

25.01.20 / Консультация / R /

I Описательная статистика

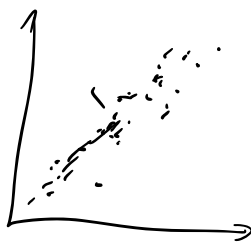
II Тестовые

III Распределения

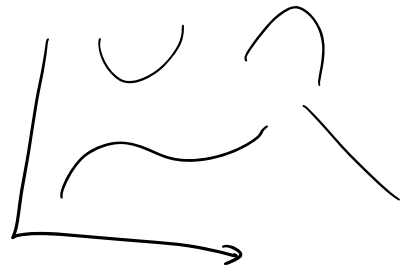
IV Линейная регрессия \rightarrow ако го направиме, минаваме

IV Лин. регрессия

- ще имаме графика



$$y = a + bx + c$$



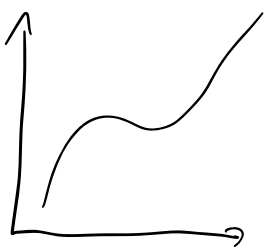
свързано с лин. регрессия

корелационен акалу
 $\text{cor} \in [-1, 1] \rightarrow$ pearson
spearman
kendall
ефективност
намалява, \uparrow стабилност
се увеличава

$|r| \in [0, 1]$

колко
коррелация
за определяне по

графика, ф-ята y
 $\rightarrow y = ax^3 + bx^2 + cx + d$



$\text{cor}(X, Y, \text{method} = 1)$

~~метод~~ $X, Y, \text{method} \leftrightarrow Y$

с корелацията взимаме връзката и въздураме най-силната променлива
таблица: Est., SE, z-value, P-value

Линейна регрессия
зависима

summary(model)

$\text{mod} = \text{lm}(y \sim x)$ \rightarrow връзка a и b
зависен от

\hookrightarrow определяме дали a и b
са значими, при $|t| > 2$
 $p\text{-value} < 0.05$

от Est. взимаме a и b

с най-голямо $|t\text{-value}| \rightarrow$ най-значим

Да определим дали модел е добър или не изва от summary
добр е когато Adj. $R^2 \rightarrow 1$

Може и с Multiple R^2

Трябва да проверим дали излизат условията

- хомоскедастичност на остатъците

погрешно
↑
остатъци с грешки

$$\begin{aligned} \Gamma &= \text{residuals}(\text{model}) \\ \downarrow \\ \Gamma &= y_i - \hat{y} \\ \downarrow \quad \quad \quad \downarrow \\ \text{реална} \quad \quad \quad &\text{прогнозирана} \\ &(\text{predict}(\text{model})) \end{aligned}$$

Условия за модел да е правилен:

1. Homoscedasticity → през графиката Γ Γ  → по тази графиката трябва да се

→ през тест $\text{bptest}(\text{model})$

→ връща p-value > 0.05

кое да изберем

→ имаме const вариация

ако я има графиката - $\text{plot}(\hat{y}, \Gamma)$

2. Да няма автокорелация

$\text{lmtest}::\text{dwtest}(\text{model})$

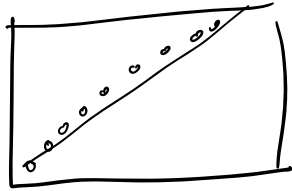
Хипотеза: липсва автокорелация
 H_0

p-value > 0.05

3. $f \sim N(0, \sigma)$

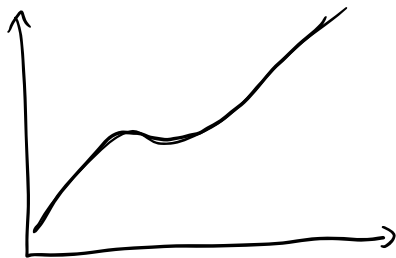
a) test - shapiro-test(f)

