



Jaulin L., Kieffer V., Didrit O., Walter E. Applied interval analysis, with examples in parameter and state estimation, robust control, and robotics. – Springer, London, 2001.



Баженов А.Н., Жилин С.И., Кумков С.И., Шарый С.П. Обработка и анализ интервальных данных. – Издательство «ИКИ»: Ижевск-Москва, 2024.

... интервального анализа данных и традиционной статистики?

традиционная
вероятностная
статистика

The diagram consists of two overlapping circles. The left circle is light red and contains the text 'традиционная вероятностная статистика'. The right circle is light green and contains the text 'интервальный анализ данных'. The circles overlap in the center, representing the intersection of the two fields.

интервальный
анализ данных

Задача восстановления зависимости

Во многих ситуациях

- мы знаем общий вид $y = f(x, c)$ функциональной зависимости вещественной величины y от величин $(x_1, \dots, x_n) = x$, но
- мы не знаем точных значений параметров $c = (c_1, \dots, c_m)$, которые определяют конкретную функцию.

Примеры

- линейная зависимость

$$y = c_0 + c_1 x_1 + \dots + c_n x_n.$$

- общая квадратичная зависимость

$$y = \sum_{i,j} c_{ij} x_i x_j.$$

- полином из экспонент (в частности, для радиоактивного распада)

$$y = \sum_{i=1}^m c_{2i-1} \exp(-c_{2i} t),$$

- и т. д.

Задача восстановления зависимости

Для нахождения c_i измеряем x_1, x_2, \dots, x_n и y несколько раз (K).

На основе результатов измерений

$$\tilde{x}^{(k)} = (\tilde{x}_1^{(k)}, \tilde{x}_2^{(k)}, \dots, \tilde{x}_n^{(k)}) \quad \text{и} \quad \tilde{y}^{(k)}, \quad k = 1, 2, \dots, K,$$

необходимо оценить значения параметров, которые «согласуются» («совместны» и т. п.) с данными, или даже наилучшим образом.

Подставляя их в выражение для искомой функции, получим

$$\begin{cases} f(\tilde{x}_1^{(k)}, \tilde{x}_2^{(k)}, \dots, \tilde{x}_n^{(k)}, c_1, c_2, \dots, c_m) = \tilde{y}^{(k)}, \\ k = 1, 2, \dots, K, \end{cases}$$

— система уравнений для определения c_i .

Задача восстановления зависимости

Измерения никогда не являются абсолютно точными.

Как следствие, необходимо принимать во внимание, что результаты измерений \tilde{v} , вообще говоря, отличаются от истинных значений v , т. е. существует ненулевая измерительная погрешность $\Delta v := \tilde{v} - v$.

Соответствующие погрешности измерений нужно учитывать при оценивании параметров c_i искомой функции f .

Модель неточностей и погрешностей в данных

Данные почти всегда неточны ...

Какую модель неопределённости данных мы принимаем?

Традиционный выбор — теоретико-вероятностная модель
(К.Ф. Гаусс, П.С. Лаплас и т. д.):

ошибки измерений и наблюдений — это
случайные величины теории вероятностей
с (более-менее) известными характеристиками

Интервальная модель погрешностей

Неопределённости и неточности в данных

часто удобнее описывать интервально

Даны интервальные оценки величин, т. е. принадлежности $x_i^{(k)}$ и $y^{(k)}$ некоторым интервалам:

$$x_i^{(k)} \in \mathbf{x}_i^{(k)} = [\underline{\mathbf{x}}_i^{(k)}, \overline{\mathbf{x}}_i^{(k)}] \quad \text{и} \quad y^{(k)} \in \mathbf{y}_i = [\underline{\mathbf{y}}^{(k)}, \overline{\mathbf{y}}^{(k)}]$$

т. е.

$$\underline{\mathbf{x}}_i^{(k)} \leq x_i^{(k)} \leq \overline{\mathbf{x}}_i^{(k)} \quad \text{и} \quad \underline{\mathbf{y}}^{(k)} \leq y^{(k)} \leq \overline{\mathbf{y}}^{(k)}$$

