



ANALÝZA A KLASIFIKACE DAT



prof. Ing. Jiří Holčík, CSc.



VIII. ANALÝZA HLAVNÍCH KOMPONENT



ZAČÍNÁME

ANALÝZA HLAVNÍCH KOMPONENT

PRINCIPAL COMPONENT ANALYSIS (PCA)

ROZKLAD PODLE VLASTNÍCH ČÍSEL

SINGULAR VALUE DECOMPOSITION (SVD)

Karhunenova-Loevova transformace

ZAČÍNÁME

- ☑ **extrakce příznaků** - hledání zobrazení (optimálního) Z , které transformuje původní m rozměrný prostor (obraz) na prostor (obraz) n rozměrný ($m \geq n$);
- ☑ **nalezení vhodné transformace** – potřeba optimalizačního kritéria:
 - obrazy v novém prostoru budou aproximovat původní obrazy ve smyslu minimální střední kvadratické odchylky;
 - obrazy v novém prostoru budou minimalizovat odhad pravděpodobnosti chyby

ZAČÍNÁME

- ☑ aby byla úloha řešitelná, hledáme zobrazení v oboru lineárních zobrazení

ZAČÍNÁME

- ☑ aby byla úloha řešitelná, hledáme zobrazení v oboru lineárních zobrazení

Jak poznáme lineární zobrazení?

ZAČÍNÁME

- ☑ aby byla úloha řešitelná, hledáme zobrazení v oboru lineárních zobrazení

Jak poznáme lineární zobrazení?

X = !

TEORIE

- ☑ předpokládejme, že je dáno K obrazů a nechť existuje m příznakových veličin, které tyto obrazy charakterizují. Tedy k -tý obraz je vyjádřen m rozměrným sloupcovým vektorem $\mathbf{y}_k \in \mathcal{Y}^m$, $k=1, \dots, K$.
- ☑ aproximujme nyní kterýkoliv obraz \mathbf{y}_k lineární kombinací n ortonormálních vektorů \mathbf{e}_i ($m \geq n$)

$$\mathbf{X}_k = \begin{matrix} & & n \\ & & \curvearrowright \\ & & \dots \\ & & \curvearrowleft \end{matrix}$$



TEORIE

- ☑ koeficienty c_{ki} lze považovat za velikost i -té souřadnice vektoru \mathbf{y}_k vyjádřeného v novém systému souřadnic s bází \mathbf{e}_i , $i=1,2,\dots,n$, tj. platí

$$\mathbf{y}_k = \sum_{i=1}^n c_{ki} \mathbf{e}_i$$

- ☑ použijeme-li jako kritérium minimální střední kvadratické odchylky, pak je

$$\hat{\varepsilon} = \|\mathbf{y}_k - \hat{\mathbf{y}}_k\|^2$$

TEORIE

- ☑ pak pomocí dříve uvedených vztahů pro \mathbf{x}_k a c_{ki} dostaneme

$$\hat{\varepsilon}_k = \|\mathbf{y}_k - \hat{\mathbf{y}}_k\|^2$$

- ☑ střední kvadratická odchylka pro všechny obrazy \mathbf{y}_k , $k=1, \dots, K$ je

$$\bar{\varepsilon} = \frac{1}{K} \sum_{k=1}^K \varepsilon_k = \frac{1}{K} \sum_{k=1}^K \|\mathbf{y}_k - \hat{\mathbf{y}}_k\|^2 = \frac{1}{K} \sum_{k=1}^K \mathbf{e}_k^T \left[\nabla \mathbf{y}_k \quad \mathbf{y}_k \right] \mathbf{e}_k$$

(je tedy závislá na volbě báze systému \mathbf{e}_i)

TEORIE

- ☑ diskrétní konečný rozvoj podle vztahu (☺) s báзовým systémem \mathbf{e}_i , optimálním podle kritéria minimální střední kvadratické chyby nazýváme diskrétní Karhunenův – Loevův rozvoj;
- ☑ aby střední kvadratická odchylka podle výše uvedeného vztahu byla minimální, musí být odečítaná hodnota na pravé straně rovnice maximální.

TEORIE

- ☑ musíme tedy maximalizovat výraz

$$\sum_{i=1}^n \mathbf{e}_i^T \mathbf{K}(\mathbf{y}) \mathbf{e}_i, \quad \text{kde } \mathbf{K}(\mathbf{y}) = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k \mathbf{y}_k^T$$

je autokorelační matice řádu m . Protože je symetrická a semidefinitní, jsou její vlastní čísla λ_i , $i=1, \dots, m$, reálná a nezáporná a vlastní vektory \mathbf{v}_i , jsou buď ortonormální, nebo je můžeme ortonormalizovat (v případě násobných vlastních čísel).

