

Физические методы исследования состава и структуры веществ

Часть II : Теория ошибок измерений

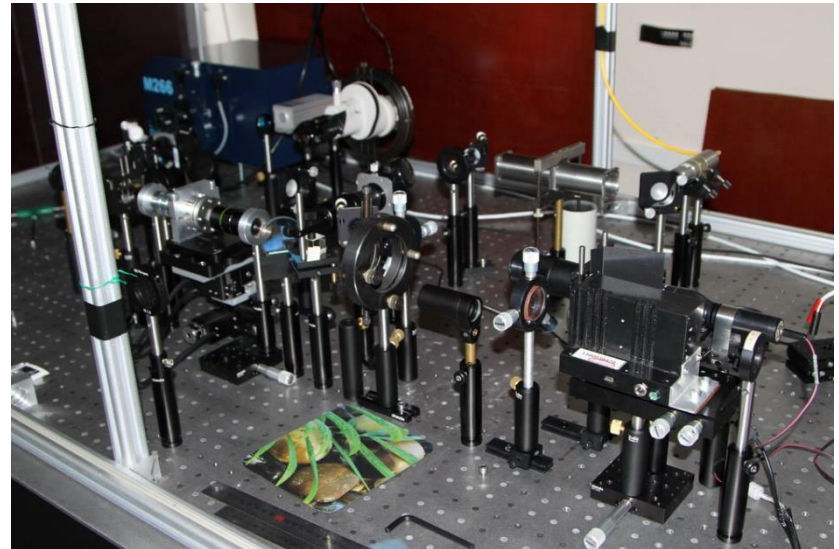
Теория ошибок и обработка результатов эксперимента.

Лекция II: Метод наименьших квадратов

Линейная регрессия, полиномиальная регрессия, нелинейная регрессия



Павел В. Зинин



Расчет среднего и доверительного интервала

1. Пусть есть n измерений величины x . Тогда среднее значение \bar{x} (mean or average) определяется по формуле:

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

2. Второй шаг – это расчет стандартного или среднеквадратичного отклонения (standard deviation) по формуле:

$$s_n = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

3. Следующий шаг - определение Δx доверительного интервала при выбранной доверительной вероятности α . Коэффициентом Стьюдента $t(\alpha, n)$ (Student coefficient) определяется из таблицы или вычисляется.

$$\Delta x = t(\alpha, n) \cdot \frac{s_n}{\sqrt{n}}$$

Тогда можно сказать, что при выбранной доверительной вероятности α (confidence), *результат измерения случайной величины x , представляется в виде*

$$\bar{x} - \Delta x \leq x \leq \bar{x} + \Delta x$$

Доверительный интервал при различных значениях доверительной вероятности

Во многих публикациях в качестве ошибки указывается только стандартное отклонение. Это означает, что значение ε в выражении для доверительного интервала выбирается равным 1. Как далеко это от значения величины $\varepsilon = t(\alpha, n)/\sqrt{n}$, которое должно использоваться для оценки доверительного интервала, при значениях доверительной вероятности или надежности 0.69, 0.9 и 0.95.

$$\Delta x = \varepsilon s_n = t(\alpha, n) \cdot \frac{s_n}{\sqrt{n}}$$

α	0.69	0.9	0.95
n	$t(\alpha, n)/\sqrt{n}$	$t(\alpha, n) \sqrt{n}$	$t(\alpha, n) \sqrt{n}$
3	0.78	1.69	2.48
4	0.61	1.18	1.59
5	0.52	0.95	1.24
6	0.46	0.82	1.05
7	0.42	0.73	0.92
8	0.39	0.67	0.84
9	0.36	0.62	0.77
10	0.34	0.58	0.72

Получение параметров из экспериментальных данных

Цель данной лекции – предложить способ обработки результатов экспериментальных работ в случае, когда неизвестные величины находятся из экспериментальных измерений, описывающихся либо линейной, либо полиномиальной, либо нелинейной зависимостями.

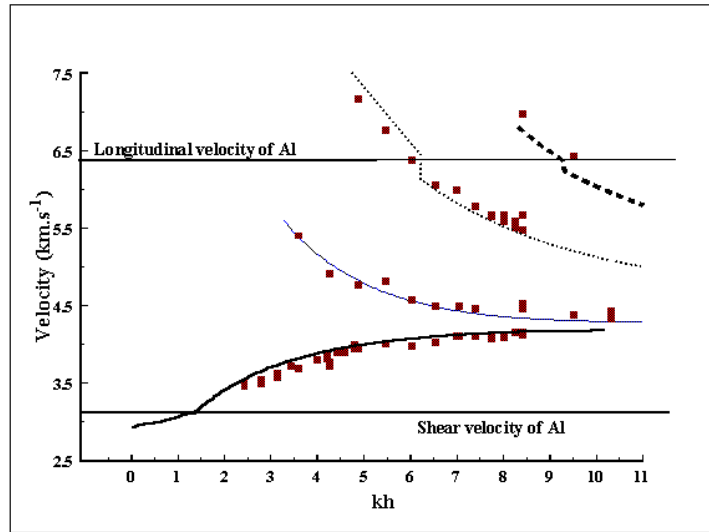


Рис. 1. Теоретические дисперсионные кривые ПАВ в оксидной пленке на алюминиевой подложке. Квадраты – экспериментально померенные данные.

Пример [1]: В тонких пленках упругие свойства можно получить, измеряя зависимость скорости поверхностных акустических волн (ПАВ) в таких пленках от частоты или дисперсионные кривые ПАВ. Дисперсионные кривые ПАВ в оксидной пленке на алюминиевой подложке, полученные путем измерения рассеяния Мандельштама – Бриллюэна (РМБ), показаны на Рис. 1. Скорости продольной и сдвиговой акустических волн в пленке были получены нелинейной подгонкой. Величина скорости сдвига v_T из барьерной пленки получается с большей точностью (± 0.6), чем значение продольной v_L скорости (± 2.1).

Метод наименьших квадратов

Метод наименьших квадратов (МНК, *Ordinary Least Squares, OLS*) — математический метод, применяемый для решения различных задач, основанный на минимизации суммы квадратов отклонений некоторых функций от искомым переменных. МНК является одним из базовых методов *регрессионного анализа* для оценки неизвестных параметров регрессионных моделей по выборочным данным.

Сущность метода наименьших квадратов

Допустим, нам известен вид функциональной зависимости физической величины y от другой физической величины x , но не известны параметры этой зависимости a_j , ($j = 1, 2, 3 \dots p$), где p — число параметров модели. В результате проведенных измерений получена таблица значений y_i при некоторых значениях x_i ($i = 1, 2, 3 \dots n$), где n — число измерений модели. Требуется найти такие значения параметров a_1, a_2, a_3, \dots при которых функция наилучшим образом описывает экспериментальные данные. Таким образом, для определения параметров a_j ($j = 1, 2, 3 \dots p$), необходимо найти минимум функции

$$S = \sum_{i=1}^n [y_i - f(x_i, \mathbf{a})]^2, \quad (1)$$

где $\mathbf{a} = [a_1, a_2, a_3, \dots, a_p]$ есть параметры модели.

Историческая справка

- Метод был предложен 1806 г. А. М. Лежандром в связи с вопросом о вычислениях кометных орбит. Ему же принадлежит название: „метод наименьших квадратов“.
- В 1809 г. К. Ф. Гаусс дал первое вероятностное обоснование метода наименьших квадратов, а в 1810 г. он же глубоко разработал вычислительную сторону вопроса и ввел символы и обозначения, сохранившиеся и поныне.
- В 1812 г. П. С. Лаплас в фундаментальном трактате по теории вероятностей получил ряд важных результатов и применил их к методу наименьших квадратов.
- Дальнейшие важные результаты были получены в теории метода наименьших квадратов в 1859 г. П. Л. Чебышевым, разработавшим теорию интерполирования по методу наименьших квадратов с помощью ортогональных полиномов, носящих его имя.
- А. А. Марков в 1898 г. в работе внес в математическую статистику ряд весьма важных идей, пояснивших суть метода наименьших квадратов.

Линник, Ю.В., Метод наименьших квадратов и основы математической теории обработки наблюдений. 1958, Москва

Метод наименьших квадратов

При интерпретации экспериментальных данных значения x_i будем считать точными. Погрешности в определении x_i приводят к дополнительному разбросу y_i и тем самым должны учитываться в отклонениях y_i от расчетной кривой. Критерий метода наименьших квадратов (МНК) требует минимальности суммы

$$S = \sum_{i=1}^n [y_i - f(x_i, \mathbf{a})]^2$$

Где $\mathbf{a} = [a_1, a_2, a_3, \dots, a_p]$ есть параметры модели.

