

# Физические методы исследования состава и структуры веществ

## Часть II : Теория ошибок измерений

### Теория ошибок и статистическая обработка результатов эксперимента. Лекция II: Метод наименьших квадратов

Линейная регрессия, полиномиальная регрессия, нелинейная регрессия



Павел В. Зинин



# Доверительная вероятность и распределение Гаусса

Непрерывная случайная величина  $x$  называется распределённой по нормальному закону, если распределение её плотности вероятности имеет вид

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-a)^2}{2\sigma^2}}$$

где  $a$  – математическое ожидание случайной величины  $x$ ,  $\sigma$  – дисперсия. Вероятность  $P$  того, что результат измерения попадет в интервал  $[x_1, x_2]$ , равна:

$$P(x_1 \leq x \leq x_2) = \frac{1}{\sigma\sqrt{2\pi}} \int_{x_1}^{x_2} e^{-\frac{(x-a)^2}{2\sigma^2}}$$

При выбранной доверительной вероятности  $P$ ,  $x$  лежит в пределах

$$x = \bar{x} \pm \varepsilon(\alpha) \cdot \sigma,$$

где  $\varepsilon$  зависит от доверительной вероятности  $\alpha$ .

	Интервал	$\alpha$	$\varepsilon$
1	$\bar{x} - \sigma \leq x \leq \bar{x} + \sigma$	0.68	1
3	$\bar{x} - 2\sigma \leq x \leq \bar{x} + 2\sigma$	0.95	2,0
4	$\bar{x} - 2,6\sigma \leq x \leq \bar{x} + 2,6\sigma$	0.99	2,6
5	$\bar{x} - 3\sigma \leq x \leq \bar{x} + 3\sigma$	0.997	3,0
6	$\bar{x} - 3,3\sigma \leq x \leq \bar{x} + 3,3\sigma$	0.999	3,3

# 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 выборочного 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  $\Delta x$  доверительного интервала при выбранной доверительной вероятности  $\alpha$ . 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  $\alpha$  (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  $t(\alpha, n)/\sqrt{n}$  в выражении для доверительного интервала выбирается равным 1. Как далеко это от значения величины  $t(\alpha, n)/\sqrt{n}$ , которое должно использоваться для оценки доверительного интервала, при значениях доверительной вероятности или надежности 0.69, 0.9 и 0.95.

$n$	$t(\alpha, n)/\sqrt{n}$			
	$\alpha$	0.69	0.9	0.95
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

# Коэффициенты Стьюдента $t(\alpha, n)$ для доверительной вероятности $\alpha$ ( $n$ - количество измерений)

$$\lim_{n \rightarrow \infty} t(\alpha, n) = \varepsilon(\alpha)$$

$$\lim_{n \rightarrow \infty} \left[ t(\alpha, n) \frac{S_n}{\sqrt{n}} \right] = \varepsilon(\alpha) \cdot \sigma$$

$$\lim_{n \rightarrow \infty} [t(\alpha, n) \cdot S_{\bar{x}}] = \varepsilon(\alpha) \sigma$$

$$* \lim_{n \rightarrow \infty} S_{\bar{x}} = \sigma$$

$n$	$\alpha$			
	0,68	0,95	0,99	0,999
2	2,0	12,7	63,7	636,6
3	1,4	4,3	9,9	31,6
4	1,3	3,2	5,8	12,9
5	1,2	2,8	4,6	8,6
6	1,2	2,6	4,0	6,9
7	1,1	2,4	3,7	6,0
8	1,1	2,4	3,5	5,4
9	1,1	2,3	3,4	5,0
10	1,1	2,3	3,3	4,8
60	1,0	2,0	2,7	3,6

## Получение параметров из экспериментальных данных: косвенные измерения

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.

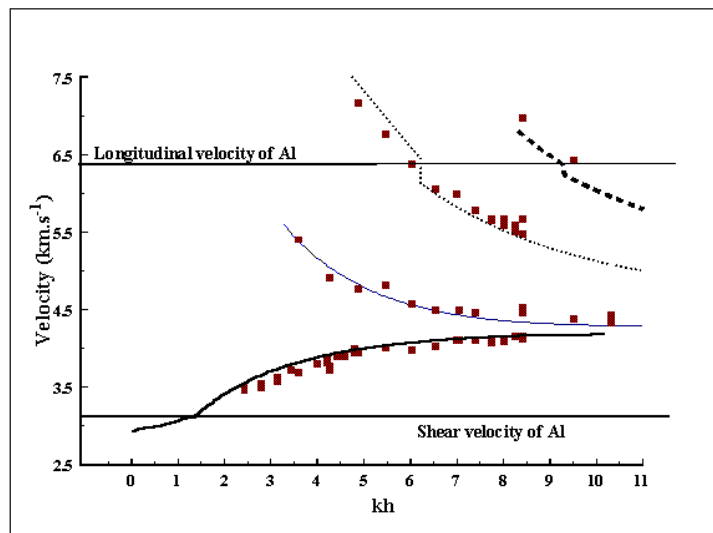


Рис. 1. Теоретические дисперсионные кривые поверхностных акустических волн в плёнке оксида алюминия на алюминиевой подложке. Квадраты – экспериментально померенные данные.

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 Манделъштама – Бриллюэна (РМБ), показаны на Рис. 1. Скорости продольной и поперечной акустических волн в пленке были получены методом наименьших квадратов (нелинейной подгонкой). Величина скорости поперечной волны  $v_T$  из барьерной пленки определяется с большей точностью ( $\pm 0.6$ ), чем значение скорости продольной волны  $v_L$  ( $\pm 2.1$ ).

# 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.

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

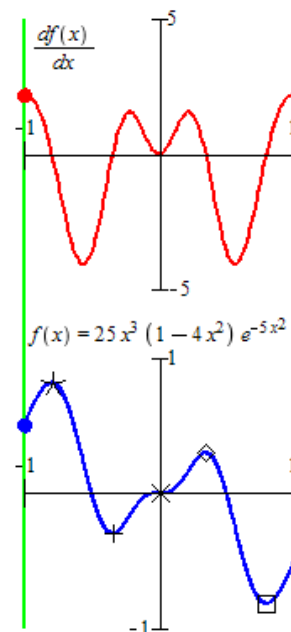


# Метод наименьших квадратов

Будем считать значения экспериментальных данных  $x_i$  ( $i = 1, 2, 3, \dots, n$ ) *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 нахождения значений параметров модели  $\mathbf{a} = [a_1, a_2, a_3, \dots, a_p]$ , при которых *сумма квадратов отклонений экспериментальных данных от теоретических минимальна*:

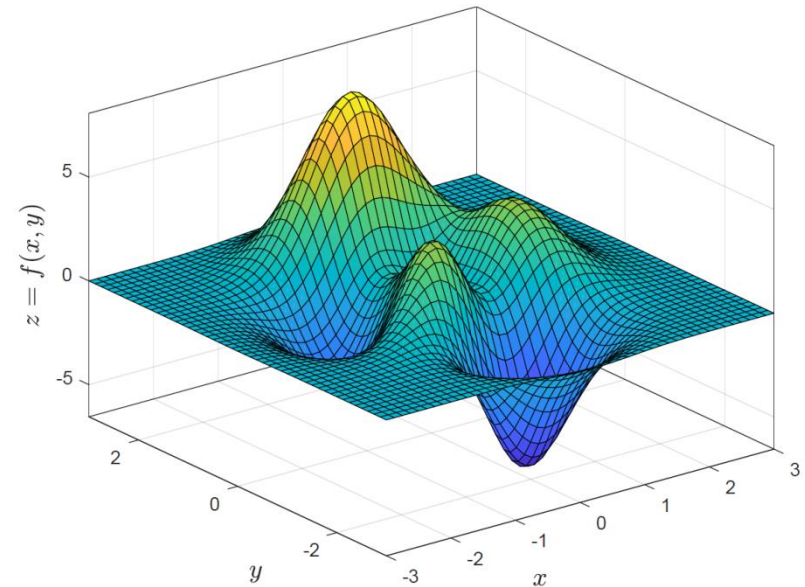
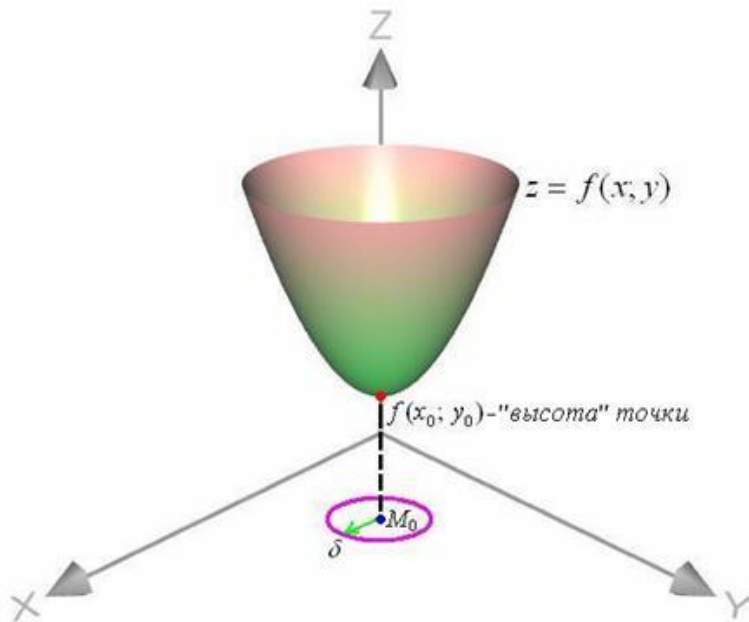
$$S = \sum_{i=1}^n [y_i - f(x_i, \mathbf{a})]^2 \quad (1)$$

**Экстремум** (*extremum* — крайний) в математике — *максимальное* или *минимальное* значение функции на заданном множестве. Точка, в которой достигается экстремум, называется *точкой экстремума*. Соответственно, если достигается минимум — точка экстремума называется *точкой минимума*, а если максимум — *точкой максимума*. В математическом анализе выделяют также понятие *локальный экстремум* (соответственно минимум или максимум).

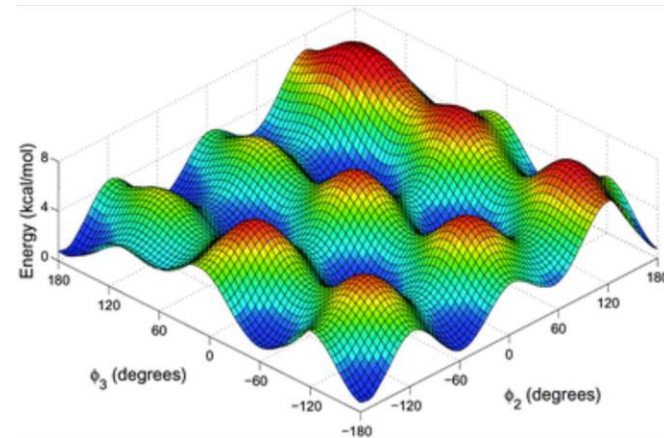


Пусть точка  $x_0$  является точкой экстремума функции  $f(x)$ , определённой в некоторой окрестности точки  $x_0$ . Тогда либо производная  $\frac{\partial f}{\partial x}$  не существует, либо  $\frac{\partial f}{\partial x}(x = x_0) = 0$ . Эти условия не являются достаточными, так, функция может иметь нуль производной в точке, но эта точка может не быть точкой экстремума, а являться, скажем, точкой перегиба, как точка  $(0,0)$  у функции  $f(x) = x^3$ .

