

Лекция 1: Экспериментальные данные. Характеристики выборки и генеральной совокупности

1. Классификация видов экспериментальных исследований.
2. Типы погрешностей измерений и их оценки.
3. Гистограмма и полигон частот.
4. Параметры распределения и их влияние на вид кривой распределения.

1 Классификация видов экспериментальных исследований

способы получения данных
через познание окружающего мира

наблюдение

воздействие наблюдателя на
объект минимально

эксперимент

наблюдение с воздействием
на объект

Результат эксперимента: выводы и рекомендации.
Информация может быть выражена в виде графиков,
чертежей, таблиц, статистических данных или словесных
описаний.

Эксперимент предполагает проведение опытов.

Опыт – воспроизведение исследуемого явления в определенных условиях проведения эксперимента при возможности регистрации его результатов.

По цели проведения и форме представления полученных результатов

эксперимент



качественный

устанавливает факт существования явления, не дает количественных характеристик объекта (словесное описание результатов эксперимента)

количественный

фиксирует существование явления, устанавливает соотношение между количественными характеристиками явления и внешнего воздействия на объект (количественное описание факторов)

Фактор – переменная величина, по предположению влияющая на результаты эксперимента.

Уровень фактора - фиксированное значение фактора относительно начала отсчета.

В отдельном конкретном опыте каждый фактор может принимать одно из значений уровня. Фиксированный набор уровней всех факторов в опыте определяет одно из возможных состояний объекта исследований.

Вид фактора	Регистрация уровня	Установка уровня
Контролируемый, управляемый	+	+
Контролируемый, неуправляемый	+	-
Неконтролируемый	-	-

количественный эксперимент



пассивный

активный

нет управляемых факторов

есть управляемые факторы

Отклик – наблюдаемая случайная величина, по предположению зависящая от факторов, т.е. некое исследуемое свойство объекта.

В количественном эксперименте необходимо:

- регистрировать уровни всех контролируемых факторов;
- количественно описать отклик – установить **функцию отклика** (зависимость между факторами и откликом).

Вид функции отклика в общем случае:

$$f(x_i, h_j) + \varepsilon_\delta$$

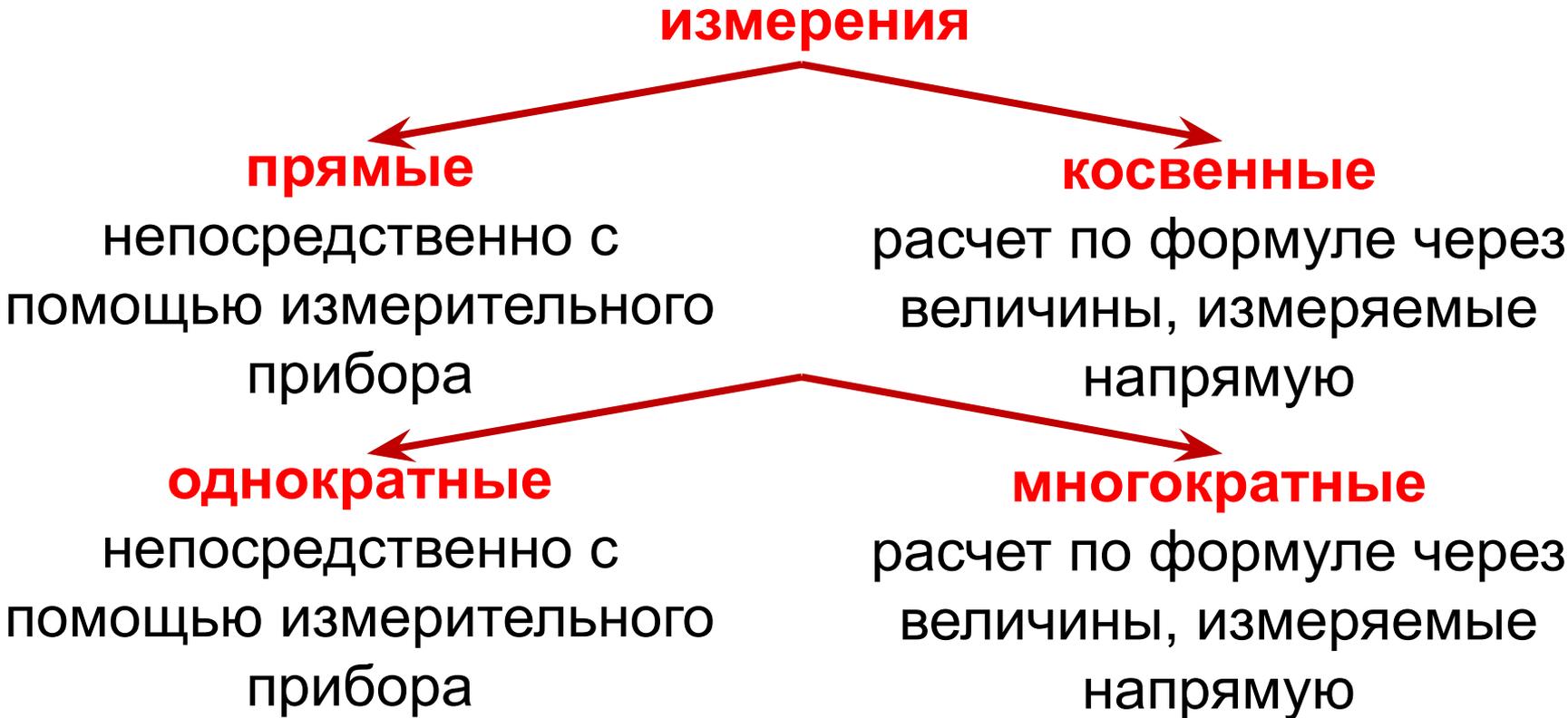
x_i – контролируемые и управляемые факторы;

