

Rekonstrukce diskrétního rozdělení psti metodou maximální entropie

Příklad

Lze nalézt ‚četnosti‘ nepozorovaných stavů tak, abychom si „vymýšleli co nejméně“?

Nechť n_i , $i = 1, 2, \dots, N$ jsou známé (absolutní) četnosti z neznámého diskrétního rozdělení pravděpodobnosti a necht' m_j , $j = 1, 2, \dots, M$ jsou neznámé četnosti z tohoto rozdělení. Hledáme hodnoty m_j tak, aby entropie celého rozdělení byla maximální.

$$J = - \sum_{i=1}^N \frac{n_i}{n} \log \frac{n_i}{n} - \sum_{j=1}^M \frac{m_j}{n} \log \frac{m_j}{n} + \lambda (n - \sum_{i=1}^N n_i - \sum_{j=1}^M m_j),$$

neznámé jsou m_j , $j = 1, 2, \dots, M$, n a λ .

Výsledek: $\frac{\partial J}{\partial n} = 0$ $\frac{\partial J}{\partial m_j} = 0$, $\frac{\partial J}{\partial \lambda} = 0$

$$n - n_0 - \sum_{j=1}^M m_j = 0$$

$$m_j = m = n_0 e^{-H_0}, \quad \text{kde} \quad n_0 = \sum_{i=1}^N n_i \quad \text{a} \quad H_0 = - \sum_{i=1}^N \frac{n_i}{n_0} \log \frac{n_i}{n_0}$$

(*)

$$H = H_0 + \log \frac{n}{n_0}, \quad \text{kde} \quad n = \underline{n_0 + M m}$$

Rekonstrukce spojitého rozdělení ze známých momentů

Příklad

Nechť μ je střední hodnota a σ rozptyl neznámého rozdělení pravděpodobnosti s hustotou $f(x)$. Jak vypadá $f(x): \mathbb{R} \mapsto [0, 1]$, která má za těchto podmínek maximální diferenciální entropii?

$$Q = - \int_{-\infty}^{\infty} f(x) \log f(x) dx + \lambda_0 \left(1 - \int_{-\infty}^{\infty} f(x) dx \right) + \lambda_1 \left(\mu - \int_{-\infty}^{\infty} x f(x) dx \right) + \lambda_2 \left(\sigma^2 - \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx \right) \quad (1)$$

Handwritten notes:
 $f_1(x) \equiv 1$
 $f_2(x) = x$
 $f_3(x) = (x - \mu)^2$
 $\frac{\partial Q}{\partial f} = 0 \quad \frac{\partial Q}{\partial \lambda_i} = 0 \quad i = 0, 1, 2$
Příklad užitečné matematiky

Výsledek

$$f(x) = \frac{1}{T} e^{a_1 x + a_2 x^2} = \dots = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (2)$$

Důkaz

1. pomocí Gibbsovy nerovnosti (viz dále)
2. pomocí variačního počtu z rovnice (1) ***1**

Zobecnění předchozího výsledku

Nechť $g_i(x)$, $i = 1, 2, \dots, K$ jsou známé funkce a $c_i = \int_{-\infty}^{\infty} g_i(x) f(x) dx$ jsou známé hodnoty. Rekonstruujeme z nich neznámou hustotu $f(x)$ tak, aby maximalizovala entropii $H(x)$.

Výsledek

$$f(x) = \frac{1}{T} e^{-\sum_{j=1}^K a_j g_j(x)}$$

a

$$H(x) = - \int_{-\infty}^{\infty} f(x) \log f(x) dx = \log T + \sum_{j=1}^K a_j c_j$$

Důkaz Nechť $\phi(x)$ je hustota rozdělení, které má větší entropii než $f(x)$ a vyhovuje podmínkám (má stejné hodnoty c_i). Platí Gibbsova nerovnost

$$H_{\phi}(x) = - \int_{-\infty}^{\infty} \phi(x) \log \phi(x) dx \leq - \int_{-\infty}^{\infty} \phi(x) \log f(x) dx =$$

$$H_{\phi}(x) > H_f(x)$$

$$= \int_{-\infty}^{\infty} \phi(x) \left(\log T + \sum_{j=1}^K a_j g_j(x) \right) dx = \log T + \sum_{j=1}^K a_j c_j = H_f(x)$$

$$\sum_{j=1}^K a_j \int_{-\infty}^{\infty} \phi(x) g_j(x) dx = \sum_{j=1}^K a_j c_j$$

což je spor.