# Поиск минимума функции двух переменных



Определение: функция достигает минимума в точке, если существует *хоть какая-то* - окрестность этой точки, в которой значение *высоты* *меньше* *всех* *остальных* значений.



# Метод наименьших квадратов: линейная регрессия

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

$$y = a \cdot x + b$$

Сумма  $S = \sum_{i=1}^n [y_i - ax_i - b]^2$

имеет минимум, когда первые производные по параметрам равна нулю:

$$\frac{\partial S}{\partial b} = -2 \sum_{i=1}^n [y_i - ax_i - b] = 0 \quad (3)$$

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

где  $n$  – число измерений. Условием минимума является равенства нулю системы  $p$  уравнений

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

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

$$\begin{aligned} nb + a \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ b \sum_{i=1}^n x_i + a \sum_{i=1}^n (x_i)^2 &= \sum_{i=1}^n x_i y_i \end{aligned} \quad (4)$$

# Метод наименьших квадратов: коэффициенты линейной регрессии

Решение системы линейных уравнений (4) имеет вид:

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

или

$$a = \frac{(\overline{xy})^2 - \bar{x} \cdot \bar{y}}{\overline{x^2} - (\bar{x})^2} \quad (6)$$
$$b = \frac{\overline{x^2 \cdot y} - \bar{x} \cdot \overline{xy}}{\overline{x^2} - (\bar{x})^2} \quad (7)$$

где

$$\overline{xy} = \frac{1}{n} \sum_{i=1}^N x_i y_i; \quad \bar{x} = \frac{1}{n} \sum_{i=1}^N x_i$$
$$\overline{x^2} = \frac{1}{n} \sum_{i=1}^N x_i^2; \quad \bar{y} = \frac{1}{n} \sum_{i=1}^N y_i$$

# Метод наименьших квадратов: коэффициенты линейной регрессии

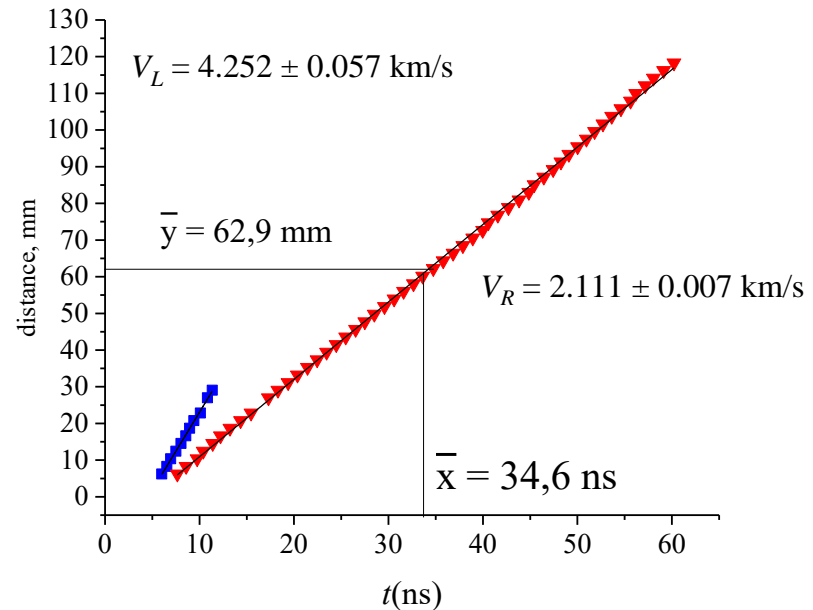
Коэффициенты регрессии можно записать в компактной форме

$$s_x = \frac{1}{n} \left( \sum_{i=1}^n x_i^2 \right) - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 = \overline{x^2} - (\bar{x})^2$$

$$a = \frac{\sum_{i=1}^n [x_i - \bar{x}] \cdot y_i}{ns_x}$$

$$b = \frac{\sum_{i=1}^n [\overline{x^2} - \bar{x}x_i] \cdot y_i}{ns_x}$$

(8)



K. 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).

# 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  $y_i$

$$a = \sum_{i=1}^n c_i \cdot y_i, \quad c_i = \frac{[x_i - \bar{x}]}{nS_x} \quad (9)$$
$$b = \sum_{i=1}^n d_i \cdot y_i, \quad d_i = \frac{[\bar{x}x_i - \bar{x}x]}{nS_x}$$

**Важно:** Статистическая теория говорит нам, что числа  $y_1, \dots, y_N$  представляют собой  $N$  результатов измерений одной и той же величины. Результат измерения каждого  $y_i$  распределен нормально около истинного значения  $a \cdot x_i + b$  с дисперсией  $\sigma_y$ . Тогда отклонения  $y_i - a \cdot x_i - b$  распределены нормально, причем все с одним и тем же центральным значением 0 и одной и той же дисперсией  $\sigma_y$ . 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  $x_i$ ,  $\sigma_{y_i}$ , используя выражение для стандартного отклонения  $S_n$ .

# Оценка коэффициентов линейной регрессии

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

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

$$\sigma_{y_i} = \sqrt{\frac{\sum_{i=1}^n [y_i - ax_i - b]^2}{n-2}} = \sqrt{\frac{S_o}{n-2}}$$

Поскольку параметр  $a$  можно представить в виде:

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

Тогда

$$\sigma_a^2 = \sum_{k=1}^N (c_k \sigma_{y_k})^2, \quad \sigma_b^2 = \sum_{k=1}^N (d_k \sigma_{y_k})^2$$

Считая что у всех измерений  $y_i$  дисперсия  $\sigma_{y_i}$  одинакова и равна  $\sigma_y$ , получим

$$\sigma_a^2 = \sigma_y^2 \sum_{k=1}^n c_k^2, \quad \sigma_b^2 = \sigma_y^2 \sum_{k=1}^n d_k^2$$

 (10)

Остаётся раскрыть суммы в выражении.

# Оценка коэффициентов линейной регрессии

$$c_i^2 = \frac{[x_i - \bar{x}]^2}{(ns_x)^2} = \frac{(x_i)^2 - 2x_i\bar{x} + (\bar{x})^2}{(ns_x)^2}$$

$$d_i^2 = \frac{[\bar{x}^2 - \bar{x}x_i]^2}{(ns_x)^2} = \frac{(\bar{x}^2)^2 - 2(\bar{x}^2)\bar{x}x_i + (\bar{x}x_i)^2}{(ns_x)^2}$$

$$= \frac{(\bar{x}^2)^2 - 2(\bar{x}^2)\bar{x}x_i + (x_i)^2(\bar{x})^2}{(ns_x)^2}$$

$$\sum_{i=1}^n c_i^2 = \frac{1}{ns_x}$$

$$\sum_{i=1}^n d_i^2 = \frac{\overline{x^2}}{ns_x}$$

$$\sum_{i=1}^n c_i^2 = \frac{\sum_{i=1}^n [(x_i)^2 - 2x_i\bar{x} + (\bar{x})^2]}{(ns_x)^2} = \frac{\overline{x^2} - 2(\bar{x})^2 + (\bar{x})^2}{ns_x^2} = \frac{\overline{x^2} - (\bar{x})^2}{ns_x^2} = \frac{1}{ns_x}$$

$$\sum_{i=1}^n d_i^2 = \frac{\sum_{i=1}^n [(\bar{x}^2)^2 - 2(\bar{x}^2)\bar{x}x_i + (x_i)^2(\bar{x})^2]}{(ns_x)^2} = \frac{n(\bar{x}^2)^2 - 2n(\bar{x}^2)(\bar{x})^2 + n(\bar{x})^2(\overline{x^2})}{(ns_x)^2} = \frac{(\bar{x}^2)s_x}{n(s_x)^2}$$



# Оценка СКО 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  $y_1, y_n$ .

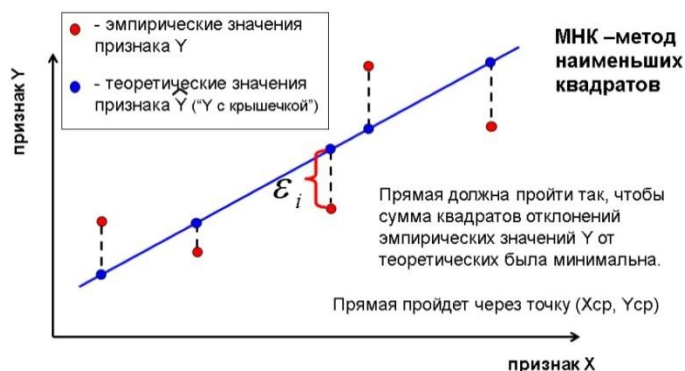
Можно показать, что фактор  $n$  в знаменателе необходимо заменить на  $(n - 2)$ . Таким образом, наш конечный ответ для погрешности в измерениях  $y_1, \dots, y_n$  есть

(10)

## Линейная регрессия

Модель – уравнение прямой –  $Y = a + b \cdot X$

Построение модели – расчет коэффициентов



$$\sigma_a = \sigma_y \sqrt{\frac{1}{ns_x}}$$

$$\sigma_b = \sigma_y \sqrt{\frac{x^2}{ns_x}}$$

$$s_x = \overline{x^2} - (\bar{x})^2$$

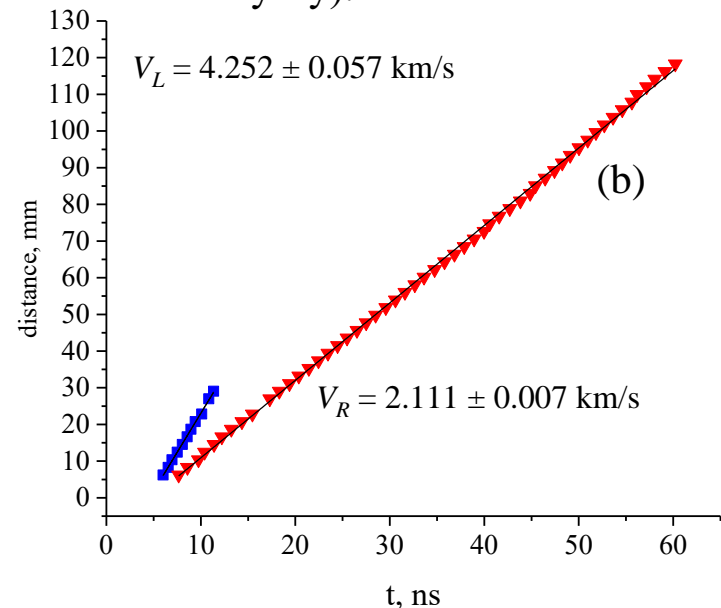
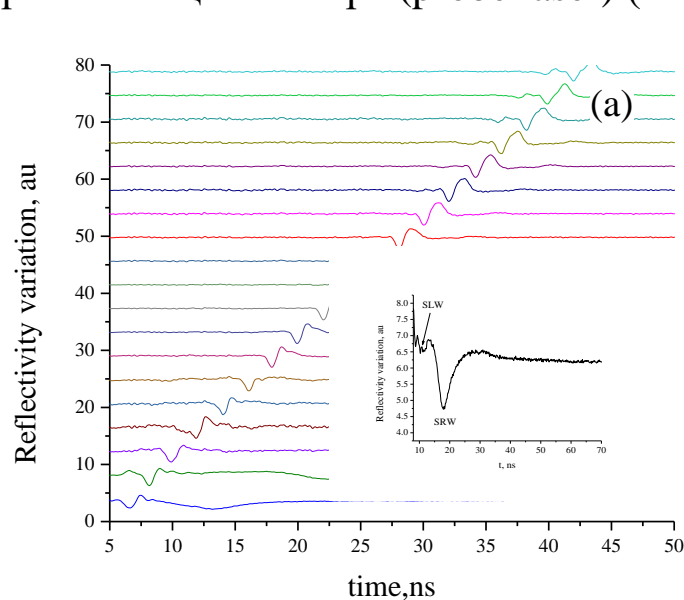
(11)

$$\sigma_y = \sqrt{\frac{\sum_{i=1}^n [y_i - ax_i - b]^2}{n - 2}}$$

(12)