Условием минимума является равенства нулю системы p уравнений

$$\frac{\partial S}{\partial a_j} = 0, j = 1, 2, \dots, p$$

Решение системы уравнений запишем в виде $\tilde{\mathbf{a}} = [\tilde{a}_1, \tilde{a}_2, \tilde{a}_3, \dots, \tilde{a}_p]$. Применим МНК к линейной зависимости ($p = 2$):

$$y = a + b \cdot x \quad (2)$$

Метод наименьших квадратов: линейная регрессия

Сумма

$$S = \sum_{i=1}^N [y_i - a - bx_i]^2$$

имеет минимум, когда первые производные по параметрам равна нулю:

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^N [y_i - a - bx_i] = 0$$

(3)

$$\frac{\partial S}{\partial b} = -2 \sum_{i=1}^N x_i [y_i - a - bx_i] = 0$$

Систему уравнений (3) можно переписать в виде системы из двух уравнений:

До начала XIX в. учёные не имели определённых правил для решения системы уравнений, в которой число неизвестных меньше, чем число уравнений; до этого времени употреблялись частные приёмы, зависевшие от вида уравнений и от остроумия вычислителей, и потому разные вычислители, исходя из тех же данных наблюдений, приходили к различным выводам. Гауссу (1795) принадлежит первое применение метода, а Лежандр (1805) независимо открыл и опубликовал его под современным названием. Лаплас связал метод с теорией вероятностей.

$$Na + b \sum_{i=1}^N x_i = \sum_{i=1}^N y_i$$

(4)

$$a \sum_{i=1}^N x_i + b \sum_{i=1}^N (x_i)^2 = \sum_{i=1}^N x_i y_i$$

Метод наименьших квадратов: коэффициенты линейной регрессии

Решение системы линейных уравнений (4) имеет вид:

$$a = \frac{\left(\sum_{i=1}^N x_i^2\right)\left(\sum_{i=1}^N y_i\right) - \left(\sum_{i=1}^N x_i\right)\left(\sum_{i=1}^N x_i y_i\right)}{N\left(\sum_{i=1}^N x_i^2\right) - \left(\sum_{i=1}^N x_i\right)^2}$$
$$b = \frac{N\left(\sum_{i=1}^N x_i y_i\right) - \left(\sum_{i=1}^N x_i\right)\left(\sum_{i=1}^N y_i\right)}{N\left(\sum_{i=1}^N x_i^2\right) - \left(\sum_{i=1}^N x_i\right)^2}$$

(5)

Коэффициенты регрессии можно записать в компактной форме

$$a = \frac{S_{x^2} \left(\sum_{i=1}^N y_i \right) - S_x \left(\sum_{i=1}^N x_i y_i \right)}{S_{xx}}$$
$$b = \frac{N \left(\sum_{i=1}^N x_i y_i \right) - S_x \left(\sum_{i=1}^N y_i \right)}{S_{xx}}$$

(6)

где

$$S_{xx} = N \left(\sum_{i=1}^N x_i^2 \right) - \left(\sum_{i=1}^N x_i \right)^2$$
$$S_{x^2} = \left(\sum_{i=1}^N x_i^2 \right) ; S_x = \left(\sum_{i=1}^N x_i \right)$$

(7)

Оценка точности определения параметров линейной регрессии

Как оценить статистическую ошибку определения коэффициентов линейной регрессии только из анализа самих данных (измеренных значений y_1, \dots, y_N)? Для этого представим коэффициенты a и b как суперпозицию слагаемых с измерениями y_i

$$a = \frac{\sum_{i=1}^N [S_{x^2} - S_x x_i] \cdot y_i}{S_{xx}} = \sum_{i=1}^N c_i \cdot y_i, \quad \text{where } c_i = \frac{[S_{x^2} - S_x x_i]}{S_{xx}} \quad (8)$$
$$b = \frac{\sum_{i=1}^N [N x_i - S_x] \cdot y_i}{S_{xx}} = \sum_{i=1}^N d_i \cdot y_i, \quad \text{where } d_i = \frac{[N x_i - S_x]}{S_{xx}}$$

Важно: Статистическая теория говорит нам, что числа y_1, \dots, y_N не представляют собой N результатов измерений одной и той же величины. Результат измерения каждого y_i распределен нормально около истинного значения $a + b x_i$ с дисперсией σ_y . Тогда отклонения $y_i - a - b x_i$ распределены нормально, причем все с одним и тем же центральным значением 0 и одной и той же дисперсией σ_y . Последнее утверждение верно не всегда. В общем случае, измерения y_i , произведенные в точке x_i может иметь собственную дисперсию.

Пример: измерения y_i в точке x_i проводились много раз. Тогда можно оценить дисперсию величины y_i в точке x_i , σ_{y_i} , используя выражение для стандартного отклонения S_{y_i} .

Оценка коэффициентов линейной регрессии

Пусть имеется функция $f(y_i)$ зависящая от N параметров y_i , измеренных с ошибками σ_{y_i} тогда ошибка вычисления функции $f(y_i)$ есть

$$\sigma_f^2 = \sum_{k=1}^N \left(\frac{\partial f(y_i)}{\partial y_i} \sigma_{y_i} \right)^2$$

Поскольку параметр a можно представить в виде:

$$a = \sum_{i=1}^N c_i \cdot y_i \Rightarrow \frac{\partial a}{\partial y_i} = c_i \quad \Rightarrow \quad \frac{\partial b}{\partial y_i} = d_i$$

Тогда

$$\sigma_a^2 = \sum_{k=1}^N (c_i \sigma_{y_i})^2, \quad \sigma_b^2 = \sum_{k=1}^N (d_i \sigma_{y_i})^2$$

Считая что у всех измерений y_i дисперсия одинакова и равна σ_y , получим

$$\sigma_a^2 = (\sigma_y)^2 \sum_{k=1}^N (c_i)^2, \quad \sigma_b^2 = (\sigma_y)^2 \sum_{k=1}^N (d_i)^2 \quad (9)$$

Остается выписать выражение для σ_y и раскрыть суммы в выражении.

Оценка коэффициентов линейной регрессии

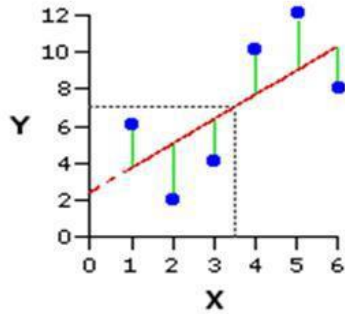
Пусть имеется функция $f(y_i)$ зависящая от параметров y_i , измеренных с ошибками σ_{z_i} тогда ошибка вычисления функции $f(y_i)$ есть

$$\begin{aligned}\sum_{i=1}^N (c_i)^2 &= \frac{\sum_{i=1}^N \left[(S_{x^2})^2 - 2S_x(S_{x^2})(x_i) + (S_x)^2(x_i)^2 \right]}{(S_{xx})^2} = \frac{N(S_{x^2})^2 - 2(S_x)^2(S_{x^2}) + (S_x)^2(S_{x^2})}{(S_{xx})^2} \\ &= \frac{N(S_{x^2})^2 - (S_x)^2(S_{x^2})}{(S_{xx})^2} = \frac{(S_{x^2}) \left[N(S_{x^2}) - (S_x)^2 \right]}{(S_{xx})^2} = \frac{S_{x^2}}{S_{xx}}\end{aligned}$$

$$\begin{aligned}\sum_{i=1}^N (c_i)^2 &= \frac{\sum_{i=1}^N [Nx_i - S_x]^2}{(S_{xx})^2} = \frac{\sum_{i=1}^N N^2(x_i)^2 - 2N(x_i)S_x + (S_x)^2}{(S_{xx})^2} \\ &= \frac{N^2 S_{x^2} - 2N(S_x)^2 + N(S_x)^2}{(S_{xx})^2} = \frac{N \left(S_{x^2} - (S_x)^2 \right)}{(S_{xx})^2} = \frac{N}{S_{xx}}\end{aligned}$$

Оценка коэффициентов линейной регрессии

Важно. Параметру a и b - точно определенные функции измеренных значений y_1, \dots, y_N . Следовательно, погрешности в a и b определяют простым расчетом ошибок в косвенных измерениях, исходя из погрешностей в y_1, y_N .



Можно показать, что фактор N в знаменателе необходимо заменить на $(N-2)$. Таким образом, наш конечный ответ для погрешности в измерениях y_1, \dots, y_N есть

