

**Московский инженерно-физический институт
(Государственный университет)**

**Учебно-методический комплекс
по дисциплине
«Математическая статистика в радиационной физике»**

**Автор: ст. преподаватель каф.№1 Морозова Н.И.
(для групп Т7-01, Т7-01М)**

**Москва
2008**

Содержание

Лекция №1

Лекция №2

Лекция №3

Лекция №4

Лекция №5

Лекция №6

Лекция №7

Лекция №8

ЛЕКЦИЯ I.

***ЭЛЕМЕНТЫ ТЕОРИИ
ВЕРОЯТНОСТЕЙ И ЗАДАЧИ
МАТЕМАТИЧЕСКОЙ
СТАТИСТИКИ***

Целью любого физического эксперимента является **определение неизвестных параметров изучаемого объекта на основании полученных результатов измерений.**

Полученные результаты физических экспериментов, представляют собой сумму двух составляющих: $\bar{x} = \mu + \sigma \sum_{j=1}^n \epsilon_j$, где μ - определяет истинное значение параметра изучаемого или измеряемого объекта, а ϵ_j - вклад j -го внешнего мешающего воздействия на результат. Влияющие на точность измерения факторы можно разделить на **две группы.**

Если группа факторов вносит в измерение **случайные по знаку и по величине, но небольшие искажения, не смещающие результат в сторону увеличения или уменьшения,** то экспериментатор имеет дело со **случайной ошибкой.**

Если появляются факторы, которые **систематически смещают центр тяжести результатов измерений в какую-либо одну сторону,** то появляется **систематическая ошибка.**

Мы будем рассматривать лишь случаи, когда математическое ожидание случайной величины X совпадает со значением неизвестного параметра θ : $M\{X\} = \theta$

Стоит задача определения θ - неизвестного параметра исследуемого объекта по данным независимых измерений этого параметра: x_1, x_2, \dots, x_n .

Оценка параметра - некоторая функция от результатов измерений $\hat{\theta}_n = \hat{\theta}(x_1, x_2, \dots, x_n)$, приближенно равная истинному значению.

Любая функция от результатов измерений называется в современной терминологии **описательной** или **дескриптивной статистикой**.

Рассмотрим каждую из этих описательных статистик и ее свойства.

Пусть у нас есть случайная величина ξ и практическая реализация этой случайной величины при измерении:

$$X_1, X_2, \dots, X_n,$$

Минимум и максимум – соответственно минимальное и максимальное значения полученных при данной выборке значений случайной величины.

Среднее (оценка среднего \bar{x}) вычисляется как среднее результатов наблюдений

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

Оценку среднего часто называют **выборочным средним**. Выборочное среднее обладает следующим **свойством**: **сумма отклонений наблюдаемых значений от среднего арифметического равна нулю**. Эта статистика единственная, которая обладает этим свойством.

Дисперсия выборки или **выборочная дисперсия** определяется следующим образом:

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

Эту описательную статистику очень часто называют **мерой рассеяния выборки случайной величины**. Эта статистика была впервые введена Фишером в 1882 г.

Стандартное отклонение (standard deviation) вычисляется как **корень квадратный из выборочной дисперсии**. На практике более удобно использовать, поскольку измеряется в тех же единицах, что и измеряемая величина, т.е. более наглядно показывает как сильно разбросаны значения выборки относительно выборочного среднего.

ЛЕКЦИЯ 2.

ОЦЕНКИ И ИХ СВОЙСТВА

ТОЧЕЧНЫЕ ОЦЕНКИ И ИХ СВОЙСТВА.

Для наглядности и простоты проанализируем следующие результаты независимых измерений случайной величины

ξ : 2.15, 2.32, 1.99, 2.26, 2.04, 2.19, 2.21, 2.14, 2.24.

Нужно *найти точечную оценку* неизвестного параметра θ

Виды оценок неизвестного параметра θ :

$$1. \hat{\theta}_n = \frac{\xi_1 + \xi_2 + \dots + \xi_n}{n} = \bar{\xi}_n \quad (\text{среднее арифметическое})$$

$$2. \hat{\theta}_n = \xi_{(1)}, \xi_{(2)}, \dots, \xi_{(n)} = \xi_{(1)}, \xi_{(2)}, \dots, \xi_{(n)}$$

$$3. \hat{\theta}_n = \frac{\xi_{(1)} + \xi_{(n)}}{2} = \frac{\xi_{(1)} + \xi_{(n)}}{2} \quad (\text{оценка по двум крайним точкам})$$

$$4. \hat{\theta}_n = \sqrt[n]{\xi_1 \cdot \xi_2 \cdot \dots \cdot \xi_n} = \sqrt[n]{\prod_{i=1}^n \xi_i} \quad (\text{среднее геометрическое})$$

На основании одних и тех же данных может быть получено по разным формулам несколько оценок, близких по величине между собой. Чтобы определить наилучшую, рассмотрим следующие свойства точечных оценок.

Состоятельность.

Оценка $\hat{\theta}_n$ является **состоятельной**, если при увеличении объема выборки n ее значение сходится к истинному значению параметра θ : $\hat{\theta}_n \xrightarrow{p} \theta$.

Здесь подразумевается **выполнение двух условий**:

1. Сходимость математического ожидания оценки к истинному значению параметра:

$$\lim_{n \rightarrow \infty} E[\hat{\theta}_n] = \theta$$

2. Сходимость дисперсии оценки к 0:

$$\lim_{n \rightarrow \infty} \text{Var}[\hat{\theta}_n] = 0$$

Если объем выборки мал, то не обязательно требовать, чтобы используемые оценки были состоятельны. При малых n сходимость не наблюдается для любых оценок.

Несмещенность.

Оценка $\hat{\alpha}_n$ является *несмещенной*, если для любого объема выборки n ее математическое ожидание совпадает с истинным значением параметра:

$$E[\hat{\alpha}_n] = \alpha$$

Эффективность.

Оценка $\hat{\alpha}_n$ называется **эффективной** для заданного объема выборки если ее дисперсия минимальна среди других видов оценок для заданного n .

Оценка $\hat{\alpha}_n$ называется **асимптотически эффективной**, если она эффективна при $n \rightarrow \infty$.

Таким образом, для того чтобы исследовать ту или иную оценку на эффективность **нужно иметь несколько видов описательных статистик**, а **эффективная оценка $\hat{\alpha}_{n, \text{эф}}$** будет выбираться из условия:

$$\frac{D[\hat{\alpha}_{n, \text{эф}}]}{D[\hat{\alpha}_n]} \leq 1$$

Кроме того, как будет показано далее для нахождения минимума дисперсии необходима информация о виде распределения случайной величины, согласно неравенства Крамера-Рао.

Достаточность. Для определения достаточности оценки $\hat{\theta}$ необходимо знать закон распределения случайной величины ξ .

Пусть $p(x, \theta)$ - плотность распределения, зависящая от параметра θ . Тогда вероятность появления выборки x_1, x_2, \dots, x_n определится соотношением:

$$\begin{aligned}
 p(x_1, \theta) \cdot p(x_2, \theta) \dots p(x_n, \theta) &= p(x_1, \theta) \cdot p(x_2, \theta) \cdot \dots \cdot p(x_n, \theta) = \\
 &= \int \dots \int p(x_1, \theta) \cdot p(x_2, \theta) \cdot \dots \cdot p(x_n, \theta) \cdot \delta(x_1 - x_1) \cdot \delta(x_2 - x_2) \cdot \dots \cdot \delta(x_n - x_n) \cdot dx_1 \cdot dx_2 \cdot \dots \cdot dx_n =
 \end{aligned}$$

Второй сомножитель в этом выражении не зависит от параметра θ и не оказывает влияния на величину вероятности. Первый же сомножитель зависит и от параметра θ , и от исходных данных и называется **функцией правдоподобия** $L(x_1, x_2, \dots, x_n, \theta)$.

Оценка $\hat{\theta}$ называется **достаточной**, если функцию правдоподобия можно представить в виде двух сомножителей:

$$L(x_1, x_2, \dots, x_n, \theta) = L_1(\hat{\theta}, \theta) \cdot L_2(x_1, x_2, \dots, x_n)$$

Пример. Попытаемся исследовать на достаточность оценку в виде среднеарифметического $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$. Далее предположим, что случайная величина ξ распределена по нормальному закону $N(\mu, 1)$. **Функция правдоподобия** для выборки x_1, x_2, \dots, x_n будет выглядеть следующим образом:

$$L(x_1, x_2, \dots, x_n) = \frac{1}{\xi^{2n}} \exp\left[-\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{2}\right],$$

тогда для оценки \bar{x} , получаем следующую функцию правдоподобия:

$$\begin{aligned} L(x_1, x_2, \dots, x_n) &= \frac{1}{\xi^{2n}} \exp\left[-\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{2}\right] \\ &= \frac{1}{\xi^{2n}} \exp\left[-\frac{\sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i \bar{x} + n \bar{x}^2}{2}\right] \\ &= \frac{1}{\xi^{2n}} \exp\left[-\frac{\sum_{i=1}^n x_i^2 - 2n \bar{x} \bar{x} + n \bar{x}^2}{2}\right] \exp\left[-\frac{\sum_{i=1}^n x_i^2}{2}\right] \end{aligned}$$

Принимая во внимание что $\frac{\sigma^2}{n} = 1$ есть $\frac{\sigma^2}{n}$, первую экспоненту можно представить в следующем виде: $\exp\left[-\frac{\sum_{i=1}^n x_i^2}{2} + \sum_{i=1}^n x_i\right]$ и в результате получаем, что **функция правдоподобия для случайной величины распределенной по нормальному закону $N(\xi, 1)$ с учетом того, что для оценивания среднего используем среднеарифметическую оценку можно представить в виде двух сомножителей:**

$$L(\xi) = \frac{1}{(2\pi)^{n/2}} \exp\left[-\frac{\sum_{i=1}^n x_i^2}{2} + \sum_{i=1}^n x_i\right]$$

таким образом среднее арифметическое является достаточной оценкой в случае если случайная величина распределена по нормальному закону.


Если аналогичные выкладки провести для оценки из примера 1.1 $\bar{x} = \frac{x_1 + x_2}{2}$, то достаточно легко можно убедиться что эта оценка не является достаточной.

Робастность.

Оценка, *свойства которой не зависят от конкретного вида распределения*, называется ***устойчивой или робастной***. Робастные оценки устойчивы к выбросам, присутствие которых в выборках обусловлено наличием так называемой «шумовой» составляющей и наложением ее распределения на основное.

На ***примере*** нашей выборки и выбранных оценок это свойство можно достаточно просто пояснить следующим образом. Заменяем значение 2.24 на значение 10.23 (ситуация, когда имеет место поломка прибора или кратковременный скачок напряжения в сети). И среди результатов выборки появляется значение, резко отличающееся от других – так называемый ***«выброс»***. Из-за этого оценка \bar{x} изменилась и стала равной 3.06, но оценка s осталась равной 2.19, т.е. не изменилась. Таким образом вторая оценка является ***более устойчивой*** по отношению к другим видам статистик из нашего примера.

Вычислительная простота.

Это свойство точечных оценок также имеет право на существование, но с развитием персональных компьютеров это свойство оценок уходит на задний план. *Из приведенных в примере оценок наиболее простой с точки зрения вычислений является* .

В итоге, понятно что найти наилучшую оценку, удовлетворяющую всем перечисленным свойствам одновременно нельзя.

Надежность оценок. (Понятие интервального оценивания.)

Точечные оценки неизвестных параметров, как правило, являются первоначальным этапом статистического оценивания результатов измерений. Поскольку далее результат конкретного эксперимента представляется в виде $\bar{x} \pm \Delta$ возникает вопрос о том, какова вероятность нахождения истинного значения μ в этом интервале. В этой связи вводится понятие ***надежности оценок***.

Надежность оценки определяется как ***вероятность нахождения истинного значения в указанном интервале с заданной вероятностью***.