TEORIE

- ☑ uspořádáme-li vlastní čísla sestupně podle velikosti, tj.

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$$

a podle toho očíslovíme i odpovídající charakteristické vektory, lze dokázat, výe uvedený výraz dosahuje maxima, jestliže platí

$$\mathbf{e}_i = \mathbf{v}_i, i=1, \dots, n$$

a pro velikost maxima je

$$\max_{\mathbf{e}} \mathbf{e}^T \mathbf{K} \mathbf{e} = \sum_{i=1}^n \lambda_i$$

TEORIE

- ☑ pro minimální střední kvadratickou odchylku tedy platí

$$\begin{aligned} \varepsilon_{\min}^2 &= \sum_{k=1}^K \mathbf{y}_k^2 - \sum_{i=1}^m \lambda_i \\ &= \text{tr}(\mathbf{Y}) - \sum_{i=1}^m \lambda_i \end{aligned}$$

TEORIE

- ☑ v některých případech je vhodnější vektory \mathbf{y}_k před aproximací centrovat se střední hodnotou

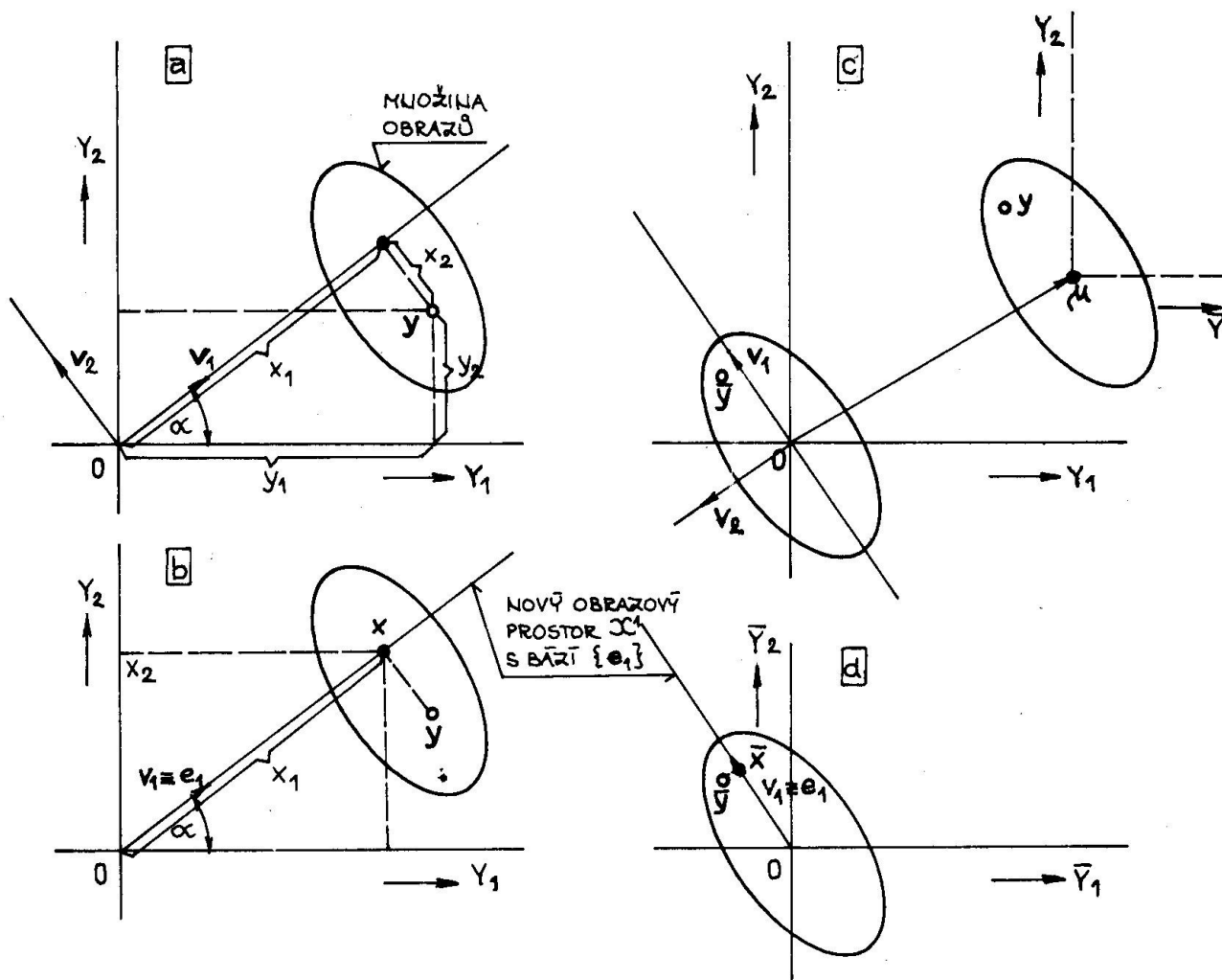
$$\bar{\boldsymbol{\mu}} = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_k$$

