

Model selection (výběr modelu)

- (i) Je náš model "vhodný"? Popisuje správně naše data?
- (ii) Máme-li více modelů, který z nich je "nejlepší"? Který si vybrat?

Grafické metody pro posouzení vhodnosti modelu

- jsou založeny na porovnání teoretického a empirického rozdělení

- vykreslení teoretické a empirické dist. funkce v jednom grafu

$$D(x) = \hat{F}_n(x) - F^*(x), x \in \mathbb{R}, \text{ kde } \hat{F}_n(x) \text{ je empirická dist. fu a } F^* \text{ je teoretická dist. fu (s odhadnutými parametry)}$$

- vykreslení teoretické hustoty a jejího odhadu do jednoho grafu (histogram, jádrový odhad, ...)

• Q-Q plot

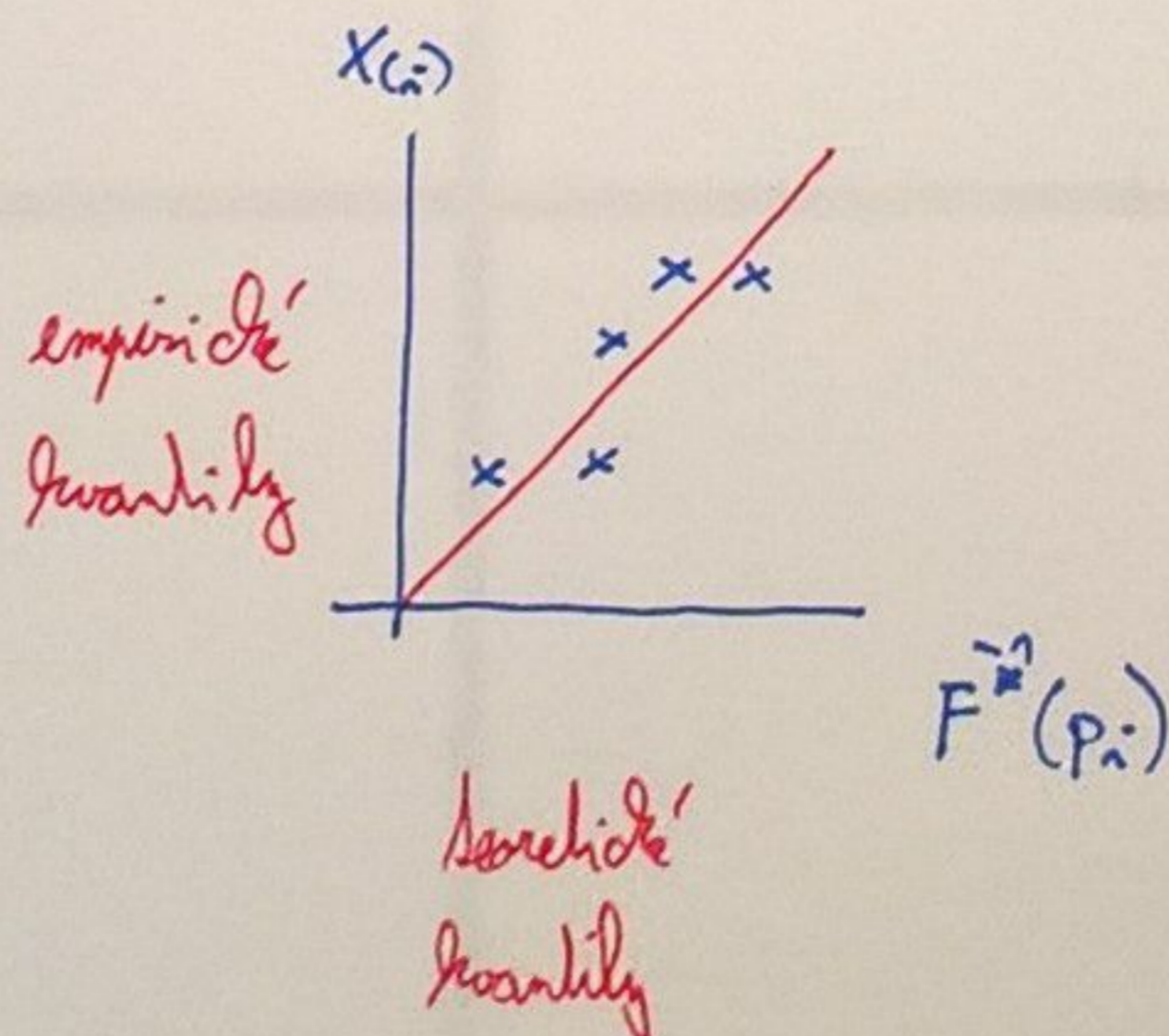
upřádané pozorování $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$

