

ОСНОВЫ МАТЕМАТИЧЕСКОЙ СТАТИСТИКИ

Задачи математической статистики

Первая задача математической статистики - указать способы сбора и группировки статистических сведений, полученных в результате наблюдений или в результате специально поставленных экспериментов.

Вторая задача математической статистики - разработать методы анализа статистических данных в зависимости от целей исследования:

а) оценка неизвестной вероятности события; оценка неизвестной функции распределения; оценка параметров распределения, вид которого известен; оценка зависимости случайной величины от одной или нескольких случайных величин и др.;

б) проверка статистических гипотез о виде неизвестного распределения или о величине параметров распределения, вид

Генеральная и выборочная совокупности

Выборочной совокупностью или просто **выборкой** называют совокупность случайно отобранных объектов. **Генеральной совокупностью** называют совокупность объектов, из которых производится выборка. **Объемом совокупности** (выборочной или генеральной) называют число объектов этой совокупности.

Пусть из генеральной совокупности извлечена выборка, причем x_1 наблюдалось n_1 раз, $x_2 - n_2$ раз, $x_k - n_k$ раз, $\sum_{i=1}^k n_i = n$ - объем выборки. Наблюдаемые значения x_i называют **вариантами**, а последовательность вариантов, записанных в возрастающем порядке - **вариационным рядом**. Числа наблюдений называют частотами, а их отношения к объему выборки - **относительными частотами**.

Статистическим распределением выборки называют перечень вариантов и соответствующих им частот или

Эмпирическая функция распределения

Эмпирической функцией распределения (функцией распределения выборки) называют функцию $F^*(x)$, определяющую для каждого значения x относительную частоту события $X < x$.

Функцию распределения $F(x) = P(X < x)$ генеральной совокупности называют теоретической функцией распределения. Из теоремы Бернулли следует, что относительная частота события $X < x$, $F^*(x)$, стремится по вероятности к вероятности $F(x)$ этого события. Отсюда следует целесообразность использования эмпирической функции распределения выборки для приближенного представления теоретической (интегральной) функции распределения генеральной совокупности.

Эмпирическая функция распределения выборки служит для оценки теоретической функции распределения генеральной совокупности.

ОЦЕНКА ПАРАМЕТРОВ РАСПРЕДЕЛЕНИЯ

Будем называть выборочной статистикой любую величину, полученную в результате обработки данных эксперимента. Любая выборочная статистика является случайной величиной.

Определение. Точечной оценкой параметра x называется его приближенное значение x^* , полученное в результате обработки выборки.

Определение. Оценка x^* параметра x называется **несмещенной**, если математическое ожидание оценки совпадает с оцениваемым параметром: $M(x^*) = x$. **Эффективной** называют статистическую оценку, которая (при заданном объеме выборки n) имеет наименьшую возможную дисперсию.

Любая величина, определенная по выборке, называется статистикой. Очевидно, что любая статистика – случайная величина.

ОЦЕНКА ПАРАМЕТРОВ РАСПРЕДЕЛЕНИЯ

Определение. Оценка параметра x называется **состоятельной**, если при увеличении объема выборки n для любого $\varepsilon \geq 0$ вероятность отклонения оценки от параметра x на величину $\geq \varepsilon$ **меньше** ε , равна 1:

Несмещенной и состоятельной оценкой математического ожидания a оцениваемого признака θ является выборочная средняя:

$$\bar{\theta} = \frac{\sum_{i=1}^n \theta_i}{n},$$

где θ_i - варианты выборки выходного параметра θ .

Несмещенной и состоятельной оценкой дисперсии D является выборочная дисперсия s^2 :

$$s^2 = \frac{\sum_{i=1}^n (\theta_i - \bar{\theta})^2}{n-1} = \frac{n}{n-1} \left(\frac{\sum_{i=1}^n \theta_i^2}{n} - (\bar{\theta})^2 \right)$$

Выборочным средним квадратическим отклонением

ОЦЕНКА ПАРАМЕТРОВ РАСПРЕДЕЛЕНИЯ

Медианой m_e называют варианту, которая делит вариационный ряд на две части, равные по числу вариант. Если число вариант нечетно, т. е. $n = 2k + 1$, $m_e = x_{k+1}$, при четном $n = 2k$ медиана $m_e = \frac{x_k + x_{k+1}}{2}$.

Размахом варьирования R называют разность между наибольшей и наименьшей вариантами: $R = x_{\max} - x_{\min}$.

Средним абсолютным отклонением θ называют среднее арифметическое абсолютных отклонений:

$$\theta = \frac{\sum_{i=1}^n |x_i - m_e|}{n}$$

Коэффициентом вариации V называют выраженное в процентах отношение выборочного среднего квадратического отклонения к x

выборочной средней:

ОЦЕНКА ПАРАМЕТРОВ РАСПРЕДЕЛЕНИЯ

Оценкой асимметрии теоретического распределения $\frac{\mu_3}{\sigma^3}$

является ее выборочное значение $E_s = \frac{n}{(n-1)(n-2)} \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{s^3}$

Оценкой теоретического эксцесса $E_k = \frac{\mu_4}{\sigma^4} - 3$ является его выборочное значение:

$$E_s = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{s^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$$

ИНТЕРВАЛЬНОЕ ОЦЕНИВАНИЕ

Доверительным интервалом для параметра θ называется интервал (θ_1, θ_2) , покрывающий истинное значение θ с заданной вероятностью $\gamma = 1 - \alpha$, то есть $P(\theta_1 < \theta < \theta_2) = \gamma$. Число $\gamma = 1 - \alpha$ называется доверительной вероятностью, а значение α - уровнем значимости. Статистики θ_1 и θ_2 , определяемые по выборке из генеральной совокупности с неизвестным параметром θ , называются нижней и верхней границами доверительного интервала (или левой и правой границами).

Чтобы найти доверительный интервал для параметра θ , необходимо знать закон распределения статистики Y , значение которой является точечной оценкой параметра θ .

Пусть существует статистика $Y = Y(\theta)$ такая, что:

а) закон распределения Y известен и не зависит от θ ,

б) функция $Y(\theta)$ непрерывна и строго монотонна по θ .