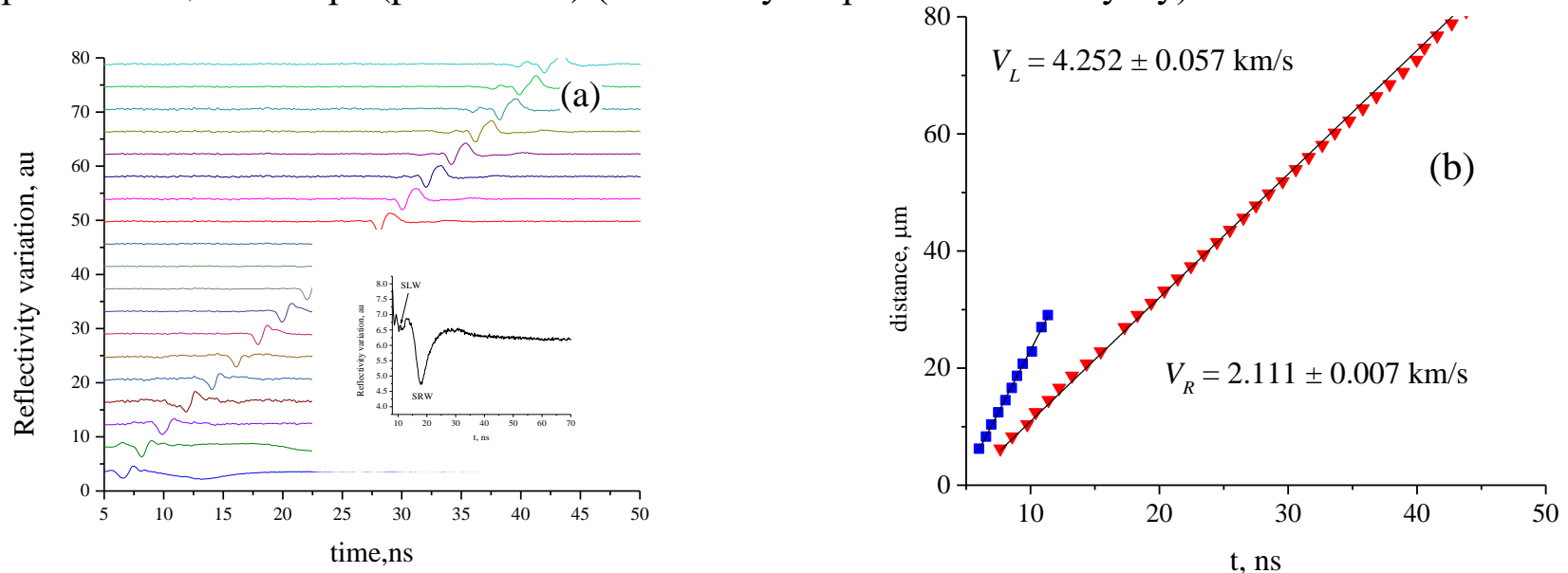
$$\sigma_a = \sigma_{yi} \sqrt{\frac{S_{x^2}}{S_{xx}}} = \sigma_{yi} \sqrt{\frac{\left(\sum_{i=1}^N x_i^2\right)}{N\left(\sum_{i=1}^N x_i^2\right) - \left(\sum_{i=1}^N x_i\right)^2}} \quad (10)$$

$$\sigma_b = \sigma_{yi} \sqrt{\frac{S_{x^2}}{S_{xx}}} = \sigma_{yi} \sqrt{\frac{N}{N\left(\sum_{i=1}^N x_i^2\right) - \left(\sum_{i=1}^N x_i\right)^2}} \quad (11)$$

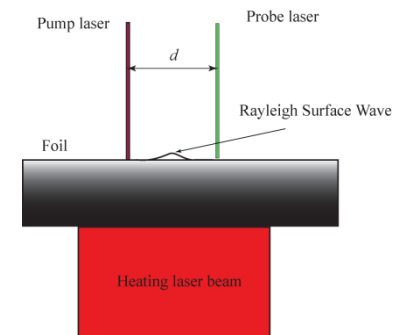
$$\sigma_{yi} = \sqrt{\frac{\sum_{i=1}^N [y_i - a - bx_i]^2}{N-2}} \quad (12)$$

Линейная регрессия: пример

В работе, [1] проводились измерения скорости поверхностной акустической волны (ПАВ) в сплаве платины при температуре 1070 К методом лазерного ультразвука (LU). В LU методе ПАВ возбуждаются импульсным лазером (pump laser), а детектируются при помощи принимающего лазера (probe laser) (См. схему в правом нижнем углу).



Принимающий лазер измеряет время прихода ПАВ по поверхности образца (Фиг. а). Путь волны определяется путем изменения расстояния между импульсным и детектирующим лазерами. График справа показывает результаты измерений: зависимость расстояния между лазерами от времени прихода ПАВ. Скорости продольных и поперечных волн определяются с использованием МНК, подгонкой экспериментальных данных линейной регрессией (Фиг. b) [1] К. Burgess, V. Prakapenka, E. Hellebrand, P. V. Zinin. “Elastic characterization of platinum/rhodium alloy at high temperature by combined laser heating and laser ultrasonic techniques”. *Ultrason.*, **54**, 963 (2014),



Метод наименьших квадратов: полином второго порядка

Запишем уравнение параболы в виде

$$y = a_2 x_i^2 + a_1 x_i + a_0$$

Минимум суммы находится из условия

$$\frac{\partial I}{\partial a_0} = -2 \sum_{i=1}^n [y_i - a_2 x_i^2 - a_1 x_i - a_0] = 0$$

$$I = \sum_{i=1}^n [y_i - a_2 x_i^2 - a_1 x_i - a_0]^2$$

$$\frac{\partial I}{\partial a_1} = -2 \sum_{i=1}^n x_i [y_i - a_2 x_i^2 - a_1 x_i - a_0] = 0$$

$$a_2 \bar{x}^4 + a_1 \bar{x}^3 + a_0 \bar{x}^2 = \overline{x^2 y}$$

$$\Rightarrow a_2 \bar{x}^3 + a_1 \bar{x}^2 + a_0 \bar{x} = \overline{xy}$$

$$\frac{\partial I}{\partial a_2} = -2 \sum_{i=1}^n x_i^2 [y_i - a_2 x_i^2 - a_1 x_i - a_0] = 0$$

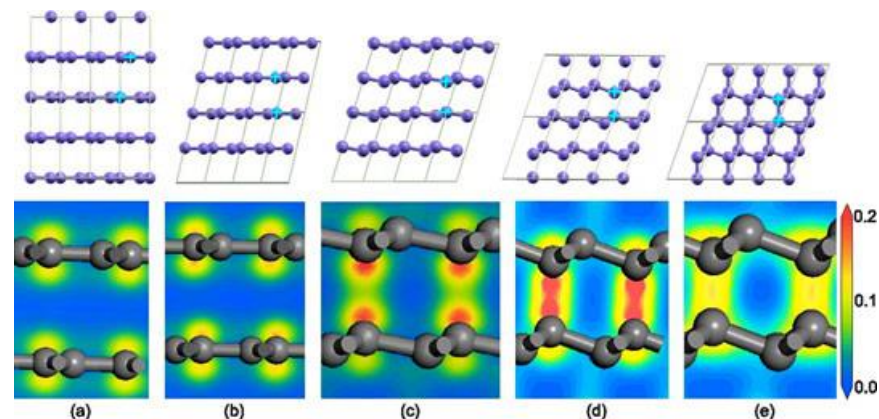
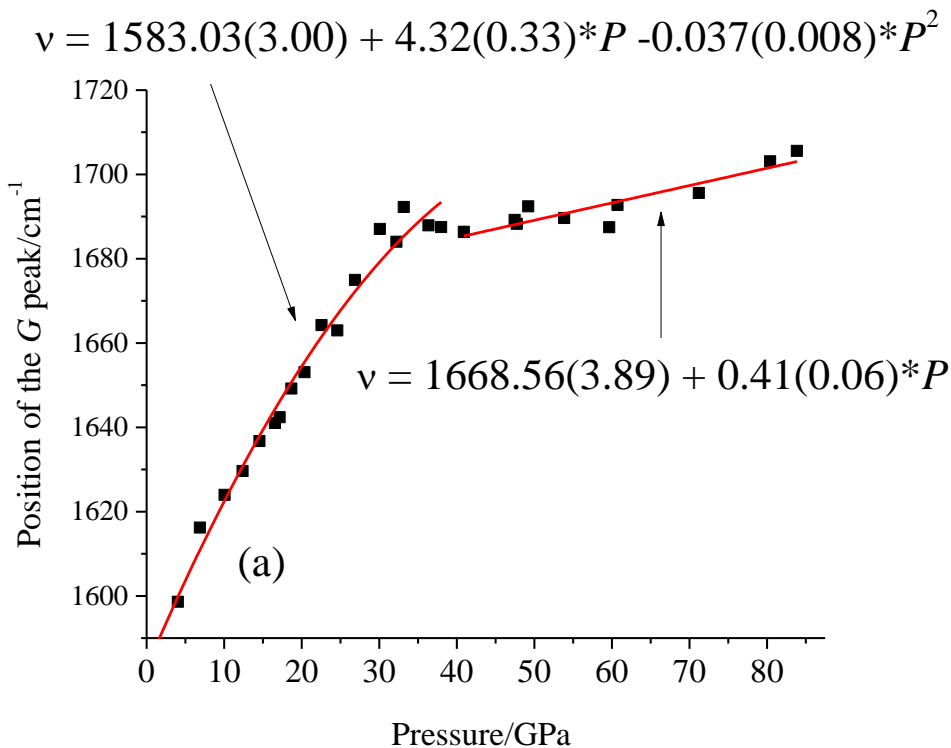
$$a_2 \bar{x}^2 + a_1 \bar{x} + a_0 = \bar{y}$$

Отсюда приходим к уравнениям

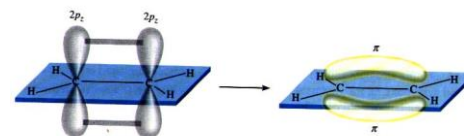
$$\begin{bmatrix} \bar{x}^4 & \bar{x}^3 & \bar{x}^2 \\ \bar{x}^3 & \bar{x}^2 & \bar{x} \\ \bar{x}^2 & \bar{x} & 1 \end{bmatrix} \begin{bmatrix} a_2 \\ a_1 \\ a_0 \end{bmatrix} = \begin{bmatrix} \overline{x^2 y} \\ \overline{xy} \\ \bar{y} \end{bmatrix}$$

Нетрудно заметить, что по мере повышения степени полинома функция аппроксимации приближается к фактическим данным, а при степени полинома, равной количеству отсчетов данных минус 1, вообще превращается в функцию интерполяции данных, что не соответствует задачам регрессии.

