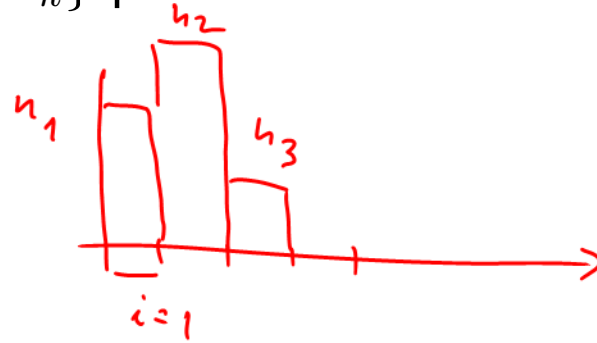


Odhad entropie z histogramu

Máme histogram $\{n_1, n_2, \dots, n_k\}$ proměnné x se šířkou přihrádky (třídy) $h > 0$.

$$\text{Platí } n = \sum_{i=1}^k n_i$$



2 případy:

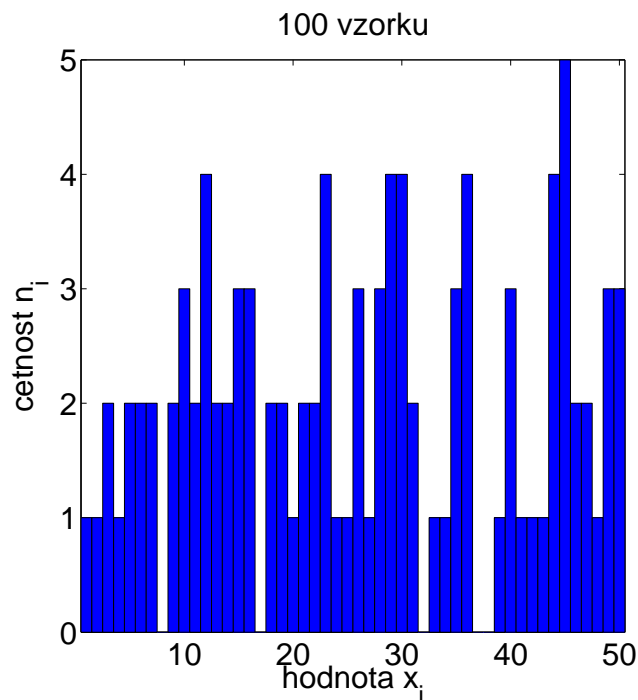
1. diskrétní náhodná proměnná: chceme entropii v přirozeném rozlišení
2. kvantizovaná spojitá náhodná proměnná: chceme entropii původní spojitě proměnné

$$\hat{H}(x) = \ln h - \sum_{i=1}^k \frac{n_i}{n} \ln \frac{n_i}{n}$$

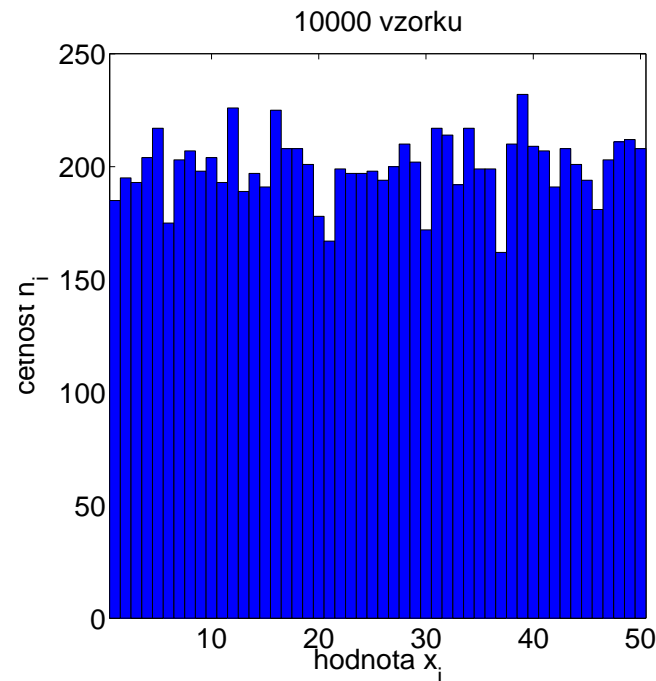
- h . . . v přirozených jednotkách oboru hodnot
- Bez členu $\ln h$ by hodnota statistiky rostla se zmenšováním rozlišení histogramu h

Je histogram kvalitním odhadem rozdělení psti?