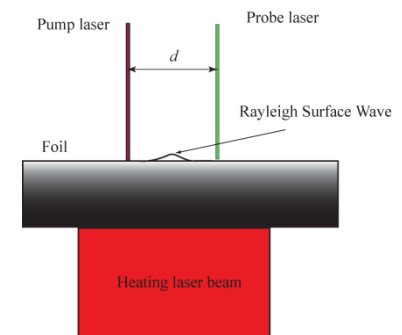
# Линейная регрессия: пример

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. a). 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] K. 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  $y = a_2 x_i^2 + a_1 x_i + a_0$

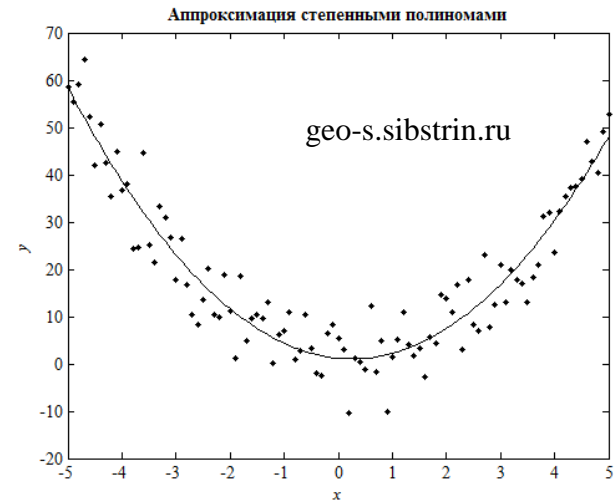
$$I = \sum_{i=1}^n [y_i - a_2 x_i^2 - a_1 x_i - a_0]^2$$

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

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

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

$\Rightarrow$



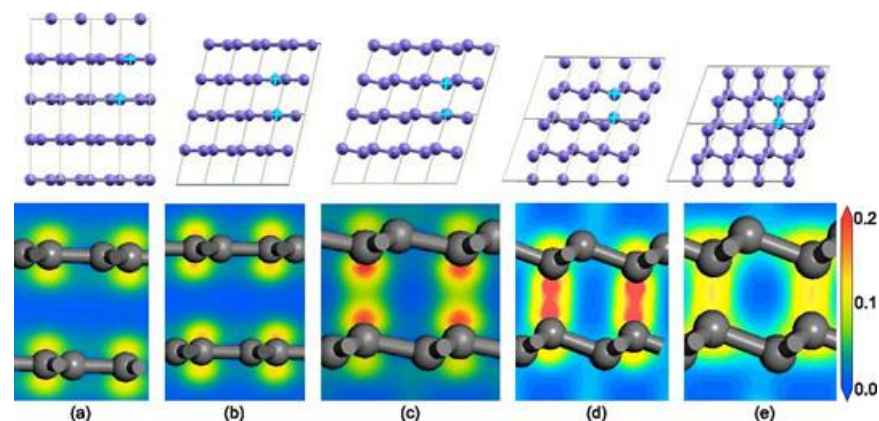
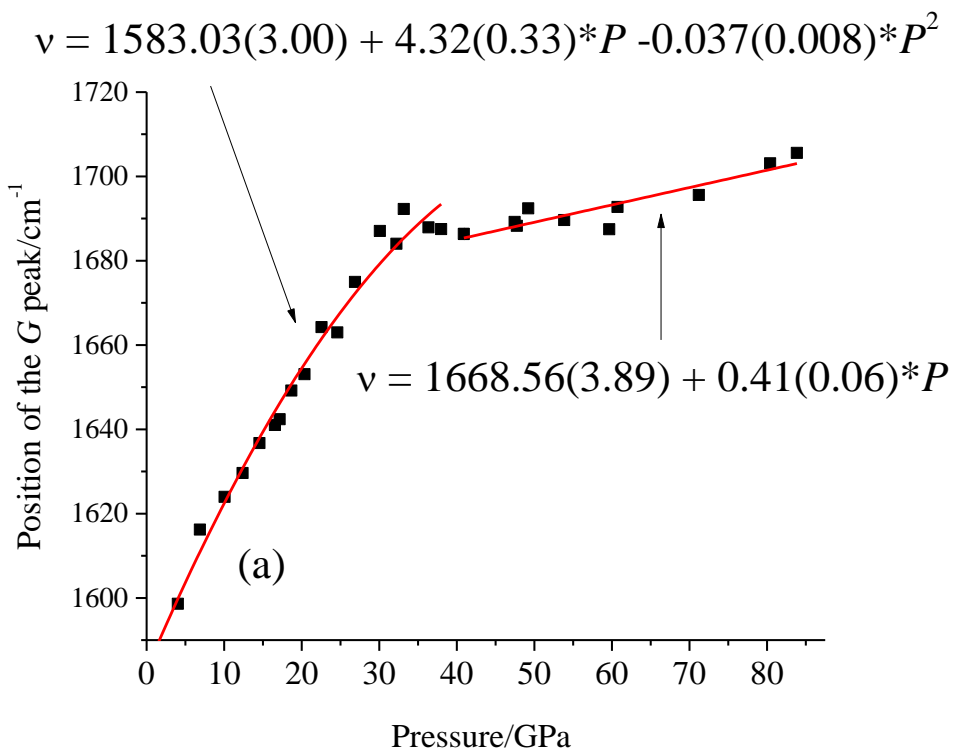
Нетрудно заметить, что по мере повышения степени полинома функция аппроксимации приближается к фактическим данным, а при степени полинома, равной количеству отсчётов данных минус 1, вообще превращается в функцию интерполяции данных, что не соответствует задачам регрессии.

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

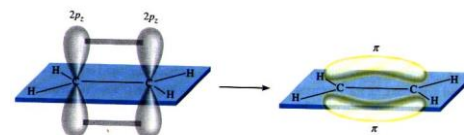
$$a_0 \bar{x} + a_1 \bar{x}^2 + a_2 \bar{x}^3 = \overline{xy}$$

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

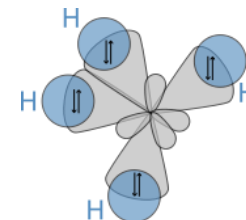
# Полиномиальная регрессия: квадратичная регрессия



$\pi$ -bonds



$\sigma$ -bonds



В работе [1] измерялась позиция пика комбинационного рассеяния  $G$  графита и фазы высокого давления  $hp$ -С графита как функция давления. Отклонение зависимости от линейной в области давлений ниже 30 ГПа связывается с трансформацией графитовых  $\pi$  связей с  $sp^2$  гибридизацией в алмазные  $\sigma$  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
$y = b \cdot x^a$	$\ln(y) = a \cdot \ln(x) + \ln(b)$
$y = b e^{ax}$	$\ln(y) = a \cdot x + \ln(b)$
$y = x / (ax + b)$	$1/y = a + b/x$
$y = b + a/x$	$y = b + a \cdot z, z = 1/x$

# 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  $\nu$  - резонаторов, испускающих и поглощающих излучение соответствующей частоты. Он предположил, что энергия вещества распределяется по резонаторам каждой частоты  $\nu$  в виде дискретных порций  $h\nu$  - квантов энергии ( $h$  - Планка постоянная).

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

$I$  – интенсивность излучений,  $\varepsilon$  - коэффициент излучения,  $\lambda$  - длина волны света,  $c_1 = 2\pi^5 h^6 c^2 / 15$ ,  $c_2 = hc/k$ , где  $h$  - постоянная Планка,  $c$  – скорость света,  $k$  – константа Больцмана:  $c_1 = 3.7410 \cdot 10^{-12}$ , вт см<sup>2</sup>,  $c_2 = 1.438$  см·град.

# Линеаризация в МНК. Закон Вина

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

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

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

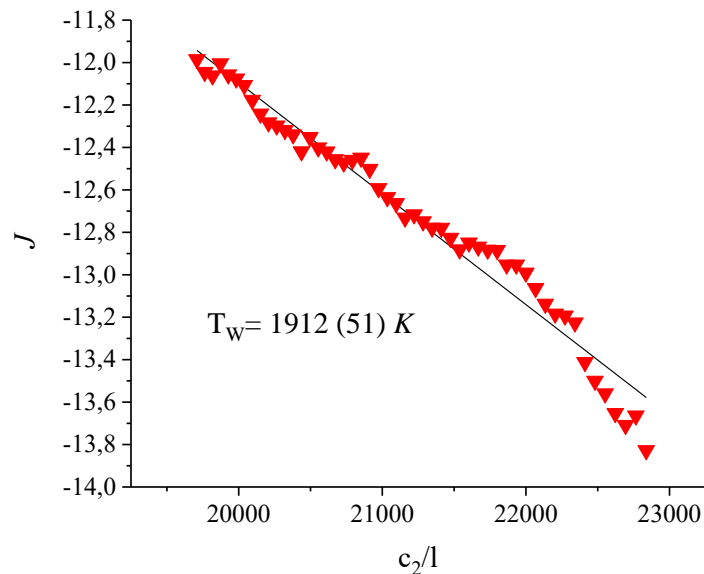
$$I(\lambda) = \frac{\varepsilon \cdot c_1}{\lambda^5 \left[ \exp\left(\frac{c_2}{\lambda T}\right) - 1 \right]}$$

Последнее уравнение называется формулой Вина и может быть линеаризовано.

Подгонка экспериментальных данных с использованием формулы Вина.

