

Přednáška 11: Bayesovská klasifikace; Úvod do grafických modelů a hierarchického učení

Ondřej Tichý

Bayesovské metody v machine learningu (BML)

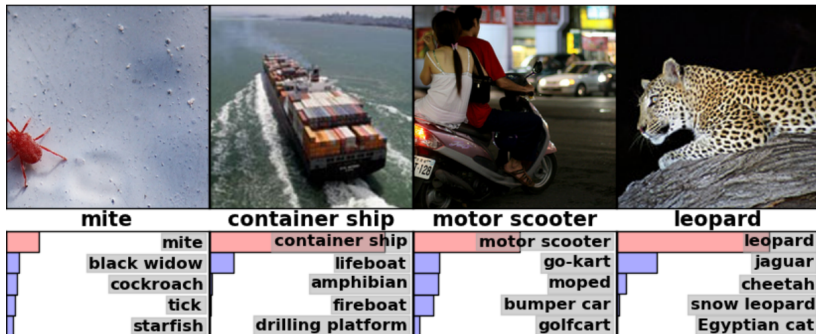
May 13, 2018

Obsah přednášky

- ▶ bayesovská klasifikace
 - ▶ principy
 - ▶ bayesian spam filtering
- ▶ grafické modely
 - ▶ hierarchické učení
 - ▶ demonstrace na maticové dekompozici

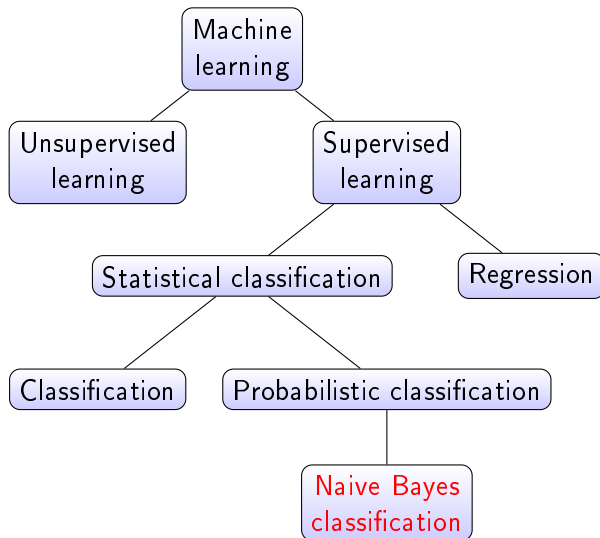
Bayesovská klasifikace

Úvod



Bayesovská klasifikace

Zařazení



Naive Bayes Klasifikátor

Princip

- ▶ mějme vektor pozorování $\mathbf{x} = (x_1, \dots, x_n)$, který chceme přiřadit ke třídě C_k ze setu $\{C_1, \dots, C_K\}$.
- ▶ teroreticky: máme tabulku pro příslušnost x_i k jednotlivým třídám a vyjádříme

$$p(C_k|\mathbf{x}), \quad (1)$$

což není vždy vhodné a možné.

Naive Bayes Klasifikátor

Princip

- ▶ mějme vektor pozorování $\mathbf{x} = (x_1, \dots, x_n)$, který chceme přiřadit ke třídě C_k ze setu $\{C_1, \dots, C_K\}$.
- ▶ teroreticky: máme tabulku pro příslušnost x_i k jednotlivým třídám a vyjádříme

$$p(C_k|\mathbf{x}), \quad (1)$$

což není vždy vhodné a možné.

- ▶ s využitím Bayesova teorému

$$p(C_k|\mathbf{x}) = \frac{p(C_k) p(\mathbf{x}|C_k)}{p(\mathbf{x})} \quad (2)$$

- ▶ $p(C_k|\mathbf{x})$ - posteriorní rozdělení
- ▶ $p(C_k)$ - apriorní rozdělení
- ▶ $p(\mathbf{x}|C_k)$ - model (věrohodnost)
- ▶ $p(\mathbf{x})$ - marginála přes \mathbf{x} (nezávisí na C_k)

Naive Bayes Klasifikátor

Princip

$$p(C_k|\mathbf{x}) \propto p(C_k) p(\mathbf{x}|C_k) = p(x_1, x_2, \dots, x_n, C_k) \quad (3)$$