Пусть задана *выборка случайной величины* $\{x_1, x_2, \dots, x_n\}$ *распределенная с плотностью вероятности* $p(x, a_1, a_2, \dots, a_k)$ Стоит **задача** оценивания неизвестных параметров распределения.

Интервальной оценкой неизвестного параметра θ называется *интервал* (a, b) *такой, что истинное значение* θ *попадает в этот интервал с заданной вероятностью* $(1-\alpha)$. Вспоминаем, приведенные в первой лекции **обозначения и определения** связанные с интервальным оцениванием:

(a, b) - *интервальная оценка параметра* θ ;

a – *нижний доверительный предел;*

b – *верхний доверительный предел;*

α – *доверительный уровень или уровень значимости;*

$(1-\alpha)$ – *доверительная вероятность.*

В некоторых прикладных задачах наоборот требуется найти для заданной интервальной оценки неизвестную доверительную вероятность. Более подробно для нескольких частных случаев вопросы, связанные с интервальным оцениванием будут рассмотрены далее в лекции 7.

ЛЕКЦИЯ 3.

ОЦЕНИВАНИЕ ПРИ ОТСУТСТВИИ ИНФОРМАЦИИ О ВИДЕ РАСПРЕДЕЛЕНИЯ

Метод подстановки.

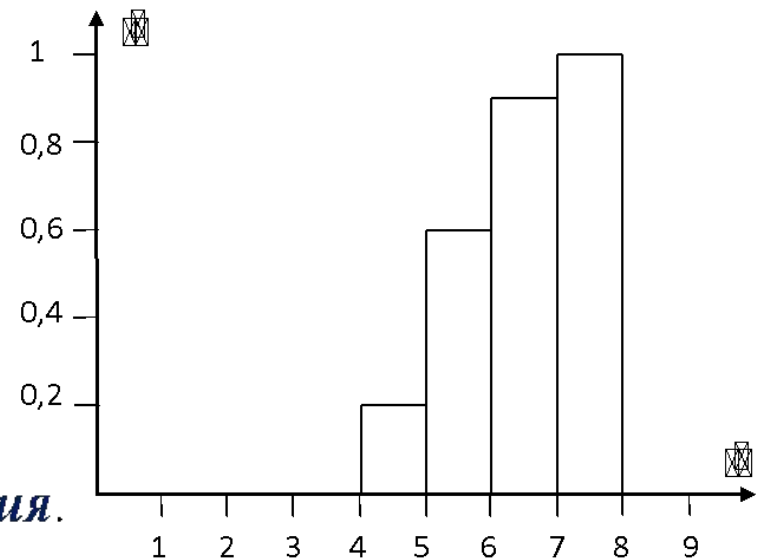
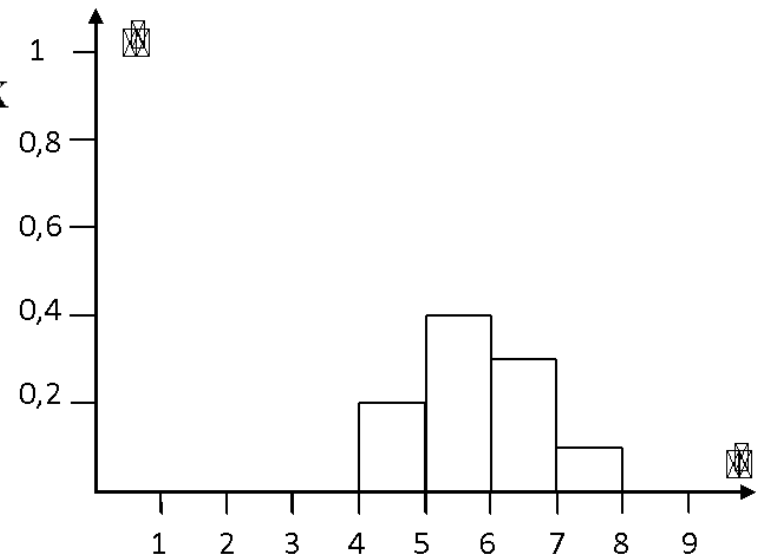
Поставим **задачу** определения оценок среднего и дисперсии по данным независимых измерений случайной **величины** x_1, x_2, \dots, x_n . Из теории вероятности известны **определения среднего** и **дисперсии $D (\sigma^2)$** :

$$\bar{x} = \int_{-\infty}^{\infty} x p(x) dx \quad D = \sigma^2 = \int_{-\infty}^{\infty} x^2 p(x) dx - \bar{x}^2 \quad (ЛЗ.1)$$

где $p(x)$ - плотность вероятности распределения случайной величины x

Принцип подстановки формулируется следующим образом. Для нахождения оценки величины x , которая зависит от функции распределения случайной величины или плотности распределения, в формулу описывающую эту зависимость вместо неизвестных $p(x)$ или $F(x)$ подставляются их эмпирические аналоги.

Эмпирическая функция
распределения $F_n(x)$ в каждой точке x
 равна *отношению числа*
выборочных значений меньших x
к размеру выборки n
 и *представляет собой*
ступенчатую функцию,
 имеющую значения от 0 до 1
 и размер ступенек кратен $1/n$,
 и определяется числом
 выборочных значений
 попавших в выделенный интервал.



Эмпирическая плотность
распределения $f_n(x)$
 равна *производной*
эмпирической функции распределения.

На рис. 1 приведены
 эмпирические плотность
 и функция распределения для следующей выборки:

Если ввести определение *единичной функции* $U(x)$ и *дельта-функции* $\delta(x)$:

$$U(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad \delta(x) = \begin{cases} \infty, & x = 0 \\ 0, & x \neq 0 \end{cases}$$

(Л3.2)

То *эмпирические функция и плотность распределения* запишутся:

$$F(x) = \frac{1}{n} \sum_{i=1}^n U(x - x_i) \quad f(x) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$$

(Л3.3)

Воспользуемся *принципом подстановки* для получения *оценок среднего и дисперсии*:

$$\bar{x} = \int_{-\infty}^{\infty} x f(x) dx = \int_{-\infty}^{\infty} x \frac{1}{n} \sum_{i=1}^n \delta(x - x_i) dx = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{Л3.4})$$

$$\begin{aligned} \sigma^2 &= \int_{-\infty}^{\infty} (x - \bar{x})^2 f(x) dx = \int_{-\infty}^{\infty} (x - \bar{x})^2 \frac{1}{n} \sum_{i=1}^n \delta(x - x_i) dx \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned} \quad (\text{Л3.5})$$

Таким образом, мы получили, что хорошо известные оценки *среднее арифметическое и дисперсия* являются *оценками, полученными методом подстановки*.

Свойства оценок среднего и дисперсии, полученных этим методом.

Проверим - является ли оценка среднего **состоятельной и несмещенной**.

Согласно определению **состоятельности** оценки мы найдем математическое ожидание и дисперсию этой оценки. При этом считаем, что все результаты измерений x_i **независимы и распределены по тому же закону, что и сама случайная величина** X другими словами $M[x_i] = \mu$, а $D[x_i] = M[(x_i - \mu)^2] = \sigma^2$.

Прикладывая оператор среднего и дисперсии к оценке \bar{x}_n (средне арифметической), и вспоминая свойства оператора математического ожидания и дисперсии получаем:

$$M[\bar{x}_n] = M\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n M[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu \quad (Л3.6)$$

$$D[\bar{x}_n] = D\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n^2} \sum_{i=1}^n D[x_i] = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n} \quad (Л3.7)$$

Из (Л3.5) следует, что оценка среднего, полученная на основе метода подстановки, является **несмещенной и состоятельной**.

Теперь проверим **состоятельность и несмещенность оценки дисперсии**, которая получена на основе метода подстановки. Для начала найдем математическое ожидание этой оценки.

$$M\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right] = \frac{1}{n} \sum_{i=1}^n M[(x_i - \bar{x})^2] = \frac{1}{n} \sum_{i=1}^n [M(x_i^2) - 2x_i \bar{x} + \bar{x}^2]$$

Рассмотрим отдельно каждое i -ое слагаемое суммы:

$$M[(x_i - \bar{x})^2] = M[(x_i - \bar{x} + \bar{x} - \bar{x})^2] = M[x_i^2 - 2x_i \bar{x} + \bar{x}^2] = M[x_i^2] - 2\bar{x} M[x_i] + \bar{x}^2 = \sigma^2 - \frac{2\bar{x}^2}{n} + \frac{\bar{x}^2}{n} = \frac{n-1}{n} \sigma^2.$$

Подставляя полученное выражение для каждого слагаемого в сумму (Л3.7) получим:

$$M\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right] = \frac{1}{n} \sum_{i=1}^n \frac{n-1}{n} \sigma^2 = \frac{n-1}{n} \sigma^2. \tag{Л3.8}$$

Этот результат говорит о том, что **для оценки дисперсии, полученной методом подстановки, выполняется первое условие состоятельности** ($\lim_{n \rightarrow \infty} M\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right] = \sigma^2$). Но условие **несмещенности не выполняется. Второе условие состоятельности** ($\lim_{n \rightarrow \infty} \frac{n-1}{n} \sigma^2 = \sigma^2$) также выполняется.

Для того, **чтобы избавиться от смещения в оценке для дисперсии** необходимо в выражение (Л3.5) вместо n подставить $(n-1)$. Тогда получим **выражение для несмещенной и состоятельной оценки дисперсии**:

$$D[\bar{x}] = \frac{\sigma^2}{n-1} \left(\frac{n}{n} - \frac{1}{n} \right) = \frac{\sigma^2}{n-1} \quad (\text{Л3.9})$$

Еще раз воспользуемся **методом подстановки** для нахождения квадрата погрешности оценки среднего (оценка дисперсии оценки среднего):

$$D[D[\bar{x}]] = \frac{\sigma^4}{(n-1)^2} = \frac{\sigma^4}{(n-1)^2} \quad (\text{Л3.10})$$

Очевидно, что **погрешность оценки среднего также обладает своей погрешностью**. В теории математической статистики доказывается, что **отношение погрешности погрешности среднего к погрешности среднего есть величина обратная пропорциональная \sqrt{n}** . Тем обстоятельством, что каждая вычисленная погрешность в свою очередь обладает своей погрешностью, и обусловлен принцип отбрасывания знаков записи результатов физического эксперимента в виде $\bar{x} \pm \Delta \bar{x}$.

ЛЕКЦИЯ 4.

ПЕРЕНОС ОШИБОК

Л4.1 Погрешность косвенных измерений (Одномерный случай).

Пусть есть некоторая случайная величина Y , которая является функцией случайных величин x_1, x_2, \dots, x_n . Стоит задача найти оценку среднего величины Y .

Например, рассмотрим случай $Y = x^2$.