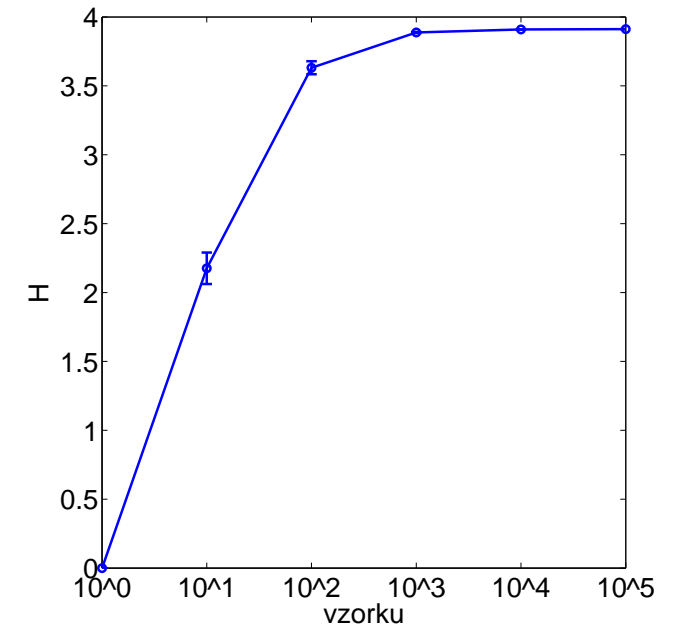
$x \in \{1, \dots, 50\}$: náhodná proměnná s rovnoměrným rozdělením



$$\hat{H} = 3.6874 \text{ nat}$$



$$\hat{H} = 3.9095 \text{ nat}$$



- teoretická entropie $\ln(50) = 3.9120 \text{ nat}$
- hodnota četnosti je náhodná proměnná

Volba šířky přihrádky histogramu

Systém $\mathbf{S} = \{s_1, s_2, \dots, s_q\}$

q dimenze histogramu

n počet měření

$\hat{\sigma}_i$ odhad rozptylu proměnné s_i

h_i šířka přihrádky pro proměnnou s_i

Předpoklad normálního rozdělení s diagonální kovarianční maticí

$$f(\underline{\mathbf{x}}) = \frac{1}{\sqrt{\det(2\pi\mathbf{S})}} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\bar{\mathbf{x}}})^\top \mathbf{S}^{-1}(\underline{\mathbf{x}} - \underline{\bar{\mathbf{x}}})}$$

$$\mathbf{S} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_q^2)$$

Scottovo pravidlo

$$\frac{h_i}{\hat{\sigma}_i} \approx \frac{3.5}{2+q\sqrt{n}}$$

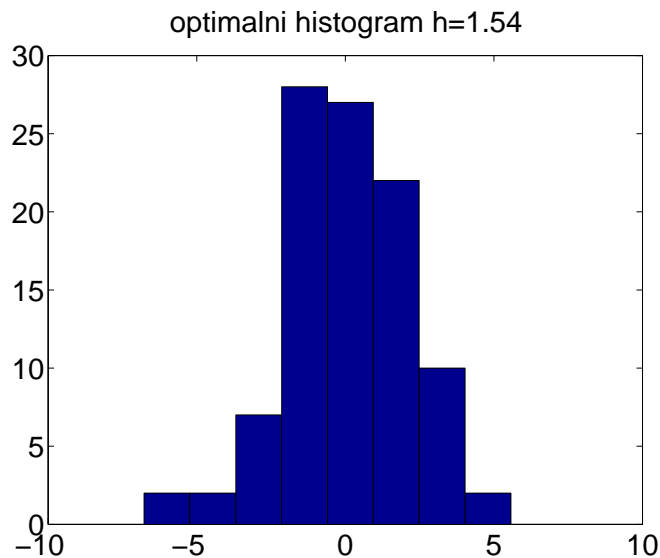
Scott, D. W. Multivariate Density Estimation: Theory Practice, and Visualization, John Wiley & Sons, Chichester 1992.

Příklad

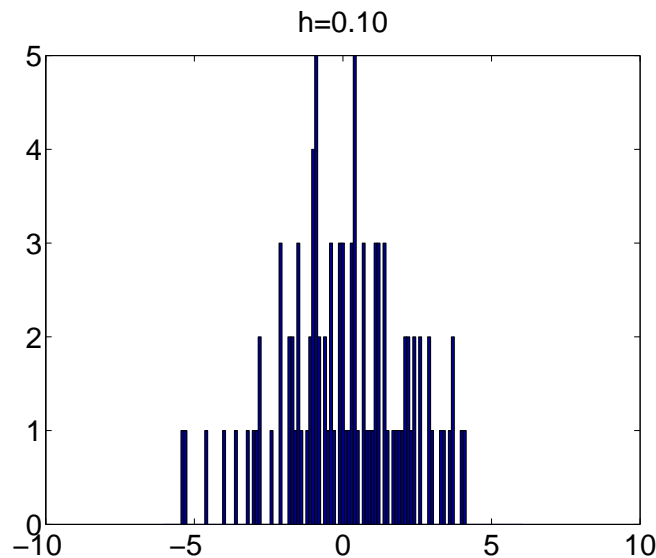
x – spojitá náhodná proměnná s rozdělením $N(0, 2)$