ИНТЕРВАЛЬНОЕ ОЦЕНИВАНИЕ

Пусть $\gamma = 1 - \alpha$ - заданная доверительная вероятность, $y_{\frac{\alpha}{2}} = y_{\frac{1-\gamma}{2}}$

$y_{1-\frac{\alpha}{2}} = y_{\frac{1+\gamma}{2}}$ - квантили статистики Y порядков $\frac{\alpha}{2}$ и $1-\frac{\alpha}{2}$ и $y_{\frac{\alpha}{2}} < Y(\theta, \theta) < y_{1-\frac{\alpha}{2}}$ с вероятностью $\gamma = 1 - \alpha$ выполняется неравенство:

Решая неравенство относительно θ , найдем границы θ_1 и θ_2

доверительного интервала для θ .

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Распределение некоторых статистик для НСВ $X \sim N(m, \sigma)$

$$N\left(m, \frac{\sigma}{\sqrt{n}}\right)$$

1. Выборочное среднее

имеет нормальное

Распределение некоторых статистик для НСВ $X \sim N(m, \sigma)$

2. Выборочная дисперсия $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ связана со случайной

величиной $\chi^2(n-1)$ соотношением: $s^2 = \frac{\sigma^2}{n-1} \chi^2(n-1)$.

3. Статистика $\frac{\bar{x} - m}{s/\sqrt{n}}$ является случайной величиной, имеющей распределение Стьюдента с $(n-1)$ степенью свободы $T(n-1)$.

3. Пусть s_1^2 и s_2^2 - выборочные дисперсии, вычисленные по независимым выборкам объема n_1 и n_2 из двух нормально распределенных генеральных совокупностей с дисперсиями σ_1^2 и σ_2^2 и

тогда отношение $\frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2}$ имеет распределение Фишера $F(n_1-1, n_2-1)$

Доверительный интервал для среднего при известной σ

Пусть x_1, \dots, x_n - выборка n наблюдений случайной величины $X \sim N(m, \sigma)$.

Рассмотрим статистику $Y = \frac{\bar{x} - m}{\sigma/\sqrt{n}} \sim N(0, 1)$. Вероятность события

$u_{\frac{\alpha}{2}} < Y < u_{1-\frac{\alpha}{2}}$ равна заданной доверительной вероятности γ , то есть

$P\left[u_{\frac{\alpha}{2}} < Y < u_{1-\frac{\alpha}{2}}\right] = \gamma$. Для квантилей стандартного (нормированного) нормального распределения $u_{\frac{\alpha}{2}} = -u_{1-\frac{\alpha}{2}}$. После решения

неравенства относительно m получим доверительный интервал в

виде:

$$\bar{x} - \frac{\sigma}{\sqrt{n}} u_{1-\frac{\alpha}{2}} < m < \bar{x} + \frac{\sigma}{\sqrt{n}} u_{1-\frac{\alpha}{2}}$$

Доверительный интервал для дисперсии

В качестве оценки дисперсии используют выборочную дисперсию

$s^2 = \frac{n}{n-1} \left(\frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2 \right)$. Статистика $Y = \frac{(n-1)s^2}{\sigma^2}$ имеет распределение $\chi^2(n-1)$.

Определим квантили $\chi_{\frac{\alpha}{2}}^2(n-1)$ и $\chi_{1-\frac{\alpha}{2}}^2(n-1)$, удовлетворяющие условию

$$P \left[\chi_{\frac{\alpha}{2}}^2(n-1) < \frac{(n-1)s^2}{\sigma^2} < \chi_{1-\frac{\alpha}{2}}^2(n-1) \right] = \gamma = 1 - \alpha$$

Решая неравенство относительно σ^2 , получим доверительный интервал для дисперсии:

$$\frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}}^2(n-1)} < \sigma^2 < \frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}}^2(n-1)}$$

ПРОВЕРКА СТАТИСТИЧЕСКИХ ГИПОТЕЗ

Статистической гипотезой H называется предположение относительно параметров или вида распределения случайной величины X . Если по выборке наблюдений необходимо проверить предположение о значении параметров известного распределения, то такие гипотезы называются **параметрическими**.

Проверяемая гипотеза называется **нулевой гипотезой H_0** . Вместе с гипотезой H_0 рассматривают одну из **альтернативных (конкурирующих)** гипотез H_1 . Правило, по которому принимается решение принять или отклонить гипотезу H_0 , называется **критерием**. Статистика Z , на основании которой принимается решение об истинности гипотезы H_0 , называют **статистикой критерия**.

Перед анализом выборки назначается некоторая малая вероятность α , называемая **уровнем значимости**.

ПРОВЕРКА СТАТИСТИЧЕСКИХ ГИПОТЕЗ

Пусть V – множество значений статистики Z , а V_k – такое

подмножество, что при условии истинности гипотезы H_0 ,

вероятность попадания статистики критерия в V_k равна α :

$$P[Z \in V_k | H_0] = \alpha$$

α – уровень значимости (малое число). Пусть z_v – выборочное значение статистики Z , вычисленное по выборке. Формулировка

критерия: отклонить гипотезу H_0 , если $z_v \in V_k$ (попадание в

критическую область); принять гипотезу H_0 , если $z_v \in V \setminus V_k$

(попадание в область принятия гипотезы).

ПРОВЕРКА СТАТИСТИЧЕСКИХ ГИПОТЕЗ

Уровень значимости α определяет размер критической области V_k . Положение критической области на множестве значений статистики Z зависит от формулировки альтернативной гипотезы H_1 . Если проверяется гипотеза $H_0: \theta = \theta_0$, а альтернативная гипотеза $H_1: \theta > \theta_0$ ($\theta < \theta_0$), то критическая область размещается на правом (левом) хвосте плотности распределения статистики Z , то есть $Z > z_{1-\alpha}$ имеет вид неравенства: $Z < z_\alpha$, где $z_{1-\alpha}$ и z_α - квантили распределения статистики Z , вычисленные при условии истинности H_0 .