Как правильно построить оценку величины Y .

$$\begin{aligned}
 Y &= x^2 \approx x_0^2 + 2x_0 \Delta x + \Delta x^2 \\
 &= x_0^2 + 2x_0 \frac{\Delta x}{\Delta x} \Delta x + \Delta x^2
 \end{aligned}$$

Тогда математическое ожидание y будет равно:

$$\begin{aligned}
 \langle Y \rangle &= \langle x^2 \rangle = \langle x_0^2 + 2x_0 \Delta x + \Delta x^2 \rangle \\
 &= \langle x_0^2 \rangle + 2x_0 \langle \Delta x \rangle + \langle \Delta x^2 \rangle = x_0^2 + 2x_0 \langle \Delta x \rangle + \langle \Delta x^2 \rangle
 \end{aligned}$$

Воспользовавшись **методом подстановки**, можем подставить в разложение оценки $\hat{\mu}$.. $\hat{\mu}$:

$$\hat{\mu} = \hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu}) + \frac{1}{2} \frac{\partial^2 \hat{\mu}}{\partial \mu^2} (\mu - \hat{\mu})^2 + \dots$$

Данная формула является **более точной** по сравнению с ранее предложенной и **содержит квадратичную поправку**. В случае, если значение этой поправки мало, простая подстановка на начальном этапе может быть удовлетворительной. Ограничимся теперь линейным членом разложения ряда и вычислим **дисперсию** $\hat{\mu}$:

$$\hat{\mu} = \hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu}) + \frac{1}{2} \frac{\partial^2 \hat{\mu}}{\partial \mu^2} (\mu - \hat{\mu})^2 + \dots$$

$$\text{Var}(\hat{\mu}) = \text{Var}\left(\hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu}) + \frac{1}{2} \frac{\partial^2 \hat{\mu}}{\partial \mu^2} (\mu - \hat{\mu})^2 + \dots\right)$$

$$\text{Var}(\hat{\mu}) = \text{Var}\left(\hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu})\right) = \text{Var}\left(\hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu})\right) = \text{Var}\left(\hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu})\right) = \text{Var}\left(\hat{\mu} + \frac{\partial \hat{\mu}}{\partial \mu} (\mu - \hat{\mu})\right)$$

Полученная формула также **приближенная**, так как нелинейные члены ряда отброшены. Но т.к. $\hat{\mu}$ уже определена неточно, то учет членов ряда более высоких порядков не приведет к существенному увеличению точности.

14.2. Формула переноса в случае зависимых измерений.

Рассмотрим более **общий случай**. Пусть **случайные величины** X_1, X_2, \dots, X_n **не являются независимыми**, в этом случае должна быть известна $\Sigma = \Sigma_{ij}$ - **ковариационная матрица** X_1, X_2, \dots, X_n . Эта матрица, также как и оценки $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ может быть получена на основе экспериментальных данных. Например, если оценки $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$ получены на основе выборок

одинакового размера:

$$\begin{matrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{matrix}$$

где первый индекс указывает на номер измерения, а второй – на номер случайной величины, то **оценка (ij)-ого элемента** ковариантной матрицы определится по формуле:

$$\Sigma_{ij} = \frac{1}{m} \sum_{k=1}^m (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)$$

где $\bar{x}_i = \frac{1}{m} \sum_{k=1}^m x_{ki}$, $\bar{x}_j = \frac{1}{m} \sum_{k=1}^m x_{kj}$ - **оценка среднего для каждого столбца**. Стоит **задача** определения Σ_{ij} и Σ_{ij} - **ковариационной матрицы** X_1, X_2, \dots, X_n , где $X_i = X_i(x_1, x_2, \dots, x_n)$ зависимые функции n случайных переменных.

Разложение в ряд Тейлора в окрестности точки (μ_1, \dots, μ_k) и последующее применение оператора математического ожидания приводит к **формуле для оценки** $\hat{\mu}_k$ аналогичной:

$$\begin{aligned} \hat{\mu}_k &= \mu_k + \frac{1}{2} \frac{\partial^2 \ln L(\mu_1, \dots, \mu_k)}{\partial \mu_k^2} \Big|_{\mu_1 = \mu_1, \dots, \mu_k = \mu_k} + \dots \\ &= \mu_k + \frac{1}{2} \frac{\partial^2 \ln L(\mu_1, \dots, \mu_k)}{\partial \mu_k^2} \Big|_{\mu_1 = \mu_1, \dots, \mu_k = \mu_k} + \dots \end{aligned}$$

Для нахождения $\hat{\mu}_k$ рассмотрим вначале **частный случай**. Пусть функция $\ln L(\mu_1, \dots, \mu_k) = \ln L(\mu_1, \dots, \mu_k)$, $j=1, k$, линейны относительно μ_k .

$$\begin{aligned} \ln L(\mu_1, \dots, \mu_k) &= \ln L(\mu_1, \dots, \mu_k) \\ &= \ln L(\mu_1, \dots, \mu_k) \end{aligned}$$

Применим оператор математического ожидания к обеим частям равенства: $E[\ln L(\mu_1, \dots, \mu_k)] = \sigma \dots$

Вычтем из первого равенства второе. Получим:

$$\begin{aligned} \ln L(\mu_1, \dots, \mu_k) - E[\ln L(\mu_1, \dots, \mu_k)] &= \dots (\mu_k - E[\mu_k]) \end{aligned}$$

Подставим в формулу для элемента ковариационной матрицы σ_{ij} , который по определению равен:

$$\begin{aligned} \sigma_{ij} &= \frac{1}{n} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j) = \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \bar{x}_i \bar{x}_j - \bar{x}_j \bar{x}_i + \bar{x}_i \bar{x}_j = \\ &= \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \bar{x}_i \bar{x}_j = \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \bar{x}_i \bar{x}_j \end{aligned}$$

Введем *векторные обозначения*:

$$\begin{aligned} \bar{x}_i &= \frac{1}{n} \sum_{k=1}^n x_{ki}, & \bar{x}_j &= \frac{1}{n} \sum_{k=1}^n x_{kj}, & \bar{x}_i \bar{x}_j &= \frac{1}{n^2} \sum_{k=1}^n \sum_{l=1}^n x_{ki}x_{lj}, \\ \sigma_{ij} &= \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \bar{x}_i \bar{x}_j, & \sigma_{ij} &= \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \frac{1}{n^2} \sum_{k=1}^n \sum_{l=1}^n x_{ki}x_{lj}. \end{aligned}$$

Тогда в векторной форме:

$$\begin{aligned} \bar{x}_i &= \frac{1}{n} \sum_{k=1}^n x_{ki}, \\ \sigma_{ij} &= \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \bar{x}_i \bar{x}_j \end{aligned}$$

Откуда, пользуясь принципом подстановки, получаем:

$$\sigma_{ij} = \frac{1}{n} \sum_{k=1}^n x_{ki}x_{kj} - \bar{x}_i \bar{x}_j$$

Обратимся к **общему случаю**, когда функции $f_{ij} = f_{ij}(x_1, \dots, x_n)$, $x_i = x_i$ **нелинейны**. Разложим каждую функцию в ряд Тейлора вокруг оценки среднего (в идеале нужно раскладывать вокруг истинного среднего, но оно все равно неизвестно, поэтому еще раз применим **принцип подстановки**) с точностью до линейных членов разложения:

$$f_{ij}(x) = f_{ij}(x_0) + \sum_k \frac{\partial f_{ij}}{\partial x_k} (x_k - x_{k0}) + \dots - f_{ij}(x_0)$$

Для такой линейной функции можно применить результат, при этом матрица A определится следующим образом:

$$A_{ij} = \frac{\partial^2 f_{ij}}{\partial x_k \partial x_l} (x_k - x_{k0}) (x_l - x_{l0}) + \dots$$

Сдвиг на const не влияет на величину дисперсии, т.к. $D[\text{const}] = 0$.

Итак, **в общем случае получим**:

$$f_{ij}(x) = \sum_k A_{ij} x_k + \dots,$$

где элементы матрицы определяются из выражения

$$A_{ij} = \frac{\partial^2 f_{ij}}{\partial x_k \partial x_l} (x_k - x_{k0}) (x_l - x_{l0}) + \dots$$

ЛЕКЦИЯ 5.

***МЕТОДЫ ПОЛУЧЕНИЯ ОЦЕНОК
ПРИ ИЗВЕСТНОМ ВИДЕ
РАСПРЕДЕЛЕНИЯ***

Л5.1. Виды распределений случайных величин, встречающиеся в радиационной физике

1. Случайная величина $\xi = \xi_0 + \sigma \sum_{i=1}^n \xi_i$ (n – велико, ξ_i – случайные факторы) – есть **результат действия большого числа случайных факторов на исходную постоянную величину**. В этом случае **применима центральная предельная теорема**, и величина

ξ распределена нормально, т.е. $\xi \sim N(\xi_0, \sigma^2 n)$, где ξ_0 , σ – неизвестные параметры нормального распределения.

2. Имеется N образцов, для которых событие A имеет место с вероятностью p (или не имеет места с вероятностью $1-p$). В эксперименте **подсчитывается случайная величина ν – количество исходов в эксперименте, когда событие A имело место**. В этом случае **действует схема независимых испытаний или схема Бернулли**, при этом распределение случайной величины

будет биномиальным: $P(\nu) = \binom{N}{\nu} p^\nu (1-p)^{N-\nu} = \frac{N!}{\nu! (N-\nu)!} p^\nu (1-p)^{N-\nu}$, где p – неизвестный параметр биномиального распределения.

3. Эксперимент состоит в том, что *в течение некоторого фиксированного промежутка времени t подсчитывается число появлений некоторого события* (например, подсчет гамма-квантов от радиоактивного источника детектирующей системой). Если можно считать, что

- для любых двух неперекрывающихся промежутков времени число зарегистрированных гамма-квантов N_1 за время t_1 и число зарегистрированных гамма-квантов N_2 за время t_2 независимы;
- вероятность одновременной регистрации двух гамма-квантов мала

то *случайное число зарегистрированных гамма-квантов за фиксированный промежуток времени подчиняется закону Пуассона:*

$$P(N) = \frac{\lambda^N}{N!} e^{-\lambda}, \text{ где } \lambda - \text{параметр распределения.}$$

4. Если число зарегистрированных гамма-квантов за фиксированный промежуток времени будет больше 20-30, то распределение зарегистрированных гамма-квантов за фиксированный интервал времени будет ближе к *нормальному распределению*.

Таким образом, *задача статистической обработки данных* в этом случае будет заключаться в *нахождении оценок параметров априорно известных видов распределений*.

Неравенство Крамера-Рао.

Рассмотрим **неравенство**, которое **позволяет установить нижнюю границу дисперсии оценки случайной величины с известным видом распределения.**

Для любого метода оценивания всегда существует предел точности. Дисперсия оценки, найденная любым методом, не может быть меньше значения, установленного этим неравенством.

Вспомним, что **оценка, которая имеет при заданном размере выборки наименьшую дисперсию** называется **эффективной**. (таким образом, зная вид распределения всегда можно задать нижний предел оценивания для дисперсии).

Пусть измерены x_1, x_2, \dots, x_n - **независимые реализации случайной величины ξ** . Считаем, что **случайная величина ξ распределена в общем случае по закону $p(x, a_1, a_2, \dots, a_k)$** , где a_1, a_2, \dots, a_k - неизвестные параметры распределения.

Для того, чтобы упростить дальнейшее изложение **рассмотрим простейший случай**, когда **закон распределения случайной величины зависит только от одного параметра**, а затем обобщим полученный результат.

Точечная оценка параметра $\theta = f(x_1, x_2, \dots, x_n)$. Закон распределения каждого выборочного значения, тот же, что и у случайной величины ξ . Если все выборочные значения независимы, то закон распределения оценки $\hat{\theta}$ определится произведением:

$$P_{\hat{\theta}}(\hat{\theta}) = \zeta_{\theta}(\hat{\theta}, a) \quad (Л4.1)$$

Будем предполагать, что оценка $\hat{\theta}$ - *несмещенная*, т.е. $M\hat{\theta} = a$. Для закона распределения $p(x, a)$ выполняется следующее *условие нормировки*:

$$\int_{-\infty}^{+\infty} p(x, a) dx = 1 \quad (Л4.2)$$

Продифференцируем обе части (Л4.2) по a . Получим:

$$\int_{-\infty}^{+\infty} p'_a(x, a) dx = 0 \quad (Л4.3)$$

Поделим и умножим обе части равенства ((Л4.3) на $p(x, a)$.

$$\int_{-\infty}^{+\infty} \frac{p(x, a)}{p(x, a)} p'_a(x, a) dx = 0$$

Получим *для случая одного неизвестного параметра распределения неравенство Крамера-Рао в виде:*

$$D\hat{\theta} \geq \frac{1}{\int_{-\infty}^{+\infty} \frac{p''_a(x, a)}{p(x, a)} p(x, a) dx} \quad (Л4.4)$$

Пример. Рассмотрим, какой вид приобретает неравенство Крамера-Рао для случайной величины, нормально распределенной. Пусть:

$$f(x, a) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right]$$

$$\ln(f(x, a)) = -\frac{(x-a)^2}{2\sigma^2} + \ln\left(\frac{1}{\sigma \sqrt{2\pi}}\right)$$

$$\frac{\partial \ln(f(x, a))}{\partial a} = \frac{x-a}{\sigma^2}$$

$$I = \int_{-\infty}^{+\infty} \left(\frac{x-a}{\sigma^2}\right)^2 \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right] dx = \frac{1}{\sigma^2}$$

Таким образом, **каким бы способом не была найдена оценка, ее дисперсия не может быть меньше величины, определяемой неравенством:**

$$D[\hat{a}] \geq \frac{\sigma^2}{n} \quad (Л4.5)$$

Ранее методом подстановки была получена **оценка среднего и дисперсия оценки среднего** $D[\hat{a}] = D\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{\sigma^2}{n}$. На основании сравнения (Л4.5) и этой оценки дисперсии оценки среднего можно сделать вывод о том, что **для нормального распределения эта оценка будет эффективной.**

ЛЕКЦИЯ 6.