Формулировка задачи

Интервальные данные —

$$\begin{array}{cccccc} \mathbf{x}_1^{(1)}, & \mathbf{x}_2^{(1)}, & \dots & \mathbf{x}_n^{(1)}, & \mathbf{y}^{(1)}, \\ \mathbf{x}_1^{(2)}, & \mathbf{x}_2^{(2)}, & \dots & \mathbf{x}_n^{(2)}, & \mathbf{y}^{(2)}, \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{x}_1^{(K)}, & \mathbf{x}_2^{(K)}, & \dots & \mathbf{x}_n^{(K)}, & \mathbf{y}^{(K)}. \end{array}$$

Подставляя их в выражение для искомой функции, получим

$$\left\{ \begin{array}{l} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2^{(k)}, \dots, \mathbf{x}_n^{(k)}, c_1, c_2, \dots, c_m) = \mathbf{y}^{(k)}, \\ k = 1, 2, \dots, K, \end{array} \right. \quad (\diamond)$$

— система уравнений для определения c_i .

Интервальная система уравнений

Но решения и множества решений для интервальных уравнений и систем уравнений могут определены очень многими способами, смысл которых — существенно разный.

Наиболее популярно объединённое множество решений, которое образовано всевозможными решениями точечных систем уравнений, содержащимися в интервальной системе:




$$\Xi_{\text{uni}} := \left\{ c \in \mathbb{R}^m \mid \begin{aligned} &(\exists x^{(k)} \in \mathbf{x}^{(k)})(\exists y^{(k)} \in \mathbf{y}^{(k)}) f(x^{(k)}, c) = y^{(k)}, \\ &k = 1, 2, \dots, K \end{aligned} \right\}$$

Интервальная система уравнений

Кроме него существуют также другие важные множества решений, в частности, допусковое множество решений:

$$\Xi_{\text{tol}} := \left\{ c \in \mathbb{R}^m \mid \begin{aligned} &(\forall x_{:}^{(k)} \in \mathbf{x}_{:}^{(k)}) (\exists y^{(k)} \in \mathbf{y}^{(k)}) f(x_{:}^{(k)}, c) = y^{(k)}, \\ &k = 1, 2, \dots, K \end{aligned} \right\}$$

Оно образовано теми решениями точечных систем уравнений, для которых значение отображения f попадает в интервалы правой части при любых значениях параметров $x_{:}^{(k)} \in \mathbf{x}_{:}^{(k)}$.

-  **Nuding E., Wilhelm J.** Über Gleichungen und über Lösungen // ZAMM. – 1972. – Bd. 52. – S. T188–T190.
-  **Neumaier A.** Tolerance analysis with interval arithmetic // Freiburger Intervall-Berichte. – 1986. – No. 6. – S. 5–19.
-  **Deif A.** Sensitivity analysis in linear systems. – Berlin: Springer, 1986.

Ранняя история исследования этого множества

— в работе А. Ноймайера: «внутренние решения».

Большое внимание уделено ему в книге А. Дейфа, но предложенные в ней методы исследования этого множества носят частный характер.



Rohn J. Inner solutions of linear interval systems // Interval Mathematics 1985 / Nickel K., ed. – New York: Springer Verlag, 1986. – P. 157–158. – (Lecture Notes in Computer Science; vol. 212).

— представление допускового множества решений
через решение системы линейных неравенств.

Можно исследовать пустоту/непустоту допускового множества решений с помощью методов линейного программирования, которые имеют полиномиальную трудоёмкость.



Shary S.P. Solving interval linear tolerance problem // Mathematics and Computers in Simulation. – 1995. – Vol. 39. – P. 53–85.

— распознающий функционал допускового множества решений интервальной линейной $m \times n$ -системы уравнений $\mathbf{Ax} = \mathbf{b}$:

$$\text{Tol}(x, \mathbf{A}, \mathbf{b}) = \min_{1 \leq i \leq m} \left\{ \text{rad } \mathbf{b}_i - \left| \text{mid } \mathbf{b}_i - \sum_{j=1}^n \mathbf{a}_{ij} x_j \right| \right\}$$



Шарый С.П. Конечномерный интервальный анализ. — Новосибирск: XYZ, 2025. Электронная книга, доступная на <http://www.nsc.ru/interval/?page=Library/InteBooks>

— глава 6 целиком посвящена допусковому множеству решений для интервальных линейных систем уравнений.

Аналогичные результаты — также в других главах.