ПРОВЕРКА СТАТИСТИЧЕСКИХ ГИПОТЕЗ

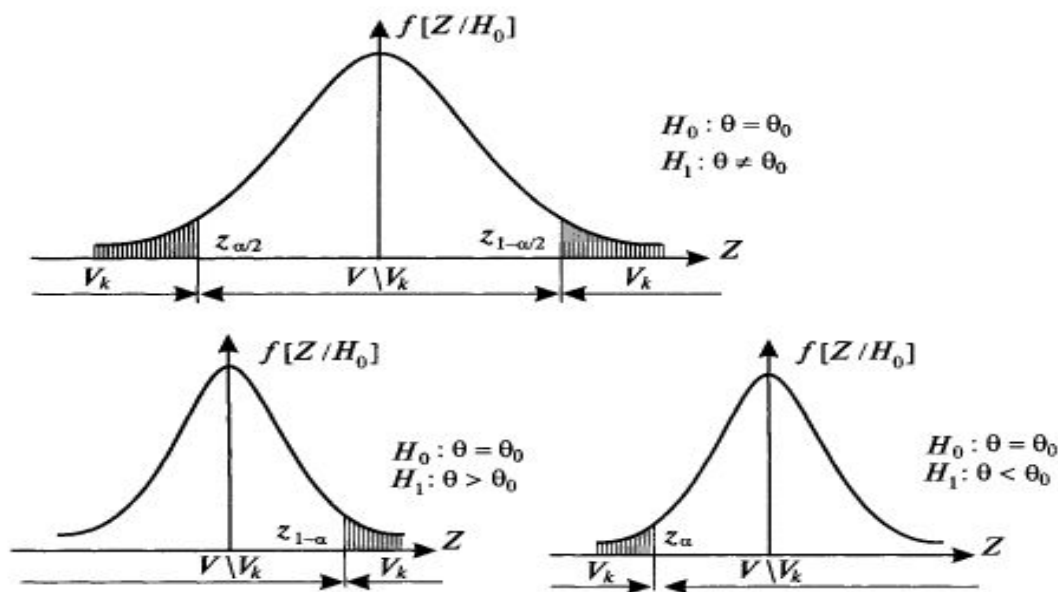
Если альтернативная гипотеза $H_1: \theta \neq \theta_0$, критическая область размещается на обоих хвостах плотности распределения Z , то есть определяется

неравенствами: $Z < -z_{\frac{\alpha}{2}}$ и $Z > z_{1-\frac{\alpha}{2}}$. В этом случае

критерий – двусторонний.

При проверке гипотез значимости α необходимо:

1. Сформулировать гипотезы H_0 и H_1 , выбрать уровень значимости α .
2. Выбрать статистику Z критерия для проверки H_0 .
3. Определить критическую область V_k , при условии, что верна H_0 .
4. Вычислить выборочное значение статистики z критерия.



Размещение критической области при различных альтернативных гипотезах

Сравнение двух дисперсий нормальных генеральных совокупностей

При заданном уровне значимости α по двум выборкам объемов n_1 и n_2 проверить нулевую гипотезу $H_0 : D(X) = D(Y)$ о равенстве генеральных дисперсий нормальных совокупностей при конкурирующей гипотезе $H_1 : D(X) > D(Y)$.

В качестве критерия выбирается отношение большей дисперсии к меньшей

распределение F с параметрами $(n_1 - 1, n_2 - 1)$, имеющей $F > F_{1-\alpha}(n_1 - 1, n_2 - 1)$

$$F_{кр} = F_{1-\alpha}(n_1 - 1, n_2 - 1)$$

Фишера (n_1 – объем выборки с большей дисперсией). Критическая область – односторонняя и определяется по правилу:

, где $F_{1-\alpha}(n_1 - 1, n_2 - 1)$ – квантиль распределения Фишера. Введем критическое число $F_{кр}$, отделяющее критическую область от не критической.

Сравнение выборочной дисперсии с гипотетической генеральной дисперсией нормальной совокупности

При заданном уровне значимости α по выборке объемом n проверить нулевую гипотезу $H_0: D(X) = \sigma_0^2$ о равенстве генеральной дисперсии $D(X)$ гипотетическому значению σ_0^2 при конкурирующей гипотезе $H_1: D(X) > \sigma_0^2$.

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

В качестве критерия выбирается статистика имеющая распределение χ^2 с $n-1$ степенью свободы.

$$\chi_{набл}^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

Наблюдаемое значение критерия вычисляется по выборке критическая область $(n-1)$ правосторонняя и определяется по правилу

, где $\chi_{кр}^2(1-\alpha, n-1)$ - квантиль распределения χ^2 с $n-1$ степенью свободы порядка $1-\alpha$.

Введем критическое число

$$\chi_{набл}^2 > \chi_{кр}^2$$

Если - нет оснований отвергнуть нулевую гипотезу.

Если - нулевая гипотеза отвергается.

Сравнение двух средних генеральных совокупностей, дисперсии которых известны

При заданном уровне значимости α по двум выборкам объемами n и m проверить нулевую гипотезу $H_0: M(X)=M(Y)$ о равенстве математических ожиданий двух нормальных генеральных совокупностей с известными дисперсиями при конкурирующей гипотезе $H_1: M(X) \neq M(Y)$. В качестве

критерия выбирается статистика
$$z = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{D(X)}{n} + \frac{D(Y)}{m}}} \sim N(0,1).$$

Наблюдаемое значение $z_{набл}$ критерия вычисляется по выборкам.