***МЕТОД НАИМЕНЬШИХ
КВАДРАТОВ***

Постановка задачи

Имеются *пары измерений*

$$x_i, y_i, \sigma_i = x_i, y_i,$$

при этом x_i измерены

без погрешности,

а y_i имеют погрешности σ_i .

Зависимость между x и y :

$$y = f(x, \theta_1, \dots, \theta_n).$$

Задача состоит в отыскании оценок неизвестных параметров

$\theta_1, \dots, \theta_n$ **и их погрешностей.** Аналитическая зависимость может

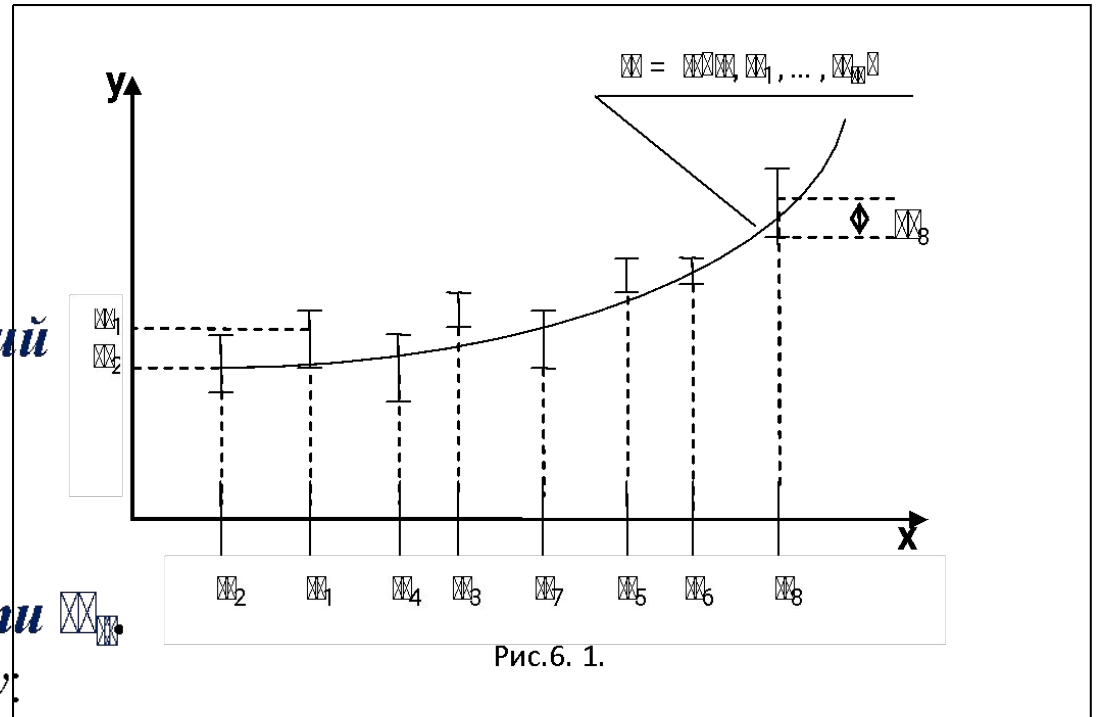
быть известна из теоретических соображений или выбрана из

практических соображений. Эту постановку задачи можно легко

распространить и на многомерный случай, когда функция $y =$

$f(x, \theta_1, \dots, \theta_n)$ зависит не от одной переменной x , а от нескольких

переменных.



Пусть случайные величины x_i подчиняются нормальному закону с параметрами $\mu, \sigma^2, \dots, \mu, \sigma^2$. В этом случае задача отыскания неизвестных параметров может быть решена **методом максимального правдоподобия**. При этом максимизируется функция правдоподобия:

$$L(x_1, \dots, x_n) = \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right]$$

(Л6.1)

где $w_i = \frac{1}{\sigma^2}$ статистический вес i -го измерения.

Очевидно, **максимум функции правдоподобия соответствует минимуму функционала**:

$$S(x_1, \dots, x_n) = \sum_{i=1}^n w_i (x_i - \mu)^2$$

(Л6.2)

Этот функционал представляет собой сумму взвешенных квадратов отклонений измеренных значений y_i от их ожидаемых значений. При этом, чем больше погрешность измерения, тем меньше вклад этого измерения в сумму.

Оценки параметров $\hat{\beta}_j$ находятся из условия минимума суммы квадратов отклонений теоретических значений от экспериментальных в точках X_j , взятых с весами обратно пропорциональными квадрату среднеквадратичной погрешности измерений. В результате **оценки параметров** могут быть найдены как решение системы уравнений:

$$\frac{\sum_{i=1}^n X_j^2 Y_i}{\sum_{i=1}^n X_j^2} = 0, \text{ при } j=1, \dots, k \quad (Л6.3)$$

В этом суть метода, получившего в дальнейшем название **метода наименьших квадратов (МНК)**.

Метод был впервые предложен Лежандром в 1805 году, дальнейшее развитие получил в работах Гаусса в 1809 году. При выводе метода было сделано предположение, что измеренные величины распределены по нормальному закону. Как показала практика и теория, этот **метод может быть применен и для других видов распределения, но только в случае нормального распределения измеренных величин**, оценки, полученные МНК совпадают с оценками полученными методом максимального правдоподобия, т.е. обладают свойством асимптотической эффективности.

Метод наименьших квадратов для случая линейной зависимости от неизвестных параметров.

Большое количество задач, решаемых МНК, может быть сведено к случаю, когда **функция** $y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_n x_n$, **зависит от неизвестных параметров линейно.** Функцию можно представить в виде:

$$y = \sum_{j=0}^n a_j x_j = \sigma_{j=0}^n a_j x_j \quad (Л6.4)$$

Введем следующее обозначение:

$$A_{ij} = x_j \quad (Л6.5)$$

*Матрица A размерности $n \times k$, элементы которой определяются выражением (Л6.5), называется **конструкционной матрицей.*** Так, например, если стоит задача аппроксимации зависимости полиномом второй степени, а число замеров равно пяти, то конструкционная матрица A размерности 5×3 запишется следующим образом.

$$\begin{vmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \\ 1 & x_4 & x_4^2 \\ 1 & x_5 & x_5^2 \end{vmatrix}$$

С учетом введенных обозначений (Л6.4) и (Л6.5) **минимизируемый функционал** переписывается в виде:

$$S(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n) = \sigma_{\mathbf{X}_1}^2 \mathbf{X}_1^T \mathbf{X}_1 + \dots + \sigma_{\mathbf{X}_n}^2 \mathbf{X}_n^T \mathbf{X}_n - \sigma_{\mathbf{X}_1}^2 \mathbf{X}_1^T \mathbf{X}_2 + \dots + \sigma_{\mathbf{X}_n}^2 \mathbf{X}_n^T \mathbf{X}_1 \quad (\text{Л6.6})$$

Минимум этого функционала находится из системы уравнений:

$$\sigma_{\mathbf{X}_1}^2 \mathbf{X}_1 + \sigma_{\mathbf{X}_2}^2 \mathbf{X}_2 + \dots + \sigma_{\mathbf{X}_n}^2 \mathbf{X}_n = \sigma_{\mathbf{X}_1}^2 \mathbf{X}_2 + \dots + \sigma_{\mathbf{X}_n}^2 \mathbf{X}_1 \quad (\text{Л6.7})$$

Линейный метод наименьших квадратов очень удобно представлять в матричной записи. Введем следующие **обозначения** векторов и матриц:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \dots \\ \mathbf{X}_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \dots \\ \mathbf{X}_n \end{bmatrix} \quad \mathbf{A} = \begin{bmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 & \dots & \mathbf{X}_1^T \mathbf{X}_n \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 & \dots & \mathbf{X}_2^T \mathbf{X}_n \\ \dots & \dots & \dots & \dots \\ \mathbf{X}_n^T \mathbf{X}_1 & \mathbf{X}_n^T \mathbf{X}_2 & \dots & \mathbf{X}_n^T \mathbf{X}_n \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} \sigma_{\mathbf{X}_1}^2 & & & \\ & \sigma_{\mathbf{X}_2}^2 & & \\ & & \dots & \\ & & & \sigma_{\mathbf{X}_n}^2 \end{bmatrix} \quad (\text{Л6.7})$$

Тогда систему (Л6.7) можно переписать в следующем виде:

$$\mathbf{A} \mathbf{X} = \mathbf{W} \mathbf{Y} \quad (\text{Л6.8})$$

Решение системы (Л6.8) запишется в виде:

$$\mathbf{X} = (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} \mathbf{Y} \quad (\text{Л6.9})$$

Ковариационную матрицу $\mathbf{V}_{\mathbf{X}}$ для нахождения погрешности параметров \mathbf{X} найдем с помощью матричной записи **формулы переноса ошибок**. Промежуточные выкладки опускаем и получаем:

$$\mathbf{V}_{\mathbf{X}} = (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \quad (\text{Л6.10})$$

Теорема Гаусса-Маркова.

Теорема позволяет оценить эффективность оценок, полученных с помощью линейного метода.

Формулируется следующим образом, *пары измерений* $\{x_i, y_i\}, i = 1, 2, \dots, n$, *при этом* x_i *измерены без погрешности (т.е. не являются случайными величинами), а* y_i *имеют погрешности* ϵ_i . Известно, что зависимость между x и y описывается формулой: $y = a_0 + a_1 x + \epsilon$. При определении оценок неизвестных параметров аппроксимирующих зависимостей *метод наименьших квадратов является оптимальным в случае линейной зависимости от параметров при любом виде распределения.*

Метод наименьших квадратов в нелинейном случае.

Метод наименьших квадратов может **применяться и в случае, когда функция $f = f(x_1, x_2, \dots, x_n)$ в функционале (Л6.2) зависит от неизвестных параметров **нелинейно**.**

Существует огромное число методов, позволяющих решить **задачу минимизации функционала** (Л6.2). В частности, иногда встречается случай, когда **можно подобрать преобразование обеих частей равенства $f = f(x_1, x_2, \dots, x_n)$ таким образом, что с вводом новых переменных $u = u(x_1, x_2, \dots, x_n)$, $v = v(x_1, x_2, \dots, x_n)$, $w = w(x_1, x_2, \dots, x_n)$, $z = z(x_1, x_2, \dots, x_n)$, ... $z = z(x_1, x_2, \dots, x_n)$ вновь полученное равенство $f = f(x_1, x_2, \dots, x_n)$ будет **содержать функцию g , зависящую от параметров линейно**. Для некоторых частных функций f при минимизации функционала (Л6.2) **удается получить нелинейную систему уравнений**, аналогичную (Л6.2), которая может быть **решена аналитически**.**

Однако, в ряде случаев нелинейности функции нельзя получить решение аналитически. В этом случае применяются специальные численные методы.