h_j – контролируемые и неуправляемые факторы;

ε_δ - ошибка эксперимента (влияние неконтролируемых факторов)

2 Типы погрешностей измерений и их оценки

Предмет количественного эксперимента - количественные величины. Для определения абсолютного значения величины ее сравнивают с единицей величины - эталоном.



Погрешность – количественная характеристика неоднозначности результата измерений.

Нельзя получить экспериментальные данные с абсолютной точностью => необходимо

- оценить значение измеряемой величины
- указать, насколько оценка близка к истинному значению
- оценить качество измерений

Пусть результат многократных измерений – случайный вектор $X = \{x_1, x_2, \dots, x_N\}$, тогда оценка значения -

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$$

Абсолютная погрешность ΔX - разность между истинным значением измеряемой величины X и его оценкой.

Абсолютная погрешность

- выражается в единицах измеряемой величины X ;
- не отражает качества измерений.

Относительная погрешность – критерий качества измерений. Безразмерная величина.

$$\varepsilon = \frac{\Delta X}{\bar{X}} \quad \text{или} \quad \varepsilon = \frac{\Delta X}{\bar{X}} \cdot 100\%$$

Высокой точности измерения соответствует малое значение относительной погрешности.

погрешность

промах

результат с аномальным
числовым значением

систематическая ошибка

постоянная составляющая,
изменяется закономерно

методологическая ошибка

неправильный выбор
метода измерения.

Максимально
учитывается введением
поправок.

инструментальная погрешность

погрешность прибора
измерения.

Средняя инструментальная
погрешность:

$$\delta = \frac{X_m}{100} K \quad \text{односторонняя шкала;}$$

$$\delta = \frac{2X_m}{100} K \quad \text{двусторонняя шкала}$$

Алгоритм обработки данных прямых многократных измерений:

1 оценить истинное значение величины $\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$

2 оценить СКО и среднеквадратическую ошибку среднего

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{X})^2}{N-1}} \quad \sigma_x = \frac{\sigma}{\sqrt{N}}$$

3 вычислить среднюю инструментальную погрешность

$$\delta = \frac{X_m}{100} K \quad \text{или} \quad \delta = \frac{2X_m}{100} K$$

4 при априорно известной доверительной вероятности P найти коэффициенты Стьюдента $t(P, W-1)$ $t(P, \infty)$

Алгоритм обработки данных прямых многократных измерений:

5 вычислить абсолютное значение случайной погрешности и абсолютное значение инструментальной погрешности

$$\Delta X_{сл} = t(P, N - 1) \cdot \sigma_x \qquad \Delta X_{np} = t(P, \infty) \cdot \frac{\delta}{3}$$

6 рассчитать абсолютную погрешность измерений

$$\Delta X = \sqrt{\Delta X_{сл}^2 + \Delta X_{np}^2} = \sqrt{(t(P, N - 1) \cdot \sigma_x)^2 + \left(t(P, \infty) \cdot \frac{\delta}{3}\right)^2}$$

7 вычислить относительную погрешность

$$\varepsilon = \frac{\Delta X}{\bar{X}} \qquad \text{или} \qquad \varepsilon = \frac{\Delta X}{\bar{X}} \cdot 100\%$$

8 записать результат в виде: $X = \bar{X} \pm \Delta X$ при доверительной вероятности P и погрешности измерений ε .