Критическая область - двусторонняя и определяется условиями

$z < u_{\frac{\alpha}{2}}$ или $z > u_{1-\frac{\alpha}{2}}$, где u_{α} - квантиль нормального нормированного

распределения порядка α . Учитывая, что $u_{\frac{\alpha}{2}} = -u_{1-\frac{\alpha}{2}}$ область принятия

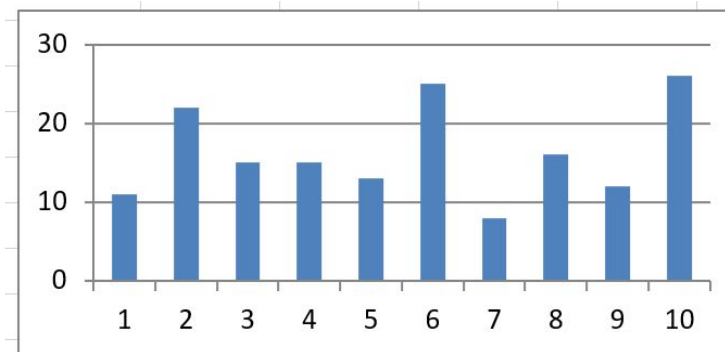
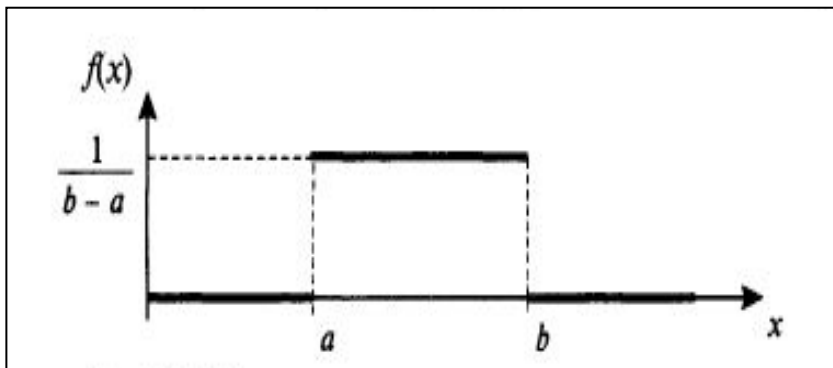
гипотезы определяется неравенством $|z| < u_{1-\frac{\alpha}{2}}$. Пусть $z_{кр} = u_{1-\frac{\alpha}{2}}$.

$$|z_{набл}| < z_{кр}$$

КРИТЕРИЙ СОГЛАСИЯ ПИРСОНА

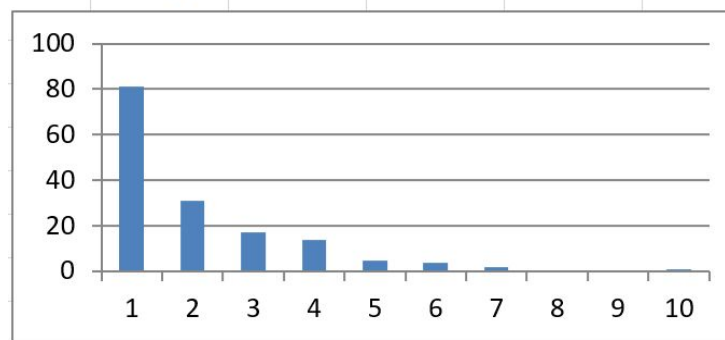
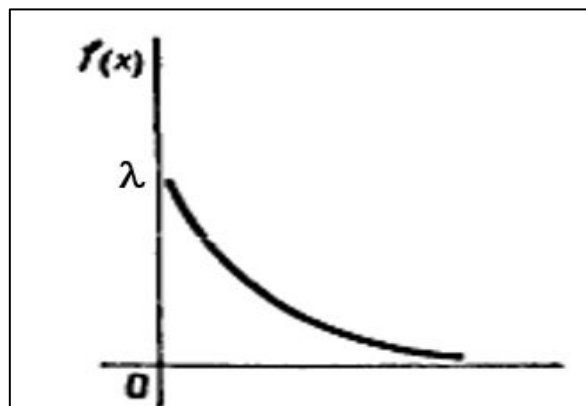
Пусть по выборке случайной величины X объема n необходимо проверить гипотезу H_0 о том, что X имеет предполагаемый закон распределения. Пусть построена гистограмма для заданной выборки случайной величины X объема n с частичными интервалами Δ_j одинакового размера h . По форме гистограммы делается предположение о типе закона распределения и по выборке вычисляются точечные оценки его параметров. Используя предполагаемый закон распределения случайной величины X , находим вероятности p_j того, что значение X принадлежит частичным интервалам Δ_j , то есть $p_j = P[X \in \Delta_j]$, где $j = 1, \dots, N$. Величины $m_j = np_j$ равны теоретическим частотам попадания случайной величины X в каждый частичный интервал. Выборочное (наблюдаемое) значение статистики критерия вычисляется
$$\chi^2_{\text{набл}} = \sum_{j=1}^N \frac{(m_j - m_j)^2}{m_j}$$

КРИТЕРИЙ СОГЛАСИЯ ПИРСОНА

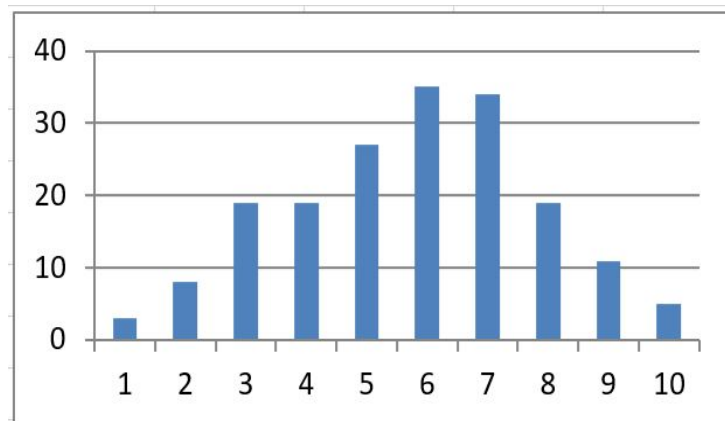
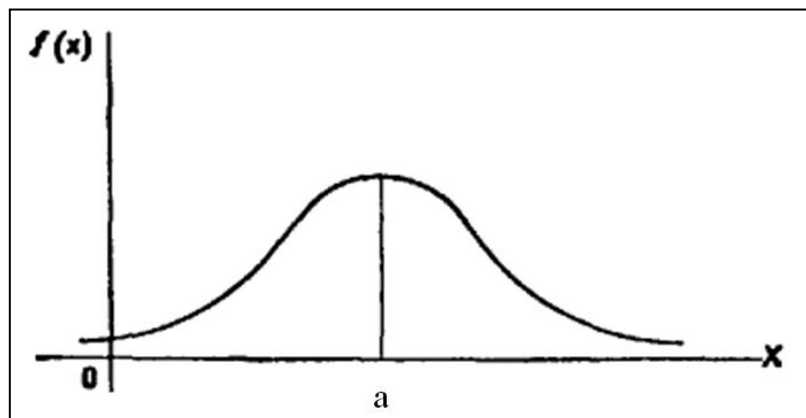