$$x_{(i)} \text{ je } p_i = \frac{i - \beta}{m + 1 - 2\beta} \quad (0 \leq \beta < 1) \text{ - } i\text{-tý kvantil}$$

opraxi: $\beta = 0,5$
 $\beta = 0,3175$

Q-Q plot je graf $[F^{\beta^{-1}}(p_i), x_{(i)}]$ pro $i = 1, \dots, m$

$F^{\beta^{-1}}$ je teoretická kvantilová funkce (s odhadnutými parametry)



Poznámka

N-P plot je Q-Q plot pro "osvětlením normality dat"

$$\beta = 0,3175 \text{ pro } m \leq 10$$

$$\beta = 0,5 \text{ pro } m > 10$$

$F^{\beta^{-1}}$ je Φ^{-1} kvantilová funkce $N(0,1)$.

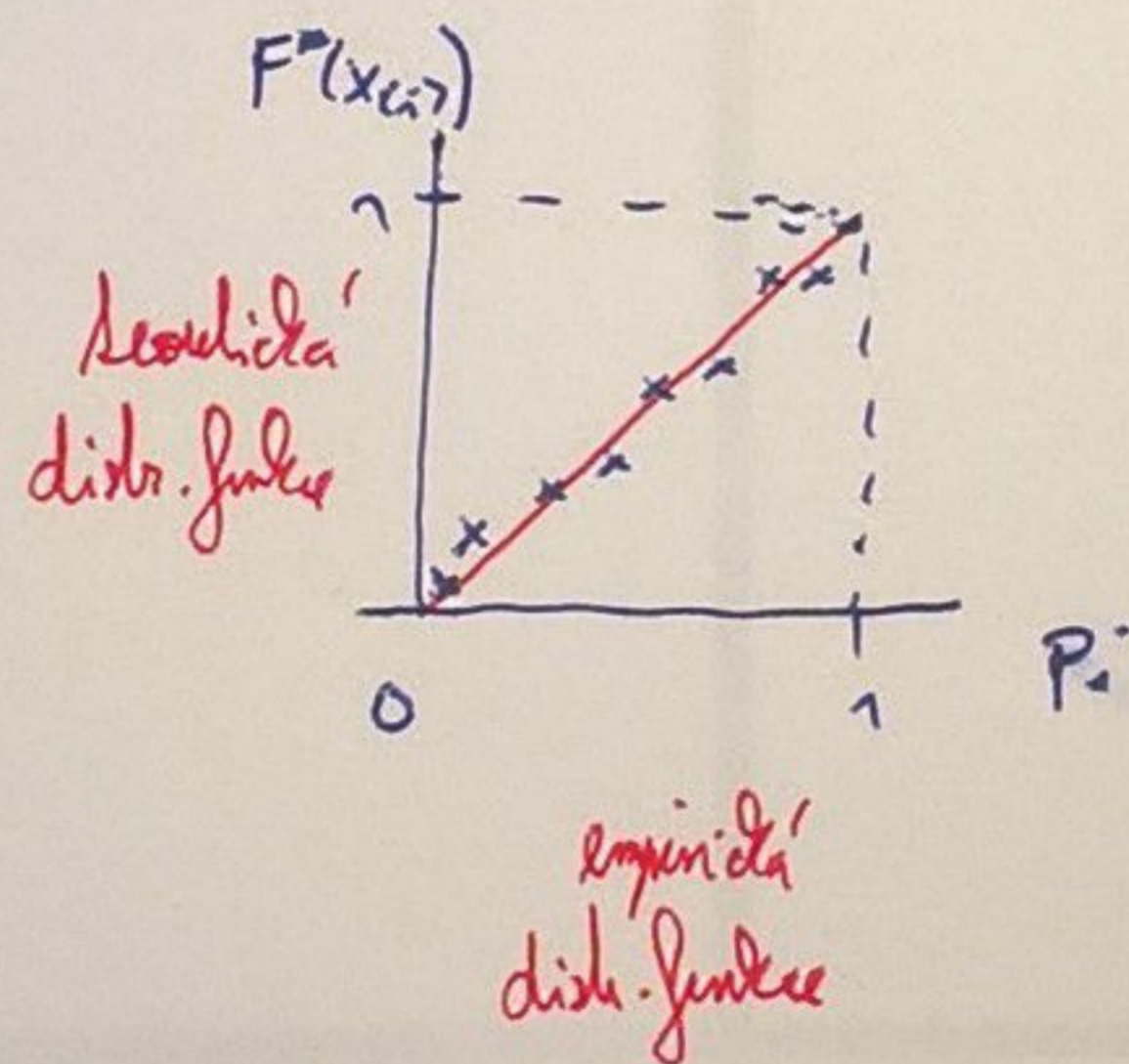
• P-P plot

- porovnání empirické a teoretické distr. funkce

$$\hat{F}_m(x) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{X_i \leq x\} \dots \text{drobná modifikace} \quad \tilde{F}_m(x) = \frac{1}{m+1} \sum_{i=1}^m \mathbb{1}\{X_i \leq x\} \dots \text{malý rozdíl} \quad P_i = \frac{i}{m+1} \text{ pro } i = 1, 2, \dots, m.$$

P-P plot je graf $[P_i, F^*(x_{(i)})]$ pro $i = 1, \dots, m$

F^* je teoretická distribuční funkce (s odhadnutými parametry)



