

# Metody pro ověření normality

- obecné metody pro ověření shody teoretického a empirického rozdělení (Dyhwajj' data mají představit; našemu modelu?)

## 1.) Grafické metody

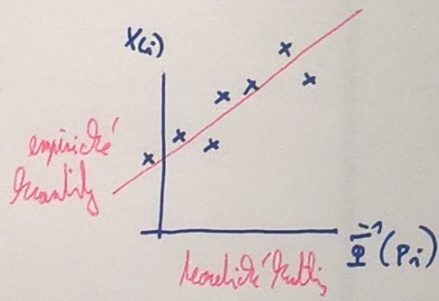
- a) histogram
  - b) jádrový odhad hustoty
  - c) boxplot
  - d) Q-Q plot = porovnání teoretické a empirické kvantily
- } do jednoho grafu + graf teoretické hustoty, odhadnutými parametry normálního rozdělení

idea: uspořádaná pozorování  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m)}$  (přes  $x_{(i)}$ ) je  $P_i = \frac{i - \beta}{m + 1 - 2\beta}$  ( $0 \leq \beta < 1$ ) - kvantil.

pro  $\beta = 0,5 \dots P_i = \frac{i - 0,5}{m}$

$\beta = 0,3175$  NP plot (normal-probability plot)

Q-Q plot je graf  $[\Phi^{-1}(P_i), x_{(i)}]$  pro  $i=1, \dots, m$ .



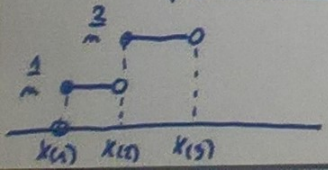
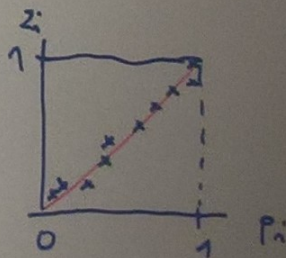
e) P-P plot (percent-percent, probability-probability) = porovnání empirické a teoretické distribuční funkce

$$\hat{F}_n(x) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{X_i \leq x\}$$

$$P_i = \frac{i}{m+1} \quad i=1, \dots, m$$

$$Z_i = \Phi\left(\frac{x_{(i)} - \hat{\mu}}{\hat{\sigma}}\right) \quad \hat{\mu} = \bar{X}, \hat{\sigma} = S$$

P-P plot je graf  $[P_i, Z_i]$  pro  $i=1, \dots, m$ .



## 2) Statistické testy

$X_1, \dots, X_n$  je máh. vzor z rozdělení s dích. fci F

$H_0$ : F je distribuční funkce normálního rozdělení ( $\mu, \sigma^2$  jsou rušivé parametry)

$H_1$ : F není . . . . .

a) testy založené na momentech

$$g_1 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^3}{\left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right\}^{3/2}} \quad \dots \text{rychlová mířička}$$

$$g_2 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^4}{\left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right\}^2} \quad \dots \text{rychlová špicálová}$$

### Teorem

Necht  $X_1, \dots, X_n$  je máh. vzor z normálního rozdělení. Pak veličiny  $g_1$  a  $g_2$  jsou asymptoticky normální a navíc asymptoticky nezávislé. Dále platí

$$Eg_1 = 0, \quad Eg_2 = \frac{3(n-1)}{n+1}, \quad D(g_1) = \frac{6(n-2)}{(n+1)(n+3)} \quad \text{a} \quad D(g_2) = \frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}$$

### Pozn.

$$\frac{g_1 - Eg_1}{\sqrt{D(g_1)}} \stackrel{H_0}{\approx} N(0,1) \quad \text{a} \quad \frac{g_2 - Eg_2}{\sqrt{D(g_2)}} \stackrel{H_0}{\approx} N(0,1)$$

testová statistika založená na mířičce

testová statistika založená na špicálovosti

Aplikace rychlové testy pro  $n$  velkých ( $n \geq 200$ , resp.  $n \geq 500$ ).