Алгоритм обработки данных косвенных многократных измерений:

Пусть $Z = f(A, B, \dots)$ - измеряемая величина; величины A, B, C, \dots - измерены прямыми многократными измерениями.

1 результаты измерения величин A, B, C, \dots обработать по алгоритму обработки прямых многократных измерений

2 оценить истинное значение измеряемой величины Z

$$\bar{Z} = f(\bar{A}, \bar{B}, \bar{C}, \dots)$$

3 найти абсолютную погрешность измеряемой величины Z

$$\Delta Z = \sqrt{\left(\frac{\partial f}{\partial A}\right)_{A=\bar{A}, B=\bar{B}, \dots}^2 (\Delta A)^2 + \left(\frac{\partial f}{\partial B}\right)_{A=\bar{A}, \dots}^2 (\Delta B)^2 + \left(\frac{\partial f}{\partial C}\right)_{A=\bar{A}, \dots}^2 (\Delta C)^2 + \dots}$$

4 найти относительную погрешность $\varepsilon = \frac{\Delta Z}{\bar{Z}}$ или $\varepsilon = \frac{\Delta Z}{\bar{Z}} \cdot 100\%$

Соотношения для расчета погрешностей косвенных измерений для простейших функций

f	Δf	$\varepsilon = \Delta f / f$
$A+B$	$\sqrt{(\Delta A)^2 + (\Delta B)^2}$	$\frac{\sqrt{(\Delta A)^2 + (\Delta B)^2}}{\bar{A} + \bar{B}}$
$A-B$	$\sqrt{(\Delta A)^2 + (\Delta B)^2}$	$\frac{\sqrt{(\Delta A)^2 + (\Delta B)^2}}{\bar{A} - \bar{B}}$
$A \cdot B$	$\sqrt{\bar{B}^2 (\Delta A)^2 + \bar{A}^2 (\Delta B)^2}$	$\sqrt{\left(\frac{\Delta A}{\bar{A}}\right)^2 + \left(\frac{\Delta B}{\bar{B}}\right)^2}$
$\frac{A}{B}$	$\frac{\sqrt{\bar{B}^2 (\Delta A)^2 + \bar{A}^2 (\Delta B)^2}}{\bar{B}^2}$	$\sqrt{\left(\frac{\Delta A}{\bar{A}}\right)^2 + \left(\frac{\Delta B}{\bar{B}}\right)^2}$

Соотношения для расчета погрешностей косвенных измерений для простейших функций

f	Δf	$\varepsilon = \Delta f / f$
$\sin A$	$\cos \bar{A} \cdot \Delta A$	$\operatorname{ctg} \bar{A} \cdot \Delta A$
$\cos A$	$\sin \bar{A} \cdot \Delta A$	$\operatorname{tg} \bar{A} \cdot \Delta A$
$\operatorname{tg} A$	$\frac{1}{\cos^2 \bar{A}} \Delta A$	$\frac{2}{\sin 2\bar{A}} \Delta A$
$\ln A$	$\frac{\Delta A}{\bar{A}}$	$\frac{1}{f} \frac{\Delta A}{\bar{A}}$

Результат косвенных измерений записывается в виде:

$Z = \bar{Z} \pm \Delta Z$ при доверительной вероятности P и погрешности измерений ε

3 Гистограмма и полигон частот.

Предварительная обработка данных начинается с определения того, какими типами переменных представлены данные.

Типы переменных (признаков) представления данных:

- **непрерывные** – представлены действительными числами (например, длина или вес);
- **дискретные** – представлены целыми, как правило, положительными числами;
- **категориальные** (например, марка кабеля, тип материала, географический регион). Значения категориальных данных не могут быть положены на числовую прямую.

Построение **гистограммы** или **полигона частот** - самый простой способ наглядного представления о распределении вероятности выпадения того или иного значения случайной величины по выборке.