Výpočet a_i , $i = 1, 2, \dots, K$

Zavedeme funkci

partition function

$$Z(a_1, \dots, a_K) = T = \int_{-\infty}^{\infty} e^{-\sum_{j=1}^K a_j g_j(x)} dx$$

Platí

$$\left[\frac{\partial Z}{\partial a_i} = \int_{-\infty}^{\infty} -g_i(x) e^{-\sum_{j=1}^K a_j g_j(x)} dx = -Z \int_{-\infty}^{\infty} g_i(x) f(x) dx = -Z c_i \right]$$

Takže a_i jsou řešením soustavy nelineárních rovnic

$$-\frac{1}{Z} \frac{\partial Z}{\partial a_i} = c_i, \quad i = 1, 2, \dots, K, \quad (3)$$

kde neznámé jsou a_i , kde c_i jsou dány a kde $g_i(x)$ jsou známé funkce.

Pozn

- soustava (3) je většinou obtížně řešitelná analyticky

⊛2 použijte pro důkaz výsledku (2) na str 28.

Identifikace struktury

Identifikace struktury

Nejlepší dekompozice systému S na množinu podsystémů $G = \{^1S, ^2S, \dots, ^qS\}$

Podproblémy

1. systematické generování rekonstrukčních hypotéz G

- neredundance $\forall i, j \in \{1, 2, \dots, q\}: ^iS \not\subseteq ^jS$
- pokrytí $\bigcup_i ^iS = S$

2. Rekonstrukce celkového systému S^* z G .

3. Vyhodnocení rekonstrukční chyby $\Delta(G) = L(S, S^*)$



(S . . množina proměnných systému)

Bez podmínky neredundance a pokrytí by byl počet hypotéz

$$2^{2^n - 1} - 1.$$

(n – počet proměnných, $n = 2 \Rightarrow 7$ hypotéz, $n = 4 \Rightarrow 32767$ hypotéz)

Zjemnění rekonstrukční hypotézy

Rekonstrukční hypotéza $G = \{^1S, ^2S, \dots, ^qS\}$ \overline{Pr}

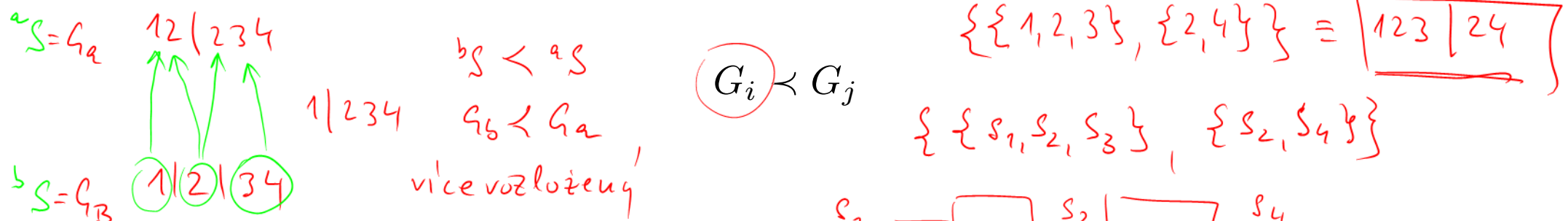
$$^1S = \{s_1, s_2, s_3\} \quad \{1, 2, 3\}$$

$$^2S = \{s_2, s_4\} \quad 123$$

$$\{2, 4\} \quad 24$$

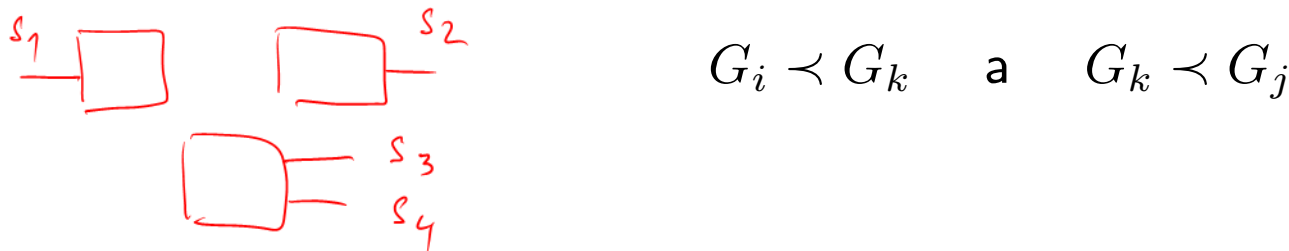
Zjemnění rekonstrukční hypotézy \prec

G_i je zjemněním G_j když pro každé $^xS \in G_i$ existuje $^yS \in G_j$ takové, že $^xS \sqsubseteq ^yS$. Značíme



Bezprostřední zjemnění \prec^*

G_i je bezprostředním zjemněním G_j jestliže neexistuje G_k takové, že



Relace bezprostředního zjemnění tvoří svaz na množině všech rekonstrukčních hypotéz.