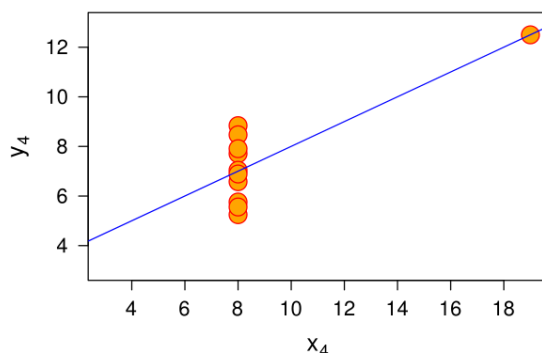
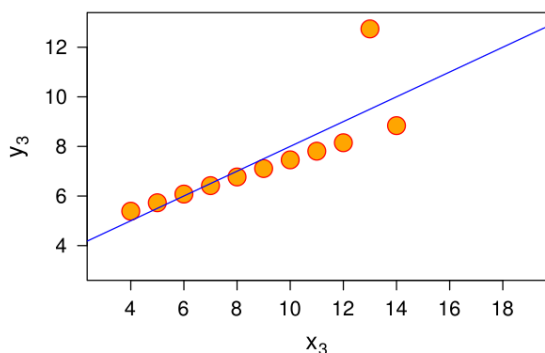
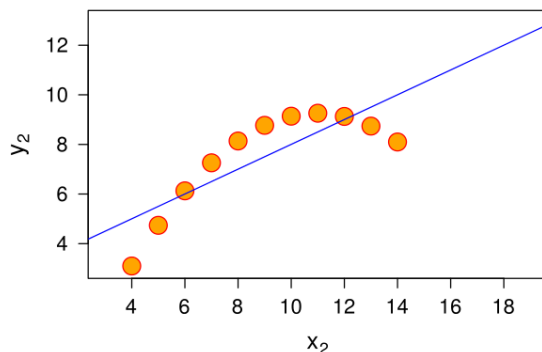
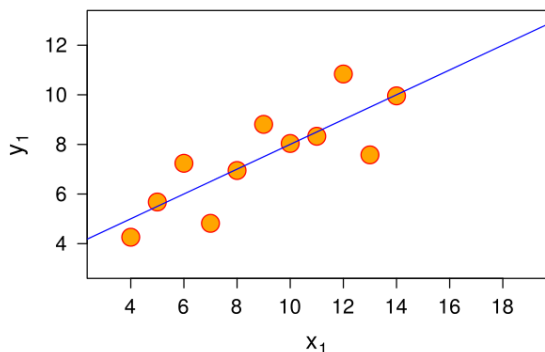
Видно, что модель не совсем адекватно описывает экспериментальные данные. Можно оценить адекватность модели с использованием коэффициента корреляции. Для хорошей модели  $R^2$  близок к 1 (0,99-0,95).

Визуальное наблюдение – хороший способ оценить адекватность модели  
визуальная оценка.



# Квартет Энскомба

Прямая линейной регрессии:  $y = 0,5x + 3$



**Квартет Энскомба** — четыре набора числовых данных, у которых простые статистические свойства идентичны, но их графики существенно отличаются. Каждый набор состоит из 11 пар чисел. Квартет был составлен в 1973 году английским математиком Ф. Дж. Энскомбом для иллюстрации важности применения графиков для статистического анализа и влияния выбросов значений на свойства всего набора данных F.J. Anscombe. “Graph in statistical analysis”. *The American Statistician*. 27(1) 17. 1973.

Среднее значение переменной  $x - 9,0$ ; среднее значение переменной  $y - 7,5$ ; дисперсия переменной  $x - 10,0$ ; дисперсия переменной  $y - 3,75$ ; корреляция между переменными  $x$  и  $y - 0,816$ .



# Нелинейная регрессия

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

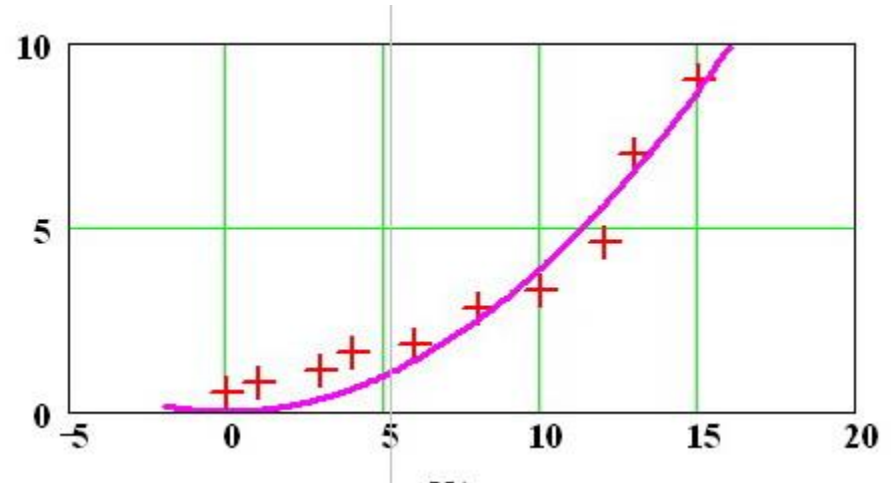
$$f(x, \tilde{a}) = a_1 \exp(-a_2 \cdot x) + a_3 \exp(-a_4 \cdot x)$$

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 искоемых значений набора параметров  $a(j=1,2,3..p)$  и оценка доверительных интервалов для этих параметров осуществляется в два этапа. В-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,  $p = 1$ ), (b) 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  $\alpha$ - 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((alpha), 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  $\chi_2$ :

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

Необходимо найти значения  $a_j$ , при которых функционал  $\chi_2$  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  $\alpha$ - 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(( $\alpha, p, n$ ))`.

# Поиск глобального минимума

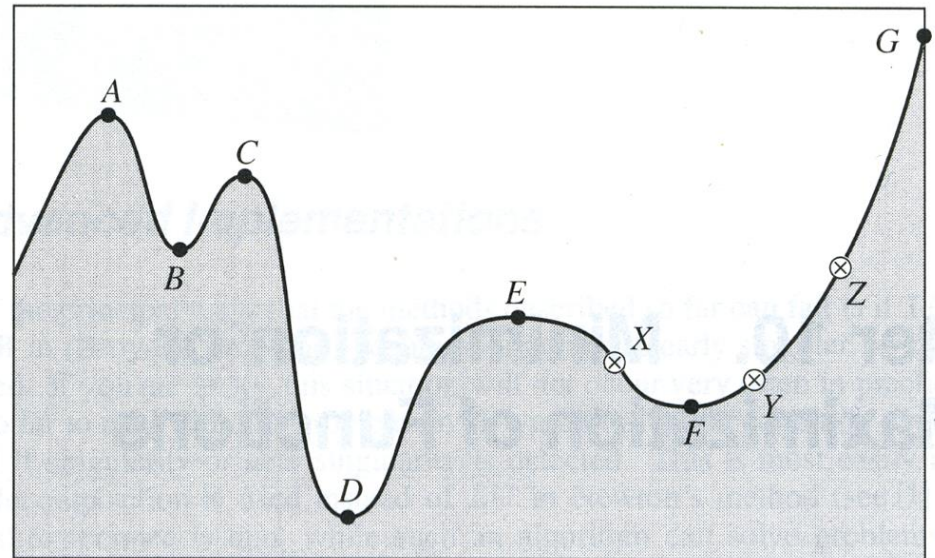
Функционалы (1) и (4) могут в общем случае иметь несколько минимумов. Решением МНК проблемы считаются значения  $a$ , при которых функционалы (1) и (4) находятся в глобальном (наименьшем) минимуме (точка D на рисунке).

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.

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



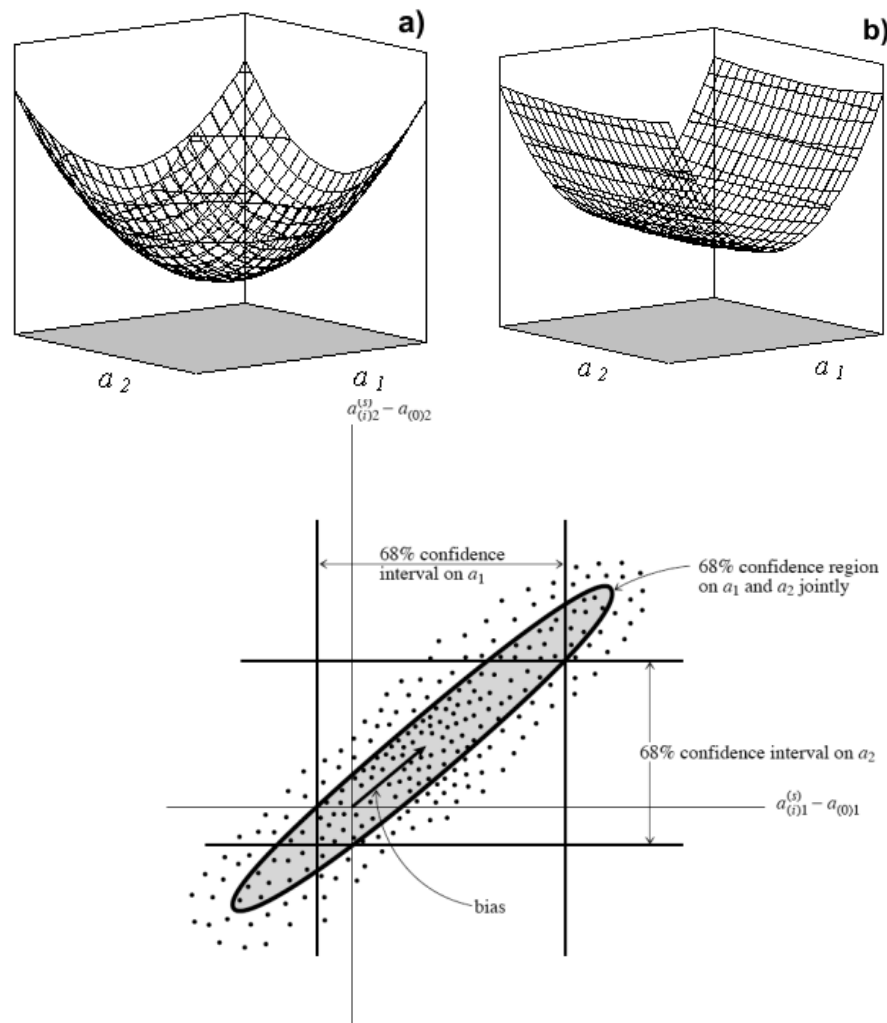
# Поиск глобального минимума и доверительной области параметров $a$

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.