Метод линеаризации Ньютона-Рафсона.

Достоинство этого метода заключается в том, что он является универсальным и не зависит от конкретного вида функции.

Пусть так же, как и ранее *имеются пары независимых измерений* $x_i, y_i, i = 1, 2, \dots, n$, при этом x_i измерены без погрешности, а y_i имеют погрешность Δy_i . Известно, что зависимость между x и y описывается нелинейной относительно параметров функцией $y = f(x, a_1, a_2, \dots, a_n)$.

1. Выбирается **начальное приближение** x_0, y_0 , на основе конкретного вида функции f . Например, для функции

$$y = \frac{a_1}{x} + a_2 x = \frac{a_1}{x} + a_2 x$$

Начальное приближение x_0, y_0 - *площадь под ломаной кривой*, соединяющей точки $(x_i, y_i), i = 1, 2, \dots, n$. Начальное приближение $x_0 = \frac{\Delta x}{2}$, где Δx - ширина на половине максимальной высоты над осью X графика ломаной соединяющей точки $(x_i, y_i), i = 1, 2, \dots, n$. Начальное приближения $y_0 = \frac{y_{\max}}{2}$ соответствующее максимальному значению $y_i, i = 1, 2, \dots, n$.

2. Функция $y = y_1, y_2, \dots, y_n$ **раскладывается в ряд Тейлора в окрестности точки** x_1, \dots, x_n с точностью до линейных членов разложения:

$$y_1, y_2, \dots, y_n = y_1(x_1, x_2, \dots, x_n), y_2(x_1, x_2, \dots, x_n), \dots, y_n(x_1, x_2, \dots, x_n) + \sigma \frac{\partial y_i}{\partial x_j} \Delta x_j, \dots, \Delta x_n,$$

где $\Delta x_j = x_j - x_j^0$.

3. Полученная функция линейна по отношению к параметрам Δx_j . Следовательно, можно применить **метод наименьших квадратов** для линейного случая. **Матричные обозначения**, аналогичные (Л6.8), будут выглядеть следующим образом:

$$\Delta y = \begin{bmatrix} y_1 - y_1^0 & \dots & y_n - y_n^0 \\ y_2 - y_2^0 & \dots & y_n - y_n^0 \\ \dots & \dots & \dots \\ y_m - y_m^0 & \dots & y_n - y_n^0 \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \dots & \dots & \dots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \Delta x = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \dots & \dots & \dots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \dots \\ \Delta x_n \end{bmatrix}$$

Л6.8

4. *Решение линейной задачи* запишется в соответствии с (Лб.10) и обозначениями (Лб.11) следующим образом:

$$\Delta \mathbf{x}^{(k)} = \mathbf{J}(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)})$$

Ковариационная матрица параметров $\Delta \mathbf{x}^{(k)}$:

$$\mathbf{C}_{\Delta \mathbf{x}^{(k)}} = \mathbf{J}(\mathbf{x}^{(k)})^{-1} \mathbf{C}_{\mathbf{F}} \mathbf{J}(\mathbf{x}^{(k)})^{-T}$$

5. *Следующее приближение значение параметров* запишется:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)}$$

6. Возвращаемся к пункту 2 и повторяем процедуру относительно первого приближения и так далее. *Итерационная формула* будет выглядеть следующим образом:

$$\Delta \mathbf{x}^{(k)} = \mathbf{J}(\mathbf{x}^{(k)})^{-1} \mathbf{F}(\mathbf{x}^{(k)})$$

Ковариационная матрица параметров:

$$\mathbf{C}_{\Delta \mathbf{x}^{(k)}} = \mathbf{J}(\mathbf{x}^{(k)})^{-1} \mathbf{C}_{\mathbf{F}} \mathbf{J}(\mathbf{x}^{(k)})^{-T}$$

7. *Очередное приближение значений параметров* запишется:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)}$$

8. Итерационный процесс продолжается до тех пор, пока не выполнятся одно из следующих предварительно выбранных условий:

а) Достигается заранее **назначенная точность вычисления параметров** ϵ_j , то есть для всех j выполняется:

$$|\Delta x_j| \leq \epsilon_j$$

б) Приращение минимизируемого функционала достигает заранее назначенной точности ϵ :

$$F(x^{k+1}) - F(x^k) \leq \epsilon$$

в) **Модуль приращения** каждого из параметров на данной итерации будет **составлять заданную долю** $\epsilon = \epsilon_{rel} \div \epsilon_{abs}$ **вычислительной погрешности**:

$$|\Delta x_j| \leq \epsilon_{rel} \cdot |x_j| + \epsilon_{abs} = \epsilon_{rel} \cdot |x_j| + \epsilon_{abs}$$

Если процесс расходится, или достигается локальный минимум вместо глобального, то существуют различные способы выхода из подобных ситуаций:

1. **Решение задачи для нескольких начальных приближений и выбор решения, которое соответствует глобальному минимуму функционала** (Л6.2). (если функционал имеет несколько локальных минимумов).

2. **Замена параметра** $\alpha_j = \frac{\alpha_j}{\alpha_j}$, если минимум находится на бесконечности, то есть начиная с некоторой итерации приращение некоторого j -ого параметра превышает заданное значение: $\Delta \alpha_j > \alpha_j$ (наблюдается **расходимость** по этому параметру).

В случае переопределенности, т.е. если два (или более) параметра α_i и α_j сильно коррелированы между собой, наблюдается расходимость.

Чтобы **проверить**, действительно ли это случай переопределенности, вычисляются элементы корреляционной матрицы: $\alpha_i \alpha_j = \frac{\alpha_i \alpha_j}{\alpha_i \alpha_j}$. Если

для каких-нибудь i и j $\alpha_i \alpha_j > \alpha_i \alpha_j \div \alpha_i \alpha_j$, то **между параметрами α_i и α_j существует сильная зависимость**. В этом случае рекомендуется **замена**: $\alpha_j = \alpha_j \alpha_i$, обоснованная конкретными условиями задачи.

Метод наименьших квадратов в случае, когда обе переменные имеют погрешность.

Пусть заданы *пары независимых измерений* x_i, y_i , $i = 1, 2, \dots, n$, при этом x_i имеют погрешность Δx_i , а y_i имеют погрешность Δy_i . Зависимость между x и y описывается нелинейной относительно параметров функцией $y = f(x, a_1, a_2, \dots, a_m)$.

Поскольку значения y_i измерены с ошибкой, то это само по себе вносит погрешность в значения a_1, \dots, a_m . **По формуле переноса ошибок** эта погрешность будет определяться:

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

Перенесем эту погрешность на значения x_i , а y_i будем считать измеренными без погрешности. Таким образом, **полная погрешность y** определится равенством:

$$\Delta y = \frac{\sum_{i=1}^n \frac{\partial f}{\partial x} \Delta x_i}{\sum_{i=1}^n \left(\frac{\partial f}{\partial x} \right)^2} + \Delta y_{\text{изм}}$$

С учетом полной погрешности функционал (Л6.2) перепишется следующим образом:

$$\sigma_{\mathbb{R}} \left[\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots \right] = \sigma_{\mathbb{R}} \left[\frac{\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots}{\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots} \right] = \sigma_{\mathbb{R}} \left[\frac{\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots}{\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots} \right] \quad (\text{Л6.12})$$

Если $\left| \frac{\partial \mathcal{L}}{\partial x_k} \right| \ll \left| \frac{\partial \mathcal{L}}{\partial x_l} \right|$, то дополнительной погрешностью, вносимой $\left| \frac{\partial \mathcal{L}}{\partial x_k} \right|$ можно пренебречь, и задача сводится к одному из случаев, рассмотренных ранее.

В противном случае минимум функционала (Л6.12) ищется с помощью **итерационной процедуры**, аналогичной рассмотренной выше. При этом *на m -ой итерации будет минимизироваться функционал:*

$$\sigma_{\mathbb{R}} \left[\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots \right] = \sigma_{\mathbb{R}} \left[\frac{\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots}{\sum_{k=1}^n \left(\frac{\partial \mathcal{L}}{\partial x_k} \right)^2 + \dots} \right] \quad (\text{Л6.13})$$

где $\left| \frac{\partial \mathcal{L}}{\partial x_k} \right|$, $\mathbb{R} = \mathbb{R}, \mathbb{R}$ – искомое $(m+1)$ -ое приближение, по которому производится минимизация функционала (Л6.12), $\left| \frac{\partial \mathcal{L}}{\partial x_k} \right|$, $\mathbb{R} = \mathbb{R}, \mathbb{R}$ – m -ое приближение, уже найденное по предыдущей итерации.

ЛЕКЦИЯ 7. I.

ИНТЕРВАЛЬНОЕ ОЦЕНИВАНИЕ

Метод доверительных интервалов.

Пусть задана выборка случайной величины \mathbb{X} x_1, x_2, \dots, x_n распределенная с плотностью вероятности $p(x, a_1, a_2, \dots, a_k)$ Стоит **задача** оценивания неизвестных параметров распределения.

Интервальной оценкой неизвестного параметра \mathbb{X} называется **интервал** (a, b) такой, что истинное значение \mathbb{X} **попадает в этот интервал с заданной вероятностью** $(1 - \mathbb{X})$. Вспоминаем, приведенные в первой лекции **обозначения и определения** связанные с интервальным оцениванием:

(a, b) - интервальная оценка параметра \mathbb{X} ;

a – нижний доверительный предел;

b – верхний доверительный предел;

\mathbb{X} – доверительный уровень или уровень значимости;

$(1 - e)$ – доверительная вероятность.

В некоторых прикладных задачах \mathbb{X} наоборот требуется найти для заданной интервальной оценки неизвестную доверительную вероятность.

Рассмотрим решение поставленной задачи для двух частных случаев, а затем обобщим решение для общего случая.

Частный случай 1.

Пусть случайная величина X распределена по нормальному закону:

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$$

При этом *стандартное отклонение σ заранее задано.* Таким образом закон распределения имеет только один *неизвестных параметр μ .* Воспользуемся *точечной оценкой среднего значения,* полученной для нормального распределения методом максимального правдоподобия:

$$\hat{\mu} = \frac{\sigma^2 \sum_{i=1}^n x_i}{n}$$

Случайная величина X , как сумма нормально распределенных случайных величин, также подчиняется нормальному распределению с параметрами $(\hat{\mu}, \sigma/\sqrt{n})$:

$$f(x) = \frac{1}{\sigma \sqrt{2\pi} (\frac{\sigma}{\sqrt{n}})} \exp\left[-\frac{(x - \hat{\mu})^2}{2(\frac{\sigma}{\sqrt{n}})^2}\right] \quad (Л7.1.1)$$

Введем **нормированную случайную величину** $t = \frac{(\bar{x} - \mu)}{s/\sqrt{n}}$,

которая согласно ЦПТ будет распределена по нормальному закону $N(0,1)$. Тогда по таблицам нормального распределения **можно подобрать такое значение t_0 , что вероятность того, что значение t не превысит t_0 равно заданной доверительной вероятности:**

$$P\left\{t \leq t_0\right\} = P\left\{\bar{x} - \mu \leq \frac{t_0 s}{\sqrt{n}}\right\} = P\left\{\mu \geq \bar{x} - \frac{t_0 s}{\sqrt{n}}\right\} = 1 - \alpha \quad (\text{Л7.1.2})$$

С учетом того, что $t = \frac{(\bar{x} - \mu)}{s/\sqrt{n}}$

$$P\left\{\bar{x} - \frac{t_0 s}{\sqrt{n}} < \mu < \bar{x} + \frac{t_0 s}{\sqrt{n}}\right\} = P\left\{-t_0 < \frac{\bar{x} - \mu}{s/\sqrt{n}} < t_0\right\} = 1 - \alpha$$

(Л7.1.3)

Таким образом, **найден доверительный интервал (a, b) :**
 $a = \bar{x} - \frac{t_0 s}{\sqrt{n}}$ и $b = \bar{x} + \frac{t_0 s}{\sqrt{n}}$, такой, что неизвестное истинное значение параметра μ попадает в него с вероятностью $1 - \alpha$

Частный случай 2.