Sharaya I.A. On unbounded tolerable solution sets // *Reliable Computing*. – 2005. – Vol. 11, No. 5. – P. 425–432.

— выведен простой критерий ограниченности
допускового множества решений.

Оказывается, что допусковое множество решений почти всегда ограничено, в отличие от объединённого множества решений.

В задачах восстановления зависимостей это означает, что множество параметров, совместных с данными, имеет ограниченную вариабельность (изменчивость).

— в задаче восстановления зависимостей по интервальным данным



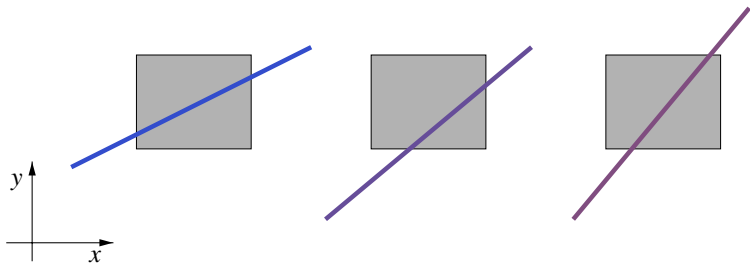
ШАРЫЙ С.П. Сильная согласованность в задаче восстановления зависимостей при интервальной неопределённости данных // Вычислительные Технологии. – 2017. – Т. 22, №2. – С. 150–172.



SHARY S.P. Weak and strong compatibility in data fitting problems under interval uncertainty // Advances in Data Science and Adaptive Analysis. 2020. Vol. 12, No. 1. Paper 2050002.

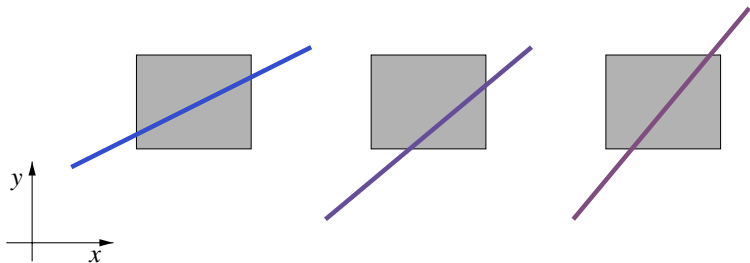
Допусковое множество решений

— в задаче восстановления зависимостей по интервальным данным



Допусковое множество решений

— в задаче восстановления зависимостей по интервальным данным



Покажем, что допустовое множество решений
также имеет вероятностный смысл.

Обработка вероятностной неопределённости

Во многих ситуациях погрешность в данных может быть адекватно описана как «случайная величина» теории вероятностей, а её распределение вероятностей более-менее известно.

Будем предполагать, что погрешности измерений, соответствующие различным переменным, независимы в смысле теории вероятностей.

Пусть распределения, описывающие случайные погрешности, абсолютно непрерывны и имеют плотности вероятности

- $p_i(\Delta x_i)$ для независимых переменных x_i
- $q(\Delta y)$ для зависимой переменной y .

Метод максимума правдоподобия

— один из наиболее популярных методов оценивания параметров:

оценкой параметров берётся значение, на котором достигается максимум «функции правдоподобия».

Обычно функция правдоподобия равна

- вероятности оценки
(для дискретных распределений),
- плотности вероятности оценки
(для абсолютно непрерывных распределений).

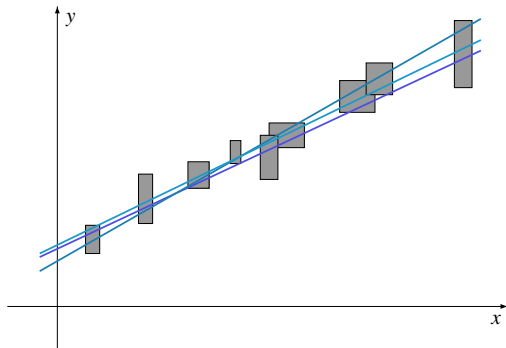


Kreinovich V., Shary S.P. Interval methods for data fitting under uncertainty: a probabilistic treatment // Reliable Computing. – 2016. – Vol. 23. – P. 105–141

Свободно доступна на веб-сайте журнала по ссылке —

<https://interval.louisiana.edu/reliable-computing-journal/volume-23/reliable-computing-23-pp-105-140.pdf>

Метод максимума совместности



Результат

Метод максимума совместности в слабой версии является не чем иным как методом максимального правдоподобия при равномерных вероятностных распределениях на интервалах данных.