$$H(x) = \ln \sigma \sqrt{2\pi e} = 2.1121$$

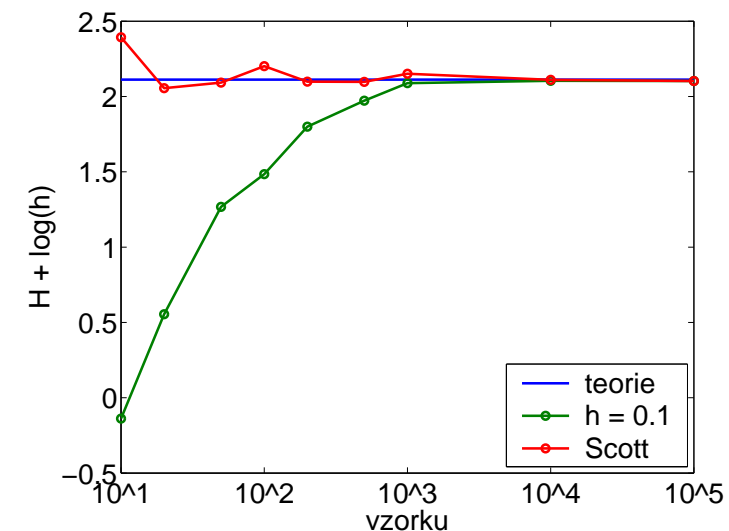
$n = 100$ vzorků, dimenze $q = 1$, takže $h = \frac{3.5 \cdot 2}{\sqrt[3]{100}} \approx 1.51$



$$\hat{H} = 2.1279 \text{ nat}$$



$$\hat{H} = 1.3026 \text{ nat}$$



Scott: pokaždé vypočteme novou hodnotu h a přepočteme histogram

- entropie z optimálního histogramu je lepším odhadem $H(x)$

Estimátor entropie bez histogramování (Kožačenko-Leoněno)

Dáno: množina vektorových měření $\{\underline{x}_i, i = 1, \dots, n\}$ z neznámého spojitého rozdělení pravděpodobnosti

Cíl: výpočet entropie bez diskretizace a histogramování

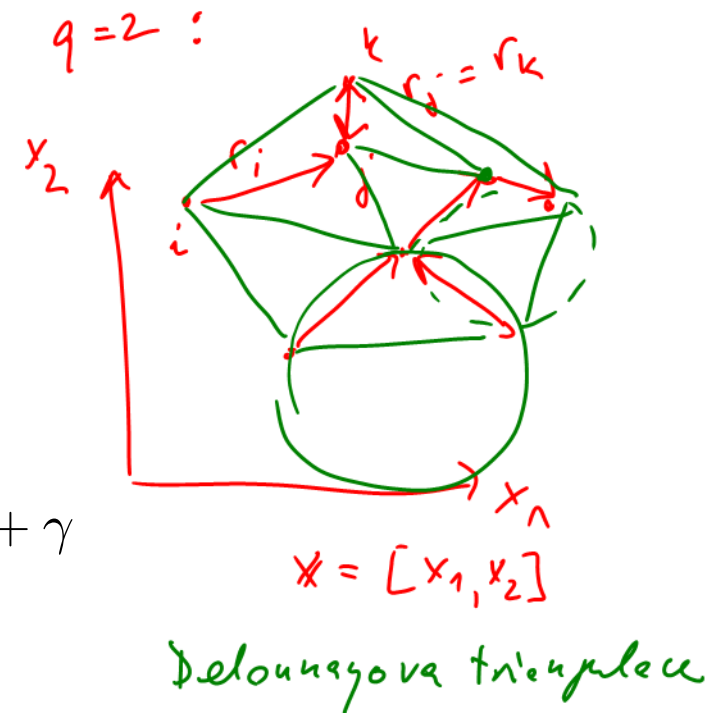
q dimenze vektoru měření

n počet měření

r_i euklidovská (L_2) vzdálenost k nejbližšímu sousedu \underline{x}_i

γ Euler-Mascheroniho konstanta ($\gamma \approx 0.5772156649$)

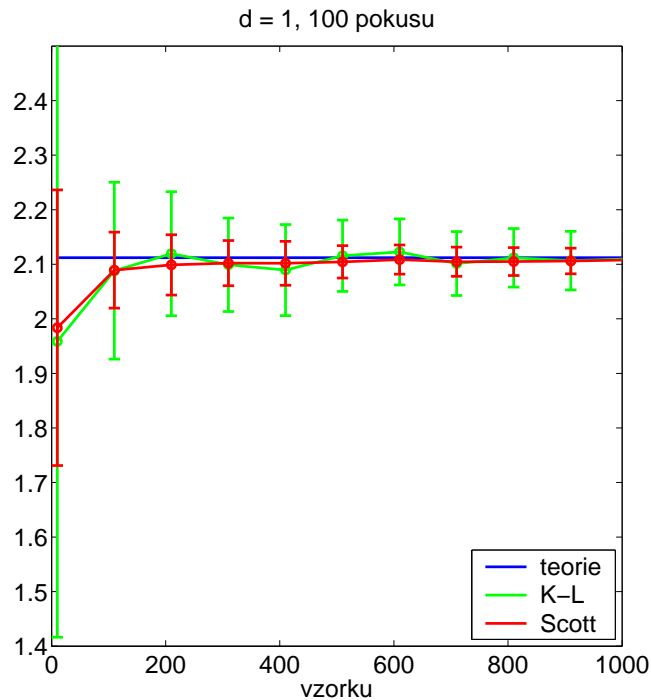
$$H = \frac{q}{n} \sum_{i=1}^n \ln r_i + \ln \frac{(n-1)\pi^{\frac{q}{2}}}{\Gamma(1 + \frac{q}{2})} + \gamma$$