Пусть случайная величина распределена по нормальному закону:

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}, \quad \text{но при этом значение } \mu \text{ неизвестно.}$$

Воспользуемся *точечными оценками среднего значения и дисперсии, полученными методом максимального правдоподобия:*

$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n}, \quad \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Случайная величина $\hat{\mu}$, также как и предыдущем случае подчиняется нормальному распределению, но нормированная величина $t = \frac{(\hat{\mu} - \mu) \sqrt{n}}{\hat{\sigma}}$, уже *не будет распределена по закону $N(0,1)$, так как в знаменателе выражения стоит величина, зависящая от исходных данных.* Согласно закона распределения функции случайных величин $t = F(x_1, x_2, \dots, x_n)$, в этом случае t будет распределено *по закону распределения Стьюдента:*

$$f(t) = \frac{\Gamma\left(\frac{m+1}{2}\right)}{\Gamma\left(\frac{m}{2}\right) \sqrt{m\pi}} \left(1 + \frac{t^2}{m}\right)^{-\frac{m+1}{2}} \quad (Л7.1.4)$$

где m – число степеней свободы, а Γ - гамма-функция, определяемая выражениями:

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt, \quad \Gamma(x) = (x-1) \Gamma(x-1); \quad (Л7.1.5)$$

$$\Gamma(x) = \frac{\Gamma(x+1)}{x}; \quad \Gamma(0) = \Gamma(1) = 1.$$

В данном случае **число степеней свободы t равно $(n-1)$** , т.к. на исходные данные накладывается одна дополнительная связь при подстановке в выражение для t вместо истинного значения μ его оценки \bar{x} .

По таблицам распределения Стьюдента для заданного числа степеней свободы $\nu = n - 1$ можно подобрать такое значение t_0 , что вероятность того, что значение t не превысит t_0 , равна заданной доверительной вероятности:

$$P\{t \leq t_0\} = P\left\{\frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \leq t_0\right\} = P\{\bar{x} - \mu \leq t_0 \frac{s}{\sqrt{n}}\} = P\{\mu \geq \bar{x} - t_0 \frac{s}{\sqrt{n}}\} \quad (Л7.1.6)$$

С учетом того, что $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$ соответствующее равенство для μ запишется:

$$P\{a \leq \mu \leq b\} = P\left\{\bar{x} - t_0 \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} + t_0 \frac{s}{\sqrt{n}}\right\} = P\left\{\bar{x} - t_0 \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} + t_0 \frac{s}{\sqrt{n}}\right\} = P\left\{\bar{x} - t_0 \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} + t_0 \frac{s}{\sqrt{n}}\right\} \quad (Л7.1.7)$$

Таким образом, для второго частного случая **найден доверительный интервал (a, b)** , $a = \bar{x} - \frac{t_0 s}{\sqrt{n}}$ и $b = \bar{x} + \frac{t_0 s}{\sqrt{n}}$, такой, что неизвестное истинное значение параметра попадает в него с вероятностью $1 - \alpha$

ЛЕКЦИЯ 7.2.

***ПРОВЕРКА СТАТИСТИЧЕСКИХ
ГИПОТЕЗ***

Критерий хи-квадрат.

Пусть на заключительном этапе обработки результатов измерений ***стоят следующие задачи:***

1. имеется выборка, состоящая из реализации случайной величины X_1, X_2, \dots, X_n ; и необходимо сделать обоснованное заключение о виде распределения случайной величины
2. измерения случайной величины X в точках x с погрешностью Δ аппроксимирующей отношение между величинами зависимости
3. несколько выборок произвольного (может быть одинакового) размера одной или нескольких случайных величин с известной или неизвестной дисперсией.

Тогда ***любое из перечисленных выше предположений экспериментатора, проверяемое на основании выборочных экспериментальных данных, будет называться статистической гипотезой.***

В первом случае это будет гипотеза о виде распределения случайной величины, во втором случае это будет гипотеза о виде аппроксимирующей зависимости между двумя случайными или в ряде случаев зависимыми величинами и третьем случае гипотеза о статистиках случайных величин о наличии между ними зависимости.

На данном занятии рассмотрим несколько ***частных случаев*** проверки статистических гипотез.

Частный случай 1. *Имеются измерения случайной величины y_i в точках x_i с погрешностью $\sigma_i = \sigma, \dots, \sigma$. Известно, что y_i распределены по нормальному закону $N(\mu, \sigma^2)$. Требуется **проверить статистическую гипотезу о том, что данная зависимость описывается кривой $y = f(x)$ на уровне значимости α** . Отсутствие параметров в формуле для $f(x)$ означает, что они заранее известны (фиксированы).*

Например: $y = f(x) = 2x + 1$. Параметры линейной зависимости здесь фиксированы. Отнормируем значения y_i :

$$z_i = \frac{y_i - f(x_i)}{\sigma}$$

Отнормированные значения будут распределяться по нормальному закону $N(0,1)$. Найдем **сумму квадратов**:

$$Q = \sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n \sigma^2 \left(\frac{y_i - f(x_i)}{\sigma} \right)^2 \quad (Л7.2.1)$$

Случайная величина Z распределена по закону χ^2_k (хи-квадрат с k степенями свободы) с плотностью распределения $f(x)$:

$$f(x) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2} \quad (Л7.2.2)$$

где $\Gamma(k/2)$ – гамма функция

Таким образом, для заданных значений n и α *по таблице хи-квадрат распределения находим значение χ^2_{α}* , соответствующее решению уравнения:

$$\chi^2_{\alpha} = \chi^2_{\alpha} - \chi^2_{\alpha} \quad (Л7.2.3)$$

В случае, *если* найденное значение $\chi^2_{\alpha} < \chi^2_{\alpha}$, то *исходная гипотеза* о том, что данные описываются кривой $f(x)$, *отвергается*. В противном случае гипотеза не отвергается. Однако, это еще не означает, что она бесспорно верна, поскольку *существует бесконечно много аналогичных гипотез*, которые не будут отвергнуты на том же уровне значимости. Поэтому, в результате проверки такой статистической гипотезы не заключают, что гипотеза верна, а *делают вывод о том, что гипотеза не противоречит данным с доверительной вероятностью $1 - \alpha$* . При этом α – вероятность отвергнуть правильную гипотезу.

Частный случай 2.

Имеются измерения случайной величины X_i в точках x_i с погрешностью Δx_i , $\Delta x_i = \Delta x, \dots, \Delta x$, известно, что X_i распределены по нормальному закону. Требуется **проверить статистическую гипотезу о том, что данная зависимость описывается кривой $f(x_1, x_2, \dots, x_n)$ на уровне значимости α при некоторых значениях неизвестных параметров $\theta_1, \dots, \theta_k$. Отнормированные значения Z_i будут функциями неизвестных параметров:**

$$Z_1, \dots, Z_n = \frac{X_i - f(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_k)}{\Delta x_i} \quad (Л7.2.4)$$

и будут как и в предыдущем случае распределены по нормальному закону. **Сумма квадратов**

$$Q = Z_1^2 + \dots + Z_n^2 = \sum_{i=1}^n \frac{(X_i - f(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_k))^2}{\Delta x_i^2} = \sum_{i=1}^n \frac{(X_i - f(x_1, x_2, \dots, x_n; \theta_1, \dots, \theta_k))^2}{\sigma^2}$$

(Л7.2.5) будет теперь **распределена по закону χ^2_{n-k} (хи-квадрат с $n-k$ степенями свободы) с плотностью распределения $f_{\chi^2_{n-k}}(Q)$. (Л7.2.2).**

Для каких значений параметров β_1, \dots, β_k следует проверить исходную гипотезу? Если найти фиксированные значения параметров $\beta_1^0, \dots, \beta_k^0$ **методом наименьших квадратов**, то они будут минимизировать значение критерия (Л7.2.5). Если гипотеза будет отвергнута при минимальном значении критерия, то она будет отвергнута и при любом другом его значении.

Число степеней свободы определяется размером выборки минус число неизвестных параметров. Для проверки исходной гипотезы при найденных значениях параметров $\beta_1^0, \dots, \beta_k^0$ **определим значение критерия** $\chi^2(\beta_1^0, \dots, \beta_k^0)$ (Л7.2.5) и по таблице распределения хи-квадрат для числа степеней свободы $n - k$ и заданного уровня значимости α **найдем** χ_{α}^2 , соответствующее решению уравнения:

$$\chi^2(\beta_1^0, \dots, \beta_k^0) = \chi_{\alpha}^2 \quad (\text{Л7.2.6})$$

В случае, если найденное значение $\chi^2(\beta_1^0, \dots, \beta_k^0) < \chi_{\alpha}^2$, то исходная гипотеза о том, что данные описываются кривой $\chi^2(\beta_1, \dots, \beta_k)$, **отвергается** (данная модельная кривая неадекватно описывает выборочные данные на заданном уровне значимости). В случае $\chi^2(\beta_1^0, \dots, \beta_k^0) > \chi_{\alpha}^2$ гипотеза **не отвергается**, т.е. предлагаемая модель не противоречит данным. Но не исключено, что **возможно существуют другие модели**, также не противоречащие данным. Поэтому **при проверке таких гипотез однозначен только отрицательный ответ.**

Частный случай 3

Пусть имеется *выборка, состоящая из независимых реализаций случайной величины* x_1, x_2, \dots, x_n . Требуется *проверить гипотезу о том, что случайная величина x подчиняется закону с плотностью распределения $f(x)$ на уровне значимости α* . Так же, как и в первом частном случае, плотность распределения может иметь параметры, но значение их заранее известно (фиксировано). Например, $f(x)$ может соответствовать плотности распределения $N(1,0)$.

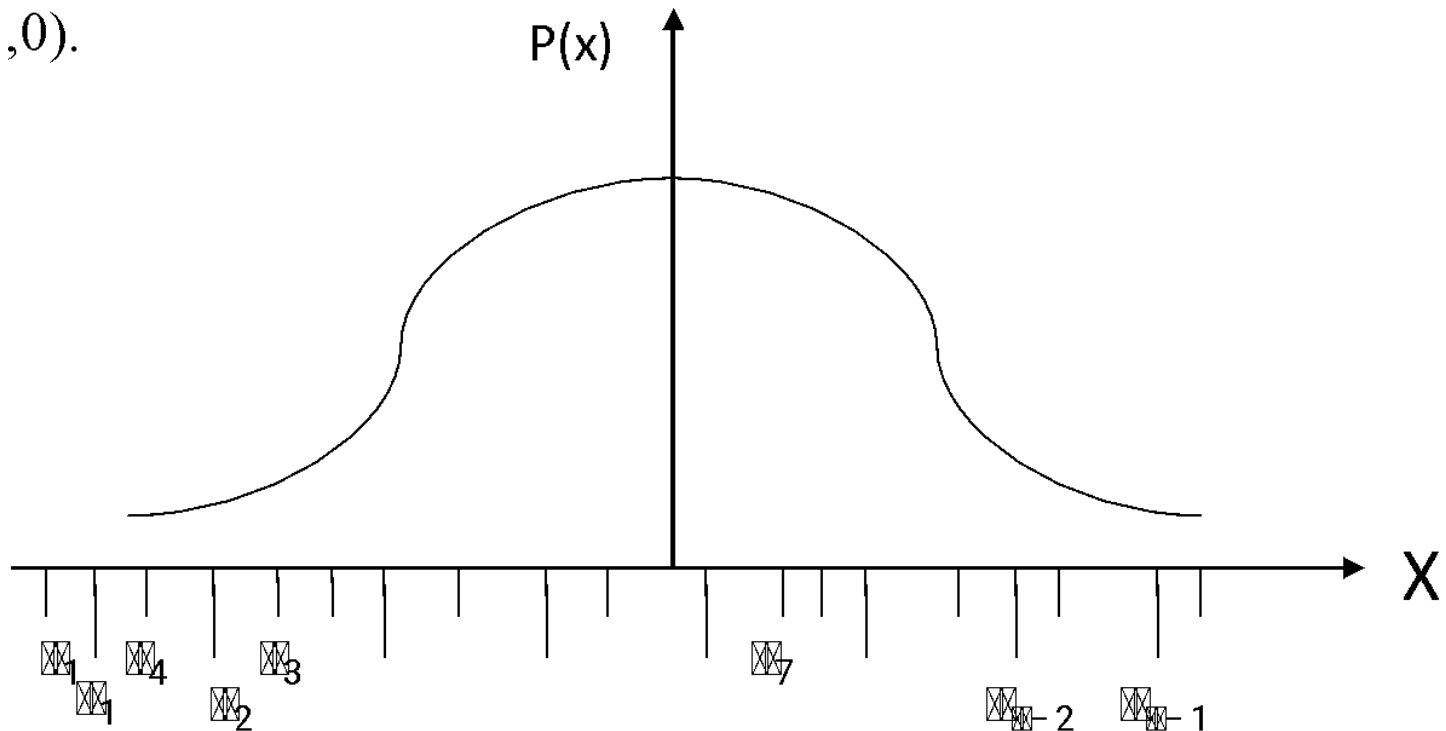


Рис. Л7.2.1

Поделим область значений случайной величины ξ (на рис. Л7.2.1) величина ξ определена от $-\infty$ до $+\infty$) на r участков:

$$\Delta_{i-1}: \xi - \infty, \Delta_i: \xi_{i-1}, \xi_i, \dots, \Delta_r: \xi_{r-1}, \xi + \infty$$

