

## Лекция 1 Баченов Мат. стат.

### Выборки (экспериментальные данные)

Опр. Генеральной совокупностью наз-ся мн-во всех исходов эксперимента

Опр. Выборочной сов-к-ой наз-ся объём эксперим-х данных из генеральной сов-ти которую мы имеем возможность наблюдать

У военных была выборочная сов-та самолётов вернувшихся из БД.

А предположение что probability равномерно распределено. Такая схема получила название схемы воливилера.

Опр. Выборки наз-ся репрезентативной если её расп-е совпадает с теоретическим. (?)

Все выборки считаем репрезентатив-и

В каждой области существуют свои способы борьбы с перепреванием.

Опр. объём выборки и наз-е кол-во экспер-х данных <sup>одна сущ. величина</sup> <sub>и набор экспериментов</sub>

Опр. А: выборкой объёма  $n$  наз-е набор экспер-х данных  $\bar{X} = (x_1, \dots, x_n)$  (апостериорное определение выборки, то что после опыта)

Опр. В: (из опыта) выборкой объёма  $n$  наз-е набор  $\bar{X} = (X_1, \dots, X_n)$  независимых, одинаково распределённых случайных величин.

Эксперт-б один в котором  $n$  незав-х случайных величин

## Взвешенные характеристики

Имеется выборка  $\vec{X} = (x_1, \dots, x_n)$   
в списке А.

Мы можем рассматривать  $\vec{x}$  как  
дискретную случайную величину

$X$	$x_1$	$\dots$	$x_n$
$P$	$1/n$		$1/n$

Математическое ожидание  $E\vec{x}$   
дискретной случайной величины называется  
среднее взвешенное

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

Это то же приближенное значение  
математического ожидания.

Точно также мы можем вычислить  
дисперсию

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

- взвешенная  
дисперсия  
(оценка неизвестной  
дисперсии)

Взвешенная функция распределения

$$F^*(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x) = \frac{\text{число } X_i \leq x}{n}$$

Теорема: Гливенко - Кантелли

Пусть  $F^*(x)$  - взвеш. ф-я распр.

$F(x)$  - теор. ф-я распр-я.

$n$  - объём выборки

$$\sup_x |F^*(x) - F(x)| \xrightarrow[n \rightarrow \infty]{P} 0$$

← сходимость по вероятности

Правила обработки стат. данных

1) Ранжирование выборки (упорядочиваем данные по возрастанию), в результате получаем вариационный ряд.

$$(X_{(1)}, X_{(2)}, \dots, X_{(n)}) \quad X_{(i)} \leq X_{(j)} \quad i \leq j$$

$X_{(i)}$  -  $i$ -я порядковая статистика

$X_{(1)} = X_{\min}$  минимальная порядковая стат-ка

$X_{(n)} = X_{\max}$

Note: Если в выборке много  
повторяющихся значений, то бывает  
удобно представлять в виде  
вариационного ряда

$X$	$X_1$	...	$X_n$
$n$	$n_1$		$n_n$

где  $n_i$  это  
число или  
частота  
значения  $X_i$ .  
Словом раз оно  
встретилось  
(удобно для  
дискретной  
величины)

Note  
в этом случае  
выборочно дисперсию  
и среднее удобно  
считать по формулам

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i \cdot n_i$$

$$\text{Дисперсия} = \text{Var} = \sigma^2$$

$$D^* = \frac{1}{n} \sum_{i=1}^n (X_i \cdot n_i - \bar{x}^2)$$

$$V_i = \frac{n_i}{n} - \text{относительная частота}$$

и можно рассматривать как теоретическую  
вероятности этих значений

$X$	$X_1$	...	вариационный ряд с относительными
$V$	$V_1$		частотами

2) Иногда бывает полезно отбросить некоторое количество самых малых и самых больших значений выборки (5, 5%). Причины

- 1) уменьшить репрезентативность
- 2) исключить выбросы (повыш. дисперсия)

Применять ко ситуации. Если распределение или нормальное, то не следует отбрасывать.