- ▶ aplikujeme řetězové pravidlo (n-krát)

$$\begin{aligned} p(x_1, x_2, \dots, x_n, C_k) = & p(x_1|x_2, \dots, x_n, C_k) \times \\ & \times p(x_2|x_3, \dots, x_n, C_k) \dots p(x_n|C_k) p(C_k) \end{aligned}$$

Naive Bayes Klasifikátor

Princip

$$p(C_k|\mathbf{x}) \propto p(C_k) p(\mathbf{x}|C_k) = p(x_1, x_2, \dots, x_n, C_k) \quad (3)$$

- ▶ aplikujeme řetězové pravidlo (n-krát)

$$\begin{aligned} p(x_1, x_2, \dots, x_n, C_k) = & p(x_1|x_2, \dots, x_n, C_k) \times \\ & \times p(x_2|x_3, \dots, x_n, C_k) \dots p(x_n|C_k) p(C_k) \end{aligned}$$

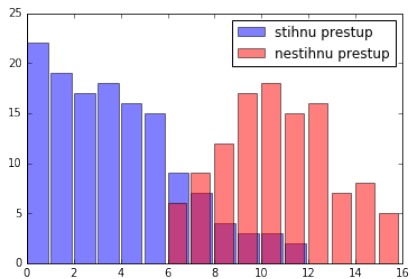
- ▶ proč “naivní” Bayes: předpokládáme podmíněnou nezávislost mezi prvky vektoru \mathbf{x} :

$$p(x_1, x_2, \dots, x_n, C_k) = p(x_1|C_k) p(x_2|C_k) \dots p(x_n|C_k) p(C_k)$$

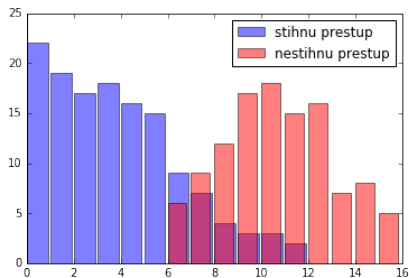
- ▶ pravděpodobnost příslušnosti k dané třídě:

$$p(C_k|\mathbf{x}) \propto p(C_k) \prod_{i=1}^n p(x_i|C_k) \quad (4)$$

Naive Bayes Klasifikátor

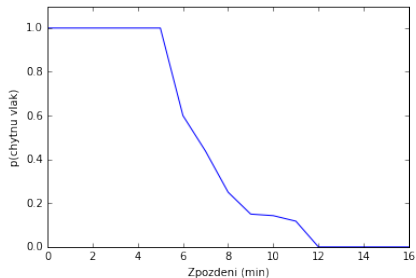
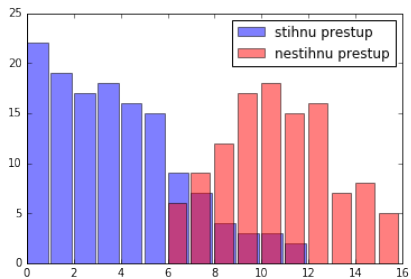


Naive Bayes Klasifikátor

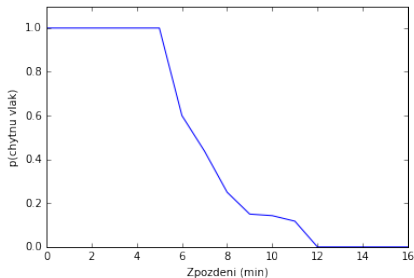
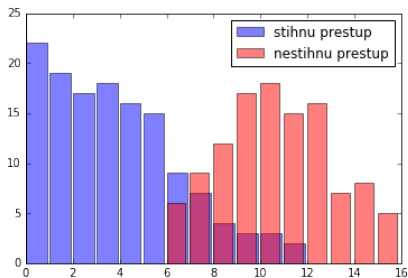


$$p(\text{stihnu} | \text{zpozdeni} = 6) = \frac{9}{9 + 6} = 0.6 \quad (5)$$

Naive Bayes Klasifikátor



Naive Bayes Klasifikátor



► MAP klasifikace: $\hat{k} = \arg \max_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i | C_k)$