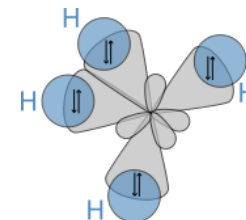
Полиномиальная регрессия: квадратичная регрессия



π -bonds



σ -bonds



В работе [1] измерялась позиция пика комбинационного рассеяния G графита и фазы высокого давления hp -C графита как функция давления. Отклонение зависимости от линейной в области давлений ниже 30 ГПа связывается с трансформацией графитовых π связей с sp^2 гибридизацией в алмазные σ связи с sp^3 гибридизацией.

[1] S. Odake, P. V. Zinin, E. Hellebrand, V. Prakapenka, *et al.* “Formation of the high pressure graphite and BC_8 phases in a cold compression experiment by Raman scattering”. *Journal of Raman Spectroscopy*, **44**, 1596 (2013).

Линеаризация в МНК

Существуют нелинейные зависимости, которые можно преобразовать в линейные.

Вид нелинейной зависимости	Получаемая линейная зависимость
$y = ax^b$	$\ln(y) = \ln(a) + b \cdot \ln(x) +$
$y = ae^{bx}$	$\ln(y) = \ln(a) + b \cdot x$
$y = x/(a+bx)$	$1/y = b + a/x$
$y = v + a/x$	$y = v + a \cdot z, z = 1/x$

Линеаризация в МНК. Закон Планка

Тепловое излучение обуславливается возбуждением частиц вещества при соударениях в процессе теплового движения или ускоренным движением зарядов (колебания ионов кристаллической решетки, тепловое движение свободных электронов и т.д.). Оно возникает при любых температурах и присуще всем телам. Характерной чертой теплового излучения является *сплошной спектр*.

Закон излучения Планка (формула Планка) - закон распределения энергии в спектре излучения равновесного при определённой температуре T . Был открыт М. Планком (M. Planck) в 1900 на основе гипотезы квантования энергии вещества. Планк моделировал вещество совокупностями гармонических осцилляторов различной частоты ν - резонаторов, испускающих и поглощающих излучение соответствующей частоты. Он предположил, что энергия вещества распределяется по резонаторам каждой частоты ν в виде дискретных порций $h\nu$ - квантов энергии (h - Планка постоянная).

$c_1 = 2\pi^5 h^6 c^2 / 15$, $c_2 = hc/k$, где h - постоянная Планка, c - скорость света, k - константа Больцмана: $c_1 = 3.7410 \cdot 10^{-12}$, вт см², $c_2 = 1.438$ см·град.

$$I(\lambda) = \frac{\varepsilon \cdot c_1}{\lambda^5 \left[\exp\left(\frac{c_2}{\lambda T}\right) - 1 \right]}$$

Линеаризация в МНК. Закон Вина

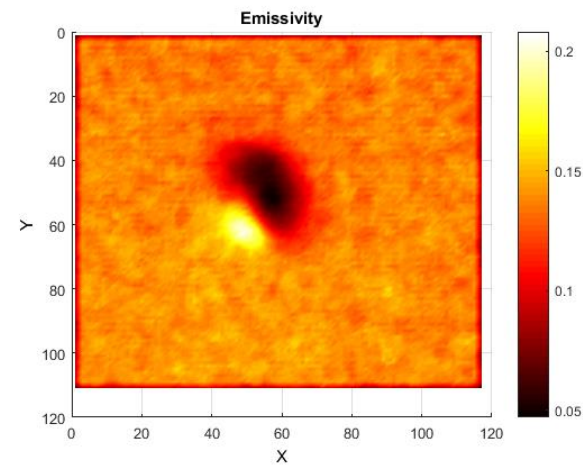
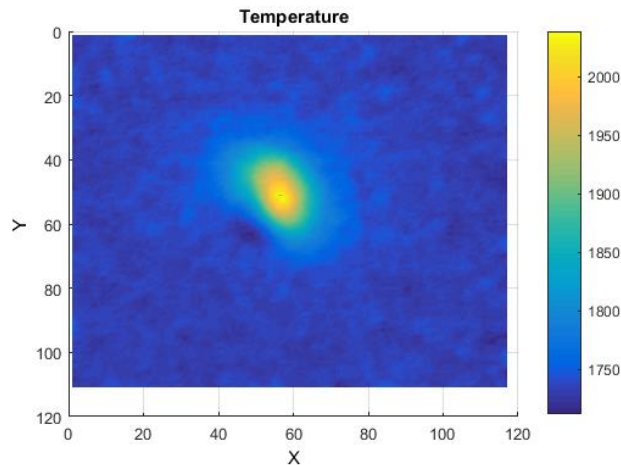
Если $\exp(c_2/\lambda T) \gg 1$, то закон Планка можно переписать

$$I(\lambda) = \frac{\varepsilon \cdot c_1}{\lambda^5 \left[\exp\left(\frac{c_2}{\lambda T}\right) - 1 \right]}$$

$$I(\lambda) = \frac{\varepsilon \cdot c_1}{\lambda^5} \exp\left(-\frac{c_2}{\lambda T}\right) \Rightarrow \frac{I(\lambda)\lambda^5}{c_1} = \varepsilon \cdot \exp\left(-\frac{c_2}{\lambda T}\right)$$

Последнее уравнение называется формулой Вина и может быть линеаризовано.

$$J = \frac{I(\lambda)\lambda^5}{c_1}, t = \frac{1}{\lambda} \Rightarrow \ln J = \ln \varepsilon - \frac{c_2}{T} t$$



Нелинейная регрессия

Применение МНК для подгонки некоторых функций требует коррекции. Таковыми являются функции, состоящие из ряда спадающих экспонент

$$f(x, \tilde{a}) = a_1 \exp(-a_2 \cdot x) + a_3 \exp(-a_4 \cdot x)$$

Рассмотрим случай одной экспоненты

$$f(x, \tilde{a}) = a_1 \exp(-a_2 \cdot x),$$

Для которой

$$S = \sum_{i=1}^n [y_i - a_1 \exp(-a_2 x_i)]^2$$

Минимизация МНК функционала обладает одним недостатком – вклад от точек с большим значением x_i экспоненциально мал по сравнению с точками с малым значением x_i . Чтобы сделать подгонку равномерной для всех значений x_i вводятся весовые коэффициенты.

$$S = \sum_{i=1}^n \left[\frac{y_i - a_1 \exp(-a_2 x_i)}{w_i} \right]^2, \quad w_i = y_i, \text{ or } \quad w_i = \sqrt{|y_i|}$$

Нелинейная регрессия

В общем случае, когда не удастся линеаризовать функционал S

$$S(\mathbf{a}) = \sum_{i=1}^n [y_i - f(x_i, \mathbf{a})]^2, \quad (1)$$

нахождение искомых значений набора параметров $a(j=1,2,3..p)$ и оценка доверительных интервалов для этих параметров осуществляется в два этапа/ Во-первых, находятся значения $a = \tilde{a}$, при которых функционал S достигает минимума.

$$S(\tilde{\mathbf{a}}) = \sum_{i=1}^n [y_i - f(x_i, \tilde{\mathbf{a}})]^2 = \min, \quad (2)$$

На втором этапе находим доверительные интервалы для значений $a = \tilde{a}$, при которых функционал S достигает минимума. Для этого строится поверхность $S(a)$ вблизи точки $a = \tilde{a}$. Тогда доверительная область: (а) интервал (одномерный случай, $p = 1$), (б) доверительная площадь (двумерный случай, $p = 2$), (с) доверительный объем ($p \geq 3$) находится внутри фигуры, определяемой уравнением:

$$S(a) = S(\tilde{\mathbf{a}}) \left[1 + \frac{p}{n-p} F(p, n-p, 1-\alpha) \right] \quad (3),$$

где α - доверительная вероятность, p - число параметров модели, n - число измерений, $F(p, n-p, 1-\alpha)$ - коэффициент F (Фишера) распределения. F коэффициент в виде $F(\alpha, p, n)$ можно посчитать в Матлабе, используя функцию `finv((α, p, n))`.

Нелинейная регрессия

В случае, когда измерения в каждой точке x_i проводились много раз, и известны нормальные отклонения в каждой точке $x_i - s_i$, тогда ищется минимум функции χ_2 :

$$\chi^2 = \sum_{i=1}^n \frac{[y_i - f(x_i, \mathbf{a})]^2}{\sigma_i^2} \quad (4)$$

Необходимо найти значения a_j , при которых функционал χ_2 достигает минимума.

$$\frac{\partial \chi^2(\mathbf{a})}{\partial a_j} = -2 \sum_{i=1}^n \left(\frac{[y_i - f(x_i, \mathbf{a})]}{s_i^2} \right) \left(\frac{\partial f(x_i, \mathbf{a})}{\partial a_j} \right), \quad j = 1, 2, \dots, p$$

$$\chi^2(\mathbf{a}) = \chi(\tilde{\mathbf{a}}) \left[1 + \frac{p}{n-p} F(p, n-p, 1-\alpha) \right],$$

где α - доверительная вероятность, p - число параметров модели, n - число измерений, $F(p, n, 1-\alpha)$ - коэффициент F (Фишера) распределения. F коэффициент в виде $F(\alpha, p, n)$ можно посчитать в Матлабе, используя функцию `finv((α, p, n))`.