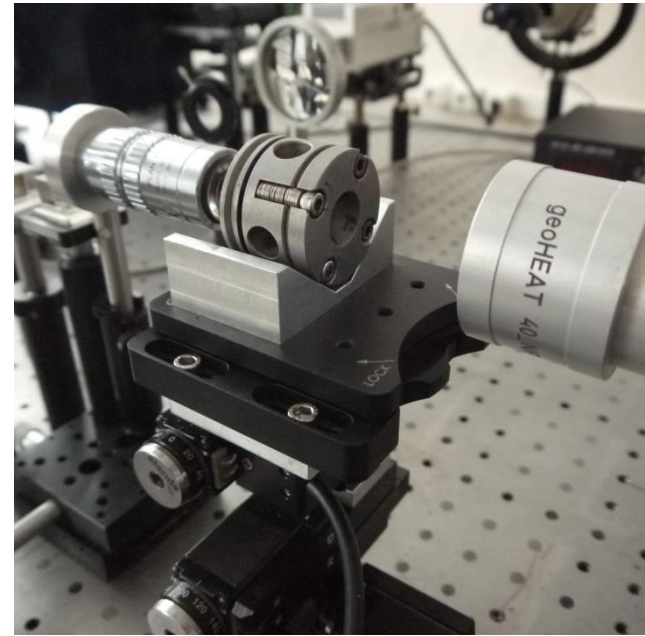
Алгоритм: (а) определяется уровень надежности или вероятности доверия (например, 68.3%, 90% и т. д.).

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

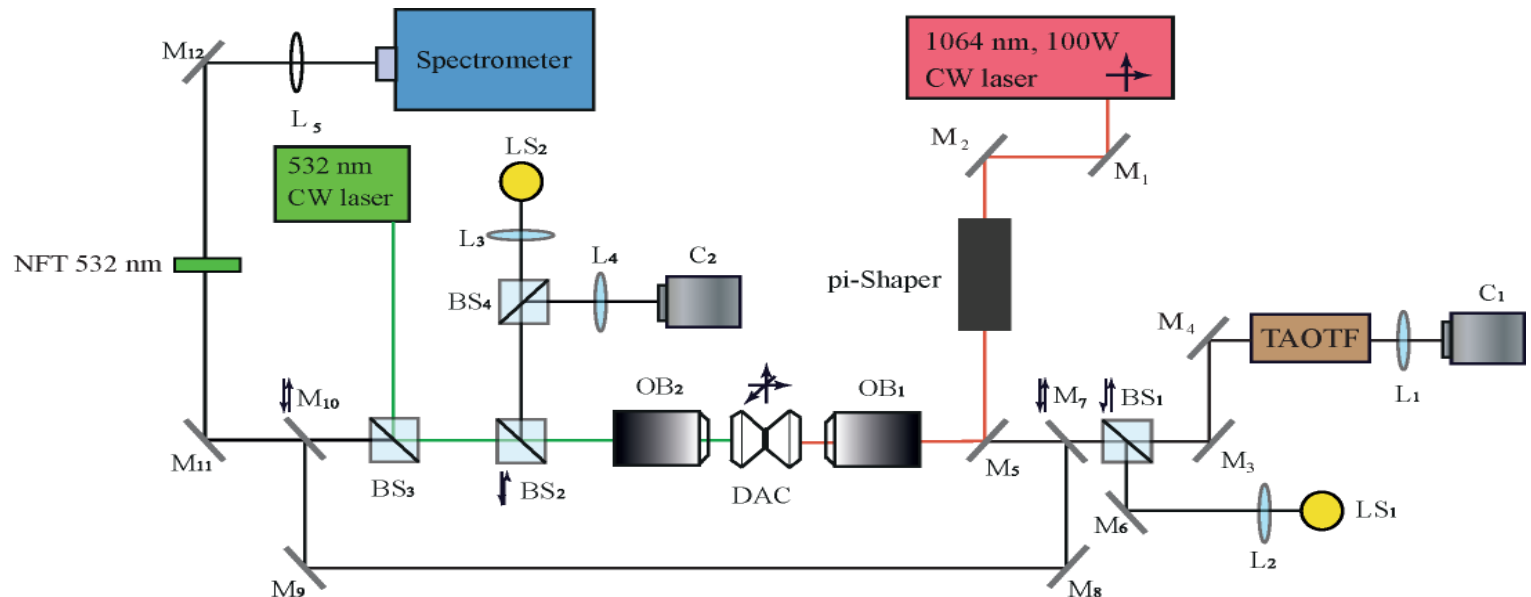
(с) Строится доверительная область по формулам (3) или (5), и определяются доверительные интервалы параметров  $a$ . В зависимости от числа параметров форма доверительной области имеет вид линии ( $p = 2$ ), эллипса ( $p = 3$ ), или эллипсоида ( $p = 2$ ).



# Photo of the laser heating in a diamond anvil cell system combined with TAOTF



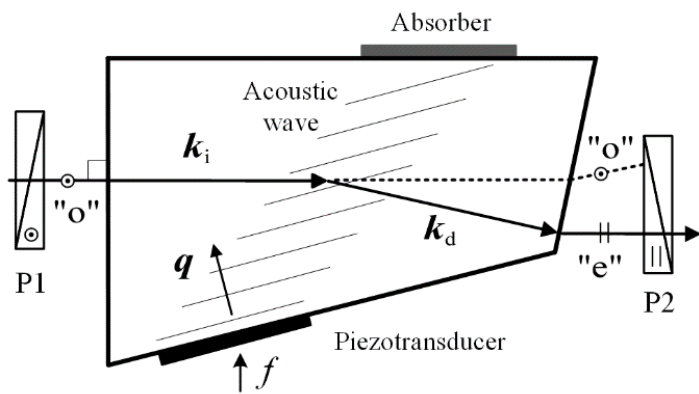
The sketch of the LH-TAOTF system.





# Imaging system based on a tandem acousto-optical tunable filter for measurements of the temperature distribution

The main component of the system is an imaging tandem acousto-optical tunable filter (TAOTF) synchronized with a video camera. A set of TAOTF spectroscopic images (up to a few hundreds) is taken by the TAOTF imaging system in order to fit the measured spectral curves in each pixel to the Planck radiation function and determine the temperature and emissivity of the sample using the gray body approximation. It was experimentally shown that this technique provides aberration-free spectral imaging suitable for precise multispectral imaging radiometry (MIR).



Schematic of non-collinear AOTF

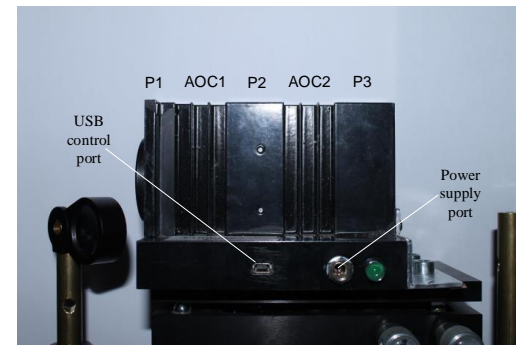
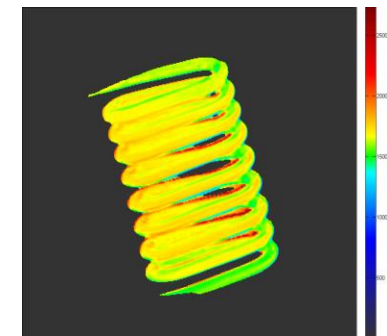


Photo of the TAOTF



A. S. Machikhin, P. V. Zinin, A. V. Shurygin, D. D. Khokhlov. *Optics Letters*, **41**(5), 901-904 (2016).

TAOTF spectroscopic image at  $\lambda = 800$  nm; (b) derived 2-D temperature distribution.

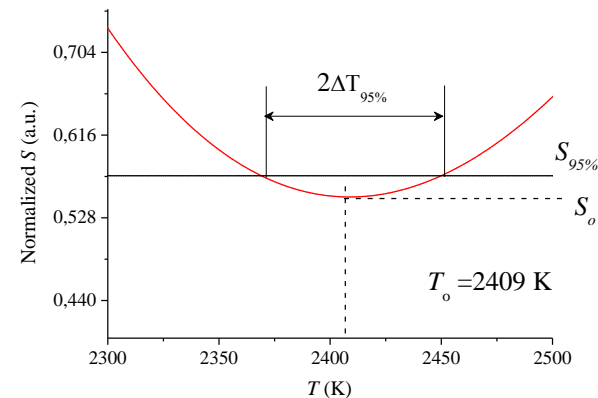
## 2-D non-linear least square fitting of the Plank's curve

The correct method of determining  $T$  from the experimentally determined  $I(\lambda_k)_{\text{corrected}}$  data is to find values of  $T_0$  and  $\varepsilon_0$  at which the function,

$$S(\varepsilon, T) = \sum_{i=1}^n [I(\lambda_i) - \varepsilon g(\lambda_i, T)]^2, \quad g(\lambda, T) = \frac{c_1}{\lambda^5 \left[ \exp\left(\frac{c_2}{\lambda T}\right) - 1 \right]}$$

has a minimum (2-D non-linear least square fitting). To decrease the effect of statistical error on temperature determination, a more stable least-squares fitting procedure was introduced. It is based on the fact that emissivity  $\varepsilon$  in eq. (1) is a linear parameter. We know that if the function  $S(T, \varepsilon)$  has a minimum at  $T_0$  and  $\varepsilon_0$ , then the following conditions should be satisfied:  $\partial S / \partial \varepsilon |_{\varepsilon=\varepsilon_0, T=T_0} = 0$ ;  $\partial S / \partial T |_{\varepsilon=\varepsilon_0, T=T_0} = 0$ . The first equation gives the value of  $\varepsilon$