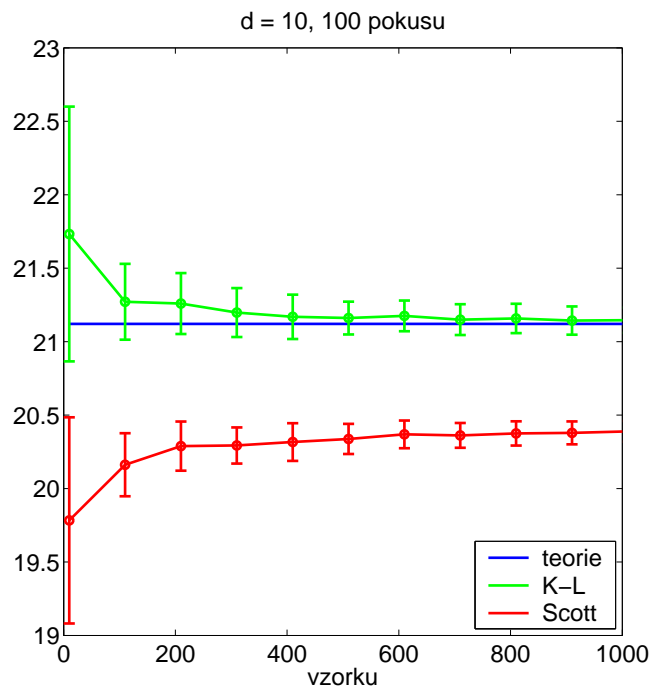
Poznámky

- množina nejbližších sousedů pro všechna \underline{x}_i lze teoreticky nalézt za dobu $O(c^q n \log n)$ pro libovolné q .
- degenerovanost: použít $\ln \max(r_i, \frac{1}{\sqrt{n}})$ místo $\ln r_i$.
- vhodné pro velká q (viz příklad)
- vhodné pro multimodální rozdělení psti
- $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt, x \in \mathbb{R}, \Gamma(k) = (k-1)!, k \in \mathbb{N}$

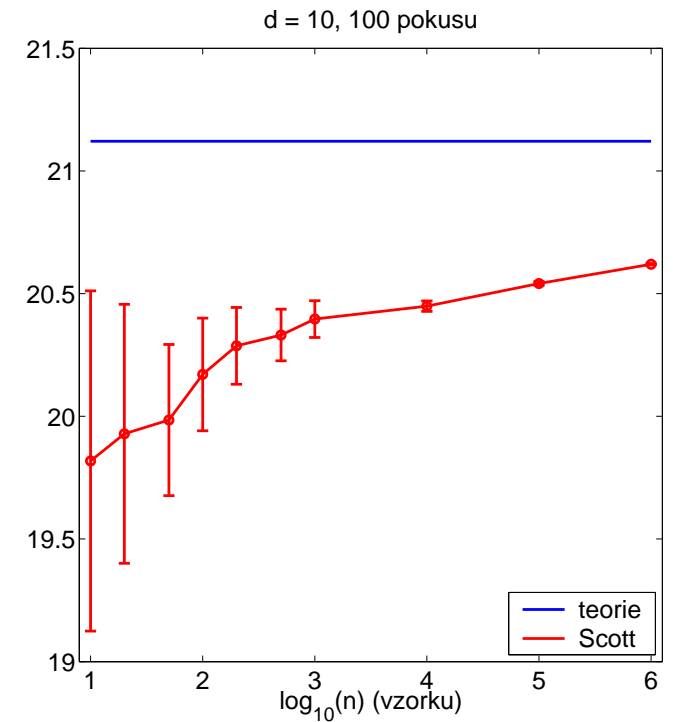
Příklad



$$\underline{\mathbf{x}} \sim N(0, 2), q = 1$$



$$\underline{\mathbf{x}} \sim N(0, 2), q = 10$$



- malá dimenze q : histogramová metoda má lepší rozptyl
- velká dimenze q : KL je méně vychýlený, rozptyly srovnatelné
- velká q : histogramová metoda konverguje (pro velmi velká n)

Odhad entropie a její chyby

Dáno: množina $\mathcal{D} = \{\underline{\mathbf{x}}_i, i = 1, 2, \dots, n\}$

Cíl: odhad entropie $\hat{H}(\mathcal{D})$ včetně chyby $\text{var}[\hat{H}(\mathcal{D})]$