Поиск глобального минимума

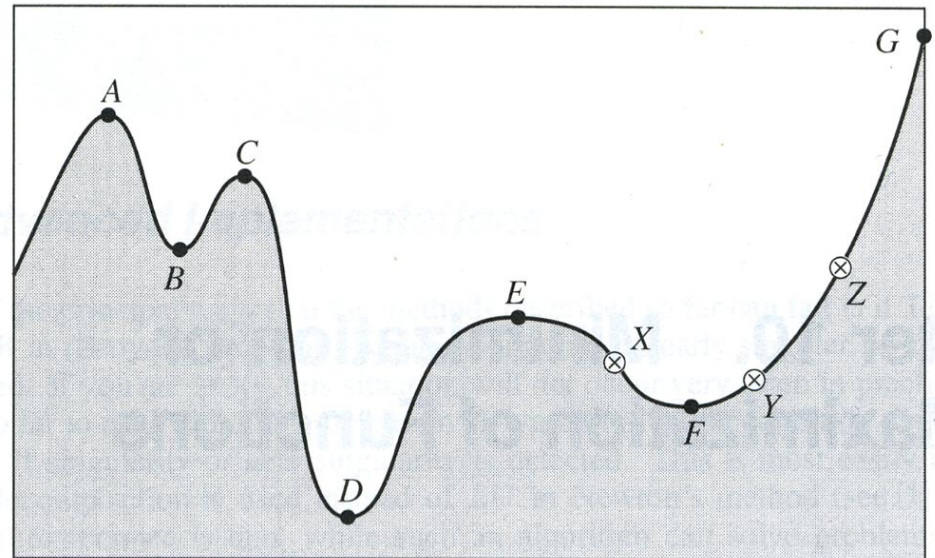
Функционалы (1) и (4) могут в общем случае иметь несколько минимумов. Решением МНК проблемы считаются значения a , при которых функционалы (1) и (4) находятся в глобальном (наименьшем) минимуме (точка D на рисунке).

Функция ошибки зависит от параметров модели МНК, и может рассматриваться как многомерная “поверхность”, на которой мы ищем минимум.

- В зависимости от сложности модели (т. е. число степеней свободы модели, M) поверхности ошибка может иметь несколько минимумов.

- Проблемой является нахождение набора параметров модели, при котором (1) и (4) находятся в глобальном, а не локальном минимуме.

$$\chi^2(a) = \chi(\tilde{a}) \left[1 + \frac{p}{n-p} F(p, n-p, 1-\alpha) \right] \quad (5)$$



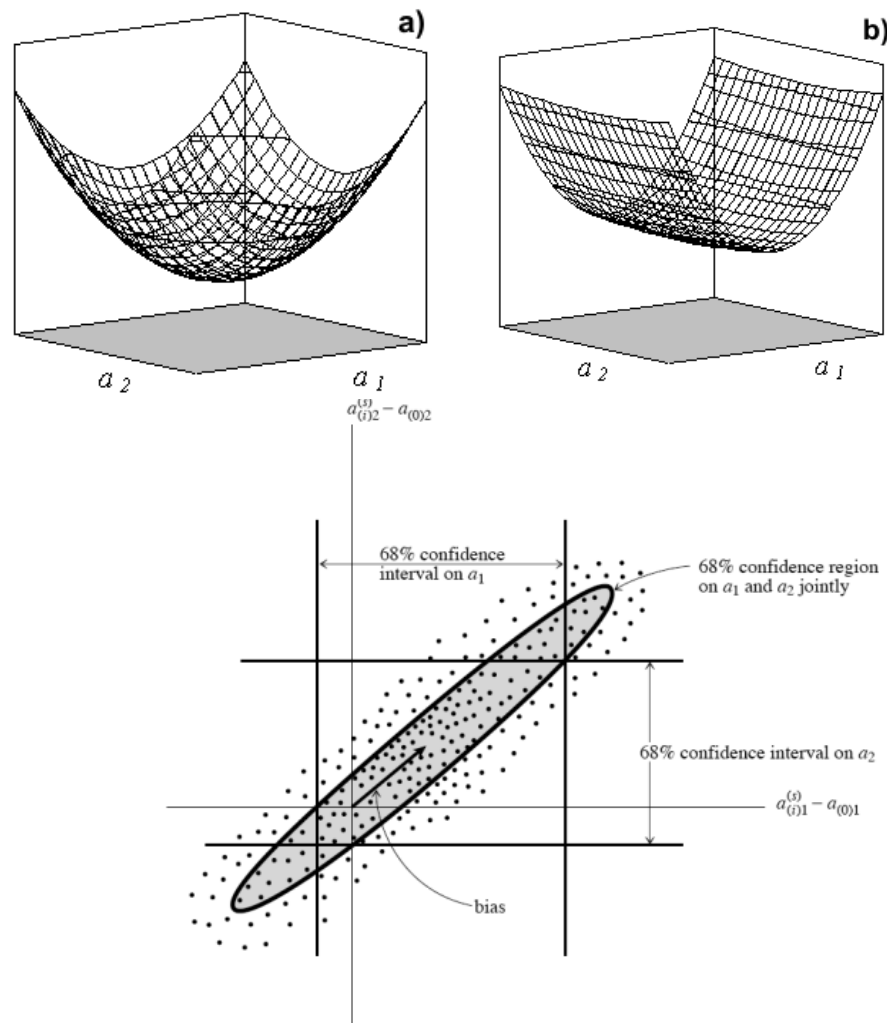
Поиск глобального минимума и доверительной области параметров a

Распределение вероятностей-это функция, определенная на p -мерном пространстве параметров a . Доверительный интервал-это регион, который содержит заданный процент от общего распределения по отношению к параметрической модели.

Алгоритм: (а) определяется уровень надежности или вероятности доверия (например, 68.3%, 90% и т. д.).

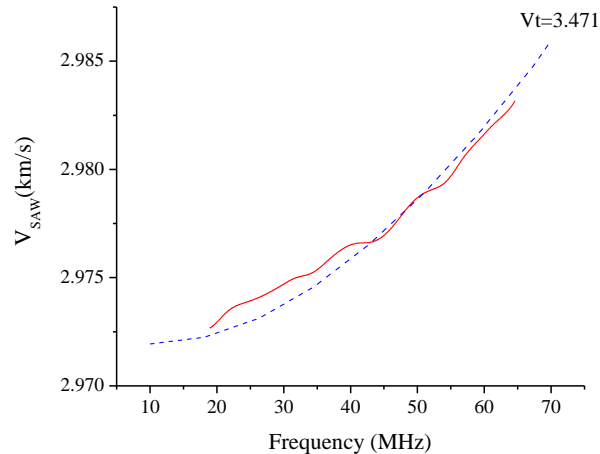
(б) Численным методом находится минимум функционала (1) или (4) и определяются оценки параметров модели a .

(с) Строится доверительная область по формулам (3) или (5), и определяются доверительные интервалы параметров a . В зависимости от числа параметров форма доверительной области имеет вид линии ($p = 2$), эллипса ($p = 3$), или эллипсоида ($p = 2$).



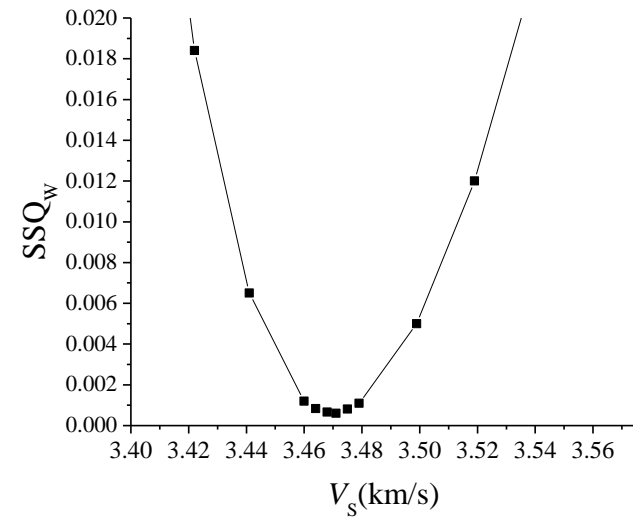
Метод наименьших квадратов. Пример: одномерная нелинейная регрессия

Пример взят из работы [1]. В этой работе измерения упругих свойств алмазоподобных пленок проводились с использованием рассеяния Мандельштама-Бриллюэна и лазерного ультразвука.



Дисперсионная кривая поверхностной акустической волны (SAW) образца с Cr-DLC покрытием. Сплошная линия результат измерения при помощи установки лазерного ультразвука, штриховая линия – получена с помощью метода наименьших квадратов.

$$SSQ = \sum_{i=1}^{i=1498} \left[(V_{SAW}(i) - V_{SAW}^{theor}(i)) / V_{SAW}(i) \right]^2$$

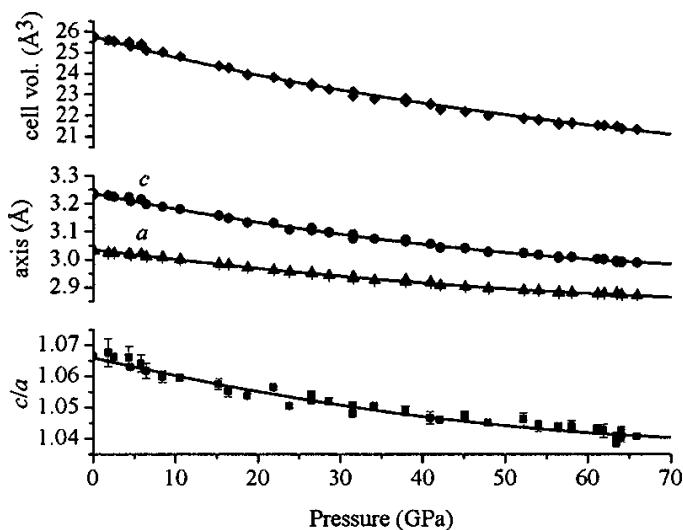


Сумму квадратов отклонений с весовыми коэффициентами как функция величины скорости сдвиговой волны в пленке V_s .