Законы
распределения:

а) равномерный;



б) экспоненциальный;



в) нормальный

КРИТЕРИЙ СОГЛАСИЯ ПИРСОНА

Гипотеза H_0 согласуется с результатами наблюдений на уровне значимости α , если $\chi_{набл}^2 < \chi_{1-\alpha}^2(N - m_p - 1)$, где $\chi_{1-\alpha}^2(N - m_p - 1)$ - квантиль порядка

$1 - \alpha$ распределения χ^2 с $N - m_p - 1$ степенями свободы, где m_p - число неизвестных параметров в распределении, оцениваемых по выборке. Если же $\chi_{набл}^2 \geq \chi_{1-\alpha}^2(N - m_p - 1)$, то гипотеза H_0 отклоняется.

При заданном уровне значимости α используется следующий расчет.

1) Для каждого частичного интервала определяются средние точки y_j частичных интервалов.

$$m_j = \frac{hn}{x_{j+1} - x_j}, \quad j = 1, \dots, N$$

2) Вычисляются теоретические частоты предполагаемого распределения: а) для равномерного закона

КРИТЕРИЙ СОГЛАСИЯ ПИРСОНА

б) для нормального закона $m_j = \frac{hn}{s} \varphi(u_j)$, $j=1, \dots, N$, где $u_j = \frac{y_j - x_0}{s}$,

$\varphi(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$ - плотность нормального нормированного распределения;

$$m_j = n \left(e^{-\lambda_0 x_j} - e^{-\lambda_0 x_{j+1}} \right), j=1, \dots, N$$

в) для показательного закона $\lambda_0 = \frac{1}{x_0}$, где x_0 - выборочный параметр показательного закона распределения:

$$\chi^2_{набл} = \sum_{j=1}^N \frac{(m_j - t_j)^2}{m_j}$$

Находится наблюдаемое значение критерия χ^2 :

3) Находится критическое значение критерия χ^2 по уровню значимости α и числу степеней свободы $\chi^2_{крит} = \chi^2_{1-\alpha, p}$, где p - число параметров распределения, равное $m_p = 2$ для нормального и равномерного законов $m_p = 1$ для показательного распределения:

4) Если $\chi^2_{набл} < \chi^2_{крит}$ - то выдвинутая гипотеза о типе закона

МОДЕЛЬ ПАРНОЙ РЕГРЕССИИ

Пусть есть результаты некоторого эксперимента $(x_i, y_i), i = 1, \dots, n$.
Надо подобрать (“подогнать”) функцию $y = f(x) = ax + b$, наилучшим образом описывающую экспериментальные точки (x_i, y_i) .

Критерием качества “подгонки” кривой выберем минимум

$$\text{функционала } F = \sum_{i=1}^n (y_i - (a + bx_i))^2 \rightarrow \min$$

Данный принцип называется методом наименьших квадратов и является одним из самых распространенных методов подбора функциональной зависимости, наилучшим образом описывающим эмпирические данные.

$$\frac{\partial F}{\partial a} = 0, \frac{\partial F}{\partial b} = 0$$

Используя необходимые условия экстремума:

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

Получим:

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

МОДЕЛЬ ПАРНОЙ РЕГРЕССИИ

Введем обозначения: $\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$, $\vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$, $\vec{s} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$

Тогда $\vec{y}^* = a \vec{s} + b \vec{x}$ - вектор "подогнанных" с помощью

регрессионного уравнения значений \vec{y} , лежащий в плоскости (гиперплоскости) векторов \vec{x} и \vec{s}

и $\vec{e} = \vec{y} - \vec{y}^*$ - вектор разности между выборочными значениями \vec{y} и "подогнанными" с помощью функции $\vec{y} = f(\vec{x})$. Наилучшее приближение к \vec{y} будет, если вектор $\vec{e} = \vec{y} - \vec{y}^*$ перпендикулярен обоим векторам \vec{x} и \vec{s} .

И $\vec{e} = \vec{y} - \vec{y}^*$ перпендикулярен обоим векторам \vec{x} и \vec{s} будет, если вектор $\vec{e} = \vec{y} - \vec{y}^*$ перпендикулярен обоим векторам \vec{x} и \vec{s} .

И $X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$, $\vec{\beta} = \begin{pmatrix} a \\ b \end{pmatrix}$. Это условие эквивалентно равенству нулю двух скалярных произведений: $\vec{e} = \vec{y} - X \vec{\beta}$ и $\vec{e} \perp \vec{x}$ и $\vec{e} \perp \vec{s}$. Обозначим:

МОДЕЛЬ ПАРНОЙ РЕГРЕССИИ

Условие ортогональности вектора \vec{e} к плоскости векторов \vec{s} и \vec{t}

$$X^T \vec{e} = \vec{0} \quad \vec{\beta} = (X^T X)^{-1} X^T \vec{y}$$

. После преобразований: . Это уравнение для коэффициентов регрессионного уравнения a и b совпадает с формулами по методу наименьших квадратов.

Пусть реальное (теоретическое) уравнение зависимости y от x имеет вид: $y = a + bx + \varepsilon$, где x - неслучайная (детерминированная) величина, y и ε - случайные величины. Считаем, что $\varepsilon \sim N(0, \sigma)$, где σ^2 - дисперсия ошибок. По данным наблюдений следует найти оценки a^* и b^* для коэффициентов уравнения регрессии и s^2 - для дисперсии ошибок σ^2 .