В нашем случае оценкой c берётся вектор значений параметров,
для которых плотность вероятности оценки, равная

$$\prod_{k=1}^K \left(q(\tilde{y}^{(k)} - f(x^{(k)}, c)) \cdot \prod_{i=1}^n p_i(\tilde{x}_i^{(k)} - x_i^{(k)}) \right),$$

достигает наибольшего значения.

Метод максимума правдоподобия

Вместо максимизации правдоподобия обычно решают
равносильную задачу максимизации его логарифма:

$$\begin{aligned} \sum_{k=1}^K \left(\log q(\tilde{y}^{(k)} - f(x^{(k)}, c)) + \sum_{i=1}^n \log p_i(\tilde{x}_i^{(k)} - x_i^{(k)}) \right) = \\ = \sum_{k=1}^K \log q(\tilde{y}^{(k)} - f(x^{(k)}, c)) + \sum_{k=1}^K \sum_{i=1}^n \log p_i(\tilde{x}_i^{(k)} - x_i^{(k)}). \end{aligned}$$

Этот шаг имеет технические причины:

с суммами гораздо удобнее работать, чем с произведениями.

Вспомним ЗБЧ, УЗБЧ, ...

Метод максимума правдоподобия

$$\sum_{k=1}^K \log q(\tilde{y}^{(k)} - f(x^{(k)}, c)) + \sum_{k=1}^K \sum_{i=1}^n \log p_i(\tilde{x}_i^{(k)} - x_i^{(k)})$$

Второе слагаемое не зависит от c ,

так что его можно опустить при поиске максимума по c .

В целом необходимо найти наибольшее значение для

$$\sum_{k=1}^K \log q(\tilde{y}^{(k)} - f(x^{(k)}, c)),$$

а также доставляющий его вектор аргумента c .

Часто мы не знаем вероятностное распределение погрешностей, а нередко даже нельзя быть уверенным в том, что погрешности подчиняются закономерностям теории вероятностей.

Предположим, что погрешности измерений Δv
локализованы на заданном интервале $[-\Delta_v, \Delta_v]$.

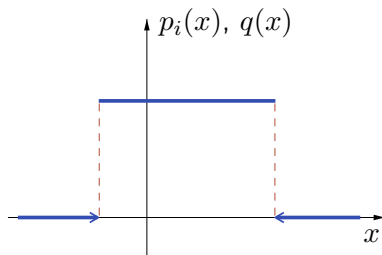
Принимаем также гипотезу, что интервалы измерений и наблюдений являются накрывающими (включающими), т. е. содержат истинные значения измеряемых величин.

Обработка интервальных данных

Чтобы выполнить вероятностный анализ обработки интервальных данных, нужно задать какое-то распределение на интервалах.

Естественно взять то, что имеет наибольшую энтропию, т.е. даёт наименьшую возможную информацию о величине.

Им оказывается равномерное вероятностное распределение.



Обработка интервальных данных

Простейший случай — точное измерение значений x_i ,

и тогда $\tilde{x}_i^{(k)} = x_i^{(k)}$ для всех i и k .

В нашей предыдущей статье показано, что метод максимума правдоподобия выбирает тогда следующее множество:

Множество всевозможных значений c , для которых справедливо $f(x^{(k)}, c) \in y^{(k)}$ при любых k .

Вероятностный подход приводит к тому же самому ответу, что и в интервальном анализе данных, где это множество называется *множеством решений*.

Но для интервальных уравнений и систем уравнений решения и множества решений могут быть определены разными способами ...

Интервальная неопределённость имеет двойственный характер

Существуют объединённое множество решений, допустовое множество решения, множества АЕ-решений, формальные решения и т. д. ...

В описанной выше ситуации главные множества решений — объединённое и допустовое — совпадают друг с другом, так что можно говорить о «множестве решений» вообще.

Метод максимума правдоподобия

В общем случае $x_i^{(k)}$ также имеют интервальную неопределённость.