Пусть выборка из экспериментальных данных: $x = \{x_1, \dots, x_N\}$.

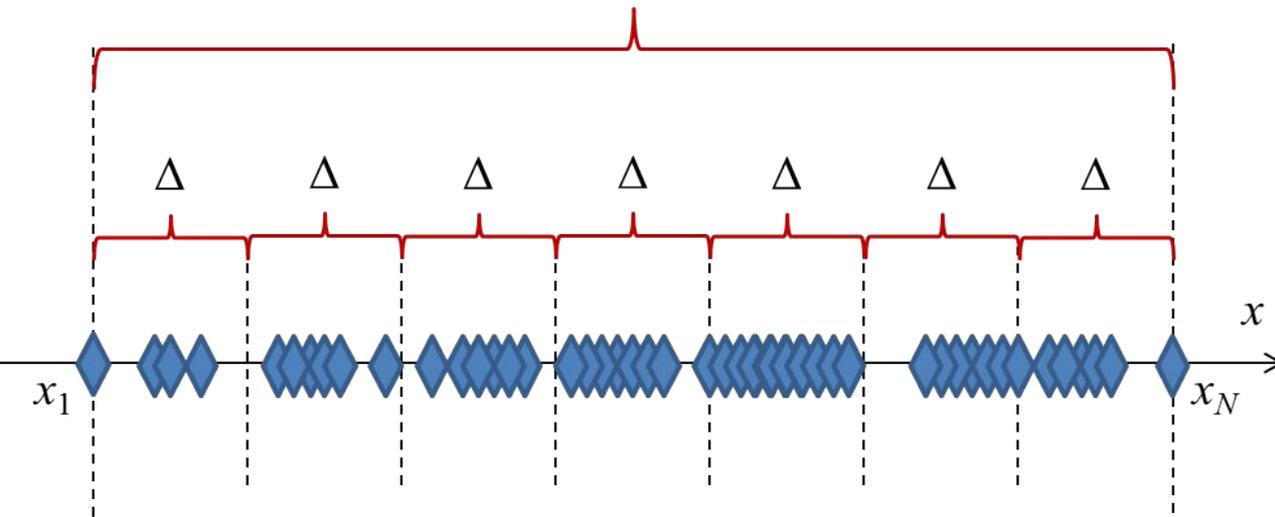
Алгоритм построения гистограммы и полигона частот

1. Построение вариационного ряда $x_1 \leq x_2 \leq \dots \leq x_N$
2. Группировка данных: разбиение отрезка $[x_1, x_N]$ на «карманы». Как и на сколько «карманов» разбивать?
Рассмотрим разбиение на «карманы» равной длины.

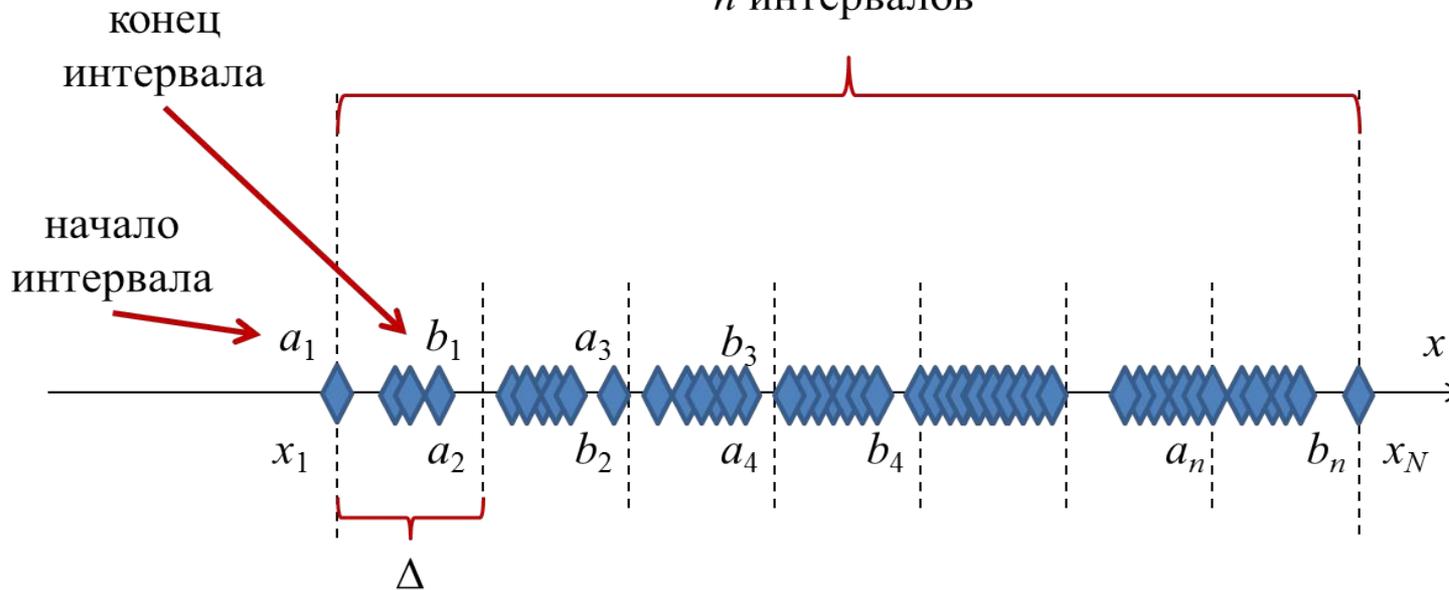
Определение числа «карманов»

