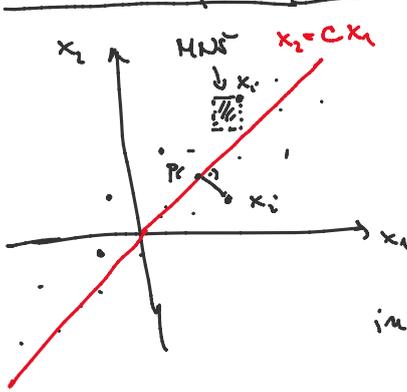


MNIST - učne napísané čísla 0-9
 .txt - 60 000 obrázkov : každý rozpredovaný do $28 \times 28 = 784$ hodnôt
 \Downarrow
 $n = 60\,000$ - odliene šeroj; \Downarrow
 $q = 784$
 - Keď každému obrázku je priradené, aké číslo to je

PCA ako optimálne kolmé projekcie

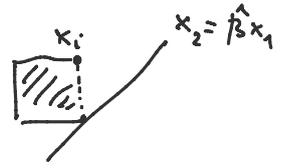


$x_1, \dots, x_n \in \mathbb{R}^2$ - centrovane! ($x_i - \bar{x}$)

$\boxed{q=2}$ - regresná priamka?

MNŠ: $x_2 = \hat{\beta} x_1$

$\hat{\beta} : \sum_{i=1}^n (x_{i,2} - \beta x_{i,1})^2 \rightarrow \min_{\beta}$



iná vzdialenosť: $\|x_i - p_i\|$

p_i - proj. x_i na priamku danú vektorom v

$p_i = \frac{v v^T}{v^T v} x_i$

$\Rightarrow \sum_{i=1}^n \|x_i - \frac{v v^T}{v^T v} x_i\|^2 \rightarrow \min_{v \in \mathbb{R}^2}$

$\boxed{q > 2}$ p_i - proj. na nadrovinu = stĺpcový priestor matice $V \in \mathbb{R}^{q \times q}$

$p_i = V(V^T V)^{-1} V^T x_i$

$\Rightarrow \sum_{i=1}^n \|x_i - V(V^T V)^{-1} V^T x_i\|^2 \rightarrow \min_{V \in \mathbb{R}^{q \times q}} \left(\begin{matrix} \text{ortogonálnu nadrovinu} \\ \text{chcem} \\ V\text{-reg.} \end{matrix} \right)$

distance:

$\sum_{i=1}^n \|x_i - V V^T x_i\|^2 \rightarrow \min_{V \in \mathbb{R}^{q \times q}} \left(\begin{matrix} V\text{-ortog. stĺpce} \end{matrix} \right)$

MNŠ:

$\hat{\beta} = (X^T X)^{-1} X^T y$

Riešenie: V : prvky q vl. vektorov matice $X^T X \approx S$

$S = \frac{1}{n-1} \sum (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{n-1} \sum x_i x_i^T = \frac{1}{n-1} \underbrace{\begin{pmatrix} x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nn} \end{pmatrix}}_{X^T} \underbrace{\begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}}_X = \frac{1}{n-1} X^T X$

$\Rightarrow V$ vrátane prvky q hlavných komponentov

Projektívne sledovanie - projection pursuit (PP)

- hľadanie štruktúry v dátach (grafickej)

→ zaujímavosť 1/2/3-rozmerné projekcie

- Kruskal (1969), Friedman & Tukey (1972)

↳ meno PP: "A projection pursuit algorithm for exploratory data analysis"

↳ často je to pripodobené im

- zvolíme index φ , ktorý vyjadruje mieru zaujímavosti projekcie

: keďže proj. bodovtina číslom φ

→ maximalizujeme tento index

⇒ PP: 1. zvolit' φ

2. nájsť maximum (loz. a glob.) cez všetky q -rozmerné projekcie

($q = 1/2/3$)

Projekčné indexy:

- je vhodné ak sú afínne invariantné: vzťahom na metáformu a škálovanie

→ rieši sa to: štandardizované dáta (unscaled): priemer 0_x

variancia I_x

↳ budeme to predpokladať $\left\{ \begin{array}{l} \text{populácie: } E(X) = 0, \text{ Var}(X) = I \\ \text{dátovo: } \bar{x} = 0, S = I \end{array} \right.$