Тогда при определённых условиях метод максимума правдоподобия выбирает, как было показано,

*множество всех таких c , для которых $f(x_i^{(k)}, c) \in y^{(k)}$
при некоторых $x_i^{(k)} \in \mathbf{x}_i^{(k)}$, $i = 1, 2, \dots, n$, и всяком k .*

Метод максимума правдоподобия

В общем случае $x_i^{(k)}$ также имеют интервальную неопределённость.

Тогда при определённых условиях метод максимума правдоподобия выбирает, как было показано,

множество всех таких c , для которых $f(x_i^{(k)}, c) \in y^{(k)}$ при некоторых $x_i^{(k)} \in \mathbf{x}_i^{(k)}$, $i = 1, 2, \dots, n$, и всяком k .

Это объединённое множество решений интервальной системы уравнений

$$\begin{cases} f(\mathbf{x}_i^{(k)}, c) = y^{(k)}, \\ k = 1, 2, \dots, K, \end{cases} \quad (\diamond)$$

построенной по данным и виду функции.

Объединённое множество решений имеет вероятностное толкование!

Более реалистичное описание практической задачи

Нередко при измерениях тех или иных величин практикуют *повторные измерения*, т. е. выполняемые неоднократно.

Для измерения величин должны быть проведены несколько измерений, по результатам которых формируется итоговое значение величины.

Это часто практикуема процедура, которую выполняют для повышения общей надёжности процесса измерения.

Примеры

- При измерении давления крови тонометром рекомендуется выполнить его 3 (три) раза, и затем результаты усредняются, если они «не сильно разнятся».
- Сверхточные атомные часы обычно состоят из нескольких независимых часов, чьи показатели выдаются пользователю, если большинство их результатов совпадают.
- Именно с помощью повторных измерений вводятся новые значения, будь это уточнённое расстояние до Луны или новое значение атомного веса элемента.

Примеры

- При измерении давления крови тонометром рекомендуется выполнить его 3 (три) раза, и затем результаты усредняются, если они «не сильно разнятся».
- Сверхточные атомные часы обычно состоят из нескольких независимых часов, чьи показатели выдаются пользователю, если большинство их результатов совпадают.
- Именно с помощью повторных измерений вводятся новые значения, будь это уточнённое расстояние до Луны или новое значение атомного веса элемента.

Результат единичного измерения не вполне надёжен, так что для повышения достоверности всегда выполняют несколько измерений, которые дают финальный результат, если «достаточно близки».

Использование повторных измерений

- Например, можно считать, что окончательный интервал результата многократного измерения сформирован операцией агрегирования индивидуальных результатов измерений:

$$\mathbf{x}_i^{(k)} = \left[\min_{\ell} x_i^{(k\ell)}, \max_{\ell} x_i^{(k\ell)} \right], \quad \mathbf{y}^{(k)} = \left[\min_{\ell} y^{(k\ell)}, \max_{\ell} y^{(k\ell)} \right].$$

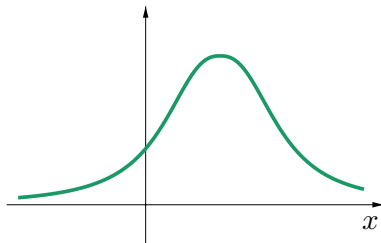
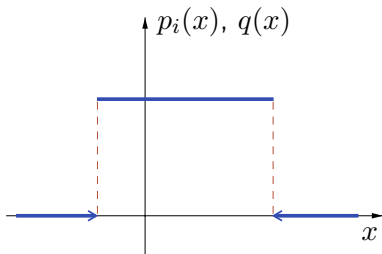
- Можно считать, что общий интервал результата берётся как среднее отдельных измерений, раздутье на несколько стандартных отклонений.
- И т. д.

Для любого измерения k вместо единственного набора $x_{\cdot}^{(k)}$ имеем несколько $x_{\cdot}^{(k\ell)} = (x_1^{(k\ell)}, x_2^{(k\ell)}, \dots, x_n^{(k\ell)})$ для различных ℓ .

Для каждой комбинации значений $x_i^{(k\ell)} \in \mathbf{x}_i^{(k)}$ можем сформировать логарифмическое правдоподобие

$$\sum_{k=1}^K \sum_{\ell} \log q(\tilde{y}^{(k)} - f(x_{\cdot}^{(k\ell)}, c)).$$

Мы не знаем точно $x_i^{(k\ell)}$,
но они лежат в интервалах $\mathbf{x}_i^{(k)}$ и равномерно распределены.