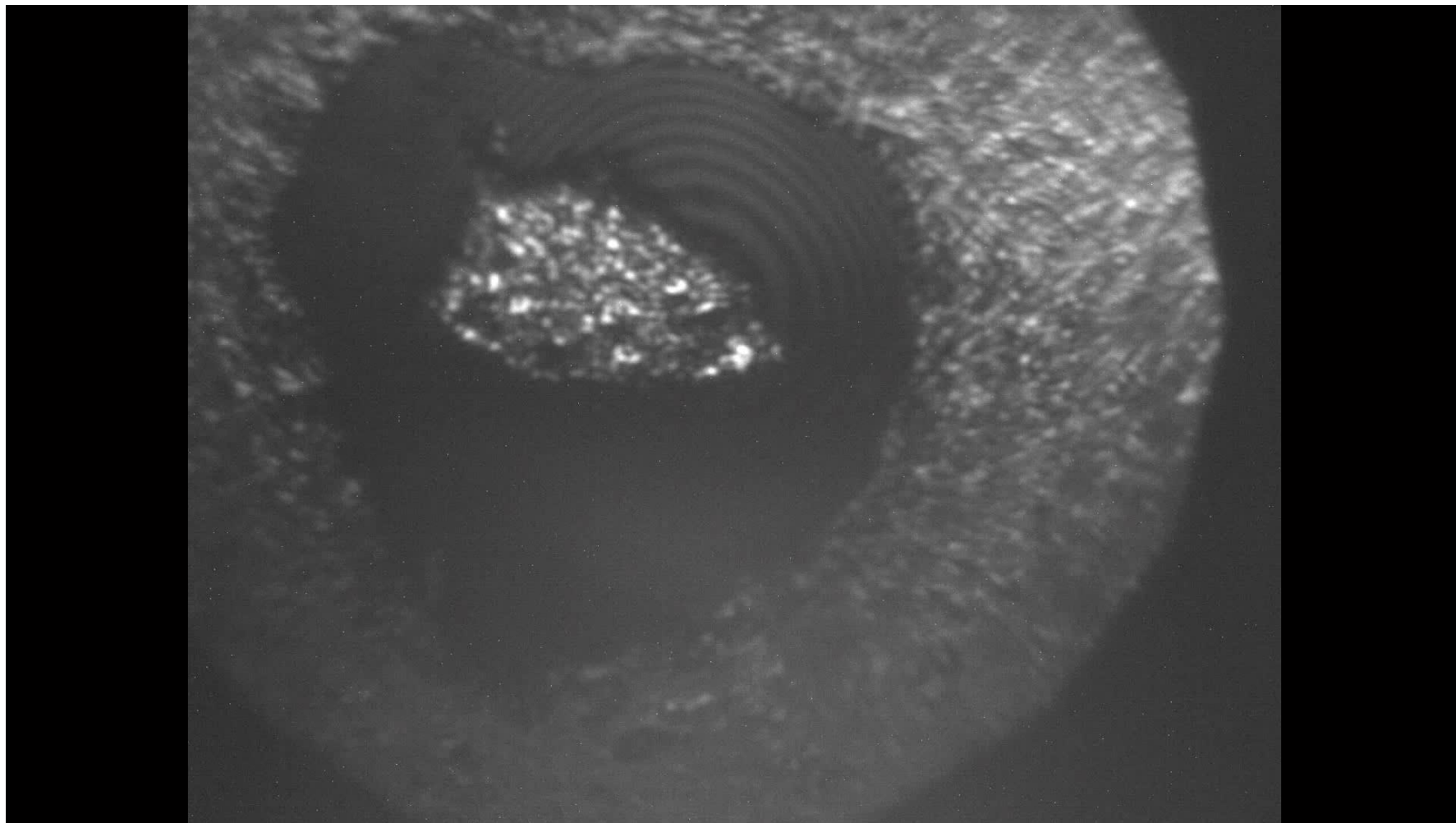
$$\varepsilon_0 = \frac{\sum_{i=1}^n [I(\lambda_i) g(\lambda_i, T)]}{\sum_{i=1}^n [g^2(\lambda_i, T)]} = \frac{\overline{I g(T)}}{g^2(T)}$$



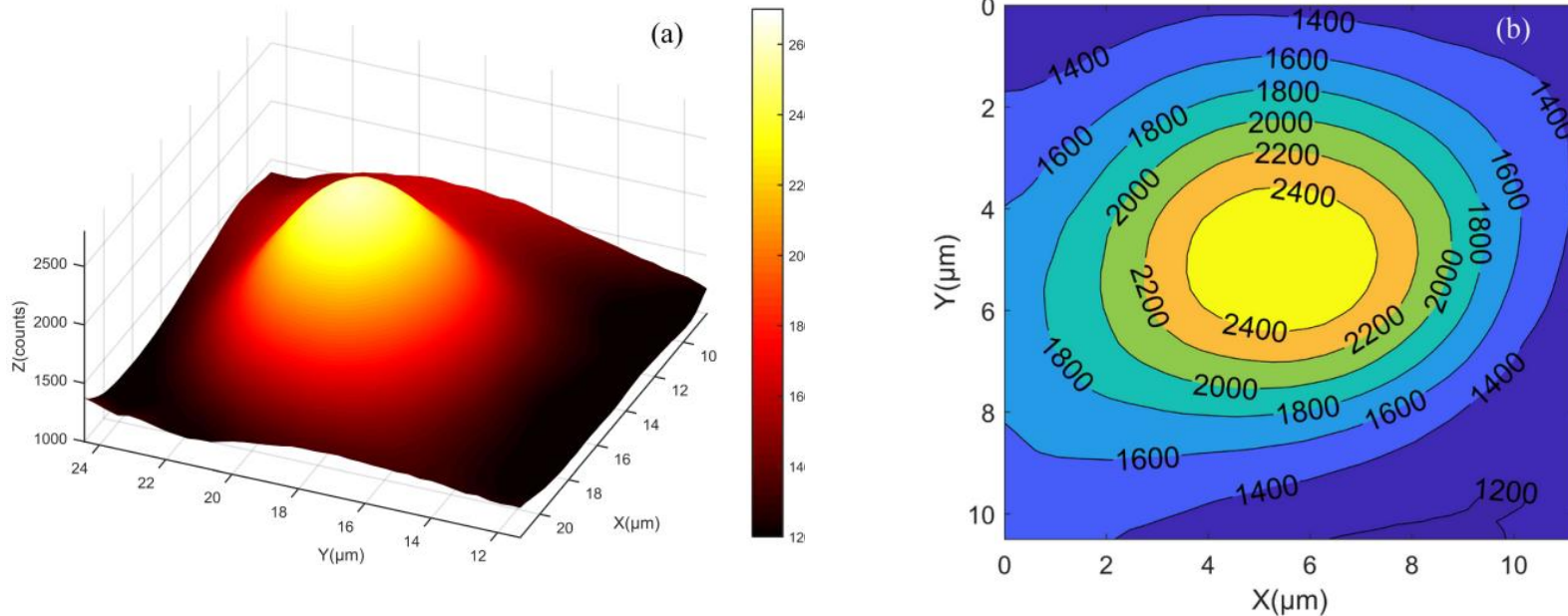
Experimental data fitting using 1-D minimization procedure for a selected point on the tungsten plate in the lamp when the laser power was 14 W. Graph shows the behavior of the normalized  $S$  (eq. 8) as a function of temperature. The function  $S$  has a minimum  $S_0$  at 2408 K. The 95% confidence interval,  $\Delta T$ , is  $\Delta T = \pm 41$  K.



# Лазерный нагрев в ячейке высокого давления



# The 2-D temperature distribution on the surface of a Fe plate heated in the DAC



(a) The 2-D on the surface of an Fe plate heated in the DAC2 at 43 GPa with 5 W laser power, 620-750 nm wavelength, and 2 s exposure time. (a) A color map and (b) a filled contour plot displaying isolines of the color map with filled areas between the isolines.

**РАСПРЕДЕЛЕНИЯ КОЭФФИЦИЕНТА ТЕПЛООВОГО ИЗЛУЧЕНИЯ И ТЕМПЕРАТУРЫ ПОВЕРХНОСТИ ВОЛЬФРАМА,  
НАГРЕТОГО ИЗЛУЧЕНИЕМ МОЩНОГО ЛАЗЕРА**

$$\varepsilon_o = \frac{\sum_{i=1}^n [I(\lambda_i)g(\lambda_i, T)]}{\sum_{i=1}^n [g^2(\lambda_i, T)]} = \frac{\sum_{i=1}^n g(\lambda_i, T)y_i}{\sum_{i=1}^n [g^2(\lambda_i, T)]} \quad \varepsilon_o = \sum_{i=1}^N [y_i \cdot G(\lambda_i, T_o)] \quad G(\lambda_i, T_o) = \frac{g(\lambda_i, T_o)}{\sum_{i=1}^N [g^2(\lambda_i, T_o)]}$$

Мы предполагаем, что все измерения  $y_i = I(\lambda_i)$  являются независимыми. Тогда стандартное отклонение величины для коэффициента излучения может быть записано

$$\sigma_{\varepsilon}^2 = \sum_{i=1}^n \left( \frac{\partial \varepsilon(y_i)}{\partial y_i} \sigma_{y_i} \right)^2$$

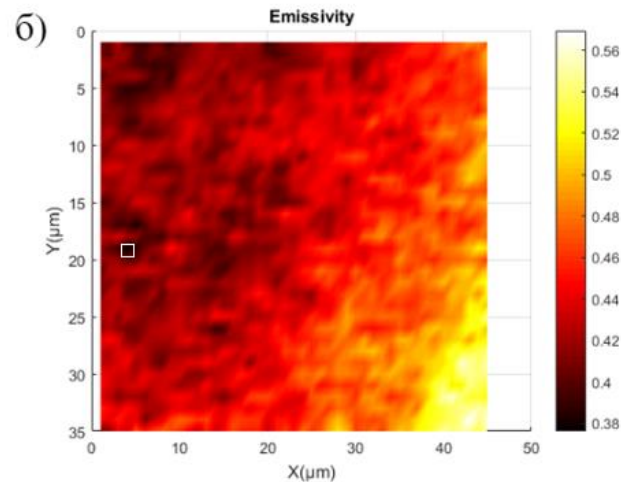
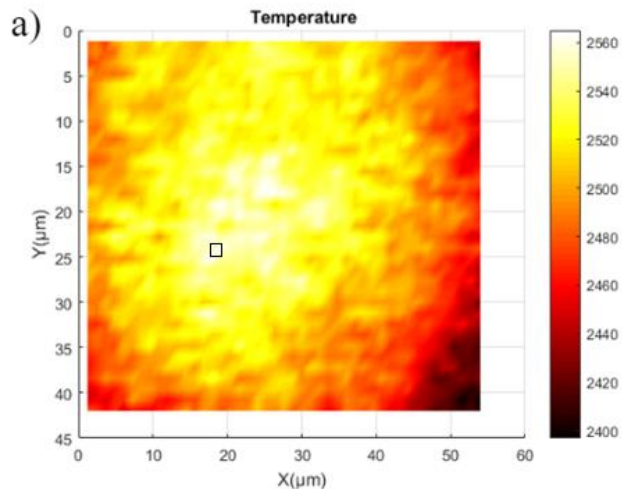
Возьмём производную  $\frac{\partial \varepsilon(y_i)}{\partial y_i}$ . Очевидно  $\frac{\partial \varepsilon(y_i)}{\partial y_i} = G(\lambda_i, T_o)$

$$\sigma_{\varepsilon}^2 = \sum_{i=1}^n \left[ \frac{\partial \varepsilon(\lambda_i)}{\partial y_i} \right]^2 = (\sigma_I)^2 \sum_{i=1}^n [G(\lambda_i, T_o)]^2 = \frac{(\sigma_I)^2 \sum_{i=1}^n [g^2(\lambda_i, T_o)]^2}{\left\{ \sum_{i=1}^N [g^2(\lambda_i, T_o)] \right\}^2} = \frac{(\sigma_I)^2}{\sum_{i=1}^N [g^2(\lambda_i, T_o)]}$$