Naive Bayes Klasifikátor

Jednoduchý příklad 2

- ▶ předpokládejme nákupní košík podle pohlaví

	masový výrobek	mléčný výrobek	pečivo
muž (M)	40%	10%	50%
žena (Z)	10%	70%	20%

- ▶ chceme určit pohlaví nakupujícího podle nákupního košíku, v němž je
 - ▶ 1 položka masový výrobek
 - ▶ 1 položka mléčný výrobek

Naive Bayes Klasifikátor

Jednoduchý příklad 2

- ▶ předpokládejme nákupní košík podle pohlaví

	masový výrobek	mléčný výrobek	pečivo
muž (M)	40%	10%	50%
žena (Z)	10%	70%	20%

- ▶ chceme určit pohlaví nakupujícího podle nákupního košíku, v němž je
 - ▶ 1 položka masový výrobek
 - ▶ 1 položka mléčný výrobek
- ▶ přímá aplikace:

$$\begin{aligned}p(M|\text{maso, mleko}) &\propto p(\text{maso, mleko}|M)p(M) = \\ &= 0.4 \times 0.1 \times 0.5 = 0.02\end{aligned}$$

$$\begin{aligned}p(Z|\text{maso, mleko}) &\propto p(\text{maso, mleko}|Z)p(Z) = \\ &= 0.1 \times 0.7 \times 0.5 = 0.035\end{aligned}$$

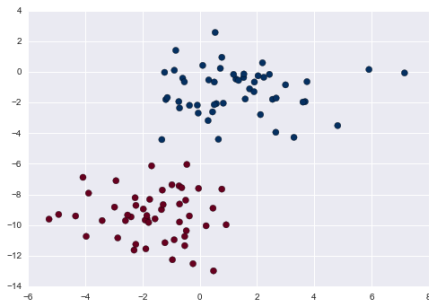
tzn. na 36% muž a na 64% žena.

Naive Bayes a další apriorní předpoklady

Apriorno

- ▶ pro třídy $p(C_k)$
- ▶ pro model $p(\mathbf{x}|C_k)$:
 - ▶ apriorní předpoklad normálního rozdělení: pro C_k mějme příslušné μ_k a σ_k^2 , pak

$$p(x_i|C_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right) \quad (6)$$

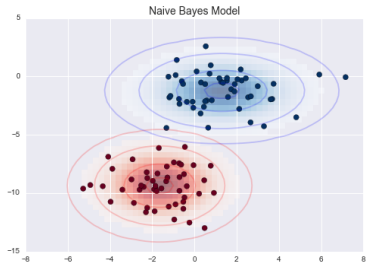
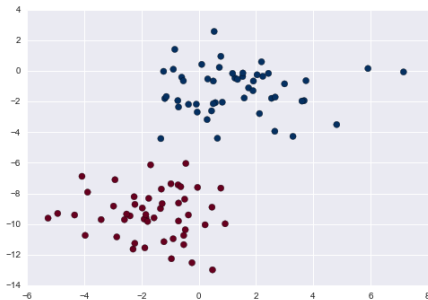


Naive Bayes a další apriorní předpoklady

Apriorno

- ▶ pro třídy $p(C_k)$
- ▶ pro model $p(\mathbf{x}|C_k)$:
 - ▶ apriorní předpoklad normálního rozdělení: pro C_k mějme příslušné μ_k a σ_k^2 , pak

$$p(x_i|C_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right) \quad (6)$$



Naive Bayes a další apriorní předpoklady

Apriorno

- ▶ multinomiální rozdělení: parametrizováno vektorem (pro k tou třídu)

$$\theta_{C_k} = (\theta_{C_k 1}, \dots, \theta_{C_k n}), \quad (7)$$

- ▶ n je počet příznaků (např. délka slovníku v klasifikaci dokumentů)
- ▶ $\theta_{C_k i}$ je pravděpodobnost $p(x_i | C_k)$, tzn. že se např. slovo objeví v dané třídě (pomocí počtu výskytů)
- ▶ Bernoulli rozdělení: jako multinomiální, ale výskyt ano-ne

$$p(x_i | C_k) = p(i | C_k) x_i + (1 - p(i | C_k))(1 - x_i). \quad (8)$$

- ▶ x_i je binární
- ▶ $p(i | C_k)$ je výskyt i tého slova

Naive Bayes

Spam filtering

- ▶ Jak poznat spam od legitimní zprávy?

