



# ANALÝZA A KLASIFIKACE DAT



prof. Ing. Jiří Holčík, CSc.



INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

# V. KLASIFIKACE PODLE MINIMÁLNÍ VZDÁLENOSTI ∞ pokračování ∞

# PRAVDĚPODOBNOSTNÍ METRIKY

- ☑ používají kompletní informaci o struktuře klasifikačních tříd danou pomocí podmíněných hustot pravděpodobnosti  $p(x|\omega_1)$  a  $p(x|\omega_2)$ ;
- ☑ metriky tohoto typu splňují následující podmínky:
  1.  $J = 0$ , pokud jsou si hustoty pravděpodobnosti rovny, tj.  
$$p(x|\omega_1) = p(x|\omega_2);$$
  2.  $J \geq 0$
  3.  $J$  nabývá maxima, pokud jsou klasifikační třídy disjunktní, tj.  $p(x|\omega_1) = 0$  a  $p(x|\omega_2) \neq 0$  a naopak.

# PRAVDĚPODOBNOSTNÍ METRIKY

## základní myšlenka:

- ☑ klasifikační chyba:

$$e = \frac{1}{2} \left\{ 1 - \int |P(\omega_1|\mathbf{x}) - P(\omega_2|\mathbf{x})| \cdot p(\mathbf{x}) d\mathbf{x} \right\}$$

integrál v tomto vztahu

$$\begin{aligned} J_K &= \int |P(\omega_1|\mathbf{x}) - P(\omega_2|\mathbf{x})| \cdot p(\mathbf{x}) d\mathbf{x} = \\ &= \int |p(\mathbf{x}|\omega_1) \cdot P(\omega_1) - p(\mathbf{x}|\omega_2) \cdot P(\omega_2)| d\mathbf{x} \end{aligned}$$

se nazývá **Kolmogorovova variační vzdálenost**

Chyba bude maximální, když integrand bude nulový, tj. když obě váhované funkce hustoty pravděpodobnosti budou totožné. Naopak chyba bude nulová, pokud se obě hustoty nebudou překrývat  $\Rightarrow$  čím větší vzdálenost mezi třídami, tím je menší chyba klasifikace a naopak.

# PRAVDĚPODOBNOSTNÍ METRIKY

Podobně jsou definovány další pravděpodobnostní míry vzdálenosti obecně vztahem

$$J(\mathbf{x}) = \int f[p(\mathbf{x}|\omega_i), P(\omega_i), i = 1,2] d\mathbf{x}$$

Chceme, aby  $J(\mathbf{x})$  byla nezáporná funkce, pro kterou je  $J(\mathbf{x})=0$ , když jsou obě hustoty pravděpodobnosti totožné a je maximální, když se obě hustoty nepřekrývají.

# PRAVDĚPODOBNOSTNÍ METRIKY

Chernoffova vzdálenost

$$J_C = -\ln \int p^s(\mathbf{x}|\omega_1) p^{1-s}(\mathbf{x}|\omega_2) d\mathbf{x}, \quad s \in \langle 0;1 \rangle$$

Bhattacharyyova vzdálenost

$$J_B = -\ln \int \sqrt{p(\mathbf{x}|\omega_1) \cdot p(\mathbf{x}|\omega_2)} \cdot d\mathbf{x}$$

Divergence

$$J_D = \int [p(\mathbf{x}|\omega_1) - p(\mathbf{x}|\omega_2)] \cdot \ln \left( \frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} \right) \cdot d\mathbf{x}$$

Patrickova-Fisherova vzdálenost

$$J_{PF} = \sqrt{\int [p(\mathbf{x}|\omega_1) - p(\mathbf{x}|\omega_2)]^2 \cdot d\mathbf{x}}$$