Jackknife

1. Pro $i = 1, 2, \dots, n$ dělej:

a. zkonstruuuj $\mathcal{D}_i = \mathcal{D} \setminus \{\underline{\mathbf{x}}_i\}$

vynecháním jednoho bodu

b. odhadni $\hat{H}_i = \hat{H}(\mathcal{D}_i)$ z \mathcal{D}_i

2. Vypočti odhad entropie \hat{H} a chyby $\text{var}[\hat{H}]$:

$$\hat{H} = \frac{1}{n} \sum_{i=1}^n \hat{H}_i, \quad \text{var}[\hat{H}] = \frac{n-1}{n} \sum_{i=1}^n \left(\hat{H}_i - \hat{H} \right)^2$$

Poznámky

- pozor na rozdíl mezi výběrovým rozptylem a rozptylem výběrového průměru
- jackknife může být použit na jakoukoliv statistiku, nejen entropii, např. na medián, . . .
- Je to metoda *Resampling Theory*

Kontingenční analýza

Jaká je pravděpodobnost, že v i -té přihrádce histogramu bude n_i hodnot když celkem udělám n měření (pokusů)?

Sekvence, jejíž prvky jsou četnosti náhodných pokusů:

$$\left\{ \mathcal{E}_i \text{ se vyskytne } n_i \text{ krát v daném pořadí} \right\}_{i=1}^k$$

pravděpodobnost takové sekvence

$$p_1^{n_1} p_2^{n_2} p_3^{n_3} \cdots p_k^{n_k} \quad \text{víme, v jakém pořadí padaly kuličky do přihrádek}$$

$$\left\{ \mathcal{E}_i \text{ se vyskytne } n_i \text{ krát v libovolném pořadí} \right\}_{i=1}^k \quad \text{nevíme, v jakém pořadí padaly do přihrádek}$$

Pravděpodobnost, že v 1. přihrádce je n_1 hodnot, ve 2. přihrádce n_2 hodnot, . . . :

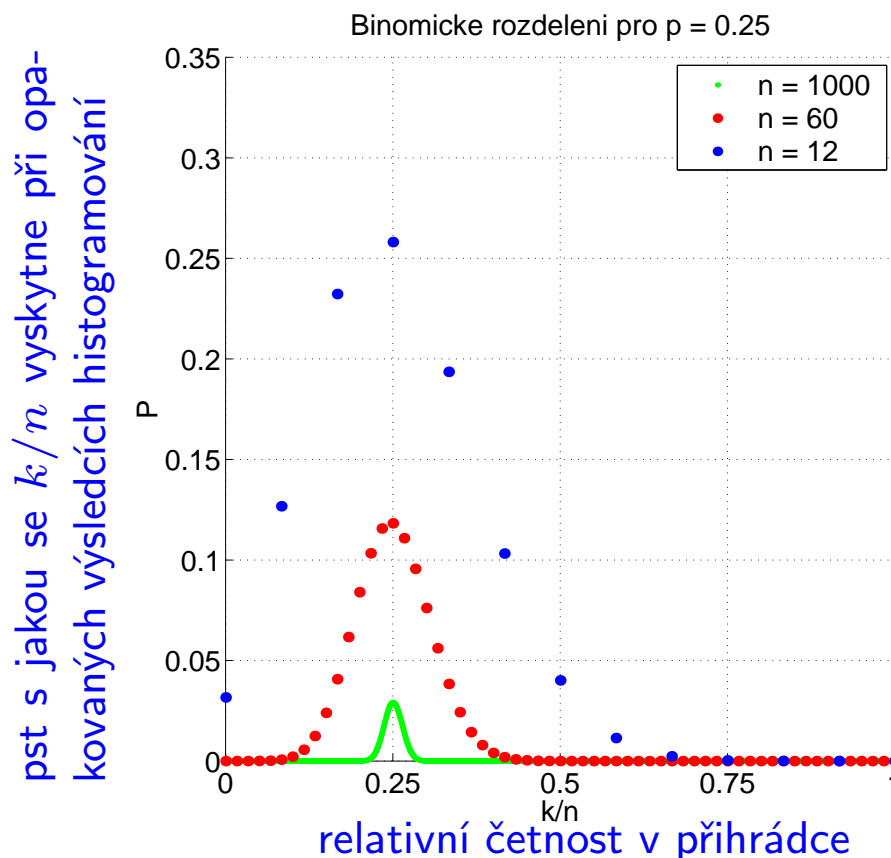
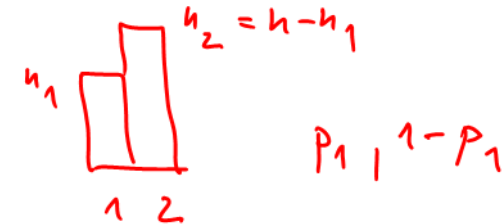
$$P(x_1 = n_1, x_2 = n_2, \dots, x_k = n_k) = \frac{n!}{n_1! n_2! \cdots n_k!} p_1^{n_1} p_2^{n_2} \cdots p_k^{n_k} \quad \text{permutace s opak.}$$

To je **multinomické rozdělení** s parametry n, p_1, p_2, \dots, p_k .