Naive Bayes

Spam filtering

- ▶ Jak poznat spam od legitimní zprávy?
- ▶ náš příklad: podle slov, které obsahuje.
- ▶ označení: S - “spam”, H - “ham”.
- ▶ pravděpodobnost, že email obsahující slovo “viagra” je spam:

Naive Bayes

Spam filtering

- ▶ Jak poznat spam od legitimní zprávy?
- ▶ náš příklad: podle slov, které obsahuje.
- ▶ označení: S - “spam”, H - “ham”.
- ▶ pravděpodobnost, že email obsahující slovo “viagra” je spam:

$$p(S|\text{"viagra"}) = \frac{p(\text{"viagra"}|S) p(S)}{p(\text{"viagra"}|S) p(S) + p(\text{"viagra"}|H) p(H)}$$

- ▶ samozřejmě, jedno slovo nestačí...

Naive Bayes

Spam filtering

- ▶ pokud předpokládáme nezávislý výskyt N slov v textu:

$$S_1, \dots, S_N$$

- ▶ potom [Graham (2002): A Plan for Spam]:

$$p(S|S_1, \dots, S_N) = \frac{p(S|S_1) \dots p(S|S_N)}{p(S|S_1) \dots p(S|S_N) + (1 - p(S|S_1)) \dots (1 - p(S|S_N))}$$

- ▶ problém s neznámými nebo s velmi málo vyskytujícími se slovy

Naive Bayes

Spam filtering

- ▶ pokud předpokládáme nezávislý výskyt N slov v textu:

$$S_1, \dots, S_N$$

- ▶ potom [Graham (2002): A Plan for Spam]:

$$p(S|S_1, \dots, S_N) = \frac{p(S|S_1) \dots p(S|S_N)}{p(S|S_1) \dots p(S|S_N) + (1 - p(S|S_1)) \dots (1 - p(S|S_N))}$$

- ▶ problém s neznámými nebo s velmi málo vyskytujícími se slovy...můžeme:

- ▶ ignorovat
- ▶ korigovat pravděpodobnost

$$\tilde{p}(S| \text{"lihovina"}) = \frac{\check{c} \cdot p(S) + n \cdot p(S| \text{"lihovina"})}{\check{c} + n}, \quad (9)$$

kde \check{c} je zvolené číslo a n je počet výskytů během učení klasifikátoru.

Naive Bayes

Spam filtering - další možnosti

- ▶ zpracování jen daného počtu slov s největší informační hodnotou
- ▶ zpracování delších řetězců
- ▶ výhody:
 - ▶ relativně jednoduchý princip
 - ▶ pro každého uživatele jiné váhy a tedy specifické chování
- ▶ nevýhody:

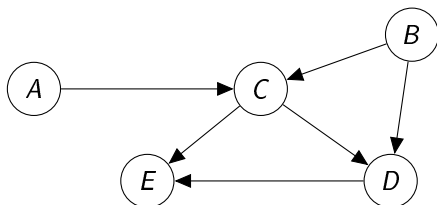
Naive Bayes

Spam filtering - další možnosti

- ▶ zpracování jen daného počtu slov s největší informační hodnotou
- ▶ zpracování delších řetězců
- ▶ výhody:
 - ▶ relativně jednoduchý princip
 - ▶ pro každého uživatele jiné váhy a tedy specifické chování
- ▶ nevýhody:
 - ▶ Bayesian poisoning (např. přidání spousty legitimních slov k textu spamu)
 - ▶ necitlivé na záměnu písmene ve slovu (“gooqle” přečteme stále jako “google”)
 - ▶ obrázky

Grafické modely

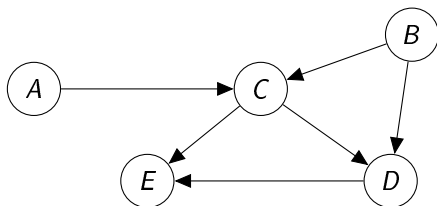
Proč? Jak?