При этом каждое выборочное значение ξ_1, \dots, ξ_n попадает в один из образовавшихся интервалов. Теоретическая вероятность попадания значения случайной величины ξ в отрезок Δ_i равна

$$P_i = \int_{\Delta_i} f(\xi) d\xi \quad (Л7.2.7)$$

ξ_i – наблюдаемое количество выборочных значений, попавших в i -ый интервал. Случайные величины ξ_i подчиняются биномиальному закону распределения с параметрами n, P_i . **Математическое ожидание** $E\xi_i = nP_i$, а **дисперсия** $D\xi_i = nP_i(1 - P_i)$. При $n \rightarrow \infty$ биномиальное распределение стремится к нормальному $N(nP_i, \sqrt{nP_i(1 - P_i)})$. Вычислим **сумму квадратов нормированных величин**:

$$\chi^2 = \sum_{i=1}^r \frac{(\xi_i - nP_i)^2}{nP_i(1 - P_i)} \quad (Л7.2.8)$$

При $n \rightarrow \infty$ случайная величина Z будет распределена по закону распределения χ^2_{r-1} (хи-квадрат с $r-1$ степенями свободы) с плотностью распределения $f_{r-1}(z)$ (см Л7.2.1). Одна дополнительная связь при определении числа степеней свободы присутствует благодаря ограничению $\sum_{i=1}^r \xi_i = n$

Таким образом, *процедура проверки гипотезы о соответствии выборки заданному распределению сводится к следующим шагам:*

1. Разбиваем область значений случайной величины на r интервалов.
2. Подсчитываем значения \bar{x}_i и теоретические вероятности \bar{p}_i (Л7.2.7).
3. Подсчитываем Z по формуле (Л7.2.8).
4. По таблице распределения хи-квадрат для заданного числа степеней свободы $r - 1$ и уровня значимости α находим χ^2_{α} , как решение уравнения $\sum_{i=1}^r \frac{\bar{x}_i^2}{\bar{p}_i - 1} = 1 - \alpha$
5. Если $\chi^2_0 < \chi^2_{\alpha}$, то исходная гипотеза отвергается, при этом вероятность отвергнуть правильную гипотезу равняется α , в противном случае гипотеза не отвергается.

Способ разбиения и количество интервалов могут повлиять на результат проверки. При разных разбиениях и одних и тех же данных одна и та же гипотеза может быть принята и отвергнута на одном и том же уровне значимости зависимости от способа разбиения. **На практике при использовании критерия хи-квадрат для проверки рассматриваемой гипотезы для разбиения области значений случайной величины на отрезки пользуются следующими эмпирическими правилами при значениях $n > 30$, $n_i = n \cdot p_i$:**

1. Разбиение должно быть таким, чтобы $n_i > 5$
2. Границы отрезков выбираются таким образом, чтобы $n_i = n \cdot p_i$ для любого i .

При выполнении этих двух правил **сходимость распределения случайной величины Z (Л7.2.8) к распределению хи-квадрат самая быстрая**, то есть получение достоверных выводов возможно при меньшем размере выборки.

Частный случай 4

Пусть имеется выборка, состоящая из независимых реализаций случайной величины X_1, \dots, X_n . Требуется *проверить гипотезу о том, что случайная величина X подчиняется закону с плотностью распределения $f(x; \theta_1, \dots, \theta_k)$ на уровне значимости α при некоторых значениях k неизвестных параметров.*

Процедура проверки гипотезы о соответствии выборки распределения сведется к следующим шагам:

1. Разбиваем область значений случайной величины на r интервалов.
2. Подсчитываем значения n_j , теоретические вероятности p_j в данном случае будут функциями от неизвестных параметров $\theta_1, \dots, \theta_k$.
3. В данном случае Z также будет функцией неизвестных параметров:

$$Z = \sum_{j=1}^r \frac{(n_j - np_j)^2}{np_j} \quad (Л7.2.9)$$

5. Если проверяемая гипотеза будет отвергнута при минимальном значении критерия (Л7.2.9), то она будет отвергнута и при любом другом его значении. Тогда, набор значений параметров находим при **решении задачи минимизации**: $x_1^*, \dots, x_n^* = \arg \min_{x_1, \dots, x_n} Z$, а минимальное значение Z^* находим при найденных значениях параметров по формуле (Л7.2.9).

6. По таблицам распределения хи-квадрат для заданного числа степеней свободы $k - 1 - k$ и уровня значимости α находим χ_{α}^2 , как решение уравнения $\sum_{i=1}^k \frac{x_i^2}{\sigma_i^2} = \chi_{\alpha}^2$

7. Если $\chi_{\alpha}^2 < \chi_{\alpha}^2$, то исходная гипотеза **отвергается**, при этом вероятность ошибки α , **в противном случае** гипотеза о соответствии выборки заданному распределению (с точностью до значений параметров) **не отвергается**.

Замечания о разбиении на интервалы остаются справедливыми и для данного случая. Задача осложняется тем, что выполнение второго эмпирического правила разбиения оказывается невозможным, по крайней мере, в общем случае. Поэтому при разбиении следует руководствоваться тем **правилом**, что **при значениях** $n_j > \frac{n \cdot \alpha_j}{k} = \frac{n \cdot \alpha_j}{k} \cdot \frac{1}{\alpha_j}$ **число точек, попадающих в каждый интервал** $n_j > \frac{n \cdot \alpha_j}{k}$

Частный случай 5

Критерий хи-квадрат применяется и в том случае, когда **нужно выяснить**, являются ли значения двух величин **коррелированными**.


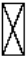

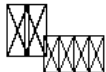
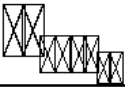
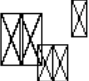







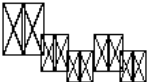
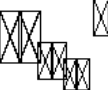

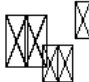


Пусть **имеются пары измерений двух случайных величин X и Y , проведенные в последовательные моменты времени:** $x_1, y_1, x_2, y_2, \dots, x_n, y_n$. Требуется **проверить гипотезу о том, что эти две случайные величины независимы**.

В основание проверки этой гипотезы положено то обстоятельство, что для **независимых случайных величин X и Y выполняется:**

$$P(X \in A, Y \in B) = P(X \in A) \cdot P(Y \in B)$$

Разобьем область значений X на k интервалов, а область Y значений на l интервалов (используя первое эмпирическое правило, упомянутое при рассмотрении третьего частного случая).

Составим следующую таблицу, содержащую в клетке n_{kj} число n_{kj} таких пар x , значения которых попадают в k -ый интервал разбиения области значений X , а значения y которых попадают в j -ый интервал разбиения области значений Y :

	<i>1</i>	...	<i>k</i>	...		
<i>1</i>				...		
...			
<i>j</i>						
...			
		...			
		...		...		<i>n</i>

В *последнем столбце* таблицы стоит количество $\square\square$ *таких пар*, значения y которых попадают в j -ый интервал разбиения области значений Y , при этом x может быть любым. Аналогично, в *последней строке* таблицы стоит количество $\square\square$ *таких пар*, значения x которых попадают в k -ый интервал разбиения области значений X , при этом y может быть любым.

Поделив каждое значение в таблице на n , получим *оценки вероятностей попадания значений пар в соответствующие интервалы*. В предположении о независимости оценки вероятности *одновременного* попадания x в k -ый интервал, а y в j -ый определяются по формуле:

$$\frac{\square\square\square}{\square\square} = \frac{\square\square}{\square\square} \times \frac{\square\square\square}{\square\square} \quad (\text{Л7.2.10})$$

Откуда, *оценочное значение количества таких пар*:

$$\square\square\square = \frac{\square\square\square \square\square\square}{\square\square} \quad (\text{Л7.2.11})$$

Поскольку величины $\chi^2_{\text{расч}}$ подчиняются биномиальному распределению, то при $n \rightarrow \infty$ случайная величина

$$\chi^2_{\text{расч}} = \sum_{i=1}^k \frac{(n_{i1} - n_{i.}n_{.1}/n)^2}{n_{i.}n_{.1}/n} + \dots + \frac{(n_{ik} - n_{i.}n_{.k}/n)^2}{n_{i.}n_{.k}/n} \quad (\text{Л7.2.12})$$

будет распределена **по закону распределения χ^2_s** (хи-квадрат с s степенями свободы) с плотностью распределения $f(x) = \frac{1}{2^{s/2} \Gamma(s/2)} x^{s/2-1} e^{-x/2}$ (см. Л7.2.2). **Число степеней свободы** в этом случае подсчитывается с учетом $n_{.1} + n_{.2} + \dots + n_{.k} - n$ дополнительных независимых связей. Таким образом

$$s = n_{.1} + n_{.2} + \dots + n_{.k} - n.$$

Если для заданного уровня значимости α и числа степеней свободы найденное по таблице хи-квадрат распределения значение $\chi^2_{\alpha, s}$, соответствующее решению уравнения, аналогичного (Л7.2.6), будет **меньше Z** , найденного по формуле (Л7.2.12), то гипотеза о независимости случайных величин X и Y **отвергается**. Вероятность отвергнуть правильную гипотезу при этом равняется величине уровня значимости α .

Пример 2

Заданные пары измерений помещены в таблице:

<i>X</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>
<i>Y</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>5</i>	<i>4</i>	<i>3</i>	<i>2</i>	<i>1</i>

Проверим гипотезу о независимости. Сами интервалы запишем в таблицу для подсчета пар, значения которых попадают в соответствующие интервалы. Оценочные значения помещены в ту же таблицу рядом с эмпирическими значениями под косой чертой:

	<i>(0,3.5)</i>	<i>(3.5,7.5)</i>	<i>(7.5,10.5)</i>	Σ
<i>(0,2.5)</i>	<i>2</i> / <i>1.2</i>	<i>0</i> / <i>1.6</i>	<i>2</i> / <i>1.2</i>	<i>4</i>
<i>(2.5,4.5)</i>	<i>1</i> / <i>1.2</i>	<i>2</i> / <i>1.6</i>	<i>1</i> / <i>1.2</i>	<i>4</i>
<i>(4.5,5.5)</i>	<i>0</i> / <i>0.6</i>	<i>2</i> / <i>0.8</i>	<i>0</i> / <i>0.6</i>	<i>2</i>
Σ	<i>3</i>	<i>4</i>	<i>3</i>	<i>1</i>

По формуле (Л7.2.12) получим $Z=5.826$. Для числа степеней $s=4$ свободы гипотеза о независимости отвергается на уровне значимости $\alpha = 0.25$ ($Z_{\alpha} = 5.39$) и не отвергается на уровне значимости $\alpha = 0.1$. ($Z_{\alpha} = 7.78$). Т.о. с 75% уверенностью можно утверждать, что X и Y зависимы, а с 90% - нельзя.