Jarqueho - Berneš test - rychlová mířička i špicálovost zároveň

$$JB = \frac{n}{6} \left( g_1^2 + \frac{(g_2 - 3)^2}{4} \right) \stackrel{H_0}{\approx} \chi_2^2$$

Použitelní pro  $n \geq 200$ .

b) regresní testy

Shapiro-Wilkův test

- je založen na porovnání 2 odhadů  $\sigma^2$  - vyšetřovací rozptylu a nejlepšího odhadu  $\sigma^2$  získaného metodou nejmenších čtverců

$$W = \frac{\left( \sum_{i=1}^m a_i X_{(i)} \right)^2}{\sum_{i=1}^m (X_i - \bar{X})^2}, \text{ kde } (a_1, \dots, a_m)' = \frac{m^T V^{-1}}{(m^T V^{-1} V^{-1} m)^{\frac{1}{2}}}, \text{ kde } m = (m_1, \dots, m_m)^T$$

$$V = (v_{ij})_{i,j=1}^m, \quad m_i = EY_{(i)}, \text{ kde } Y_1, \dots, Y_m \text{ je máh. vektor } \sim N(0, 1)$$

$$v_{ij} = C(Y_{(i)}, Y_{(j)}) \quad Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(m)}$$

je us. máh. vektor

$X_1, \dots, X_m$  máh. vektor  $N(\mu, \sigma^2)$

$$X_{(i)} = \mu + \sigma \cdot Y_{(i)} \quad \text{pro } i = 1, \dots, m$$

$(\sigma m_i - \sigma m_i)$

$$E X_{(i)} = \mu + \sigma m_i$$

$$X_{(i)} = \mu + \sigma m_i + \varepsilon_i, \text{ kde } \varepsilon_i = \sigma (Y_{(i)} - m_i) \text{ jsou chyby modelu, varianční matice } \sigma^2 V.$$

$$X_{(i)} = \begin{pmatrix} X_{(1)} \\ \vdots \\ X_{(m)} \end{pmatrix}$$

$$\text{odhadneme } \sigma \text{ metodou nejmenších čtverců } \hat{\sigma} = \frac{m^T V^{-1} X_{(i)}}{m^T V^{-1} m}.$$

$$W = \frac{(m^T V^{-1} m)^2 \cdot \hat{\sigma}^2}{m^T V^{-1} V^{-1} m \cdot \sum_{i=1}^m (X_i - \bar{X})^2}$$

Pozn.

- pro  $V$  jsou známy pro  $m \leq 20$ , pro  $m > 20$  se používají aproximace.
- $W \leq 1$ . Pro alternativní měřící malé hodnoty  $W$ . Rozdělení  $W$  za platnosti  $H_0$  je tabelováno.
- test se hodí pro malé rozsahy vzátek ( $m \leq 50$ ).

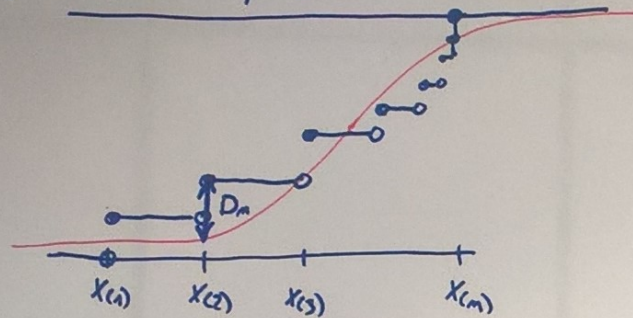
c) test závislosti na empirické distribuční funkci

Kolmogorovův - Smirnovův test

$H_0: F = F^*$ , kde  $F^*$  je distribuční funkce  $N(\mu, \sigma^2)$ ,  $\mu, \sigma^2 > 0$  známe.

$H_1: F \neq F^*$

$$\hat{F}_m(x) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{X_i \leq x\}$$



$F^*$

$$D_m = \max_{x \in \mathbb{R}} \left\{ \left| \hat{F}_m(x) - F^*(x) \right| \right\} = \max_{i=1, \dots, m} \left\{ \left| \frac{i}{m} - \Phi\left(\frac{x_{(i)} - \mu}{\sigma}\right) \right| \right\}$$

$\sqrt{m} D_m \stackrel{H_0}{\rightsquigarrow} \sup_{t \in [0,1]} |B(t)|$ , kde  $B(t)$  je Brownův most v  $C(0,1)$ .

Rozdělení n.o.  $Y = \sup_{t \in [0,1]} |B(t)|$  je známe, ale nemá se vyjádřit v uzavřené formě.

$$F_Y(y) = 1 - 2 \sum_{j=1}^{\infty} (-1)^{j+1} e^{-2j^2 y^2}, \quad y > 0$$

$$\approx 1 - 2e^{-2y^2}$$

(1- $\alpha$ )-kvantil je přibližně  $\sqrt{\frac{1}{2} \log \frac{2}{\alpha}}$ , spec. pro  $\alpha = 0,05$  je 95% - kvantil 1,36.