- ▶ každý bod grafu reprezentuje náhodnou proměnnou
- ▶ hrany reprezentují statistickou závislost mezi proměnnými

Grafické modely

Proč? Jak?



- ▶ každý bod grafu reprezentuje náhodnou proměnnou
- ▶ hrany reprezentují statistickou závislost mezi proměnnými
- ▶ ekvivalentní zápis:

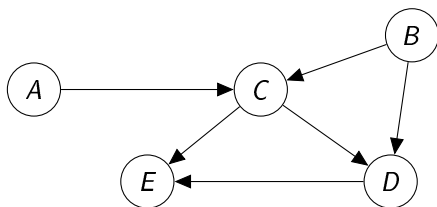
$$p(A, B, C, D, E) = p(A) p(B) p(C|A, B) p(D|B, C) p(E|C, D)$$

- ▶ formálně:

$$p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i | X_{\text{parents}(i)}) \quad (10)$$

Grafické modely

Proč? Jak?

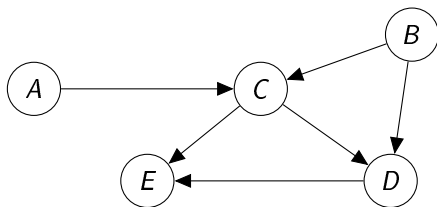


Jaké to má výhody?

- ▶ intuitivní reprezentace a vizualizace vztahů a vazeb mezi proměnnými
- ▶ vazby jsou přímo viditelné
 - ▶ např. otázka “Je A závislé na B (s předpokladem znalostí hodnot uzlu C)?”
- ▶ možnost tzv. message-passing přístupu
 - ▶ určení $p(A|C = c)$?

Grafické modely

Proč? Jak?



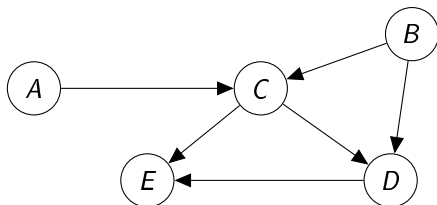
$$p(A|C=c) = \frac{p(A, C=c)}{p(C=c)} \quad (11)$$

► naivně:

$$p(A, C=c) = \sum_{B,D,E} p(A, B, C=c, D, E) \quad (12)$$

Grafické modely

Proč? Jak?



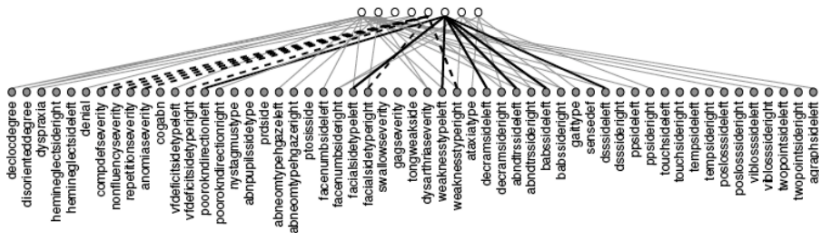
$$p(A|C=c) = \frac{p(A, C=c)}{p(C=c)} \quad (13)$$

► efektivně:

$$\begin{aligned} p(A, C=c) &= \sum_{B,D,E} p(A) p(B) p(C=c|A,B) p(D|B, C=c) p(E|C=c, D) \\ &= \sum_B p(A) p(B) p(C=c|A,B) \sum_D p(D|B, C=c) \sum_E p(E|C=c, D) \\ &= \sum_B p(A) p(B) p(C=c|A,B) \end{aligned}$$

