

# Principal component regression (PCR)

- met. extrakcie prem.: aj ako vstupy do prediktívnych modelov

↳ PCA pre lin. reg.: PCR ..., regulácia na hlavných komponentoch

data:  $y, X$

1.) PCA: vycentrujeme  $X \Rightarrow$  PCA:  $Z = XQ$   
 (+ vyčistíme)

- zoberieme prvých  $q$  hl. komp. (stĺpcov  $Z$ )  $\Rightarrow \tilde{Z} = \begin{bmatrix} | & | & & | \\ 1 & z_1 & \dots & z_q \\ | & | & & | \end{bmatrix}$

2.) lin. reg.:  $y = \tilde{Z}\beta + \varepsilon$  (\*)

$\rightarrow$  robíme s (\*)

• môže zabrániť preháňaniu (zharňuje sa výt. súmou)

$\rightarrow$  nepoznáme sa pri PCA na  $y \Rightarrow$  nemusí dobre fungovať

$\rightarrow$  CR: predpokladá, že smery  $X$  s najv. rozptýlením sú najviac informatívne pre  $y$

• voľba  $q$ : rozsah, val.

• výpočet (\*): jednoduchý, lebo stĺpce  $Z$  sú ortogonálne

• typicky  $X$  standardizujeme pred PCA

• vstiah s ridge: tá znižuje  $\beta$  viac pre nekoreš. PC

: PCR<sup>-1</sup> úplne vymaže  $\beta$  od niektorého PC

## Deo-Hitters

• PC 1: výška - dätini a starí

• PC 2: výška - menej dätini a starí

- SV: starí, JV: dätini

## Autoencoders (Autoassociative neural networks)

- autoassociatívne neuronové siete

### Neuronové siete (NS)

lin. reg.:  $y_i = x_i^T \beta + \varepsilon_i$

$$\hat{y}_i = x_i^T \hat{\beta}$$

pre nové  $x$ :  $\hat{y}(x) = x^T \hat{\beta}$

$$\hat{\beta} = \underset{\beta}{\text{min}} \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

$\rightarrow$  Machine learning pohľad:

(1) štruktúrna modeler:

vstup  $x \in \mathbb{R}^2 \rightarrow$  výstup:  $\hat{y}(x) = x^T \beta$  (pre dané  $\beta$ )

(2) fitovanie:

data:  $x_1, \dots, x_n \rightarrow \hat{y}_i = x_i^T \beta$

$y_1, \dots, y_n$

$$\Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - x_i^T \beta)^2 \rightarrow \underset{\beta}{\text{min}}$$

$\Rightarrow \beta^*$ : maturované parametre (valky)

vNS: komplikované (1.)

↳ pomocou  $x_1, \dots, x_n$  predpovedáme  $y_i$

(1.) model:  $\hat{y} = f_{\theta}(x)$   
 $\uparrow$  parametre (valky) NS

aktivácia  $\sigma(\cdot)$  (nem ... 0 ...)

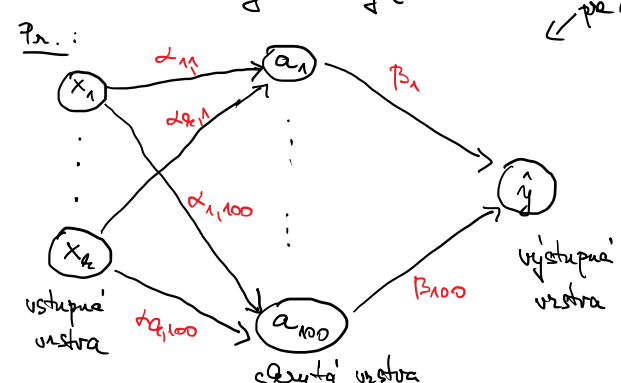
1. model:  $\hat{y} = f_{\theta}(x)$   
 - ako neuróny v mozgu  
 - parametre (váhy) NS  $\leftarrow \text{param. } \in \mathbb{R}$

aktivácia funkcia (niez lineárna)