Поскольку количество точек измерения, измеренных методом LU, велико - 1498, ошибка измерений величины сдвиговой волны, полученная из подгонки дисперсионной кривой, мала ($\sim 0,13$ м/с), и мала по сравнению с точностью метода Мандельштама-Бриллюэна.

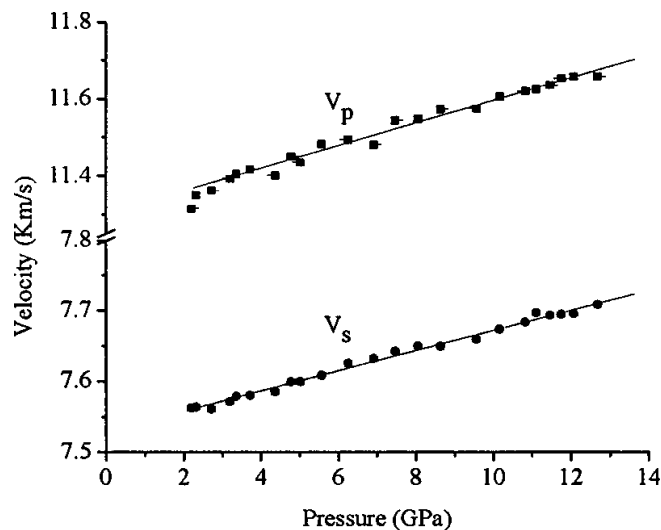
[1]. S. Berezina, P. V. Zinin, D. Schneider, *et al.* "Combining Brillouin spectroscopy and laser-SAW technique for elastic property characterization of thick DLC films". *Ultrasonics*. **43**, 87 (2004).

Метод наименьших квадратов. Пример: Двумерная нелинейная регрессия



Изменение параметров решетки TiB_2 с давлением.

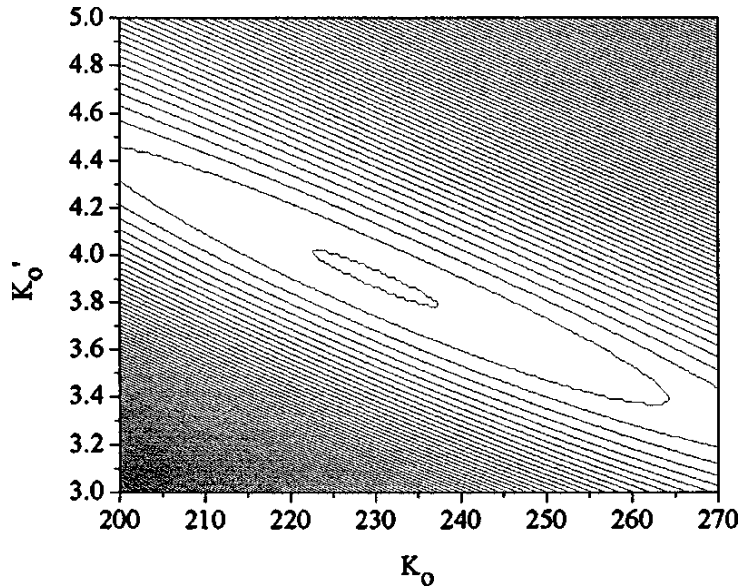
В работе [1] сжимаемость TiB_2 определялась отдельно из экспериментов по рентгеновской дифракции и на основе ультразвуковых измерений на образцах, загруженных в ячейки с алмазными наковальнями (до 65,9 ГПа) и помещенных в аппараты высокого давления (до 13,9 ГПа).



Зависимости скоростей продольных и сдвиговых волн в TiB_2 от давления.

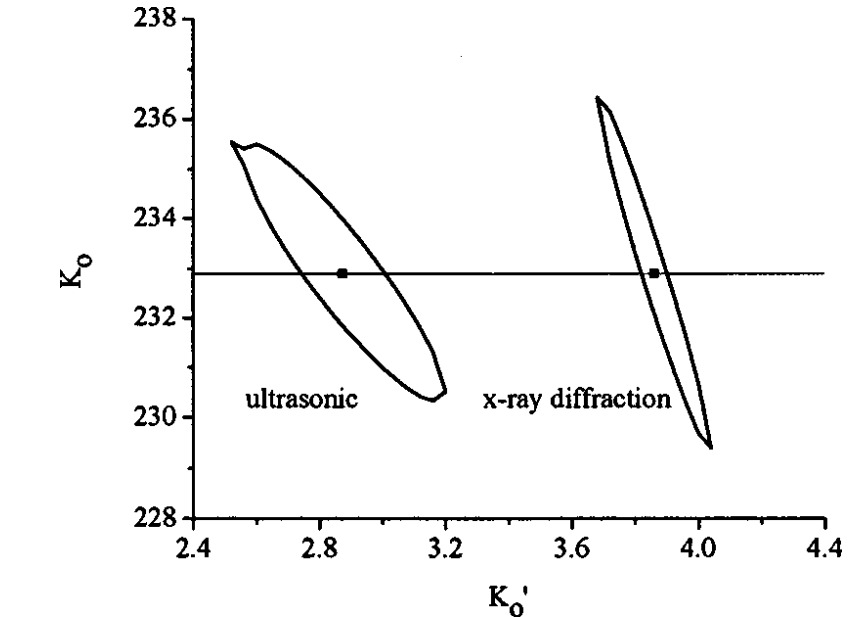
Целью статьи было получение величины объемного модуля K_0 и его производной по давлению, K'_0 , входящих в уравнение состояния третьего порядка, Берча-Марнагама.

Метод наименьших квадратов. Пример: Двумерная нелинейная регрессия



Контуры поверхности функции $S = \sum [V_i(\text{measured}) - V_i(\text{theoretical})]^2$, построенные в пространстве параметров K_o и K'_o , и показывающие доверительную область этих параметров. Данные были взяты из рентгенологических измерений.

Уравнение состояния третьего порядка Берча-Марнагама



Доверительные эллипсы, построенные для доверительной вероятности 95% и полученные на основании подгонки экспериментальных ультразвуковых и рентгеновских данных уравнением Берча-Марнагама.

$$P = \frac{3}{2} K_o f (1 + 2f)^{\frac{5}{2}} \left[1 + \frac{3}{2} (K'_o - 4) f \right], \quad f = \left[\left(\frac{V_o}{V} \right)^{\frac{2}{3}} - 1 \right]$$

Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия

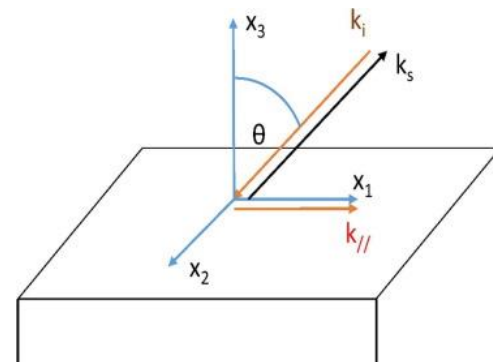
Рассеяние Манделъштама – Бриллюэна широко используется для определения упругих констант тонких пленок и веществ, находящихся при высоких давлениях. Описание процесса рассеяния света поверхностными тепловыми фононами возможно с использованием функции Грина, в которую упругие константы входят неявным образом [1]. Связь функции Грина с колебаниями поверхностных фононов следует из флуктуационно-диссипационной теоремы. Если амплитуда колебаний поверхностного фотона с частотой ω и волновым числом \mathbf{k}_{\parallel} , параллельным поверхностным образцу, обозначить как $u_3(\mathbf{k}_{\parallel}, \omega)$, то средний квадрат флуктуации термодинамической величины $\langle |u_3(\mathbf{k}_{\parallel}, \omega)|^2 \rangle$ пропорционален мнимой части функции Грина G_{33} [1]:

$$\langle |u_3(\mathbf{k}_{\parallel}, \omega)|^2 \rangle \propto \frac{T}{\omega} \text{Im} G_{33}(\mathbf{k}_{\parallel}, \omega)$$