3) Часто если объем данных большой и распределение равномерное вариационный ряд заменяют интервальным рядом. Данные разбиваются на  $k$  интервалов  $[a_i, a_{i+1})$  - интервал

$n_i$  - частота данного интервала

Note:

число интервалов выбирается по формуле Стержеса

$$k \approx 1 + \log_2 n$$

$k \times \sqrt{n}$

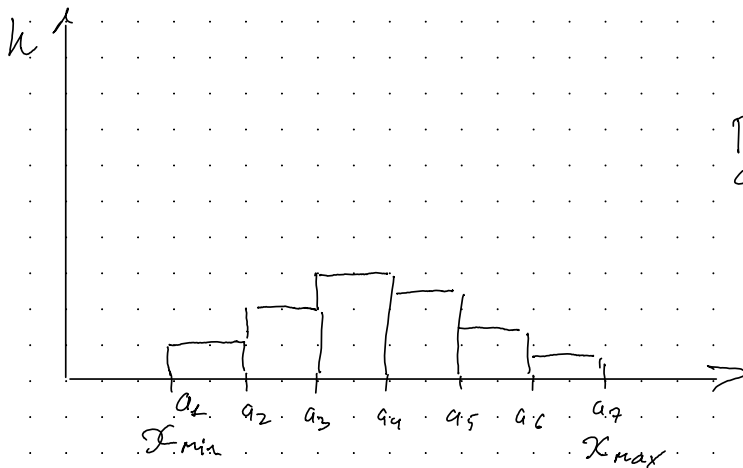
а) Можно разбивать на интервалов  
отрезки  $[x_{\min}, x_{\max}]$  или сегмент  
 $[x_{\min} - \frac{l}{2}, x_{\max} + \frac{l}{2}]$  где  $l$  длина  
интервала

б) Интервал берется либо одинаковой  
длины  $l \approx (x_{\max} - x_{\min})/k$  или  
равноканонический, в каждом интервале  
приблизительно, равное число  
элементов.

Число интервалов может зависеть  
от размаха интервала

## Гистограмма

Состоит совокупная фигура из  
прямоугольников основание которых  
длины интервалов групп  $h$  а  
высота прямоугольников  $h_i = \frac{n_i}{n \cdot L}$



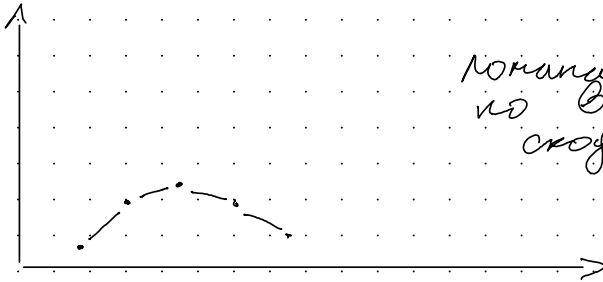
При  $n \rightarrow \infty$  гистограмма  
сподитсе по вероятности к  
плотности распредел-



## Поллигон

На числовой прямой отмечены точки  
уже  $c_i$   $k_i$   $c_i = \frac{q_i + q_{i+1}}{2}$  - середина  
интервала

$$h_i = \frac{k_i}{n}$$



показание координаты  
по вершине будет  
сходиться к  
плотности распределения  
вершин

Число с помощью компьютеров  
изображается в виде функции  
распределения.

можно специализироваться в  
экспертные системы и анализ данных

Нужно разобраться, в курсе что  
будет

$$\frac{6+3+9+7+5+9}{6} = \frac{10+10+14}{6} = \frac{34}{6} = \frac{17}{3}$$

$$\approx 5,67$$

$$\frac{1}{5} (0,333^2 + 2,667 + 1,667 + 1,333 + 0,667 + 3,333) \approx 4,667$$