- по правилу Стерджесса: $n = 1 + 3,322 \cdot \lg N$, 
- по формуле Брукса и Каррузера: $n = 5 \cdot \lg N$
- по формуле: $n = \sqrt{N}$

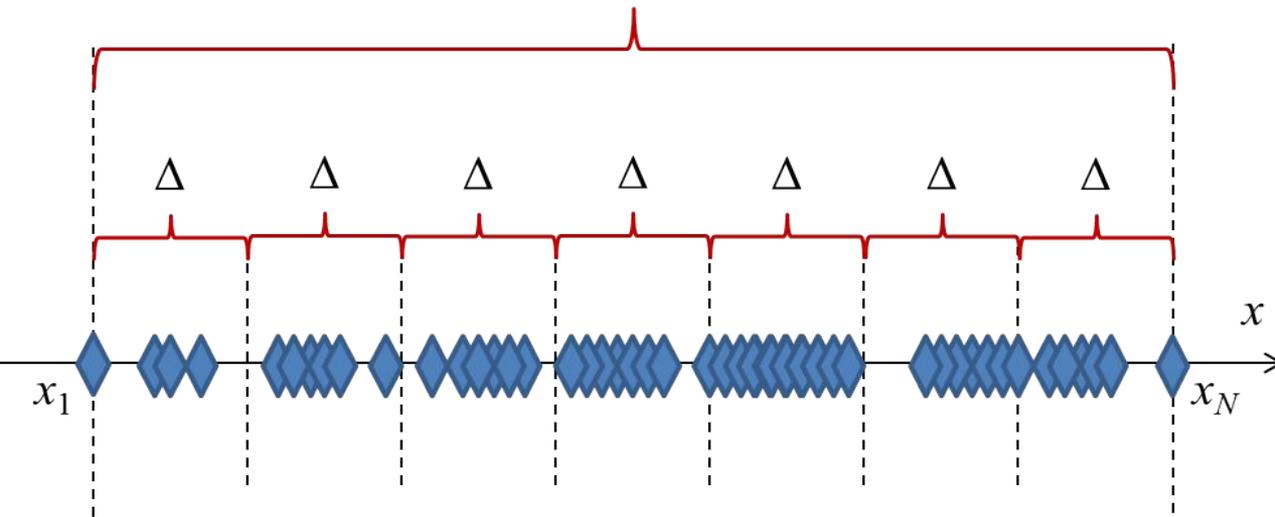
n интервалов



n интервалов



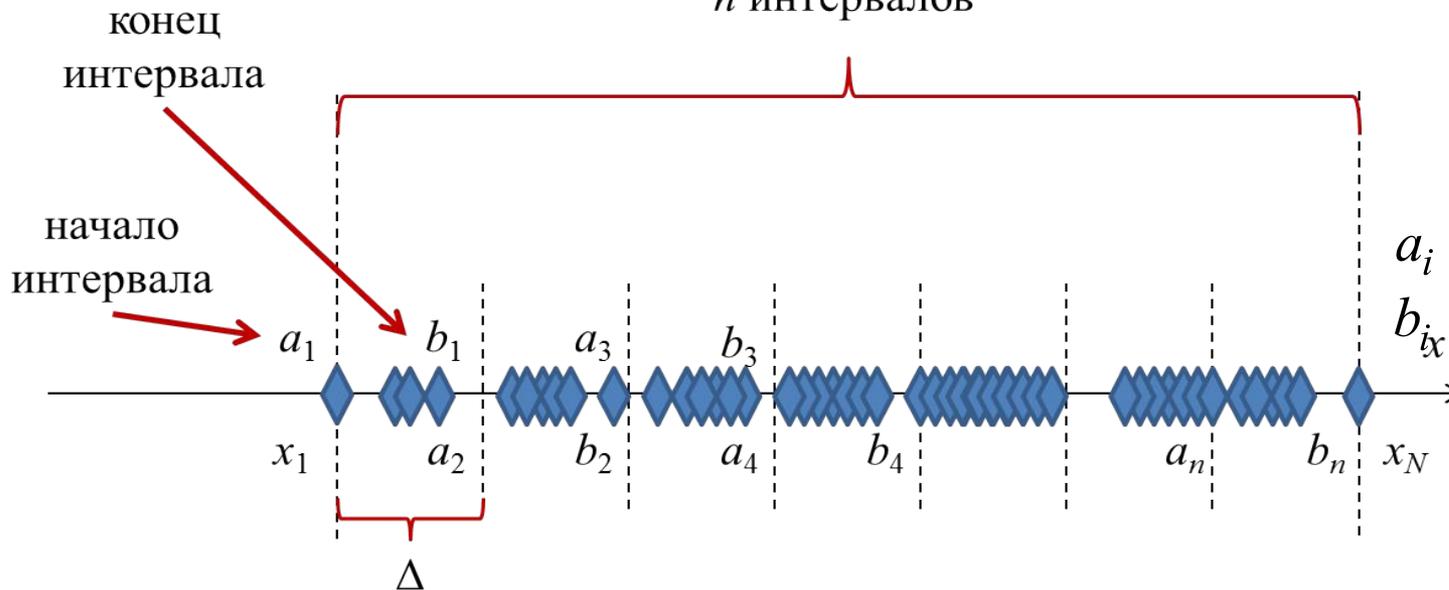