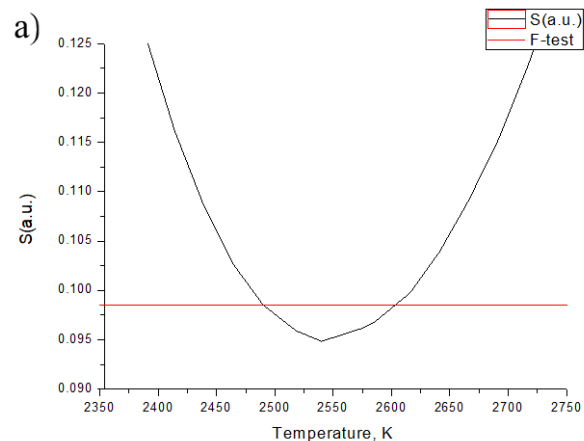
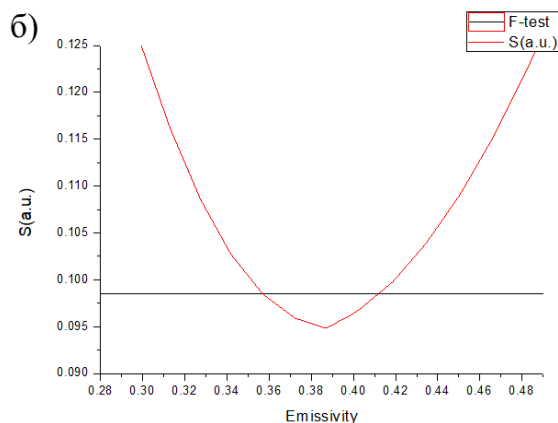
**Окончательно**

$$\sigma_{\varepsilon} = \frac{\sigma_I}{\sum_{i=1}^N [g^2(\lambda_i, T_o)]}, \quad \sigma_I = \sqrt{\frac{\sum_{i=1}^N [I(\lambda_i) - \varepsilon_o g(\lambda_i, T_o)]^2}{N - 2}}$$

# РАСПРЕДЕЛЕНИЯ КОЭФФИЦИЕНТА ТЕПЛОВОГО ИЗЛУЧЕНИЯ И ТЕМПЕРАТУРЫ ПОВЕРХНОСТИ ВОЛЬФРАМА, НАГРЕТОГО ИЗЛУЧЕНИЕМ МОЩНОГО ЛАЗЕРА

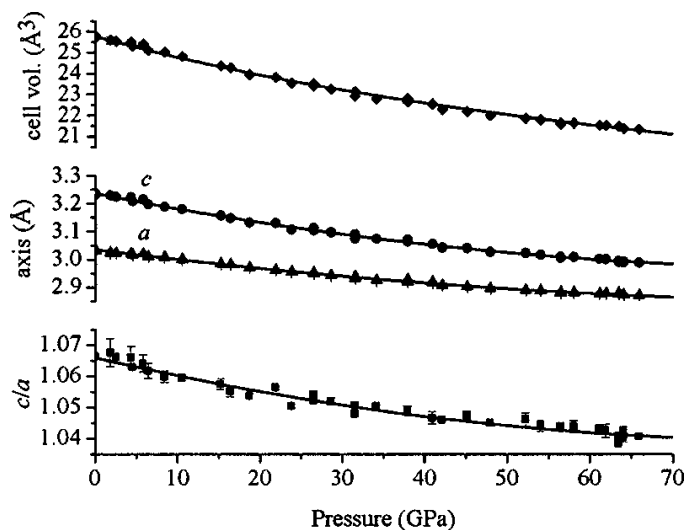


Распределение эффективного коэффициента теплового излучения(а) и температуры (б) при нагреве вольфрамовой лампы излучением лазера 8 W, экспозиция  $\frac{1}{4}$  s.



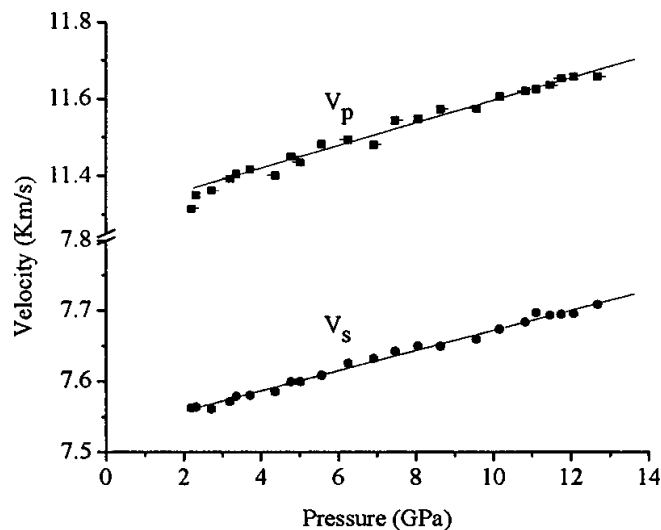
Поведение суммы квадратов отклонений  $S$  как функции (а) от температуры и (б) коэффициента теплового излучения, моделируемой в точке в центре пятна нагрева. Пересечения горизонтальной линии с функцией  $S$  определяют доверительный интервал в пределах которого лежат значения температуры с доверительной вероятностью 0,95.

# Метод наименьших квадратов. Пример: Двумерная нелинейная регрессия



Изменение параметров решетки  $\text{TiB}_2$  с давлением.

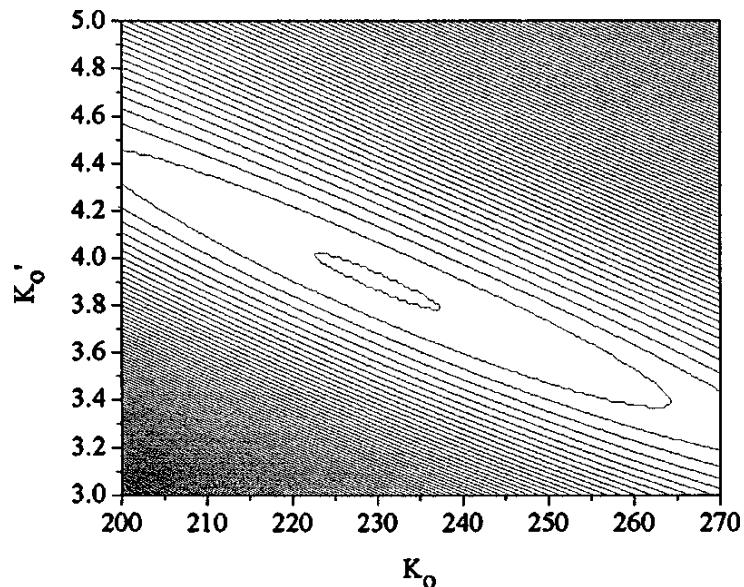
В работе [1] сжимаемость  $\text{TiB}_2$  определялась отдельно из экспериментов по рентгеновской дифракции и на основе ультразвуковых измерений на образцах, загруженных в ячейки с алмазными наковальнями (до 65,9 ГПа) и помещенных в аппараты высокого давления (до 13,9 ГПа).



Зависимости скоростей продольных и сдвиговых волн в  $\text{TiB}_2$  от давления.

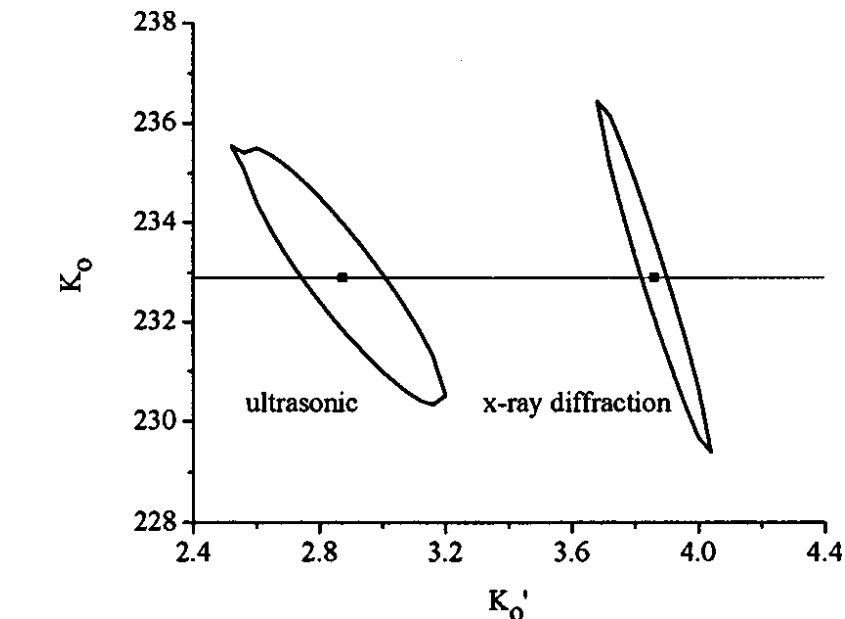
Целью статьи было получение величины объемного модуля  $K_0$  и его производной по давлению,  $K'_0$ , входящих в уравнение состояния третьего порядка, Берча-Марнагама.

# 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

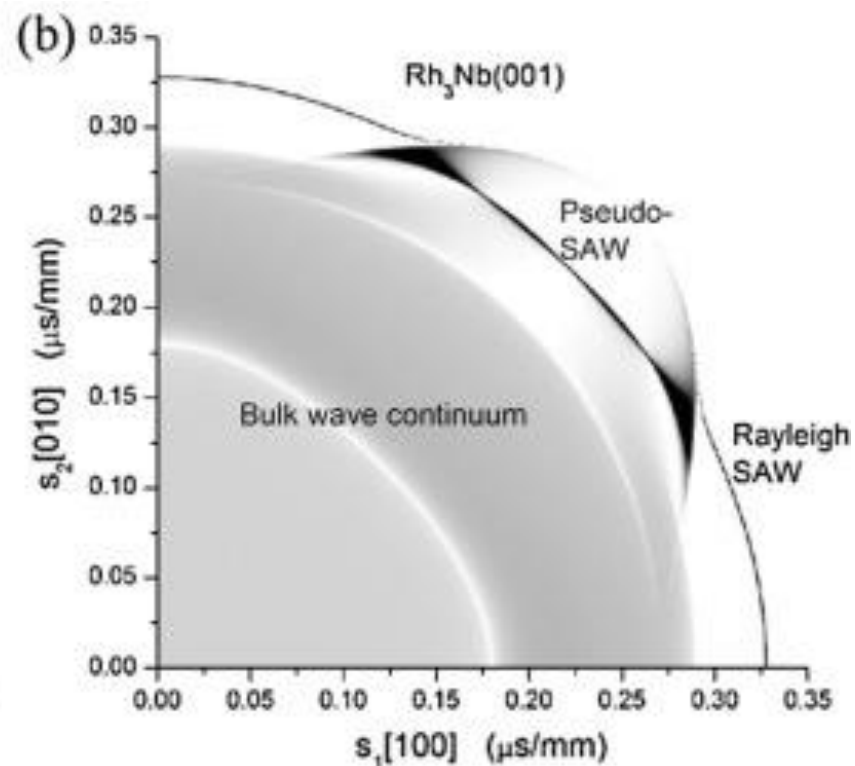
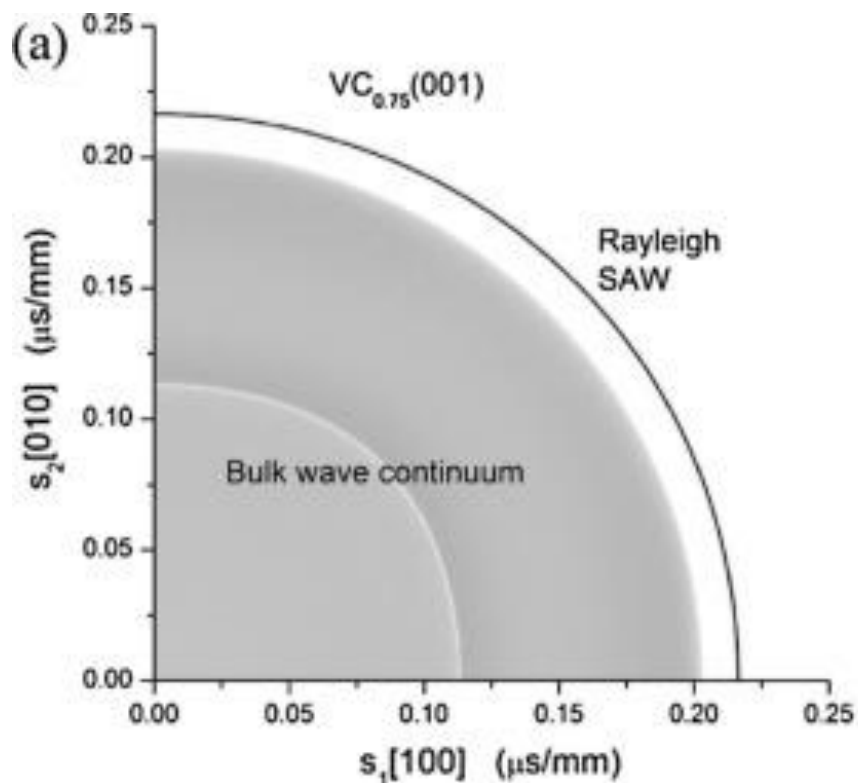


Доверительные эллипсы, построенные для доверительной вероятности 95% и полученные на основании подгонки экспериментальных ультразвуковых и рентгеновских данных уравнением Берча-Марнагама.