Statistické testy pro posouzení vhodnosti modelu

X_1, \dots, X_m je náh. výběr s distribuční funkcí F

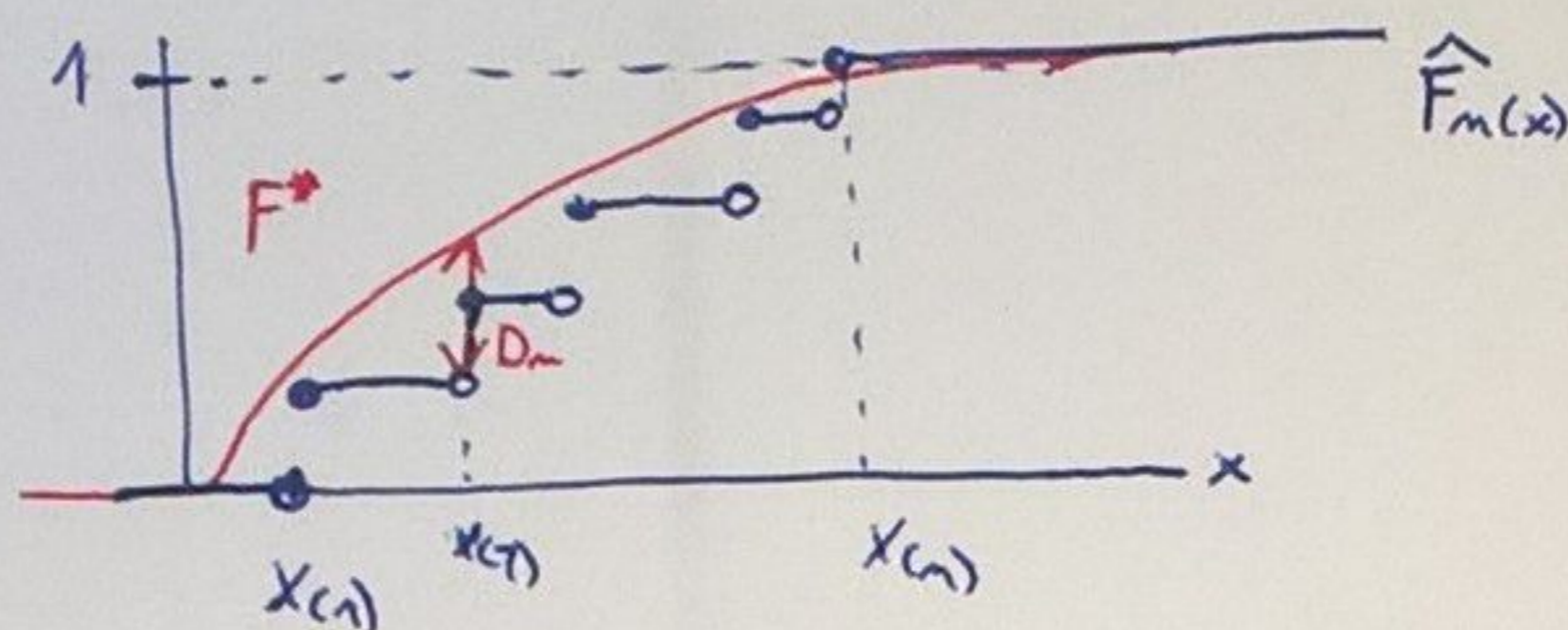
$$H_0: F = F^*$$

F^* je nějaká známá distribuční funkce

$$H_1: F \neq F^*$$

• Kolmogorov - Smirnov test

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}$$
 je empirická distribuční funkce



$$D_n = \max_{x \in \mathbb{R}} |\hat{F}_n(x) - F^*(x)| = \max_{i=1, \dots, n} \left| \frac{i}{n} - F^*(x_{(i)}) \right|$$

Teorem

je-li F^* spojitá, pak na hladině H_0 $\sqrt{n} D_n$ má asymptoticky při $n \rightarrow \infty$ rozdělení $\sup_{t \in [0,1]} |B(t)|$, kde $B(t)$ je Brownův most v $C(0,1)$.

Poznámka

Rozdělení máh. veličiny $Y = \sup_{t \in [0,1]} |B(t)|$ je známe, ale nemá se vyjádřit v uzavřené podobě. Její distribuční funkce je

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

$$\rightarrow \text{přibližná kvantilová funkce } \bar{F}_Y^{-1}(\alpha) = \sqrt{\frac{1}{2} \log \frac{2}{1-\alpha}}$$

$$\approx 1 - 2e^{-2y^2}, \quad y > 0$$

Poznámka:
 použijeme $\sqrt{n} D_n$. je-li $\sqrt{n} D_n > \bar{F}_Y^{-1}(1-\alpha)$... zamítáme H_0 na hladině významnosti α