n интервалов



$$\Delta = \frac{x_N - x_1}{n},$$

$$a_1 = x_1, \quad b_n = x_N, \quad a_i = b_{i-1}, \quad \text{для } i = 2 \dots n$$

n интервалов



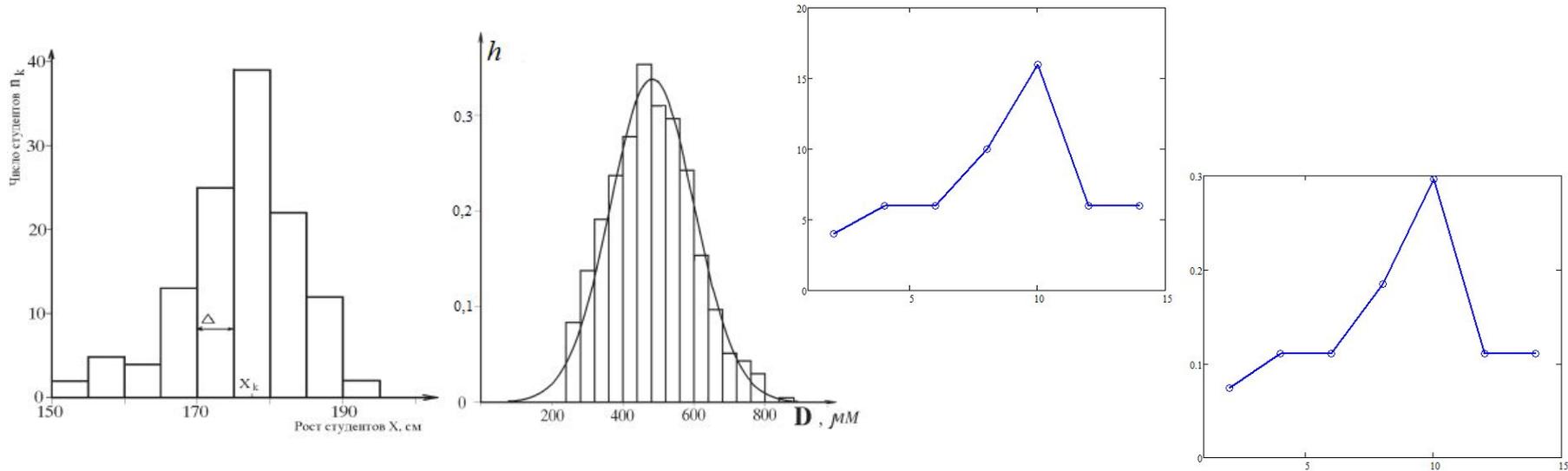
$$a_i = x_1 + (i-1)\Delta,$$
$$b_{i_x} = x_1 + i\Delta$$