Grafické modely

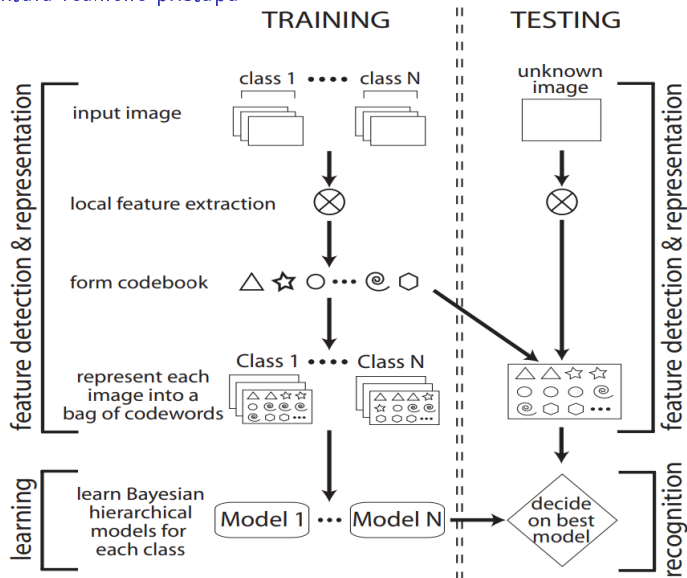
- ▶ rozpoznání a analýza řeči
- ▶ rozpoznání a analýza obrazu
- ▶ aplikace v medicíně (závislosti atd.)



[Ghahramani et al.]

Naive Bayes

Struktura reálného přístupu



[Fei-Fei et al.]

Grafické modely a hierarchické učení

Příklad na maticové dekompozici

- ▶ Mějme pozorování (data) D . Jak se “naučit” parametry θ z D ?
- ▶ Příklad na maticové dekompozici:

$$f(D|A, X, \omega) = \mathcal{N}\left(AX^T, \omega^{-1}I_p \otimes I_n\right),$$

$$f(\omega) = \mathcal{G}(\vartheta_0, \rho_0),$$

$$f(A) = \mathcal{N}(\mathbf{0}, I_p \otimes I_r),$$

$$f(X) = \mathcal{N}(\mathbf{0}, I_n \otimes I_r).$$

Grafické modely a hierarchické učení

Příklad na maticové dekompozici

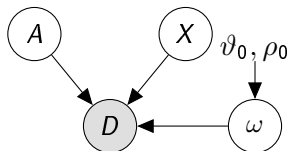
- ▶ Mějme pozorování (data) D . Jak se “naučit” parametry θ z D ?
- ▶ Příklad na maticové dekompozici:

$$f(D|A, X, \omega) = \mathcal{N}(AX^T, \omega^{-1}I_p \otimes I_n),$$

$$f(\omega) = \mathcal{G}(\vartheta_0, \rho_0),$$

$$f(A) = \mathcal{N}(\mathbf{0}, I_p \otimes I_r),$$

$$f(X) = \mathcal{N}(\mathbf{0}, I_n \otimes I_r).$$



Grafické modely a hierarchické učení

Příklad na maticové dekompozici

- Zkoumali jsme odhad počtu komponent - váhy v_1, \dots, v_r na obrázky

$$f(D|A, X, \omega) = \mathcal{N}\left(AX^T, \omega^{-1}I_p \otimes I_n\right),$$

$$f(\omega) = \mathcal{G}(\vartheta_0, \rho_0),$$

$$f(A|V) = \mathcal{N}(\mathbf{0}, I_p \otimes V^{-1})$$

$$f(v_k) = \mathcal{G}(\alpha_0, \beta_0),$$

$$f(X) = \mathcal{N}(\mathbf{0}, I_n \otimes I_r).$$

Grafické modely a hierarchické učení

Příklad na maticové dekompozici

- Zkoumali jsme odhad počtu komponent - váhy v_1, \dots, v_r na obrázky

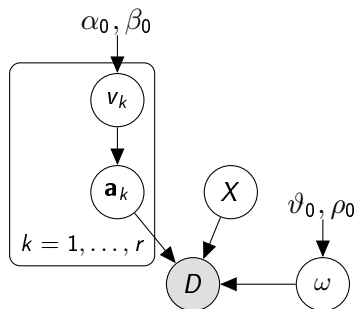
$$f(D|A, X, \omega) = \mathcal{N}(AX^T, \omega^{-1}I_p \otimes I_n),$$

