

Основные понятия  
математической статистики:  
оценки параметров  
распределения, проверка  
гипотез, системы случайных  
величин: корреляция,  
регрессия

Лекция 17

# Способы организации выборки

**1. Вариационный ряд** – элементы выборки упорядочивают по величине:  $x^{(1)} < x^{(2)} < \dots < x^{(n)}$   $\longrightarrow$   $x^{(1)} = \min \{x_n\}$

$$x^{(n)} = \max \{x_n\}$$

**2. Размах выборки** - разность между максимальным и минимальным элементами выборки  $w = x^{(n)} - x^{(1)}$

**3. Пусть выборка содержит  $k$  различных элементов.**

**Частота элемента выборки  $n_i$**  - число раз, которые данный элемент встречается в выборке

**4. Мода** – элемент выборки с наибольшей частотой

**5. Статистический ряд** – таблица:  $\{x_i \rightarrow n_i\}$

			...	
			...	

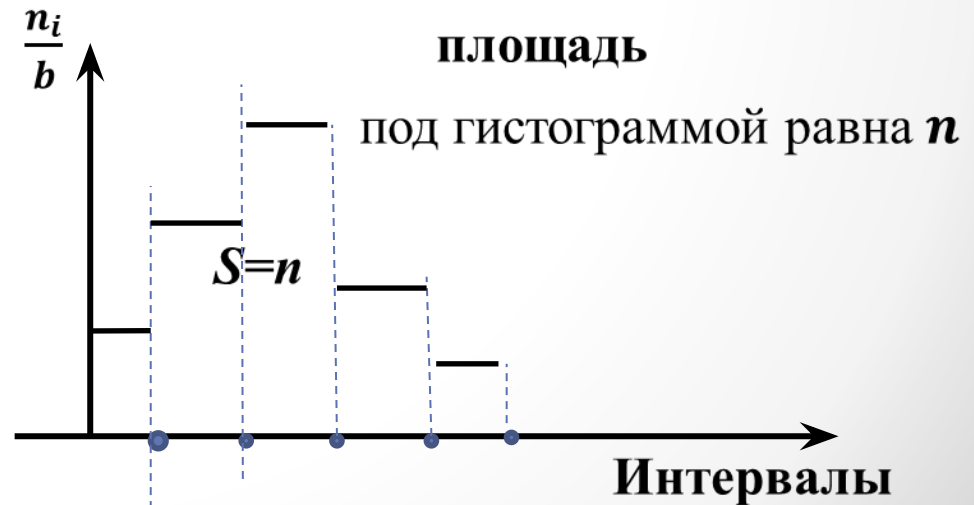
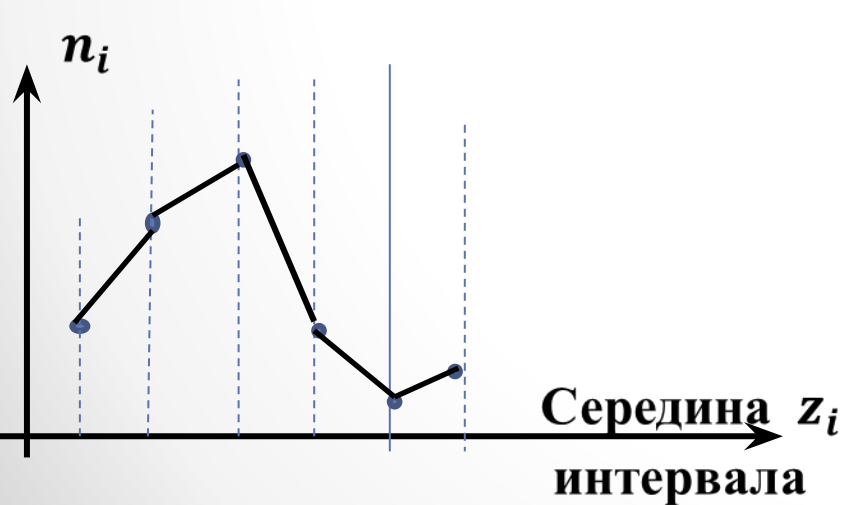
$$\sum_{i=1}^k n_i = n$$

**сумма частот всех элементов  
равна объему выборки**

# Способы описания выборки

При **большом** объеме выборки ее элементы объединяют в группы (разряды, карманы): **выбирают ширину интервала**  $b = \frac{w}{k}$ , где  $k \approx \sqrt{n}$ , или  $b = \frac{w}{1+3,2 \lg n}$ , а частота  $n_i$  - количество элементов выборки, попавшее в  $i$ -й интервал (элемент, совпадающий с внешней границей интервала считают в последующем). Кроме того вычисляю **середину каждого интервала**  $z_i = \frac{x_i + x_{i+1}}{2}$  и относительную частоту  $\frac{n_i}{n}$  - *оценку вероятности* попадания значения случайной величины в данный интервал:  $\sum_{i=1}^k \frac{n_i}{n} = 1$  (Образец табл. Стр.