3. Вычисление числа значений, попавших в каждый интервал и построение (нормированной) *гистограммы*

$$T_i = \sum_{j=1}^N t_{j,i}, \quad t_{j,i} = \begin{cases} 1, & \text{если } x_j \in [a_i, b_i], \\ 0, & \text{если } x_j \notin [a_i, b_i]. \end{cases} \quad h_i = \frac{T_i}{N \cdot \Delta} \quad \text{- нормировка } T_i$$

ИЛИ

4. Определение координат центров отрезков c_i и построение *полигона (относительных) частот* – ломанной по точкам (c_i, T_i) или (c_i, h_i)



$h_i \cdot \Delta$ - вероятность попадания результата отдельно измерения в данный интервал. Полная вероятность равна 1, значит

$$\sum_{i=1}^N h_i \Delta = 1$$

При увеличении числа измерений в пределе получаем вместо гистограммы **кривую распределения** – график **функции плотности вероятности $f(x)$** .

Следовательно,

$$\int_{-\infty}^{+\infty} f(x) dx = 1$$

Вероятность попадания измеряемой величины в интервал $(-\infty, x]$ называют **функцией распределения** или **интегральной функцией распределения**:

$$F(x) = \int_{-\infty}^x f(z) dz$$

Исходя из определения,

$$F(-\infty) = 0 \quad F(+\infty) = 1 \quad P(x_1 < x < x_2) \equiv \int_{x_1}^{x_2} f(x) dx = F(x_2) - F(x_1)$$

4 Параметры распределения и их влияние на вид кривой распределения

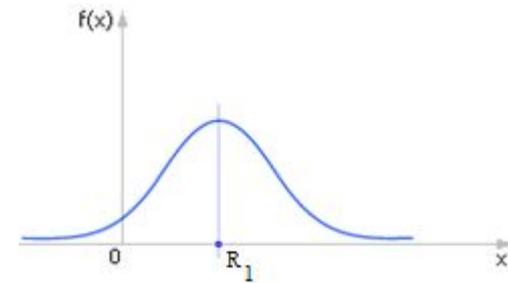
Характер кривой распределения описывается специальными мерами.

Центр распределения характеризуется *средним значением μ* , *медианой Me* и *модой Mo* .

Среднее значение (первый начальный момент) равно математическому ожиданию случайной величины:

$$\mu = R_1 = \frac{1}{N} \sum_{i=1}^N x_i = \int_{-\infty}^{+\infty} x f(x) dx$$

R_1 - центр тяжести в геометрии распределения.



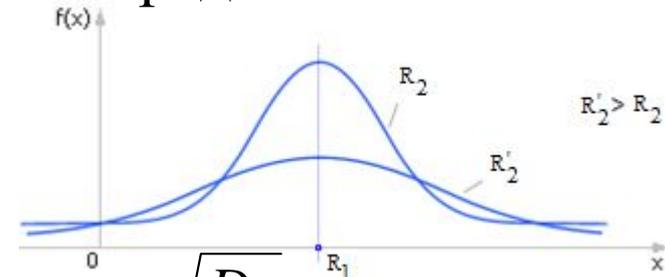
Медиана делит площадь, ограниченную функцией плотности вероятности, на две равные части $P(X \leq Me) = F(Me) = 0,5$

Мода является наиболее вероятным значением случайной величины, то есть соответствует значению x , для которого $f(x) = \max$

Рассеяние случайных величин вокруг центра группирования оценивается дисперсией, стандартным отклонением, коэффициентом вариации и размахом.