a místo s obrazem \mathbf{y}_k počítáme s jeho centrovanou verzí $\tilde{\mathbf{y}}_k = \mathbf{y}_k - \bar{\boldsymbol{\mu}}$.

Postup výpočtu se nemění, ale místo autokorelační matice používáme disperzní matici ve tvaru

$$D(\mathbf{y}) = \frac{1}{K} \sum_{k=1}^K \tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^T. \text{ Platí } D(\tilde{\mathbf{y}}) = D(\mathbf{y}) + \bar{\boldsymbol{\mu}} \bar{\boldsymbol{\mu}}^T$$

GEOMETRICKÁ INTERPRETACE



VLASTNOSTI

- ☑ při daném počtu n členů rozvoje poskytuje ze všech možných aproximací nejmenší střední kvadratickou odchylku;
- ☑ při použití disperzní matice jsou transformované souřadnice nekorelované; pokud se výskyt obrazů řídí normálním rozložením zajišťuje nekorelovanost i jejich nezávislost;
- ☑ vliv každého členu uspořádaného rozvoje se zmenšuje s jeho pořadím;
- ☑ změna požadavků na velikost střední kvadratické odchylky nevyžaduje přepočítávat celý rozvoj, nýbrž jen změnit počet jeho členů.

ROZDĚLENÍ DO TŘÍD

Jak se změní podmínky, když obrazy \mathbf{y} budou platit, které budou vymezeny jako části spojitého obrazového prostoru γ^m ?

- ✓ Výskyt obrazů v jednotlivých klasifikačních třídách bude popsán podmíněnými hustotami pravděpodobnosti $p(\mathbf{y}|\omega_r)$, $r=1,2,\dots,R$ a apriorní pravděpodobnost klasifikačních tříd bude $P(\omega_r)$.

V tom případě autokorelační matice bude

$$K(\mathbf{y}) = \sum_{r=1}^R P(\omega_r) \int_{\gamma^m} \mathbf{y} \cdot \mathbf{y}^T p(\mathbf{y}|\omega_r) d\mathbf{y} = \int_{\gamma^m} \mathbf{y} \cdot \mathbf{y}^T p(\mathbf{y}) d\mathbf{y}$$

ROZDĚLENÍ DO TŘÍD

- ☑ disperzní matice

$$D(\mathbf{y}) = \int_{\gamma^m} \mathbf{R}(\omega) \cdot (\mathbf{y} - \boldsymbol{\mu})^T (\mathbf{y} - \boldsymbol{\mu}) p(\mathbf{y} | \omega) d\mathbf{y}$$

kde

$$\boldsymbol{\mu} = \int_{\gamma^m} \mathbf{y} p(\mathbf{y} | \omega) d\mathbf{y}$$

nebo vztahem

$$\begin{aligned} D(\mathbf{y}) &= \int_{\gamma^m} \mathbf{R}(\omega) \cdot (\mathbf{y} - \boldsymbol{\mu})^T (\mathbf{y} - \boldsymbol{\mu}) p(\mathbf{y} | \omega) d\mathbf{y} \\ &= \int_{\gamma^m} (\mathbf{y} - \boldsymbol{\mu})^T (\mathbf{y} - \boldsymbol{\mu}) p(\mathbf{y}) d\mathbf{y} \end{aligned}$$

ROZDĚLENÍ DO TŘÍD

kde střední hodnota μ je vážený průměr středních hodnot všech tříd, tj.

$$\mu = \frac{1}{R} \sum_{r=1}^R P(\omega_r) \cdot \int_{\gamma_m} y \rho(y | \omega_r) dy = \int_{\gamma_m} y \rho(y) dy$$