181). **Графическое представление** – *полигон частот* (или относительных частот) и *гистограмма* – статистические аналоги функции распределения  $f(x)$



Полигон частот

Гистограмма

# Числовые характеристики выборки

Выборочное среднее  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$  или  $\bar{x} = \frac{\sum_{i=1}^n n_i x_i}{n}$

Выборочная дисперсия  $S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$  или  $S^2 = \frac{\sum_{i=1}^n n_i (x_i - \bar{x})^2}{n}$ .

Для выборок малого объема ( $n < 30$ ) вводят исправленную дисперсию

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

Пример. Для выборки из 5 чисел 3, 5, 5, 8, 4 получаем

$$\bar{x} = \frac{3+5+5+8+4}{5} = 5; \quad \tilde{S}^2 = \frac{(3-5)^2 + (5-5)^2 + (5-5)^2 + (8-5)^2 + (4-5)^2}{4} = 3,5$$

• Excel надстройки «Пакет анализа» Описательная статистика

• Среднее 5

• Стандартная ошибка 0,836660027  $\varepsilon_1 = \frac{s}{\sqrt{n}}$

• Медиана 5

• Мода 5

• Стандартное отклонение 1,870828693 S

• Дисперсия выборки 3,5 S<sup>2</sup>

• Эксцесс 2

• Асимметричность 1,145405322

• Интервал 5

• Минимум 3

• Максимум 8

• Сумма 25

• Счет 5

• Уровень надежности(95,0%) 2,322940635 -  $\varepsilon = \frac{s}{\sqrt{n}} t_{0,05}(n-1)$

# Статистическое оценивание. Точечные оценки

Точечной оценкой  $\tilde{\theta}$  неизвестного параметра  $\theta$  называют приближенное значение этого параметра, полученное по выборке  $\tilde{\theta} = \tilde{\theta}(x_1, x_2, \dots, x_n)$  или «статистика».

## Качество оценок.

**1. Состоятельность.** Оценка параметра сходится по вероятности к самому параметру при  $n \rightarrow \infty$   $\lim_{n \rightarrow \infty} P(|\tilde{\theta}_n - \theta| < \varepsilon) = 1$ .

Или чем больше объем выборки, тем точнее оценка

**Пример.**  $M[X] = \bar{x}$  выборочное среднее – состоятельная оценка математического ожидания (теорема Чебышева)

**2. Несмещенность.** Математическое ожидание оценки параметра равно самому параметру :  $M[\tilde{\theta}] = \theta$ . **Пример 1.** Выборочное среднее является несмещенной оценкой математического ожидания. **Пример 2.** Выборочная дисперсия  $S^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$  является смещенной оценкой ( $M[S^2] = \frac{n}{n-1} S^2$ ), а оценка  $\tilde{S}^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$  является несмещенной

**3. Эффективность.** Оценка должна обладать наименьшей дисперсией

# Интервальные оценки. Уровень значимости

Интервальные оценки или. доверительные интервалы вводятся с целью определения *точности* оценки.

Доверительным интервалом для параметра  $\theta$  называют интервал  $(\theta_1, \theta_2)$ , содержащий истинное значение параметра

с заданной вероятностью  $P = 1 - \alpha$  :

$$P(\theta_1 < \theta < \theta_2) = 1 - \alpha.$$

$(1 - \alpha)$  – доверительная вероятность;

$\alpha$  – число - вероятность, которую называют *уровнем значимости*, характеризует точность оценивания. Обычно выбирают  $\alpha = 0,1; 0,05$

**Пример.** Доверительный интервал для математического ожидания при неизвестной заранее дисперсии:

$$P(\bar{x} - \varepsilon < M[X] < \bar{x} + \varepsilon) = 1 - \alpha, \quad \text{где} \quad \varepsilon = \frac{s}{\sqrt{n}} t_{1-\frac{\alpha}{2}}(n-1)$$