Сечение рассеяния света поверхностными тепловыми фононами имеет вид

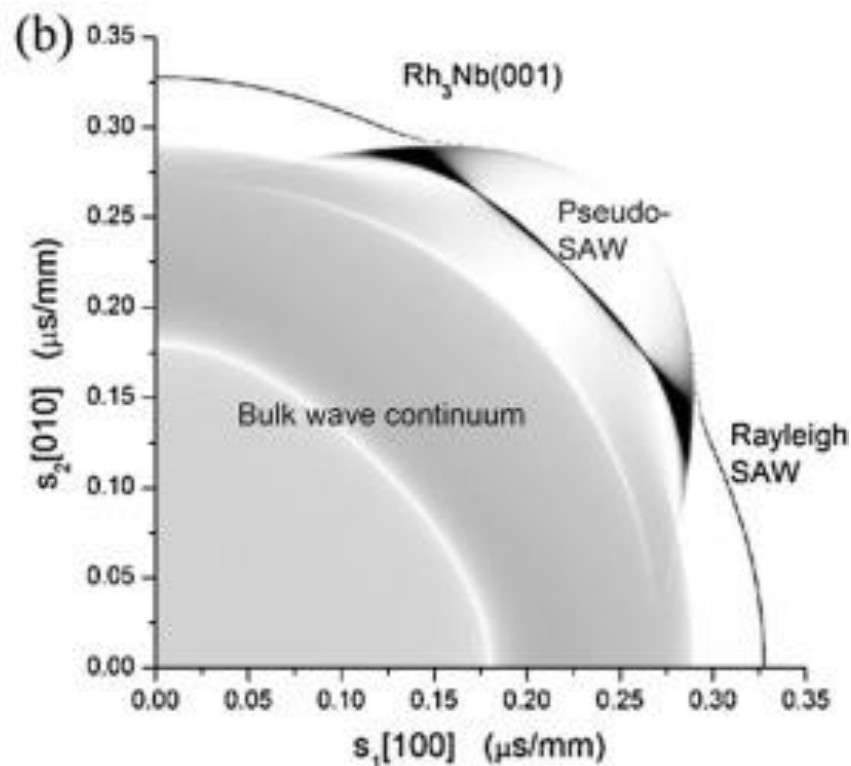
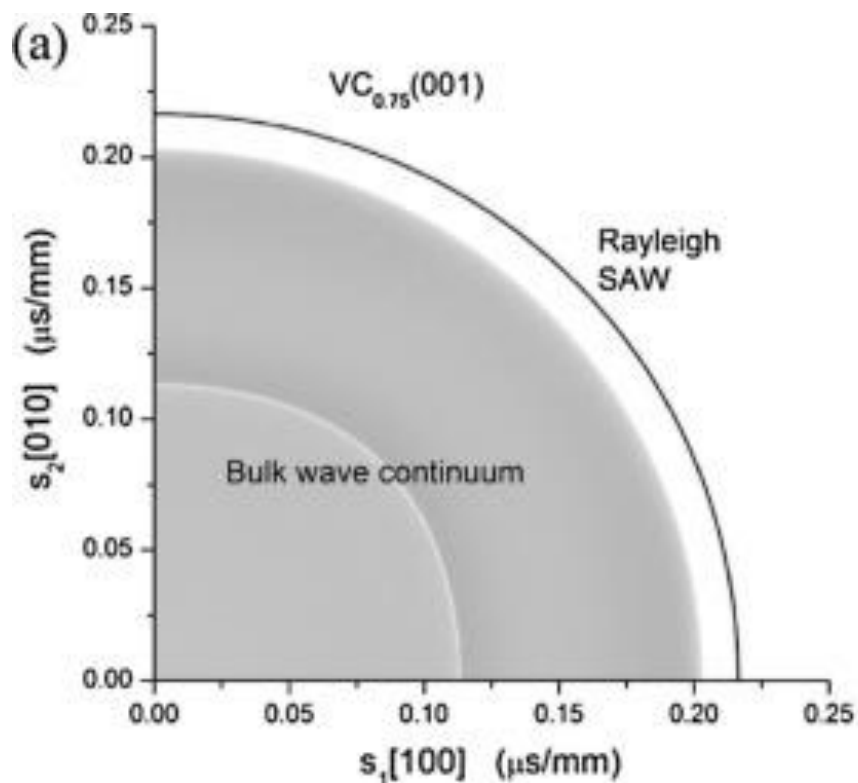
$$\frac{d^2 \sigma}{d\Omega d\omega} = \frac{AT}{\omega} \text{Im} G_{33}(\mathbf{k}_{\parallel}, \omega)$$

Векторная схема обратного поверхностного рассеяние Манделъштама – Бриллюэна.



M. G. Beghi, A. G. Every, V. Prakapenka and P. V. Zinin. “Measurements of the Elastic Properties of Solids by Brillouin Spectroscopy”, in T. Kundu ed., *Ultrasonic Nondestructive Evaluation: Engineering and Biological Material Characterization*. Taylor & Francis, N.Y., chapter 10, second edition, 540 (2012).

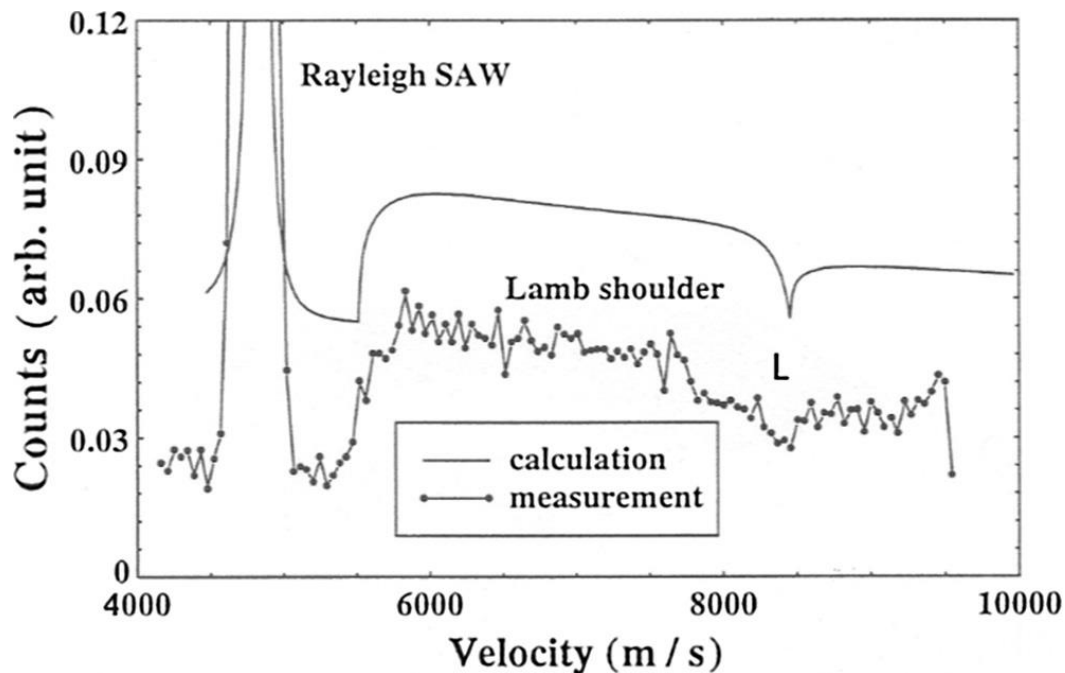
Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



Расчет функции $\text{Im}G_{33}(s_{\parallel})$ для (001) поверхности (a) $VC_{0.75}$ и (b) Rh_3Nb .

A.G. Every, C. Sumanya, B.A. Mathe, X. Zhang, J.D. Comins Optimized determination of elastic constants of crystals and their uncertainties from surface Brillouin scattering. *Ultrasonics*, **69**, 273 (2016).

Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



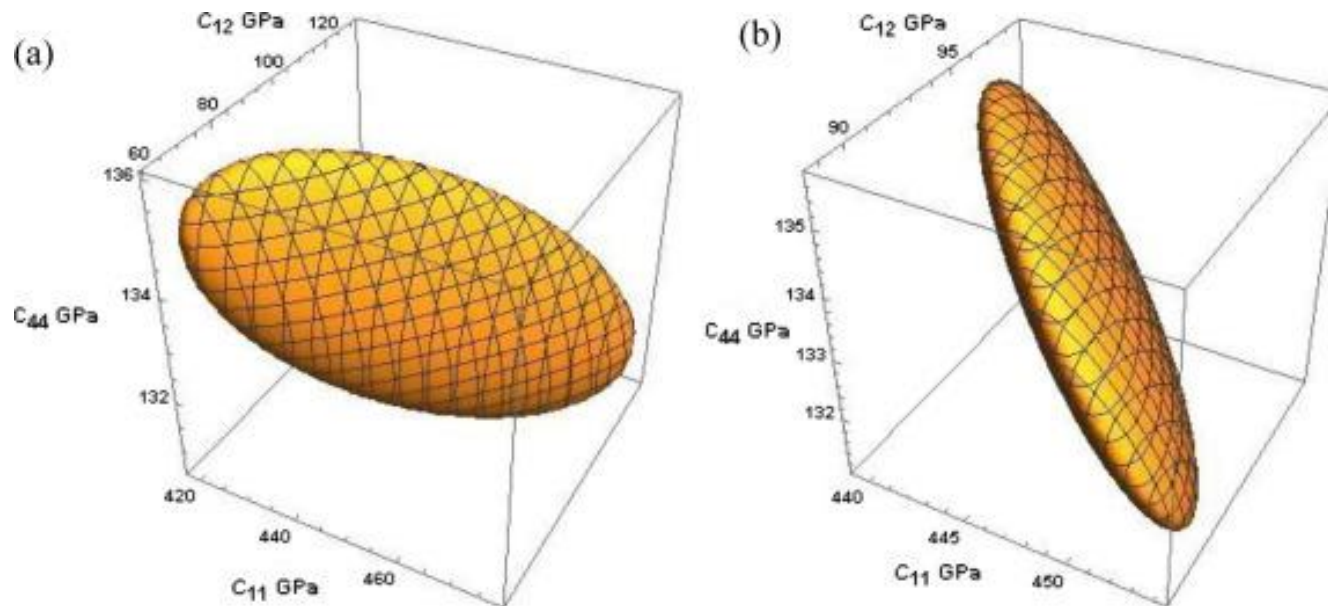
Спектр поверхностного рассеяние Манделъштама – Бриллюэна для $VC_{0.75}$ в направлении $[0\bar{1}1]$ на поверхности (110).

$$\chi^2(C_{11}, C_{22}, C_{33}) = \sum_{SAW_i} (V_i^{meas} - V_i^{calc})^2 + W \sum_{Li} (V_i^{meas} - V_i^{calc})^2$$

В работе [1] была разработана стратегия для оптимального определения трех упругих констант C_{11} , C_{12} и C_{44} кубического кристалла и их экспериментальных ошибок.