Test lze použít jen pro  $\mu$  a  $\sigma^2$  známe.

Lillieforsova modifikace K-S testu

- zde již testujeme předání  $H_0$  (F je normální s rušivými parametry)

$\mu$  a  $\sigma^2$  odhadneme z dat  $\rightarrow \hat{\mu}, \hat{\sigma}^2$

$$D_m^* = \max_{i=1, \dots, m} \left\{ \left| \frac{i}{m} - \Phi\left(\frac{x_{(i)} - \hat{\mu}}{\hat{\sigma}}\right) \right| \right\}$$

Testová statistika  $\sqrt{m} D_m^*$  již nemá za  $H_0$  nějaké uvedené rozdělení  $\rightarrow$  musíme upravit "kvantily"  
pro  $n \geq 30$  a  $\alpha = 0,05$  se používá 0,886.

dobrá kritéria:  $m \cdot \int_{-\infty}^{\infty} (\hat{F}_m(x) - F^*(x))^2 \underbrace{\Psi(F^*(x))}_{\text{váhová funkce}} f^*(x) dx$

### Gramérův-von Misesův test

$\psi(y) = 1$   
 $CVM = m \int_{-\infty}^{\infty} (\hat{F}_m(x) - F^*(x))^2 f^*(x) dx = \frac{1}{12m} + \sum_{i=1}^m \left( p_{(i)} - \frac{2i-1}{2m} \right)^2$ , kde  $p_{(i)} = \Phi\left(\frac{x_{(i)} - \hat{\mu}}{\hat{\sigma}}\right)$ .

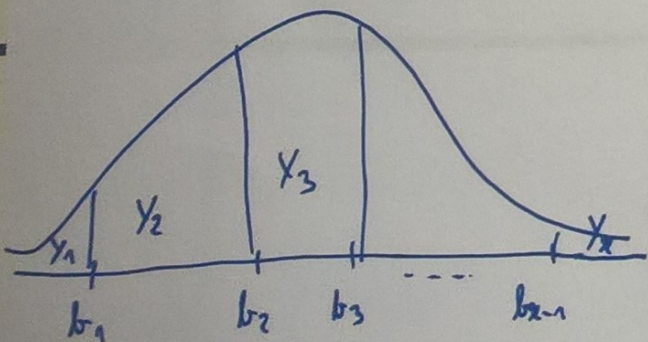
### Andersonův-Darlingův test

$\psi(y) = \frac{1}{y(1-y)}$   $0 < y < 1$  ... dáva větší váhu pozorování na chodech

$AD = m \int_{-\infty}^{\infty} \frac{(\hat{F}_m(x) - F^*(x))^2}{F^*(x)(1-F^*(x))} f^*(x) dx = -m - \frac{1}{m} \sum_{i=1}^m (2i-1) (\log p_{(i)} + \log(1-p_{(m-i+1)}))$ .

d) testy dobré shody

### Pearsonův $\chi^2$ -test dobré shody



označme  $Y_i =$  počet pozorování, které padnou do intervalu  $(b_{i-1}, b_i]$  pro  $i=1, \dots, k$   
 $b_0 = -\infty, b_k = \infty$

$p_i = P(\text{dané pozorování padne do intervalu } (b_{i-1}, b_i]) = p_i(\mu, \sigma) = P(X_j \in (b_{i-1}, b_i]) = \int_{b_{i-1}}^{b_i} f(x, \mu, \sigma) dx$

$m p_i =$  očekávaný počet pozorování, které padnou do intervalu  $(b_{i-1}, b_i]$   
 $m p_i(\mu, \sigma)$

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - n p_i(\mu, \sigma))^2}{n p_i(\mu, \sigma)} \quad H_0^* \chi_{k-1}^2$$

- Testujeme  $H_0^*$ , tedy  $\mu$  a  $\sigma^2$  známé.

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - n p_i(\hat{\mu}, \hat{\sigma}))^2}{n p_i(\hat{\mu}, \hat{\sigma})} \quad H_0 \chi_{k-3}^2$$

- Testujeme  $H_0$ , tedy  $\mu$  a  $\sigma^2$  různé parametry.

Poznámka

Jaké volit intervaly  $(b_{i-1}, b_i]$ ?

$$p_i(\hat{\mu}, \hat{\sigma}) = \frac{1}{k}, \quad \forall i = 1, \dots, k$$

$$k = 2 \cdot n^{\frac{2}{5}}$$

$$= 15 \left( \frac{n}{100} \right)^{\frac{2}{5}}$$