b) $ggnorm(f)$; $ggline(f)$



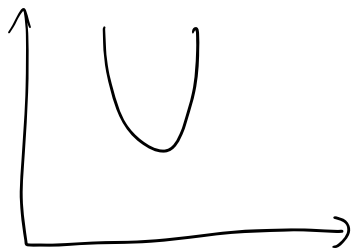
→ когато опашките са дълго го линията

По-сложни връзки



$$y = ax^3 + \dots + nx + b$$

$$y \sim I(x^3) + I(x^2) + \cancel{Ix}$$



$$y \sim I(x^2) + x$$

Тестове

1, 2 променливи - t test, wilcoxon

> 3

- ANOVA, ~~Kruskal~~ ^{Kruskal}

$$Y = \begin{bmatrix} \end{bmatrix}$$

$$X = \begin{bmatrix} \end{bmatrix}$$

$$Z = \begin{bmatrix} \end{bmatrix}$$

1. Дали са $N(\quad)$ → shapiro-test($\overset{Z}{Y}$)

ако всички са $p\text{-value} > 0.05$
⇒ ANOVA

H_0 - показвателите са равни/различни

ANOVA: $\text{aov}(\text{value} \sim \text{group}, \text{DF}) \Rightarrow \text{p_value}$ от него
називане
дали са различни
или не
(зависи от H_0)

↳ първо трябва $\text{DF} = \text{data.frame}(\text{group} = 1, 2, 3), \text{value} = c(x, y, z))$
 > 0.05 са еднакви

ако има разлики от Shapiro \Rightarrow използваме Kruskal
 т.е. поне едно $\text{p_value} < 0.05$

III Разпределения

1. Дискретни разпределения

B_i, P_0, Ge
 Ber

Ber - използваме за "нещо дали ще стане или не"
 - параметър p

$\text{rbinom}(n, p, size=1, p1)$

B_i - n -ка броя независими опита

Примерна задача

$n=100$ (100 забвучени цети)

$p=0.1$ (плащане на 10% от сумата)

\Rightarrow задачата е p \Rightarrow тогава е d

$P(k > 3) = ?$

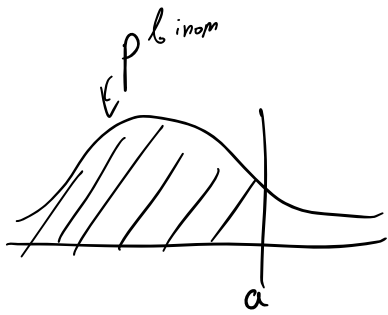
$\text{rbinom}(3, 100, 0.1) \Rightarrow P(X \leq 3)$

\Rightarrow трябва $1 - \text{rbinom}(3, 100, 0.1)$
 или

$\text{rbinom}(3, 100, 0.1, \text{low.tail} = \text{TRUE})$

Ако имам $P(k \geq 2, k \leq 4) \rightarrow \Delta \text{binom с } 2 \text{ и } 3$

$$\rightarrow P(k \leq 4) - P(k \leq 2)$$



$\Delta \text{binom} \rightarrow$ знаем вероятността, търсим a

$\Gamma \text{binom} \rightarrow$ симуляция

$\Delta \text{binom} \rightarrow$ в точки

- Теоретично - го първи член

- Поасонов - времеви интервал
 \rightarrow колко събития във времеви интервал

2. Непрерывни

- равномерно - хвърляне на зар
 гандом генератор
 тотото

- експоненциално - време или две събития

- нормално разпределение - природата

~~изразява~~

Описателна статистика

← непрекъснати

	параметрична	непараметрична
показва	mean	median
дисперсия	sd var	IQR → 3-ти quartile - 1-ви quartile mad

за показване и визуализация
hist()
geom

дискретно → mode
категорични с matplotlib
→ barplot()

как се правят ф-ции

$$x^2 + x^4 + x^6 + \dots + x^{16}$$

```
foo ← function(x) {
  s ← 0
  for(i in 1:8) {
    s = s + x^(2*i)
  }
  return(s) // връща s
}
```

X = []

```
mean = function(x) {
  return (sum(x)/len(x));
}
```

дисперсия, стандартно отклонение