Příklad

Jaké hodnoty relativní četnosti k/n mohu očekávat v první přihrádce dvoupřihrádkového ‚histogramu‘, když se hodnota vyskytuje s pravděpodobností $p_1 = 0.25$ (a druhá s $p_2 = 0.75$)?

$$P(x_1 = k) = \binom{n}{k} p_1^k (1 - p_1)^{n-k}$$



- nejistotu musíme brát v úvahu, když činíme nějaký závěr z relativních četností

Vlastnosti multinomického rozdělení

Nechť $H = \{n_1, n_2, \dots, n_k\}$ má multinomické rozdělení a $k - 1$ je počet nezávislých prvků v H . Pak:

$$E\left(\frac{n_i}{n}\right) = p_i \quad \text{var}\left(\frac{n_i}{n}\right) = \frac{p_i(1-p_i)}{n}$$
$$\text{cov}\left(\frac{n_i}{n}, \frac{n_j}{n}\right) = -\frac{p_i p_j}{n} \quad q \stackrel{\text{def}}{=} \sum_{i=1}^k \frac{(n_i - n p_i)^2}{n p_i} \quad (\text{Pearson})$$

Veličina q má při $n \rightarrow \infty$ asymptoticky rozdělení χ_{k-1}^2 s hustotou:

$$f_m(x) = \frac{x^{\frac{m}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{m}{2}} \Gamma\left(\frac{m}{2}\right)}$$

$m = k - 1$: počet nezávislých prvků v $\{n_1, \dots, n_k\}$

Distribuční funkce je neúplná gamma funkce

Matlab: `gammainc(x/2,m/2)`

$$Q\left(\frac{x}{2}, \frac{m}{2}\right) = \frac{1}{\Gamma\left(\frac{m}{2}\right)} \int_0^{\frac{x}{2}} e^{-t} t^{\frac{m}{2}-1} dt, \quad m > 0$$

Pearsonova Statistika

$$q = \sum_{i=1}^k \frac{(n_i - n p_i)^2}{n p_i}$$

n_i – nahistogramované četnosti

p_i – model

tj., co v histogramu očekáváme

Př:

1. 2-D histogram nezávislých veličin x, y , pak $p_{ij} = p_i \cdot p_j = p(x = x_i) \cdot p(y = y_j)$ a my použijeme odhad p_{ij} modelu

$$p_{ij} = \frac{n_i}{n} \cdot \frac{n_j}{n} \quad q = \sum_i \sum_j \frac{(n_{ij} - n p_{ij})^2}{n p_{ij}}$$

(možné, ale nepraktické)

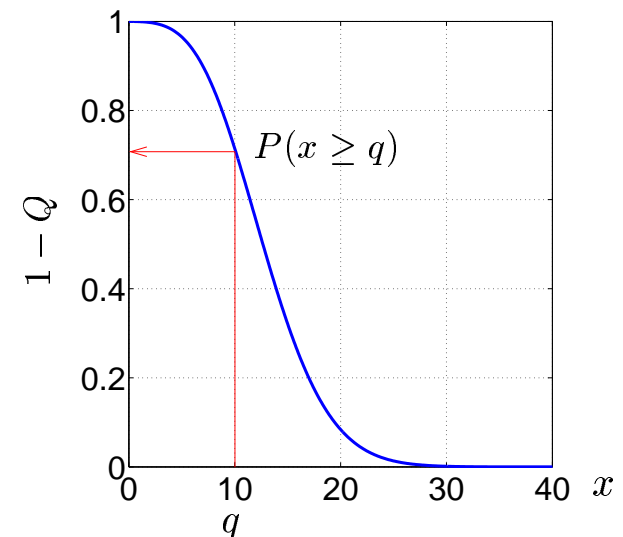
2. $p_i = p(x_i | \Theta)$

Pak

$$1 - Q\left(\frac{q}{2}, \frac{m}{2}\right) = P(x \geq q)$$

je pravděpodobnost, že změřená hodnota statistiky je ve skutečnosti větší než q , za předpokladu platnosti modelu.

m – počet přihrádek minus počet dodatečných podmínek, které musí splňovat soubor $\{n_i\}$ a které jsou potřeba k výpočtu hodnoty p_i (např. $\sum_{i=1}^k n_i = n$ a $\sum_{i=1}^k n_{ij} = n_i$ pro Př. 1).



Standardní kontingenční test

Nulová hypotéza H_0 : tvrzení X platí

H_0 je náš „model“

Chyba: „Odmítnu H_0 , a (ale) H_0 ve skutečnosti platí“

chyba 1. druhu

Cíl: $P(\text{error}) \leq \alpha$

α : hladina významnosti

α : penále, které musím zaplatit, když udělám chybu.