- dosiahneme to: 1.) centrovanie: $\tilde{x}_i = x_i - \bar{x}$

2.) $S = \frac{1}{n-1} \tilde{X}^T \tilde{X} = U \Lambda U^T$: spät. rozklad

$$\Rightarrow S^{-1} = U \Lambda^{-1} U^T = \underbrace{U \Lambda^{-1/2}}_V \underbrace{\Lambda^{-1/2} U^T}_{V^T}$$

$$\tilde{x}_i = V^T \tilde{x}_i \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{maticovo: } \begin{pmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{pmatrix} = \begin{pmatrix} (V^T \tilde{x}_1)^T \\ \vdots \\ (V^T \tilde{x}_n)^T \end{pmatrix} = \begin{pmatrix} \tilde{x}_1^T V \\ \vdots \\ \tilde{x}_n^T V \end{pmatrix} = \tilde{X} V = \tilde{X} U \Lambda^{-1/2}$$

teoreticky: $\Sigma = \text{Var}(X) = U \Lambda U^T$

$$\Rightarrow \text{Var}(V^T X) = V^T \text{Var}(X) V =$$

$$= \Lambda^{-1/2} U^T \cdot U \Lambda U^T \cdot U \Lambda^{-1/2} = \Lambda^{-1/2} \Lambda \Lambda^{-1/2} = I$$

→ teda presnejšie $\tilde{X} = \tilde{X} U \Lambda^{-1/2}$ (tie budeme nazývať X)

$$\text{PCA: } Y = \tilde{X} U$$

\Rightarrow standardizované je faktický PCA, ale ještě potřebujeme pomocou $\Lambda^{-1/2}$ variancie na 1.

normálne dáta sú "nudné"

\rightarrow väčšina projekcií mnohorozmerných dát vyzerá približne normálne

\rightarrow normálne rozdelenie má ^{všetky} marginálne znomy normálne

\Rightarrow indexy často vyjadrujú mieru nenormality

\rightarrow chceme $\varphi(Y) = 0$ pre $Y \sim N()$

$$\varphi(Y) \geq 0 \quad \forall Y$$

ideálne: $\varphi(Y) > 0$ pre $Y \not\sim N$

1-rozmerný: ($q=1$): ^{teoretický} projekcia: $Y = h^T X$, $h \in \mathbb{R}^q$

\uparrow 1-rozmerný \uparrow q -rozmerný

$$Y = h_1 X_1 + \dots + h_q X_q$$

výberovo: $y_i = h^T x_i$, $i=1, \dots, n$

$$\underbrace{(y_1, \dots, y_n)}_{y^T} = (h^T x_1, \dots, h^T x_n) = h^T \underbrace{(x_1, \dots, x_n)}_{X^T}$$

$$y^T = h^T X^T$$

$$y = X h$$

φ : 1.) variancia: $\varphi(Y) = D(Y) = D(h^T X) \rightarrow \max_h$: to maximalitý je PC1

\Rightarrow pre toto φ je PP = PCA

2.) momenty

a) špicatost: $\varphi(Y) = |kurt(Y) - 3| \rightarrow \max$

$$kurt(Y) = E \left[\left(\frac{Y - E(Y)}{\sqrt{D(Y)}} \right)^4 \right] : \text{špicatost}$$

$$kurt(N(0,1)) = \underline{3}$$

• Standardizované dáta: $\mu_1 = E(Y) = 0$

$$\mu_2 = E(Y^2) = 1$$

- pracuje sa s kumulantami (nie s momentami) - niečo ako momenty

- pracuje sa s kumulantami (nie s momentami) - niečo ako momenty

↳ lebo kumulanty sú viateľky nulové pre $N(0,1)$

: κ_j - j-tý kumulant pre $\mu_1=0, \mu_2=1; \kappa_3=\mu_3$
: μ_j - j-tý moment: $E(Y^j)$ $\kappa_4=\mu_4-3$

⇒ maximalizujeme $\varphi(Y) = |\kappa_4| = |\mu_4 - 3|$

b) zložitejšie kombinácie: napr. $\varphi(Y) = \frac{\kappa_3(Y)^2}{12} + \frac{\kappa_4(Y)^2}{48} \rightarrow \max$

výberovo: odhady: $\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n Y_i^j$, alebo iný odhad momentov

⇒ $\hat{\kappa}_3 = \hat{\mu}_3, \hat{\kappa}_4 = \hat{\mu}_4 - 3$

- výpočtovo náročný, tradičný postup

- veľmi nerobustný: outliery prudko ovplyvňujú $\hat{\kappa}_4$

maximalizácia φ - nekondičná f., zvyšuje veľa lok. maxím ⇒ zložité

- treba použiť algoritmy, kt. sa nezaseďujú v lok. optime - napr.: genetický alg.,
particle swarm optimization,
...

→ Stoch. optim. metódy
(Hartman)

praktické použitie:

→ každé lok. max. môže predstavovať zaujímavú projekciu ↗ výhoda

→ spustíme alg. viackrát, zapamätáme si výsledky

→ vyberáme čo najrozdielnejšie proj. z tých získaných { výrazné pohľady $\varphi(Y)$
ubly