# PRAVDĚPODOBNOSTNÍ METRIKY

resp. tzv. zprůměrněné verze, zahrnující i apriorní pravděpodobnost jednotlivých klasifikačních tříd

zprůměrněná Chernoffova vzdálenost

$$J_C = -\ln \int [p(\mathbf{x}|\omega_1)P(\omega_1)]^s [p(\mathbf{x}|\omega_2)P(\omega_2)]^{1-s} d\mathbf{x}, \quad s \in \langle 0;1 \rangle$$

zprůměrněná Bhattacharyyova vzdálenost

$$J_B = -\ln \int \sqrt{p(\mathbf{x}|\omega_1).P(\omega_1).p(\mathbf{x}|\omega_2)P(\omega_2)}.d\mathbf{x}$$

zprůměrněná divergence

$$J_D = \int [p(\mathbf{x}|\omega_1).P(\omega_1) - p(\mathbf{x}|\omega_2)P(\omega_2)] \ln \left( \frac{p(\mathbf{x}|\omega_1).P(\omega_1)}{p(\mathbf{x}|\omega_2).P(\omega_2)} \right) .d\mathbf{x}$$

zprůměrněná Patrickova-Fisherova vzdálenost

$$J_{PF} = \sqrt{\int [p(\mathbf{x}|\omega_1).P(\omega_1) - p(\mathbf{x}|\omega_2).P(\omega_2)]^2 .d\mathbf{x}}$$

# PRAVDĚPODOBNOSTNÍ METRIKY

uvedené výrazy se liší zejména pracností výpočtu a vazbou k hodnotám chybné klasifikace. Tato vazba je vyjádřena hodnotami  $D(\mathbf{x})$  a  $H(\mathbf{x})$  – dolním a horním odhadem pravděpodobnosti chybného zařazení.

$$H_C(\mathbf{x}) = \min_{0 \leq s \leq 1} J_C(s)$$

$$H_B(\mathbf{x}) = J_B$$

V případě, že známe dichotomické pravděpodobnostní míry a je třeba řešit problém klasifikace do více tříd, lze definovat kritérium, např. podle vztahu

$$J(\mathbf{x}) = \sum_{r=1}^R \sum_{q=r+1}^R P(\omega_r) \cdot P(\omega_q) \cdot J_{rq}(\mathbf{x})$$



# PRAVDĚPODOBNOSTNÍ METRIKY

základní nevýhodou pravděpodobnostních metrik je požadavek na znalost hustot pravděpodobnosti a jejich integrace (numerické?)

za určitých předpokladů o typu rozložení mohou být tyto vztahy integrovány analyticky

# PRAVDĚPODOBNOSTNÍ METRIKY

za předpokladu normálního rozložení ( $\mu_i$  jsou střední hodnoty a  $\Sigma_i$  kovarianční matice)

Chernoffova vzdálenost

$$J_C = \frac{1}{2}s(1-s)(\mu_2 - \mu_1)^T \Sigma_s^{-1}(\mu_2 - \mu_1) + \frac{1}{2} \ln \left( \frac{|\Sigma_s|}{|\Sigma_s|^{1-s} |\Sigma_s|^s} \right), \quad s \in \langle 0;1 \rangle$$

$$\text{kde } \Sigma_s = (1-s) \cdot \Sigma_1 + s \cdot \Sigma_2$$

Bhattacharyyova vzdálenost je pro  $s=0,5$

Divergence

$$J_D = \frac{1}{2}(\mu_2 - \mu_1)^T (\Sigma_1^{-1} + \Sigma_2^{-1})(\mu_2 - \mu_1) + \text{Tr}(\Sigma_1^{-1}\Sigma_2 + \Sigma_1\Sigma_2^{-1} - 2I)$$

Patrickova-Fisherova vzdálenost  $J_{PF} = \frac{1}{\sqrt{(2\pi)}} [ |2\Sigma_1|^{-\frac{1}{2}} + |2\Sigma_2|^{-\frac{1}{2}} -$

$$- 2|\Sigma_1 + \Sigma_2|^{-\frac{1}{2}} \cdot \exp \left\{ -\frac{1}{2}(\mu_2 - \mu_1)^T (\Sigma_1 + \Sigma_2)^{-1}(\mu_2 - \mu_1) \right\} ]$$

# PRAVDĚPODOBNOSTNÍ METRIKY

pokud  $\Sigma_1 = \Sigma_2 = \Sigma$ , pak se vztahy pro Bhattacharyyovu a divergenční vzdálenost zjednoduší na

$$J_M = J_D = 8J_B = (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1),$$

což je výraz pro Mahalanobisovu vzdálenost.

Příprava nových učebních materiálů  
oboru Matematická biologie

je podporována projektem ESF

č. CZ.1.07/2.2.00/07.0318

# „VÍCEOBOROVÁ INOVACE STUDIA MATEMATICKÉ BIOLOGIE“



INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