ЛЕКЦИЯ 8

ЭЛЕМЕНТЫ ДИСПЕРСИОННОГО АНАЛИЗА

Дисперсионный анализ широко используется при обработке данных в радиационной физике и биофизике. **Основное назначение метода** – выявление факта однородности или неоднородности выборки вследствие разных причин.

Рассмотрим совокупность результатов измерений (выборку размера n): x_1, x_2, \dots, x_n . Понятно, что все выборочные значения отличаются друг от друга. Возможны **две причины отличия**:

- **статистическая**, т.е. значения различаются как практическая реализация одной и той же физической случайной величины;
- **детерминистическая**, т.е. значения являются реализацией разных случайных величин, различающихся между собой.

Выборка, **состоящая из одной и той же случайной величины**, называется **однородной**. Выборка, **состоящая из разных случайных величин**, называется **неоднородной**.

Целью дисперсионного анализа является выяснение характера различий между отдельными значениями выборки, т. е. необходимость решения следующей задачи: носят ли различия случайный (стохастический) характер или детерминированы.

В некоторых случаях **детерминированные различия** носят **мешающий** характер. По существу *в этом случае имеем дело со смесью распределений с различными средними значениями*. Наличие детерминированных различий в такой ситуации затрудняет интерпретацию результатов измерений. В других случаях **детерминированные различия вызываются специально**, с целью изучения эффекта каких-то воздействий. **Дисперсионный анализ** представляет собой существенное *обобщение проверки гипотезы о равенстве средних двух выборок*.

В дальнейшем будем считать, что *все измеряемые величины подчиняются нормальному закону с плотностью вероятности*:

$$f(x, a, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{(x-a)^2}{2\sigma^2}\right] \quad (Л8.1)$$

где параметры **a** и **σ** неизвестны. Если выборка **неоднородна**, то считается, что *она состоит из подвыборок, подчиняющихся нормальным распределениям с разными средними значениями, но с одинаковым средне-квадратическим отклонением σ*

Проверка однородности выборки.

Рассмотрим следующий *пример из биофизической практики.* Для проведения исследования биолог выращивает культуры клеток в нескольких чашках Петри. Затем он берет несколько проб из каждой чашки и проводит необходимые измерения. Он предполагает, что пробы, взятые из разных чашек, одинаковы в том смысле, что соответствующие результаты измерений являются реализациями одной и той же случайной величины. На самом же деле условия выращивания культур в разных чашках могут различаться. В свою очередь это различие приводит к различию значений средних и может исказить конечный результат исследования. Поэтому, *прежде всего следует убедиться, что выборка однородна.* Сделать это можно с помощью дисперсионного анализа, который излагается далее применительно к данной задаче.

Пусть выборка объемом n состоит из r подвыборок (групп): x_{ij} , где $i=1, \dots, r$ – номер подвыборки (чашки), $j=1, \dots, n_i$ – номер измерения внутри подвыборки, $\sum_{i=1}^r n_i = n$. Величины x_{ij} подчиняются нормальному распределению (Л8.1).

Имеется подозрение, что средние значения $\bar{x}_{i\cdot}$ могут различаться между собой. Таким образом, следует проверить гипотезу о равенстве средних значений $\bar{x}_{i\cdot}$ разных подвыборок друг другу. Если эта гипотеза будет отвергнута, то это будет означать, что выборка неоднородна.

Введем следующие обозначения:

$$\bar{x}_{i\cdot} = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \quad (\text{Л8.2}) \quad - \text{выборочное среднее (оценка}$$

среднего значения) *в i -ой подвыборке;*

$$\bar{x} = \frac{1}{n} \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij} \quad (\text{Л8.3}) \quad - \text{выборочное среднее}$$

(оценка среднего значения) *по всей выборке.*

Рассмотрим *величину*
$$Q = \sigma_{\bar{x}=\bar{x}}^2 \sigma_{\bar{x}-\bar{x}}^2 - \bar{x}^2 \quad (Л8.4)$$

Если выборка *однородна*, то величина $\frac{\bar{x}}{\bar{x}-\bar{x}}$ является *оценкой величины дисперсии распределения* (Л8.1). Преобразуем величину (Л8.4) следующим образом:

$$\begin{aligned} Q &= \sigma_{\bar{x}=\bar{x}}^2 \sigma_{\bar{x}-\bar{x}}^2 - \bar{x}^2 = \sigma_{\bar{x}=\bar{x}}^2 \sigma_{\bar{x}-\bar{x}}^2 - \bar{x}^2 + \bar{x}^2 - \bar{x}^2 = \\ &= \sigma_{\bar{x}=\bar{x}}^2 \sigma_{\bar{x}-\bar{x}}^2 - \bar{x}^2 + 2\sigma_{\bar{x}=\bar{x}}^2 \bar{x} - \bar{x} \sigma_{\bar{x}=\bar{x}}^2 - \bar{x}^2 + \\ &+ \sigma_{\bar{x}=\bar{x}}^2 \sigma_{\bar{x}-\bar{x}}^2 - \bar{x}^2 \end{aligned}$$

Второе слагаемое в этой сумме равно нулю, т.к. $\sigma_{\bar{x}=\bar{x}}^2 \bar{x} - \bar{x}^2 = 0$ в соответствии с определением (Л8.2). В третьем слагаемом нет зависимости от индекса j , поэтому $\sigma_{\bar{x}=\bar{x}}^2 \bar{x} - \bar{x}^2 = \bar{x} \sigma_{\bar{x}=\bar{x}}^2 - \bar{x}^2$.

Окончательно имеем: $Q = Q_1 + Q_2$, где $Q_1 = \sigma_{\bar{x}=\bar{x}}^2 \sigma_{\bar{x}-\bar{x}}^2 - \bar{x}^2$, $Q_2 = \sigma_{\bar{x}=\bar{x}}^2 \bar{x} - \bar{x}^2$ (Л8.5)

Как видно, величина Q_1 характеризует *отклонение экспериментальных значений от подвыборочных средних*, а Q_2 - *отклонение подвыборочных средних от общего выборочного среднего*.

Выдвигаем гипотезу о том, что выборка *однородна*, т.е. все x_{ij} *подчиняются распределению* (Л8.1). В этом случае можно показать, что величины $\frac{\sum x_{ij}^2}{n}$, $\frac{\sum x_{i.}^2}{r}$, $\frac{\sum x_{.j}^2}{n}$ *независимы* и *подчиняются распределению хи-квадрат* с числом степеней свободы $(n-1)$, $(n-r)$, $(r-1)$ соответственно. Можно также показать, что величины: $\frac{\sum x_{ij}^2}{n} = \frac{\sum x_{i.}^2}{n-r}$, $\frac{\sum x_{ij}^2}{n} = \frac{\sum x_{.j}^2}{n-r}$ (Л8.6), являются *несмещенными оценками величины* σ^2 . Однако т.к. истинное значение σ^2 неизвестно, применить критерий хи-квадрат для проверки гипотезы непосредственно нельзя. Поэтому в качестве статистики для проверки гипотезы выбирается соотношение:

$$F = \frac{\frac{\sum x_{i.}^2}{n-r}}{\frac{\sum x_{.j}^2}{n-r}} = \frac{\sum x_{i.}^2}{\sum x_{.j}^2} \quad (\text{Л8.7}),$$

которое называется *дисперсионным отношением*. Если умножить числитель и знаменатель дисперсионного отношения на масштабный множитель σ^2 , то получим отношение двух величин, распределенных по закону хи-квадрат с соответствующим числом степеней свободы.

Английским математиком Р.Фишером показано, что величина:

$$F = \frac{\frac{\sum_{i=1}^m x_i^2}{m}}{\frac{\sum_{j=1}^k y_j^2}{k}}$$

(Л8.8),

где две независимые выборки: x_1, x_2, \dots, x_m и y_1, y_2, \dots, y_k являются реализациями одной и той же случайной величины, распределенной по нормальному закону с параметрами $N(0, \sigma^2)$, **подчиняются закону распределения с плотностью вероятности:**

$$f(z) = \frac{\Gamma\left(\frac{m+k}{2}\right)}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{k}{2}\right)} \frac{1}{\left(\frac{m}{2}z^2 + \frac{k}{2}\right)^{\frac{m+k}{2}}}$$

(Л8.9),

где Γ – гамма-функция. Это распределение называется **распределением Фишера с числом степеней свободы m и k** . Значения интегралов табулированы для многих значений z_0, m, k .

Если фактическое значение Z , вычисленное на основании экспериментальных данных по формулам (Л8.1), (Л8.3), (Л8.7), (Л8.8), удовлетворяет неравенству $\bar{a} \leq \bar{a}_0$, то **гипотеза о равенстве всех средних друг другу не отвергается**. В противном случае, когда $\bar{a} > \bar{a}_0$, **гипотеза отвергается, и выборка признается неоднородной**. При этом вероятность ошибки равна α .
 Заметим, что **гипотеза отвергается в двух случаях:**

- Дисперсионное отношение велико ($F \gg F_{\alpha}$)
- Дисперсионное отношение мало ($F \ll F_{\alpha}$)

Первый случай понятен; он соответствует ситуации, когда разброс значений подвыборочных средних относительно общего среднего велик по сравнению с разбросом внутри каждой подвыборки. Это как раз и свидетельствует о различии истинного среднего значения в каждой подвыборке по сравнению с другими подвыборочными средними.

Второй случай требует комментария.

Исходную выборку можно представить в виде:

$x_{11}, x_{12}, \dots, x_{1n_1}$

$x_{21}, x_{22}, \dots, x_{2n_2}$

\dots, \dots, \dots

$x_{r1}, x_{r2}, \dots, x_{rn_r}$

Здесь строки являются **подвыборками**, число элементов в каждой строке в общем случае **различно**. *Формально можно считать подвыборками не строки, а столбцы.* Тогда число элементов в каждом столбце в первых подвыборках будет равно r , а в последних оно будет меньше, т.к. различны. Может оказаться так, что средние в различных новых, соответствующих столбцам, подвыборках существенно отличаются друг от друга, в то время как построчные средние, напротив, могут быть близки по значениям. Легко показать, что *именно этой ситуации соответствует малое значение дисперсионного отношения R (Л8.7).*

На практике подобную ситуацию можно представить следующим образом.

Обратимся вновь к примеру выращивания клеток в чашках Петри. Допустим, последовательность действий биолога следующая: выращивается культура в одной чашке, берется готовая проба, проводится измерение, берется вторая проба, проводится измерение и т.д. Далее выращивается культура во второй чашке, и вновь пробы и измерения проводятся по той же схеме. Затем все то же проделывается с третьей чашкой и так далее. Допустим теперь, что культура клеток неустойчива (например, она постепенно разлагается). Тогда свойства пробы зависят от времени между моментом окончания процесса выращивания и моментом измерения. Тогда истинные значения средних проб, взятых из одной чашки в разные моменты времени, могут различаться, и наоборот, можно ожидать равенства истинных средних у всех первых проб, у всех вторых проб и т.д.

Понятно, что во втором случае, (R мало) неприятие гипотезы об однородности выборки также физически обосновано, как и в первом случае, когда R велико.

И в заключении, *если гипотеза об однородности выборки отвергается, то можно взять любые две подвыборки (например, с наиболее различающимися подвыборочными средними) и получить интервальную оценку разности средних.*