[1] A.G. Every, C. Sumanya, B.A. Mathe, X. Zhang, J.D. Comins Optimized determination of elastic constants of crystals and their uncertainties from surface Brillouin scattering. *Ultrasonics*, **69**, 273 (2016).

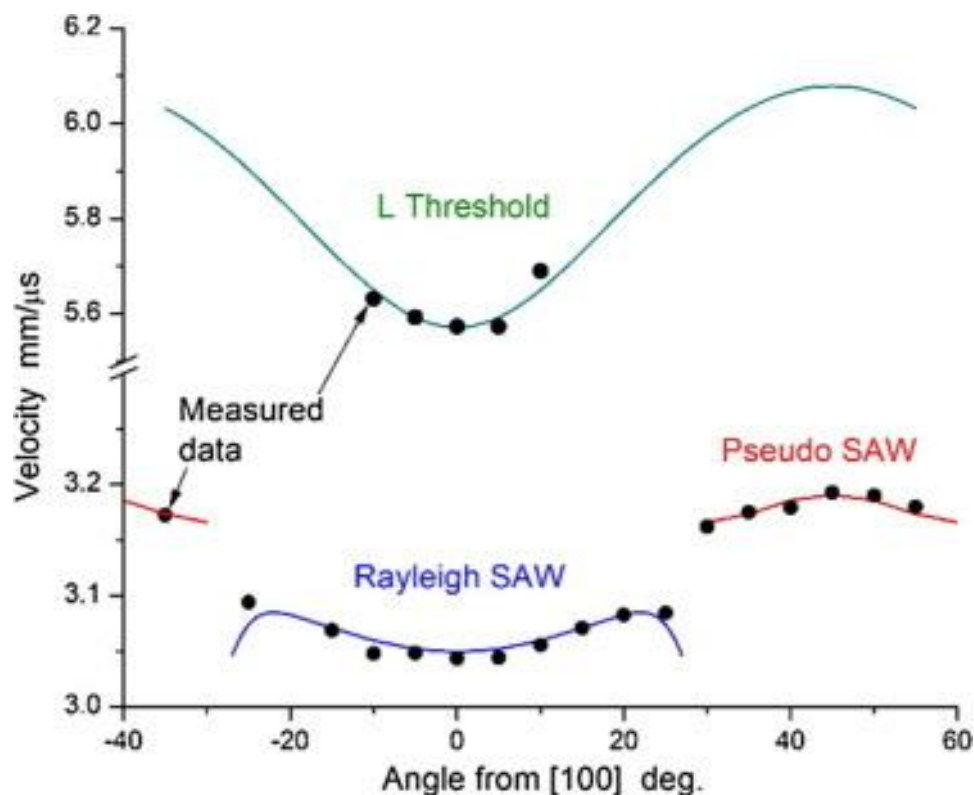
Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



Эллипсоиды доверительных областей для $VC_{0.75}$, (a) $W = 0.1$, вытянутый эллипсоида с основными осями в (C_{11}, C_{12}) плоскости, (b) $W = 50$, сплюснутый эллипсоид перпендикулярный направлению $(C_{11}, C_{12}, 2C_{44})$ [1].

Для $W = 0.1$, $C_{11} = 448.4 \pm 30$ ГПа и $C_{12} = 95.6 \pm 35$ ГПа очень слабо связаны. Принимая $W = 50$, величина $C_L = 403.5 \pm 1$ ГПа жестко ограничено. В работе [1] была выбрана стратегия, заключающаяся в том, чтобы варьировать значение W для получения неопределенность для $C_L \pm 4,2$ ГПа связанной с ошибкой прибора. Значение W , которое обеспечивает такой уровень точности является $W = 5$.

Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



Скорости Рэлеевской, псевдо-ПАВ, продольной волн как функции направления на (0 0 1) поверхности Rh_3Nb . Упругие константы, полученные с использованием МНК: $C_{11} = 368.5$ ГПа, $C_{12} = 186.0$ ГПа, $C_{44} = 161.4$ ГПа.

A.G. Every, C. Sumanya, B.A. Mathe, X. Zhang, J.D. Comins Optimized determination of elastic constants of crystals and their uncertainties from surface Brillouin scattering. *Ultrasonics*, **69**, 273 (2016).

Информационный критерий Акаики

Информационный критерий Акаике (AIC) — критерий, применяющийся исключительно для выбора из нескольких статистических моделей. Разработан в 1971 как информационный критерий Хироцугу Акаике и был предложен им в статье 1974 года [1]. Выражение для (AIC) имеет вид:

$$AIC = 2p + n \left[\ln \frac{(2\pi \cdot S_m)}{n} + 1 \right],$$

где n - число наблюдений в эксперименте, S - остаточная сумма квадратов и p - число параметров модели.

$$S_m = \sum_{i=1}^n [y_i - f(\tilde{a}, x_i)]^2 \quad \tilde{a} = [\tilde{a}_1, \tilde{a}_2, \tilde{a}_3, \dots, \tilde{a}_p]$$

Критерий не только вознаграждает за качество приближения, но и штрафует за использование излишнего количества параметров модели. Считается, что наилучшей будет модель с наименьшим значением критерия AIC.

[1] Akaike, H. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*. **19**. 716 (1974).

http://en.wikipedia.org/wiki/Akaike_information_criterion

Application of Akaike Information Criterion Statistics to High Pressure Equation of State Models

$$P = \frac{3}{2} K_o f (1 + 2f)^{\frac{5}{2}} \left[1 + \frac{3}{2} (K_o' - 4) f \right], \quad f = \left[\left(\frac{V_o}{V} \right)^{\frac{2}{3}} - 1 \right]$$

Model (Ultrasonic data,)	AIC	K_o	K_o'
Murnaghan equation	-127.8	226.1 ± 1.0	3.64 ± 0.19
3 rd order Birch-Murnaghan equation	-128.8	225.8 ± 1.0	3.73 ± 0.19

Model (x-ray data)	AIC	K_o	K_o'
Murnaghan equation	51.2	237.9 ± 12.1	2.62 ± 0.73
3 rd order Birch-Murnaghan equation	51.9	235.9 ± 11.6	2.92 ± 0.66

Model (Combined ultrasonic and x-ray data)	AIC	K_o	K_o'
Murnaghan equation	73.8	232.9 ± 5.5	2.91 ± 0.36
3 rd order Birch-Murnaghan equation	74.0	231.7 ± 5.4	3.15 ± 0.34

The Akaike information criterion (AIC) is calculated for different models so as to give the optimum result from the refinements for both the ultrasonic and x-ray measurements. The lower the value of AIC, the better the chosen model and adjustable parameter number. In all cases the pressure-density data was used in the fitting so as to give a justifiable comparison between models, both in the ultrasonic and x-ray fitting.

George M. Amulele,
Application of Akaike Information Criterion Statistics to High Pressure Equation of State Models,
 private communication, 2003

Работа над проектами

Алексей	Применение МНК для данных, описываемых законом Вина. Показать как пишется программа в Матлабе.
Александр	Рассчитать среднеквадратичное отклонения величины сопротивления, полученного Ван дер Пау методом.
Камиль	1. Применение Akaike information criterion для линейной и нелинейной регрессии. 2. Методы нахождения минимума функции нескольких переменных: Simplex method.
Лев	Применение МНК для данных, описываемых зависимости вида $y = a \exp(b \cdot x)$. (Матлаб)
Александра	Нахождения аналитического выражения ошибки измерения расстояния до сферического зеркала.
Демид	МНК метод для нахождения параметров нелинейной регрессии.
Юля	Применение МНК для данных, описываемых законом Вина. Получение аналитического выражение для среднеквадратичного отклонения величин излучения и температуры, полученных МНК методом.

Домашнее чтение

1. **Дрейпер Н., Г. Смит.** Прикладной регрессионный анализ, 3-е издание. В 2-х кн. М.: Финансы и статистика, 366 с., 1986.
2. **Диденко, Л.Г., Керженцев, В.В.** Математическая обработка и оформление результатов эксперимента Издательство: М.: МГУ Переплет: мягкий; 110 страниц; 1977.
3. **Яворский В. А.** Планирование научного эксперимента и обработка экспериментальных данных. Методические указания к лабораторным работам. Москва: Московский физико-технический институт (государственный университет). Факультет молекулярной и биологической физики. 2011.
4. **Сквайрс, Д.,** *Практическая физика.* Москва: Мир. 1971.
5. **Costa, K.D., S. Kleinstein, U. Hershberg** *Systems Biology: Biomedical Modeling Model Fitting and Error Estimation.* 2019.
http://clip.med.yale.edu/courses/brdu/Costa_ODE.pdf