$$f(\omega) = \mathcal{G}(\vartheta_0, \rho_0),$$

$$f(A|V) = \mathcal{N}(\mathbf{0}, I_p \otimes V^{-1})$$

$$f(v_k) = \mathcal{G}(\alpha_0, \beta_0),$$

$$f(X) = \mathcal{N}(\mathbf{0}, I_n \otimes I_r).$$



Grafické modely a hierarchické učení

Příklad na maticové dekompozici

- Říkali jsme si o konvolučním modelu křivek

$$f(D|A, X, \omega) = \mathcal{N}\left(AX^T, \omega^{-1}I_p \otimes I_n\right),$$

$$f(\omega) = \mathcal{G}(\vartheta_0, \rho_0),$$

$$f(A|V) = \mathcal{N}(\mathbf{0}, I_p \otimes V^{-1})$$

$$f(v_k) = \mathcal{G}(\alpha_0, \beta_0),$$

$$X = BU,$$

$$f(U) = \mathcal{N}(\mathbf{0}, I_n \otimes I_r),$$

$$f(b|\sigma) = \mathcal{N}(\mathbf{0}, \sigma^{-1}I_n),$$

$$f(\sigma) = \mathcal{G}(\zeta_0, \eta_0).$$

Grafické modely a hierarchické učení

Příklad na maticové dekompozici

- Říkali jsme si o konvolučním modelu křivek

$$f(D|A, X, \omega) = \mathcal{N}\left(AX^T, \omega^{-1}I_p \otimes I_n\right),$$

$$f(\omega) = \mathcal{G}(\vartheta_0, \rho_0),$$

$$f(A|V) = \mathcal{N}(\mathbf{0}, I_p \otimes V^{-1})$$

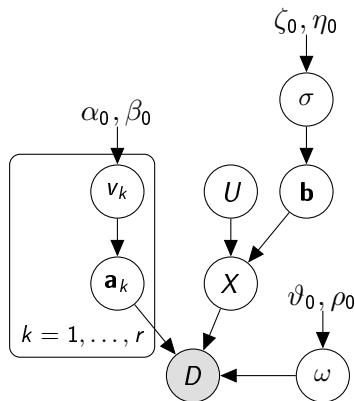
$$f(v_k) = \mathcal{G}(\alpha_0, \beta_0),$$

$$X = BU,$$

$$f(U) = \mathcal{N}(\mathbf{0}, I_n \otimes I_r),$$

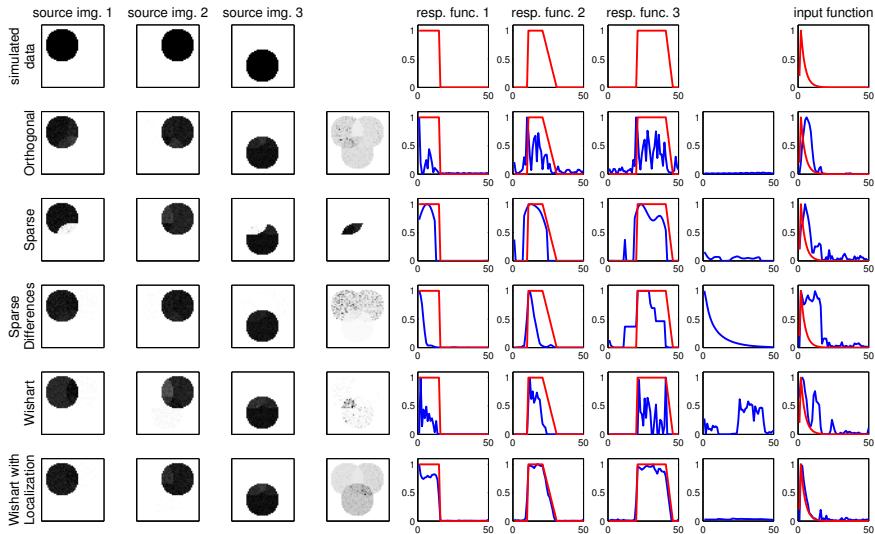
$$f(b|\sigma) = \mathcal{N}(\mathbf{0}, \sigma^{-1}I_n),$$

$$f(\sigma) = \mathcal{G}(\zeta_0, \eta_0).$$



Grafické modely a hierarchické učení

Aplikace na cvičná data



Grafické modely a hierarchické učení

Aplikace na data z renální scintigrafie