Имеет место **теорема Гаусса-Маркова**, утверждающая, что при использовании метода наименьших квадратов для определения оценок a^* и b^* коэффициентов уравнения регрессии по формуле

МОДЕЛЬ ПАРНОЙ РЕГРЕССИИ

В соответствии с теоремой Гаусса-Маркова оценки $a^* \sim N(a, \sigma_a)$, $b^* \sim N(b, \sigma_b)$. Имеют место следующие формулы для дисперсий

оценок:

$$D(a^*) = \sigma_a^2 = \sigma^2 \frac{\sum_{i=1}^n x_i^2}{n \left(\sum_{i=1}^n x_i^2 - n(\bar{x})^2 \right)}, \quad D(b^*) = \sigma_b^2 = \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}.$$

Оценки дисперсий:

$$s^2 = \frac{\sum_{i=1}^n (y_i - a^* - b^* x_i)^2}{n-2}, \quad s_a^2 = s^2 \frac{\sum_{i=1}^n x_i^2}{n \left(\sum_{i=1}^n x_i^2 - n(\bar{x})^2 \right)}, \quad s_b^2 = \frac{s^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}.$$

При проверке гипотезы $H_0: b = b_0$ с альтернативой $H_1: b \neq b_0$ используется статистика $t = \frac{b^* - b_0}{s_b}$, имеющая распределение Стьюдента с $(n - 2)$ степенями свободы. Доверительный интервал для

$$b \in (b^* - t_\gamma s_b; b^* + t_\gamma s_b)$$

МОДЕЛЬ ПАРНОЙ РЕГРЕССИИ

Доверительный интервал для коэффициента a с надежностью γ имеет вид: $(-t_\gamma s_a; a^* + t_\gamma s_a)$.

Коэффициентом детерминации (или долей объясненной дисперсии) называется величина: $R^2 = \frac{TSS - ESS}{TSS} = \frac{RSS}{TSS}$ где

$$ESS = \sum_{i=1}^n (y_i - y_i^*)^2 \quad RSS = \sum_{i=1}^n (y_i^* - \bar{y})^2$$

. Очевидно, что $0 \leq R^2 \leq 1$. Если, $R^2 = 0$, то регрессия не улучшает качество прогноза y по сравнению с тривиальным: \bar{y} . Если $R^2 = 1$, то это означает точную "подгонку": все наблюдаемые точки лежат на регрессионной кривой.

Проверяется гипотеза $H_0: b = b_0$ с альтернативой $H_1: b \neq b_0$.
Статистика $F = \frac{\sum_{i=1}^n (y_i - y_i^*)^2}{n-2}$

$\sim F(1; n - 2)$. Если при заданном уровне значимости α критическое значение $F_{\alpha; 1; n - 2}$ меньше наблюдаемого $y = y$

МОДЕЛЬ МНОЖЕСТВЕННОЙ РЕГРЕССИИ

Рассматривается следующее уравнение модели:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2 + \dots + \beta_k x_k^k + \varepsilon,$$

где $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_k)^T$ - вектор коэффициентов модели, x_1, x_2, \dots, x_k - регрессоры (независимые переменные), ε - случайная компонента, $\varepsilon \sim N(0, \sigma)$, где σ^2 - дисперсия ошибок. По результатам

экспериментов $(x_i; y_i), i = \overline{1, n}$, надо определить значения коэффициентов $\beta_0, \beta_1, \beta_2, \dots, \beta_k$.

$$X = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^k \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & x_n^2 & \dots & x_n^k \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}, \quad y^* = \begin{pmatrix} y_1^* \\ y_2^* \\ \dots \\ y_k^* \end{pmatrix}$$

Пусть $e = y - y^*$, где $y^* = (y_1^+ \ y_2^+ \ \dots \ y_k^+)^T$ -

прогноз значений объясняемой переменной. Пусть

β - вектор прогнозируемых значений y^* . Тогда e - вектор разности между выборочными и прогнозируемыми значениями y . Будем

МОДЕЛЬ МНОЖЕСТВЕННОЙ РЕГРЕССИИ

Для оценки коэффициентов множественной регрессии необходимо использовать формулу, совпадающую со случаем парной регрессии:

Формула для оценки дисперсии ошибок σ^2 :
$$s^2 = \frac{\vec{e}^T \vec{e}}{n - (k + 1)}$$

Коэффициент детерминации определяется так же, как и

$$R^2 = \frac{RSS}{TSS} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\vec{y}^T \vec{y} - n(\bar{y})^2}{\vec{y}^T \vec{y} - n(\bar{y})^2}$$

для парной регрессии:

$$V(\vec{\beta}^*) = s^2 (X^T X)^{-1}$$

Для определения дисперсий оценок коэффициентов регрессии вычисляется матрица ковариаций:

дисперсий оценок: $V(\vec{\beta}^*) : s_i^2 = s^2 q_{ii}$ берутся диагональные элементы q_{ii} матрицы $\beta_i \in (\beta_i^* - t_\gamma s_i; \beta_i^* + t_\gamma s_i)$. Доверительный интервал

МОДЕЛЬ МНОЖЕСТВЕННОЙ РЕГРЕССИИ

Для проверки значимости коэффициентов регрессии β_i можно использовать гипотезу $H_0: \beta_i = 0$ при альтернативе $H_1: \beta_i \neq 0$. Двусторонняя критическая область формируется из условия $\left| \frac{\beta_i^*}{s_i} \right| > t_{кр}$,

где $t_{кр}(\alpha; n - (k + 1))$ - критическая точка распределения Стьюдента для уровня значимости α и числа степеней свободы $n - (k + 1)$.