 měření na F^*

Poradímky

- keď hce použiť jin ro prípadě, že H_0 je plně specifikovaná modelem. Pokud jsou data použita nejprve použita k odhadu parametrů, následně krit. hodnota je měřena (keď je použitá konzervativní)

- kritická hodnota ro tomto prípadě závisí na testovanej distribucii

- pro normální rozdělení je použitá - Lillieforsova test (95% kvantil je 1,36 míří na 0,886)

- pro ostatní rozdělení se p-hodnota testu měří pomocí MC simulací:

X_1, \dots, X_n je náh. vzor z distribuční funkci $F(x, \theta)$

$H_0: F = F(x, \theta)$ pro nějaké θ

(i) odhadneme parametru θ z dat $\rightarrow \hat{\theta}$

(ii) měříme hodnota testov. statistiky pro hypotetickou distribuční funkci $F(x, \hat{\theta})$, označíme ji T_0

(iii) regenerujeme nový náhodný vzor o rozsahu n z rozdělení z distr. fci $F(x, \hat{\theta})$

(iv) odhadneme parametru θ z těchto dat $\rightarrow \tilde{\theta}$

(v) měříme hodnota testov. statistiky pro hypotetickou distr. fci $F(x, \tilde{\theta})$, označíme ji T

(vi) když (iii) \rightarrow (v) měřící statistiku opakuje; p-hodnota odhadneme jako podíl případů, když $T > T_0$.