Дисперсия (второй момент) – это математическое ожидание квадрата отклонения случайной величины от их среднего арифметического значения.

$$D_x = R_2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$$



Среднее квадратическое отклонение, СКО: $\sigma = \sqrt{D_x}$

Стандартное отклонение:

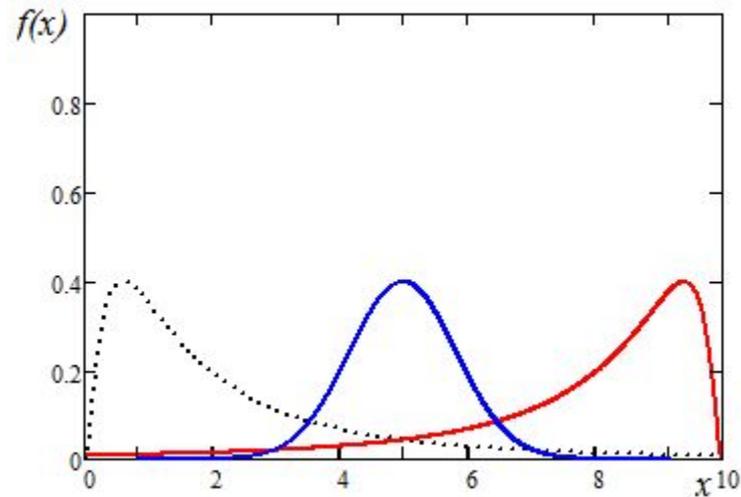
$$\sigma_{ст} = \sigma \begin{cases} 1/\sqrt{N}, & \text{для выборки} \\ 1/\sqrt{N-1}, & \text{для генеральной совокупности} \end{cases}$$

Коэффициент вариации – отношение стандартного отклонения к математическому ожиданию случайной величины.

Размах является разностью между большим и меньшим элементом выборки, то есть он равен $w = x_{\max} - x_{\min}$

Скошенность распределения, когда один хвост кривой распределения крутой, а другой - пологий, характеризует **коэффициент асимметрии**, a_3 .

$$a_3 = \frac{R_3}{\sigma_{cm}^3} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - M_x)^3}{\sigma_{cm}^3} = \frac{1}{\sigma_{cm}^3} \int_{-\infty}^{\infty} (x - M_x)^3 f(x) dx$$



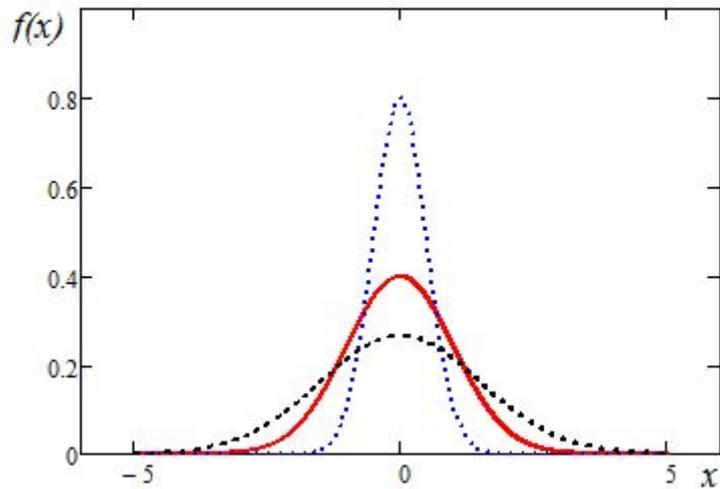
Синим – симметричное ($a_3=0$).

Черным - положительная асимметрия ($a_3 < 0$).

Красным - отрицательной асимметрия ($a_3 > 0$).

Протяженность распределения описывается **коэффициентом эксцесса (куртозиса) a_4** .

$$a_4 = \frac{R_4}{\sigma_{cm}^4} - 3 = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - M_x)^4}{\sigma_{cm}^4} - 3 = \frac{1}{\sigma_{cm}^4} \int_{-\infty}^{\infty} (x - M_x)^4 f(x) dx - 3$$



Красным – нормальное распределение ($a_3=0$)

Синим – менее протяженное распределение ($a_3<0$).

Черным – более протяженное распределение ($a_3>0$).

Квантиль - значение, которое заданная случайная величина не превышает с фиксированной вероятностью. Т.е. квантиль можно рассматривать как обратную величину функции $F(x)$.