Для оценки значимости уравнения регрессии, то есть проверки факта улучшения прогноза значения y по сравнению с тривиальным: используется критерий Фишера. Гипотеза $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ (свободный член β_0 - не учитывается) проверяется с использованием

статистики $F \sim F(k; n - (k + 1))$. Если наблюдаемое значение F больше критического $F_{кр}(\alpha; k; n - (k + 1))$, то гипотеза

"ПОДГОНКА" КРИВОЙ

Пусть данные некоторого эксперимента представлены в виде таблицы значений переменных x и y . Ставится задача об отыскании аналитической зависимости между x и y , то есть некоторой формулы

$y = f(x)$, явным образом выражающей y , как функцию x . Естественно требовать, чтобы график искомой функции $y = f(x)$ не слишком уклонялся от экспериментальных точек (x_i, y_i) . Эту задачу называют "сглаживанием" экспериментальных данных или "подгонкой" кривой. Для решения можно использовать метод наименьших квадратов (МНК). Согласно МНК указывается вид эмпирической формулы $y = Q(x, a_0, a_1, \dots, a_n)$, где a_0, a_1, \dots, a_n – числовые параметры. Наилучшими значениями параметров a_0, a_1, \dots, a_n считаются те, для которых сумма квадратов отклонений функции

$$S(a_0, a_1, \dots, a_n) = \sum_{i=1}^m (Q(x_i, a_0, a_1, \dots, a_n) - y_i)^2 \rightarrow \min$$

"ПОДГОНКА" КРИВОЙ

Используя необходимые условия экстремума функции нескольких переменных, получаем систему уравнений для определения параметров $\frac{\partial S}{\partial a_0}, \frac{\partial S}{\partial a_1}, \dots, \frac{\partial S}{\partial a_n}$. Если система имеет единственное решение, то оно

является искомым и аналитическая зависимость между экспериментальными данными определяется формулой: $y = f(x)$. В общем случае система уравнений для определения оптимальных значений параметров нелинейна.

Пусть для переменных x и y соответствующие значения экспериментальных данных (x_i, y_i) не располагаются вблизи прямой. Тогда можно выбрать новые переменные $X = \phi(x, y)$ и $Y = \psi(x, y)$ так, чтобы преобразованные экспериментальные данные $X_i = \phi(x_i, y_i)$ и

"ПОДГОНКА" КРИВОЙ

Числа k и b определяются по встроенной функции **ЛИНЕЙН**, где вместо x_i и y_i подставляют соответствующие значения X_i и Y_i .
Функциональная зависимость $y = f(x)$ определена неявно уравнением