• Andersonova - Darlingova test

patří do třídy testů dané statistikou: $m \cdot \int_{-\infty}^{\infty} (\hat{F}_m(x) - F^*(x))^2 \cdot w(F^*(x)) f^*(x) dx$ pro nějakou váhovou funkci w .

$w(y) = \frac{1}{y(1-y)}$ $0 < y < 1$... dáva větší váhu pozorováním na chvostech

$$A^2 = m \cdot \int_{-\infty}^{\infty} \frac{(\hat{F}_m(x) - F^*(x))^2}{F^*(x)(1-F^*(x))} f^*(x) dx = -m - \frac{1}{m} \sum_{i=1}^m (2i-1) [\log F^*(X_{(i)}) + \log(1-F^*(X_{(m+1-i)}))]$$

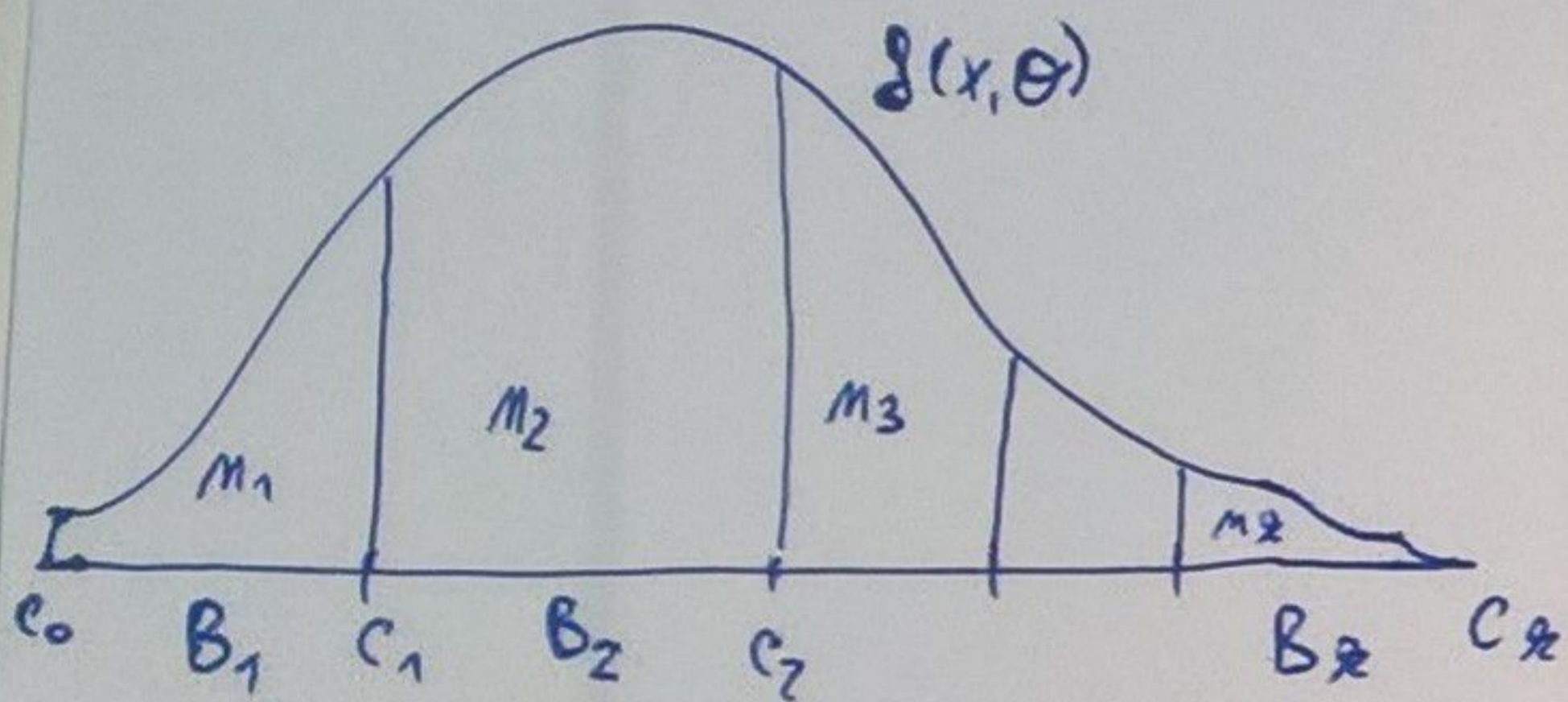
Porovnání

Pro F^* známé hodnoty A^2 pářím na levé distribuci F^* \Rightarrow mezikvily "univerzální" kvantily.

p-hodnota se učí pomocí simulací Monte Carlo.

Gramérus-von Misesův test - používá náhodou funkci $\rho(y) = 1$.

• Pearsonův χ^2 - test dobré shody



• dva hodnot n.o. X_i rozdělíme na k intervalů B_1, \dots, B_k

• označíme m_j počet pozorování, které padnou do intervalu B_j

• označíme $\mu_j(\theta) = P(X_i \in B_j) = \int_{c_{j-1}}^{c_j} f(x, \theta) dx = F(c_j, \theta) - F(c_{j-1}, \theta)$

• celkem máme m pozorování; tedy v intervalu B_j by mělo být $m \cdot \mu_j(\theta)$ pozorování (očekávané četnosti).

• by pozorování \triangleright pozorováními četnostmi m_j :

$$\chi^2 = \sum_{j=1}^k \frac{[m_j - m \cdot \mu_j(\theta)]^2}{m \cdot \mu_j(\theta)}$$

Testování

že platí H_0 má levá statistika χ^2 při $m \rightarrow \infty$ asymptoticky χ^2 rozdělení \triangleright $(k-1)$ stupni volnosti

je-li $\chi^2 > \chi_{1-\alpha}^2 (k-1)$... zamítáme H_0 .

Pravidla