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

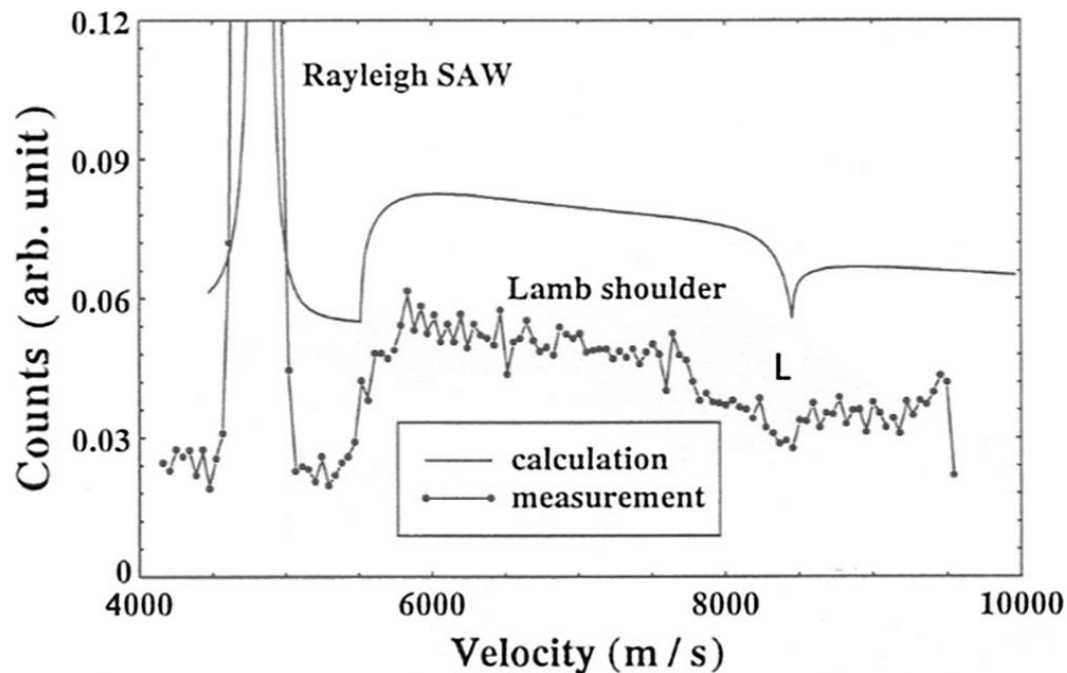


# Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



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

# 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).

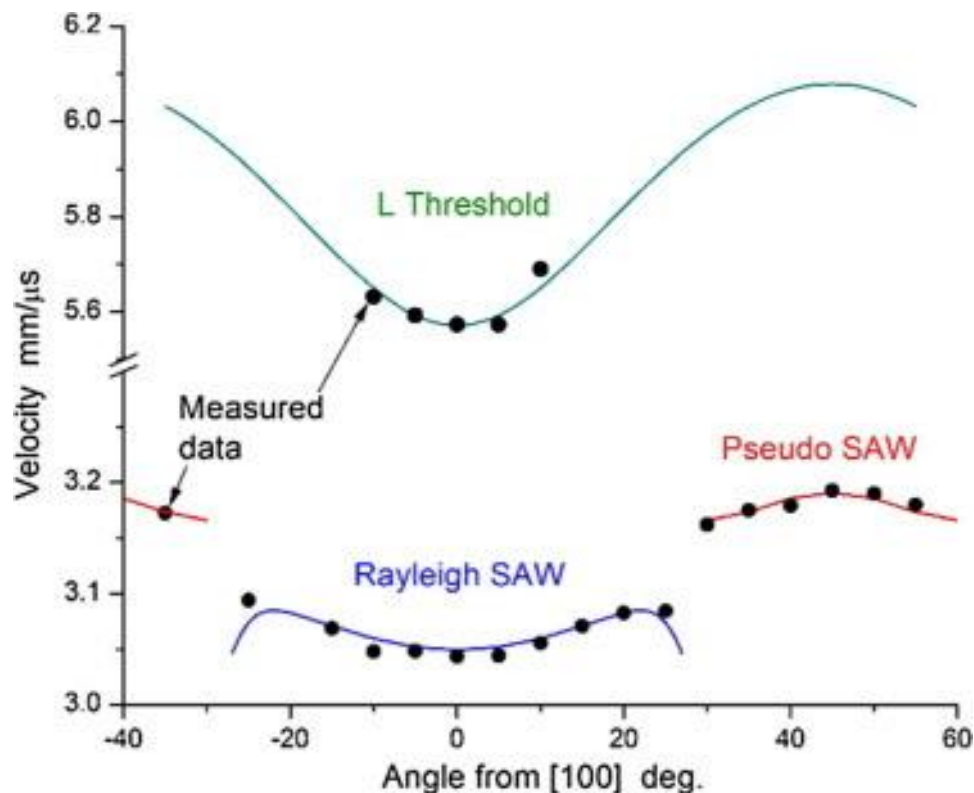
$$\chi^2(C_{11}, C_{22}, C_{33}) = \sum_{SAW_i} (V_i^{meas} - V_i^{calc})^2 + W \sum_{Li} (V_i^{meas} - V_i^{calc})^2$$

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).



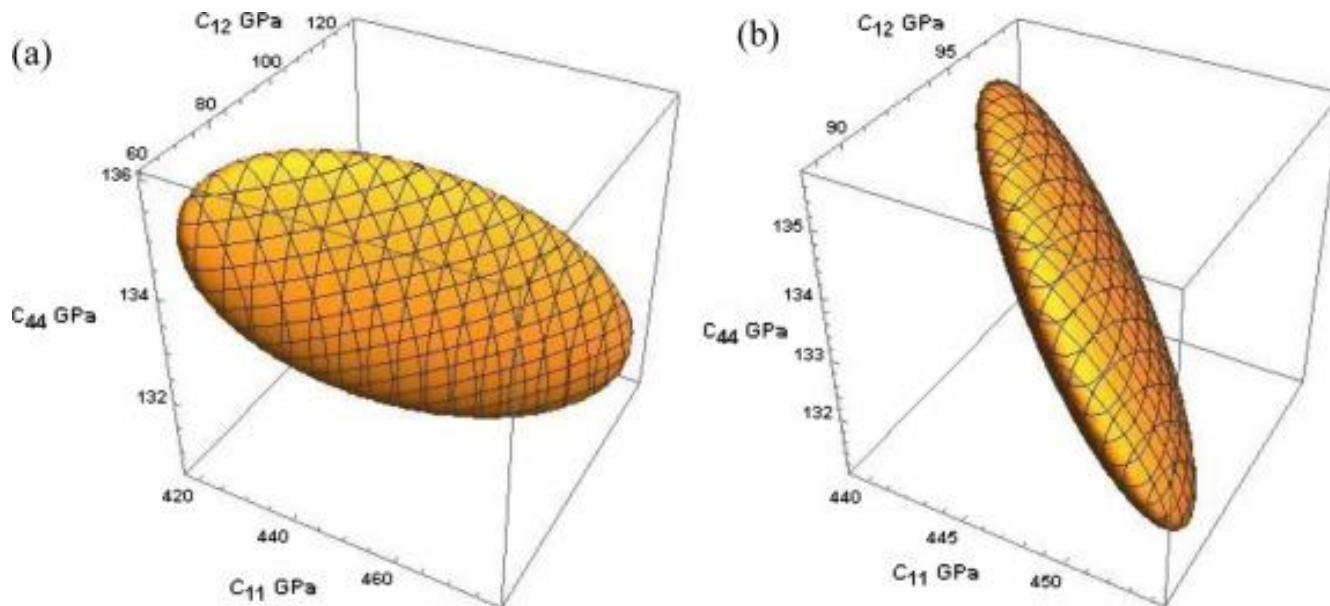
# Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



Скорости Рэлеевской, псевдо-ПАВ, продольной волн как функции направления на (0 0 1) поверхности  $\text{Rh}_3\text{Nb}$ . 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).

## Метод наименьших квадратов. Пример: Трехмерная нелинейная регрессия



Эллипсоиды доверительных областей для  $VC_{0.75}$ , (a)  $W = 0.1$ , вытянутый эллипсоид 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  $W = 5$ .

# 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](http://en.wikipedia.org/wiki/Akaike_information_criterion)

# Application of Akaike Information Criterion 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	$226.1 \pm 1.0$	$3.64 \pm 0.19$
3 <sup>rd</sup> order Birch-Murnaghan equation	-128.8	$225.8 \pm 1.0$	$3.73 \pm 0.19$

Model (x-ray data)	AIC	$K_o$	$K_o'$
Murnaghan equation	51.2	$237.9 \pm 12.1$	$2.62 \pm 0.73$
3 <sup>rd</sup> order Birch-Murnaghan equation	51.9	$235.9 \pm 11.6$	$2.92 \pm 0.66$

Model ( <b>Combined ultrasonic and x-ray data</b> )	AIC	$K_o$	$K_o'$
Murnaghan equation	73.8	$232.9 \pm 5.5$	$2.91 \pm 0.36$
3 <sup>rd</sup> order Birch-Murnaghan equation	74.0	$231.7 \pm 5.4$	$3.15 \pm 0.34$

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

# 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. Худсон Д. *Статистика для физиков*. Мир. Москва. 1967
4. Тейлор, Д., Введение в теорию ошибок. Москва. Мир. 1985.
5. Сквайрс, Д., *Практическая физика*. Москва: Мир. 1971.
6. 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](http://clip.med.yale.edu/courses/brdu/Costa_ODE.pdf)