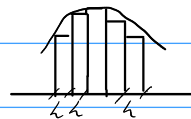


Хистограм - прва процена гушине



ида је да функцију гушине апроксимирамо
'степенастим' функцијама

↳ 'део по део' кантилативне функције

$$P\{x \leq X \leq x+h\} \approx f(x) \cdot h \Rightarrow f(x) \approx \frac{P\{x \leq X \leq x+h\}}{h}$$

↓
мало

Ако смо извукли ПСУ (x_1, x_2, \dots, x_n) , $P\{x \leq X \leq x+h\}$ можемо
оценити као $\frac{\text{број параметара које су унутар } [x, x+h]}{n}$.

Гушину сваке поједине $t \in [x, x+h]$ ћемо оценити са

$$\hat{f}(t) = \frac{\text{број параметара које су унутар } [x, x+h]}{n \cdot h}$$

Проблем: како изабрати h , интервале на којима оцењујемо f ,
колико тих интервала да изберемо?

Пошто је препорука за изборова параметара, важно је напоми-
нути да се облик хистограма може значајно променити ако само мало
променимо нпр. h , ако не усмањемо $[]$, него $[)$ или $(]$ итд.

Дакле, хистограм треба схватити као првоначну процену гушине
радијене из које су датим нашим подацима.

Постоји прецизнија оцена гушине, а то је оцена језира.

За познавање: Silverman: Density estimation for statistics
and data analysis

$$h \cdot k \approx R$$

Поступак за цртање хистограма

1. изабрати број интервала k
 често се користи Sturges-ово правило: $k = \lceil \log_2 n \rceil + 1$
 $\lceil \cdot \rceil$ - заокруживање на већи цео број
 2. одређити узорачки распон $R = x_{(n)} - x_{(1)}$
 3. h треба да је такво да $h \cdot k \approx R$
 4. одређујемо интервале над којима цртамо хистограм:
 (обично су то конјугирани интервали)
 $(x_0, x_0 + h], (x_0 + h, x_0 + 2h], \dots, (x_0 + (k-1)h, x_0 + kh],$
 x_0 је мао мање од $x_{(1)}$
- Sturges-ово правило је настало из претпоставке да подаци долазе из нормалне расподеле, а самим широкимзакључавањем нормалне расподеле симолном.
 - Није теоријски поткрепљено, али у пракси даје прихватљиве резултате
 - За разгознаре: Scott's rule, Freedman and Diaconis's rule
- правила за одређивање h
 која имају теоријску основу

Нешто лакше за цртање је хистограм фреквенција:

у хистограму фреквенција над сваким интервалом $(x_0 + ih, x_0 + (i+1)h]$

повучемо цртну линију која је једнаке броју појава из узорка које су у њему наше.

Boxplot - vizuelni prikaz raspodele podataka po kvantilima

p -ти квантил расподеле: $t \in \mathbb{R}$ ш.г. $F(t) = P\{X \leq t\} = p$, $0 < p < 1$

Примери: медијана расподеле је 0,5-ти или 50%-квантил
1. квантил := 25%-квантил
3. квантил := 75%-квантил

Статистике којима одањујемо квантилне расподеле се зову
узорачки квантили

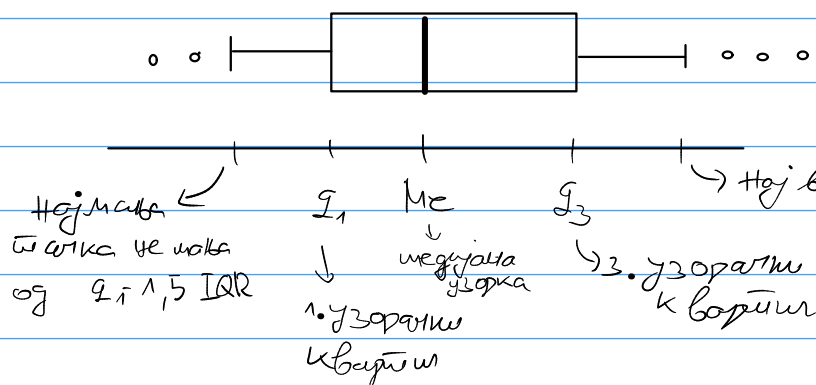
Постоји 9 дефиниција узорачких квантила, све су јакo сличне

Интуйтивно, 1. узорачки квантил је број од којег је 25% података из узорка мање, а 3. узорачки квантил је број од којег је 75% података из узорка мање, односно од којег је 25% података из узорка веће.

Како цртамо boxplot?

$$IQR = z_3 - z_1$$

↓ интерквartilно
расподеље



Патке ван $z_1 - 1,5 IQR$ и $z_3 + 1,5 IQR$ зовемо аутопајери

Интуйтивно, аутопајери су патке које одступају од очекиваних података, а ово је само један начин на који се могу дефинисати.