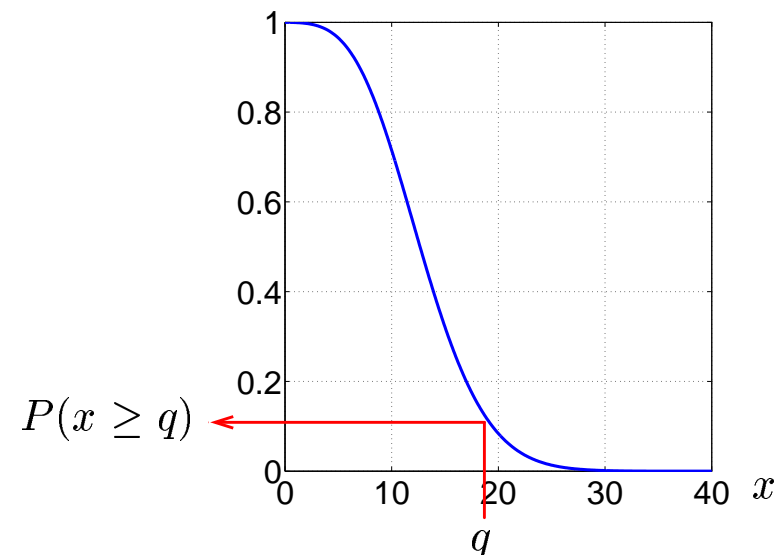
Řešení: Procedura statistického testu

1. vyslov H_0
2. změř n hodnot $\mathcal{D} = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n\}$
3. vypočti z \mathcal{D} statistiku q , která měří nesoulad s H_0
4. zvol (malé) α
5. odmítني H_0 když $P(x \geq q) > \alpha$

např. složky \underline{x} jsou statisticky nezávislé

např. Pearsonovu statistiku
typicky $\alpha = 0.01$ nebo $\alpha = 0.05$
 H_0 platí, takže to je $P(\text{error}) > \alpha$

- malé α dovolí tolerovat velké q
- 'Jsem velmi tolerantní a odmítnu H_0 jen, když je ve zřejmém rozporu s daty.'



Náš problém

Například:

- $H_0: p(\mathbf{a}, \mathbf{b}) = {}^1p(\mathbf{a}) \cdot {}^2p(\mathbf{b})$ pro test nezávislosti subsystémů
- $H_0: p(\mathbf{a}, \mathbf{b}, \mathbf{c}) = {}^1p(\mathbf{a}, \mathbf{b}) \cdot {}^2p(\mathbf{c} | \mathbf{b})$ pro test statistické významnosti rekonstrukce struktury systému

Pozn: $\{\mathbf{a}, \mathbf{b}\}$, $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ jsou rozklady množiny (vzorkovacích) proměnných systému. Můžeme si představit, že \mathbf{a} , \mathbf{b} , \mathbf{c} jsou vektorové proměnné.

Procedura testu

1. vyslov H_0
2. změř n hodnot $\mathcal{D} = \{\underline{\mathbf{x}}_1, \underline{\mathbf{x}}_2, \dots, \underline{\mathbf{x}}_n\}$
3. vypočti Pearsonovu statistiku q z \mathcal{D}
4. vypočti $\alpha = P(x \geq q)$
5. je-li dáno \mathcal{D} , pak H_0 platí s pravděpodobností alespoň α
 - odmítnutím H_0 udělám chybu α : $P(\text{odmítnu} \wedge \text{platí}) = \alpha$
 - malé $\alpha \Rightarrow$ mohu odmítnout
 - velké $\alpha \Rightarrow$ nemohu odmítnout = musím přijmout
 - $P(H_0 \text{ platí}) = P(\text{přijmu} \wedge \text{platí}) + \underbrace{P(\text{odmítnu} \wedge \text{platí})}_{\alpha} \geq \alpha$

Postup pro $p(\mathbf{a}, \mathbf{b}) = {}^1p(\mathbf{a}) \cdot {}^2p(\mathbf{b})$