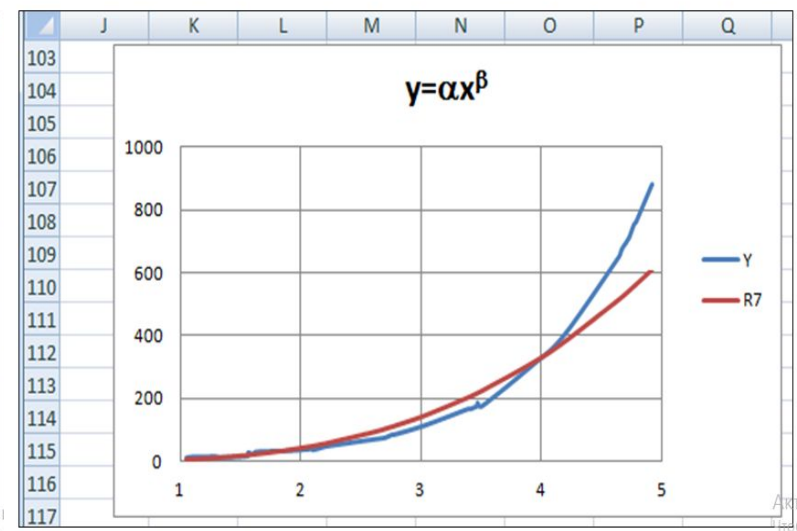
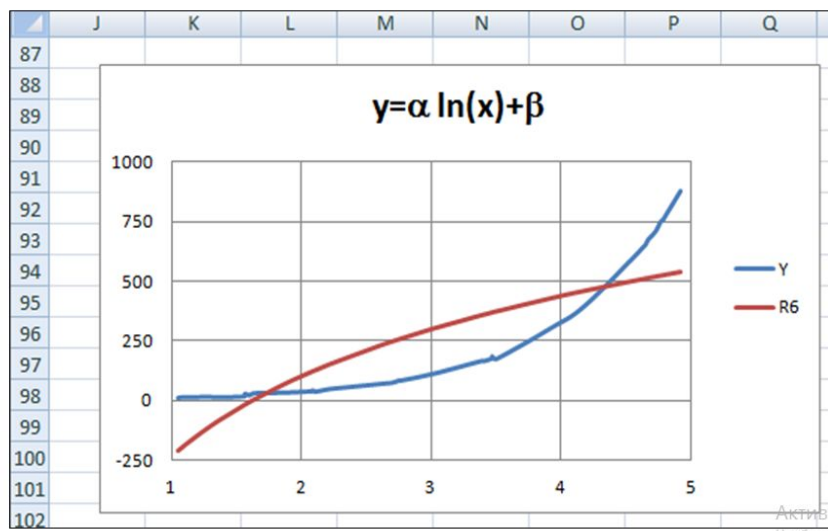
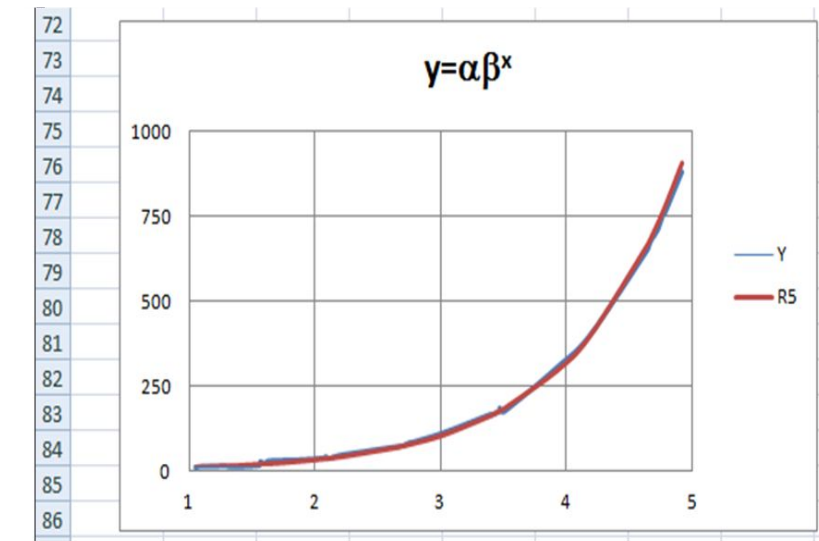
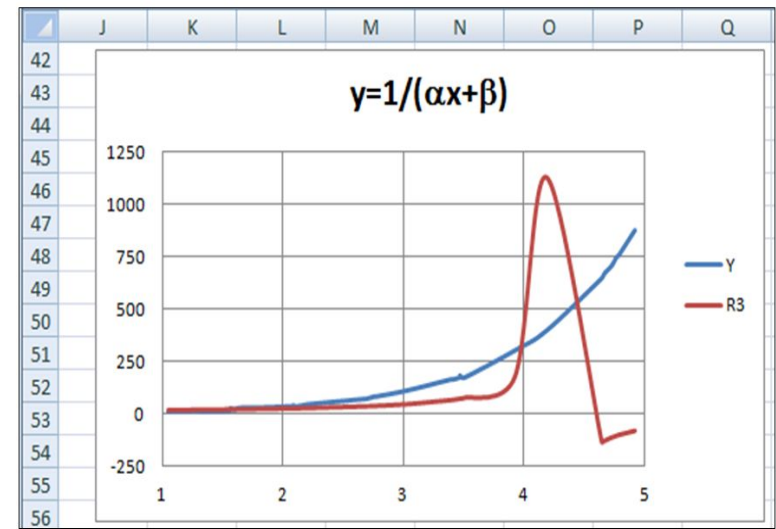
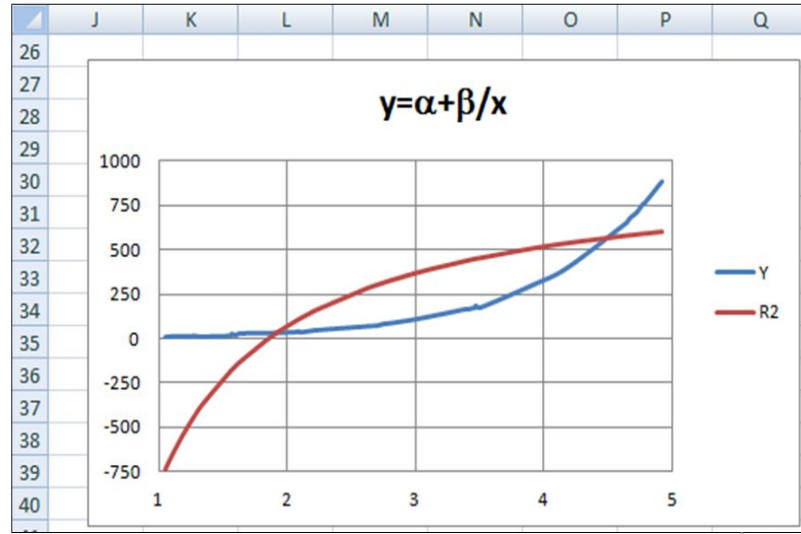
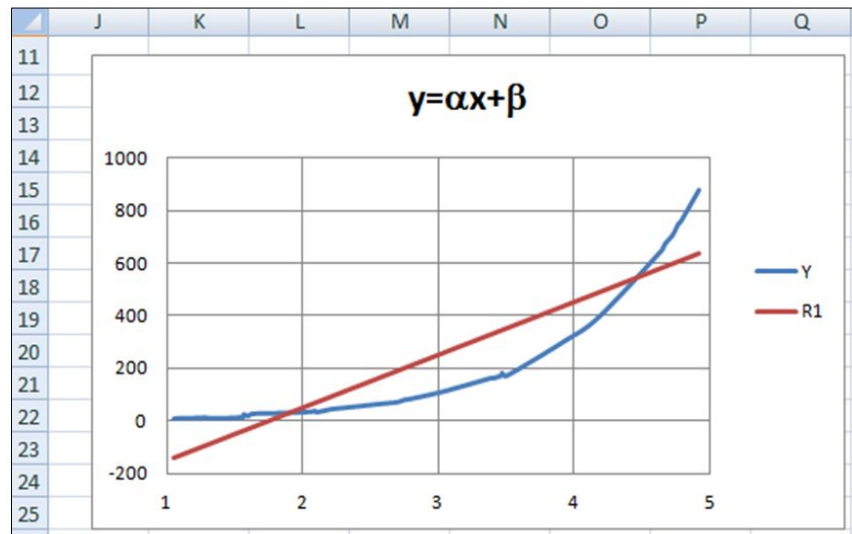
$\psi(x, y) = k\phi(x, y) + b$, которое разрешимо относительно y в частных случаях.

Рассмотрим семь вариантов преобразования переменных, в которых возможно явное выражение переменной y из уравнения $\psi(x, y) = k\phi(x, y) + b$. В таблице приведены формулы преобразования переменных, явная эмпирическая формула $y = f(x)$, зависящая от двух параметров α и β , выражение этих параметров через коэффициенты, полученные с помощью МНК.

"ПОДГОНКА" КРИВОЙ

№	Эмпирическая формула	Выравнивание данных (преобразование переменных)
1.	$y = \alpha x + \beta, \alpha = k, \beta = b$	$X = x, Y = y$
2.	$y = \alpha + \frac{\beta}{x}, \alpha = k, \beta = b$	$X = x, Y_1 = \underline{xy}$
3.	$y = \frac{1}{\alpha x + \beta}, \alpha = k, \beta = b$	$X = x, Y_2 = \frac{1}{y}$
4.	$y = \frac{x}{\alpha x + \beta}, \alpha = k, \beta = b$	$X = x, Y_3 = \frac{x}{y}$
5.	$y = \alpha \beta^x, \alpha = e^b, \beta = e^k$	$X = x, Y_4 = \ln y$
6.	$y = \alpha \ln x + \beta, \alpha = k, \beta = b$	$X_1 = \ln x, Y = y$
7.	$y = \alpha x^\beta, \alpha = e^b, \beta = k$	$X_1 = \ln x, Y_4 = \ln y$

"ПОДГОНКА" КРИВОЙ



Экспериментальные данные наилучшим образом аппроксимируются регрессионной кривой