Pozn: nejde o zjemňování rozkladu, rekonstrukční hypotézy netvoří rozklad (proměnné se opakují).

Generátor rekonstrukčních hypotéz

Vstup: Rekonstrukční hypotéza $G = \{{}^1\mathbf{S}, {}^2\mathbf{S}, \dots, {}^q\mathbf{S}\}$

Výstup: Všechna bezprostřední zjemnění G

Procedura (prohledávání svazu (G, \prec^*))

1. Pro $i = 1, 2, \dots, q$ dělej kroky 2,3.
2. Jestliže $|{}^i\mathbf{S}| \geq 2$ potom nahraď ${}^i\mathbf{S}$ množinou všech podmnožin ${}^i\mathbf{S}$ o velikosti $|{}^i\mathbf{S}| - 1$.
3. Po odstranění redundantních podmnožin v každém rozkladu dostaneme seznam bezprostředních zjemnění G . Pozn: může se opakovat na stejné úrovni zjemnění.

123 | 234 | 14

$$G = \{\{1, 2, 3\}, \{2, 3, 4\}, \{1, 4\}\} = \{{}^1\mathbf{S}, {}^2\mathbf{S}, {}^3\mathbf{S}\}$$

~~23~~ | 13 | 12 | 234 | 14

$$G_1 = \{\{1, 2\}, \{1, 3\}, \{2, 3, 4\}, \{1, 4\}\}$$

123 | 34 | 24 | ~~23~~ | 14

$$G_2 = \{\{1, 2, 3\}, \{2, 4\}, \{3, 4\}, \{1, 4\}\}$$

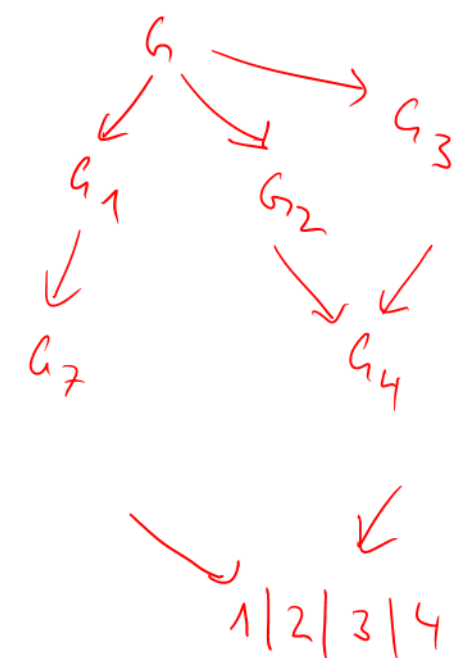
123 | 234 | ~~14~~ | ~~14~~

$$G_3 = \{\{1, 2, 3\}, \{2, 3, 4\}\}$$

$$G_4 = \{\{1, 2\}, \{1, 3\}, \{2, 3\}, \{3, 4\}, \{2, 4\}, \{1, 4\}\}$$

$$G_5 = \{\{1, 2\}, \{1, 3\}, \{2, 3, 4\}\}$$

$$G_6 = \{\{1, 2, 3\}, \{2, 4\}, \{3, 4\}\}$$



Prohledávání svazu zjemnění

- Rekonstrukční chyba Δ monotonně roste podél každé větve svazu zjemnění: Je-li $G_i \succ G_k \succ G_j$ potom $\Delta(G_i) \leq \Delta(G_k) \leq \Delta(G_j)$

• **Nejjednodušší postup:** Hledáme cestu, podle které Δ roste nejpomaleji

- **Složitější postup:** Na každé úrovni nás zajímá množina řešení s nejmenším Δ
 - Není nutno expandovat všechny větve: postačí metoda větví a mezí

[Narenda&Fukunaga 1977]