Равномерные распределения имеют ограниченный носитель, хотя у типичных вероятностных распределений они неограниченны.

Предложение

Всевозможные оценки параметров, полученные методом максимума правдоподобия, образуют

*Множество всех таких c , что $f(x^{(k\ell)}, c) \in y^{(k)}$
для любых $x_i^{(k\ell)} \in x_i^{(k)}$, $i = 1, 2, \dots, n$, и всяком k .*

Доказательство

Если для некоторого c условие

$$f(x_{\cdot}^{(k\ell)}, c) \in \mathbf{y}^{(k)}$$

не удовлетворено при каких-то $x_{\cdot}^{(k\ell)} \in \mathbf{x}_{\cdot}^{(k)}$, то

$$q(\tilde{y}^{(k)} - f(x_{\cdot}^{(k\ell)}, c)) = 0,$$

и тогда $\log 0 = -\infty$.

Поэтому в сумме

$$\sum_{k=1}^K \sum_{\ell} \log q(\tilde{y}^{(k)} - f(x_{\cdot}^{(k\ell)}, c)).$$

появляются слагаемые $(-\infty)$, тогда как остальные слагаемые конечны.

Доказательство

Если для некоторого c условие

$$f(x_{\cdot}^{(k\ell)}, c) \in y^{(k)}$$

не удовлетворено при каких-то $x_{\cdot}^{(k\ell)} \in \mathbf{x}_{\cdot}^{(k)}$, то

$$q(\tilde{y}^{(k)} - f(x_{\cdot}^{(k\ell)}, c)) = 0,$$

и тогда $\log 0 = -\infty$.

Поэтому в сумме

$$\sum_{k=1}^K \sum_{\ell} \log q(\tilde{y}^{(k)} - f(x_{\cdot}^{(k\ell)}, c)).$$

появляются слагаемые $(-\infty)$, тогда как остальные слагаемые конечны.

⇒ Вся сумма $= -\infty$, так что её конечный максимум не достигается.

Доказательство (продолжение)

Для всякого c , выбираемого методом максимума правдоподобия, в самом деле должны иметь $f(x_{:}^{(k\ell)}, c) \in y_{:}^{(k)}$ для всех $x_{:}^{(k\ell)} \in x_{:}^{(k)}$.

Так как рассматриваем равномерные распределения на $y^{(k)}$, $k = 1, 2, \dots, K$, то всем значениям в пределах этих интервалов соответствуют одни и те же значения плотности.

Тогда для всех таких кортежей c будем иметь одинаковые значения слагаемых $\log q(\tilde{y}^{(k)} - f(x_{:}^{(k\ell)}, c))$.

Поэтому все такие кортежи c будут выделены методом максимума правдоподобия.

В целом формула

*Множество всех таких c , что $f(x_i^{(k\ell)}, c) \in y^{(k)}$
при любых $x_i^{(k\ell)} \in \mathbf{x}_i^{(k)}$, $i = 1, 2, \dots, n$, и всяком k*

определяет *допусковое множество решений* для интервальной системы уравнений (\blacklozenge), построенной по функции f и интервальным данным для отыскания оценок параметров c .

Допусковое множество решений

также получает вероятностный смысл.

Дополнительные бонусы —

меньшая трудоёмкость и конечная переменность:

- в отличие от объединённого множества решений, для которого распознавание и оценивание является NP-трудным даже если $f(x, c)$ — линейная функция от c ,
- вычисление допустимого множества решений для линейной функции $f(x, c)$ сводится к решению системы линейных неравенств и потому имеет полиномиальную трудоёмкость.

Дополнительные бонусы —

меньшая трудоёмкость и конечная вариабельность:

- в отличие от объединённого множества решений, для которого распознавание и оценивание является NP-трудным даже если $f(x, c)$ — линейная функция от c ,
- вычисление допускового множества решений для линейной функции $f(x, c)$ сводится к решению системы линейных неравенств и потому имеет полиномиальную трудоёмкость.
- оценки параметров с помощью допускового множества решений почти всегда имеют конечную вариабельность в силу критерия ограниченности И.А. Шарой.

Спасибо за внимание!