Podud jsou data nejprve považována k odhadu nezávislého p -rozměrného parametru, pak $\chi^2 \stackrel{H_0}{\approx} \chi^2(k-1-p)$. (metoda minimálního χ^2)

Volba k ? - heuristická pravidla: $k \approx 2 \cdot n^{\frac{2}{5}}$ či $k \approx 15 \cdot \left(\frac{n}{100}\right)^{\frac{2}{5}}$

Intervaly B_1, \dots, B_k se pak volí „stejně pravděpodobné“, tedy C_1, C_2, \dots, C_{k-1} jsou $\frac{1}{k}, \frac{2}{k}, \dots, \frac{k-1}{k}$ - hranice rozdělení s hustotou $f(x, \hat{\theta})$.

Typy vhodného modelu (z několika kandidátů)

- princip Occamovy břitvy (princip logické úspornosti) - vybereme co nejjednodušší model (rozhněme se přecem, lepší interpretovatelnost)

- judgement-based přístup (záložný na subjektivním úsudku analytika)
 - rozhodnutí založené na vizuálních grafech a tabulkách (tail vs. mod fit)
 - rozhodnutí založené na důvěrnější zkušenosti (Paretovo rozdělení pro výši příjmů, Benfordovo pro rozdělení číselných promích čísel)
 - model je plně určen situací, kterou má popisovat (alternativní rozdělení pro házení mincí)

• score-based přístup (záložný na číselných charakteristikách; objektivní)

- nejmenší hodnota testové statistiky vybraného testu
 - nejvyšší p -hodnota vybraného testu
 - nejvyšší hodnota věrohodnosti (logaritmičeská věrohodnosti)
 - nejmenší hodnota nějakého penalizačního kritéria, např.
- } neberou do úvahy složitost modelu (počet parametrů)

$$AIC = -2 \log L(\hat{\theta}) + 2 \cdot p$$

$L(\hat{\theta})$ věrohodnost
 p počet odhadovaných parametrů

$$BIC = -2 \log L(\hat{\theta}) + p \cdot \log n$$