• Число  $t_{1-\frac{\alpha}{2}}(n-1)$ - квантиль распределения Стьюдента находим по статистическим таблицам или в Excel (функции → статистические → Стьюдентраспобр)

**Чем больше уровень значимости, тем выше точность оценивания**

# Проверка статистических гипотез

Статистическая гипотеза  $H$  – это предположение относительно параметров или вида распределения (проверяемая гипотеза называется нулевой):

**Пример 1.**  $H_0: M[X] = t \Rightarrow H_1: M[X] \neq t$  (альтернативная гипотеза)

**Пример 2.**  $H_0$ : случайная величина распределена по нормальному закону

$H_1$  : случайная величина не распределена по нормальному закону

**Критерий** - правило, согласно которому принимается решение принять или отвергнуть нулевую гипотезу.

Перед проверкой задается малая вероятность  $\alpha$  – уровень значимости, которая определяет размер критической области  $V_k$  статистики критерия  $Z$ .

*Если выборочное значение статистики критерия попадает в критическую область  $Z_{\text{выб}} \in V_k$ , гипотеза  $H_0$  отклоняется, то есть  $\alpha$  – вероятность совершить ошибку, отвергнув правильную гипотезу.*

**О достоверности выводов, полученных при заданном уровне значимости:**

$\alpha \geq 0,1$     **высокий** уровень значимости  $\rightarrow$  данные согласуются с  $H_0$

$\alpha = 0,05$     значимость  $H_0$  возможна, но есть сомнения в истинности

$\alpha = 0,02$     имеют место сильные доводы против  $H_0$

$\alpha \leq 0,01$     основная гипотеза  $H_0$  наверняка ложная

# Выборочный коэффициент корреляции. Оценка

Для системы случайных величин  $\{X, Y\}$  вводится характеристика – **ковариация (корреляционный момент)**:

$$\text{cov}[X, Y] = K_{XY} = M[(X - M[X])(Y - M[Y])] = M[XY] - M[X]M[Y]$$

Для независимых  $X, Y$  ковариация  $\text{cov}[X, Y] = 0$ .

**Коэффициент корреляции**  $r_{XY} = \frac{K_{XY}}{\sigma_X \sigma_Y}$  – безразмерный коэффициент, который определяет степень линейной корреляционной зависимости между случайными величинами.

**Свойства  $r_{XY}$** : 1) если  $Y = AX + B$ , то  $|r_{XY}| = 1$

2)  $|r_{XY}| \leq 1$  3) если  $r_{XY} = 0$ , то случайные величины  $X, Y$

называют **некоррелированными**. Независимые случайные величины являются некоррелированными.

**Оценкой коэффициента корреляции** является выборочный

коэффициент корреляции  $\widetilde{r}_{XY} = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2 \sum_{k=1}^n (y_k - \bar{y})^2}}$

**Excel функции статистические корреляция**



# Регрессионные модели

Пусть коэффициент корреляции между двумя случайными величинами **значимо** отличается от нуля и близок к единице.

Выдвигаем гипотезу: случайные величины связаны линейной корреляционной зависимостью  $Y = AX + B + \varepsilon$

Это уравнение называют **уравнением линейной регрессии**.

**Регрессия** – оптимальная зависимость, которая обеспечивает аппроксимацию опытных данных с наибольшей точностью, то есть с минимальной случайной ошибкой  $\varepsilon$ . Наилучшие оценки для коэффициентов регрессии  $A$ ,  $B$  получают по **методу наименьших квадратов** :

$$S(\tilde{A}, \tilde{B}) = \sum_{k=1}^n (y_k - (\tilde{A}x_k + \tilde{B}))^2 - \min$$

Excel  $\longrightarrow$  точечная диаграмма  $\longrightarrow$  линия тренда с указанием уравнения и

качества аппроксимации - коэффициент детерминации  $R^2 = \frac{\sum_{k=1}^n ((\tilde{A}x_k + \tilde{B}) - \bar{y})^2}{\sum_{k=1}^n (y_k - \bar{y})^2}$

(*правой кнопкой на точку*). На заключительной стадии обязательно проверяют **статистическую значимость** (можно доказать, что в доверительный интервал для коэффициента  $A$  не содержит  $A = 0$ ) и **адекватность модели** (случайные ошибки наблюдений – **остатки** распределены с нулевым средним  $\bar{\varepsilon} = 0$ ).

Excel «Анализ данных» Регрессия