1. Z kontingenční tabulky vypočteme

skutečná četnost $n(a_i, b_j)$

četnost predikovaná modelem $n \cdot p(a_i) \cdot p(b_j)$

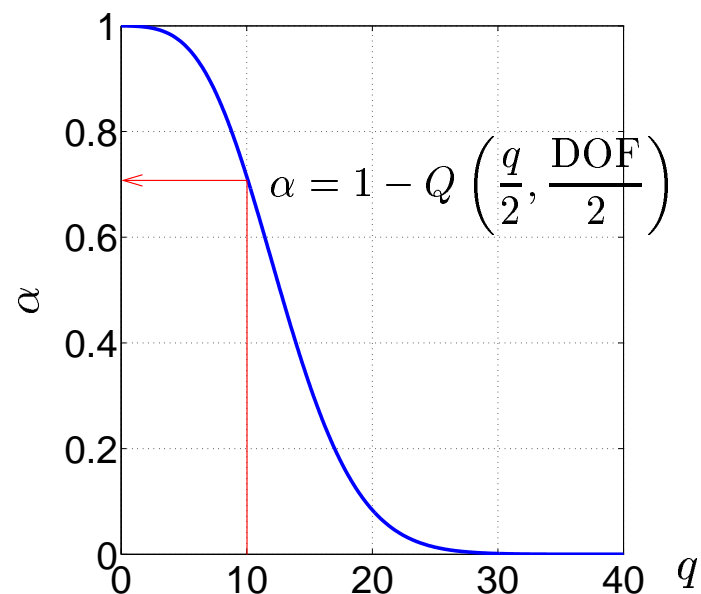
$$q = \sum_i \sum_j \frac{\left(n_{ij} - \frac{n_{i \cdot} \cdot n_{\cdot j}}{n} \right)^2}{\frac{n_{i \cdot} \cdot n_{\cdot j}}{n}}$$

2. stupně volnosti: $\text{DOF} = (r - 1)(c - 1)$

3. vypočteme $\alpha = 1 - Q\left(\frac{q}{2}, \frac{\text{DOF}}{2}\right)$

4. vyjde-li malé α , pak tvrdím, že \mathbf{a} a \mathbf{b} závislé

5. vyjde-li velké α , pak tvrdím, že \mathbf{a} a \mathbf{b} jsou nezávislé s pravděpodobností alespoň α



poznámky

Počet stupňů volnosti $\text{DOF} = rc - (r + c) + 1 = (r - 1)(c - 1)$

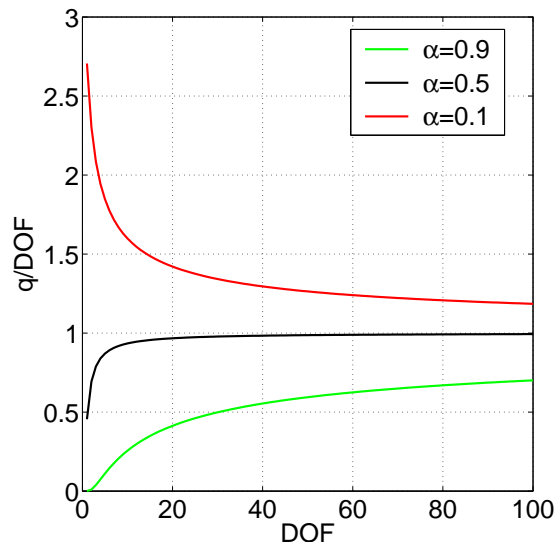
- Máme rc prvků v tabulce, ale použili jsme dodatečné vztahy

$$n_{.j} = \sum_i n_{ij}, \quad n_{i.} = \sum_j n_{ij}, \quad \text{kterých je dohromady } r + c.$$

- Ale tyto podmínky nejsou nezávislé, protože $\sum_j n_{.j} + \sum_i n_{i.} = 2n$, odečteme 1.

$$q = 2Q^{-1}\left(1 - \alpha, \frac{m}{2}\right), \quad m = \text{DOF}$$

rychlost růstu prahu přijatelnosti
s rostoucím rozlišením tabulky



$\frac{q}{\text{DOF}}$:

- velké rozlišení \Rightarrow velký počet DOF $\Rightarrow q \approx \text{DOF}$
- α přestává mít vliv
- téměř vše začíná být nezávislé
- **ale:** redukce rozlišení kvantizací zachová nezávislost
- \Rightarrow nemusíme se bát redukce rozlišení

Příklad z parlamentu: která dvojice hlasuje nezávisle na ostatních?

1. Nalezení nezávislé dvojice $(i, j) \in \binom{N}{2}$

použijeme vzájemnou informaci I

$$\arg \min_{i,j} I(\{s_i, s_j\}, \{s_k, s_l, s_m\})$$

2. Ověření statistické významnosti

$$H_0: p(s_1, s_2, s_3, s_4, s_5) = p(s_i, s_j) \cdot p(s_k, s_l, s_m)$$

s_i	s_j	I_{ij}											
			s_4			0	0	1	1				
			s_5			0	1	0	1				
s_1	s_2	0.0277	s_1	s_2	s_3	0	0	0	2	4	2	5	13
s_1	s_3	0.0299				0	0	1	3	3	9	3	18
s_1	s_4	0.0308				0	1	0	3	5	4	1	13
s_1	s_5	0.0303				0	1	1	3	3	4	5	15
s_2	s_3	0.0294				1	0	0	1	3	3	3	10
s_2	s_4	0.0263				1	0	1	1	3	0	1	5
s_2	s_5	0.0300				1	1	0	3	3	3	4	13
s_3	s_4	0.0273				1	1	1	1	4	3	5	13
s_3	s_5	0.0293											
s_4	s_5	0.0249											
						17	28	28	27				100

$$q = 15.07 \Rightarrow \text{nezávislé s } p \geq 0.82$$

Konec