$$a_1 = \varphi_1(\alpha_{1,0} + \alpha_{1,1}x_1 + \dots + \alpha_{q,1}x_q)$$

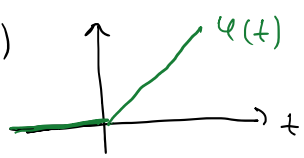
$$\vdots$$

$$a_{100} = \varphi_1(\alpha_{100,0} + \alpha_{100,1}x_1 + \dots + \alpha_{q,100}x_q)$$



$$f_{\theta}(x) = \hat{y} = H_2(\beta_0 + \beta_1 a_1 + \dots + \beta_{100} a_{100})$$

napr.  $\varphi(t) = \max(0, t)$   
 $\theta = (\alpha, \beta)$     ReLU



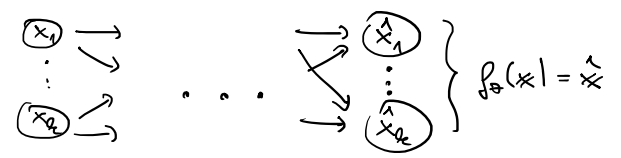
- šifry: ako neuróny poskytujú informáciu

2. klasifika: napr.  $x_1, \dots, x_n \in \mathbb{R}^q$   
 $y_1, \dots, y_n \in \mathbb{R}$   
 odhad:  $f_{\theta}(x_1), \dots, f_{\theta}(x_n)$

$$\Rightarrow \theta^* \in \arg \min_{\theta} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

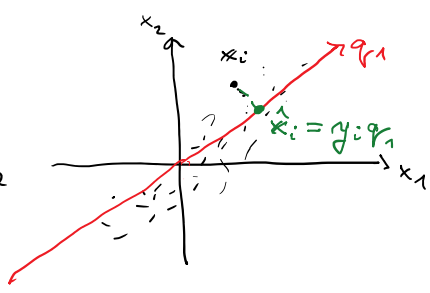
autoencoders:

1. pomocou  $x \in \mathbb{R}^q$  modeluje znova  $x \in \mathbb{R}^q$   
 -> "auto", "autoasociatívne"

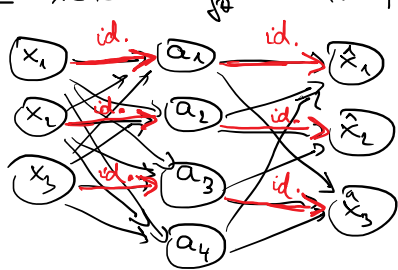


2.  $\hat{x}_i$  bližšie  $x_i$   
 $\Rightarrow$  chyba =  $\sum_{i=1}^n \underbrace{\|x_i - f_{\theta}(x_i)\|}_{\in \mathbb{R}^q}^2 \rightarrow \min_{\theta}$

• podobne PCA: vstup:  $x \in \mathbb{R}^q$   
 výstup:  $y \in \mathbb{R}^q$  ...  $\hat{x}$  bližšie  $x$ , a  $y$  sú jeho súradnice



Pr.:  $x \in \mathbb{R}^3 \mapsto f_{\theta}(x) \in \mathbb{R}^3$ , 3 vrstvy



optimálne:  $\hat{x}_1 = x_1$   
 $\hat{x}_2 = x_2$   
 $\hat{x}_3 = x_3$

$$a_1 = 1 \cdot x_1 + 0 \cdot x_2 + 0 \cdot x_3$$

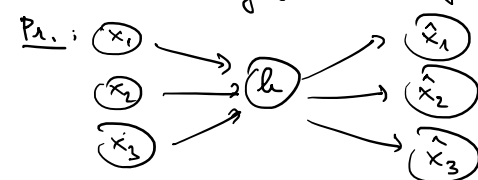
$$\hat{x}_1 = 1 \cdot a_1 + 0 \cdot a_2 + 0 \cdot a_3 + 0 \cdot a_4$$

atd.

-> zbytočné!

$\Rightarrow$  cieľ: zahaňiť, aby sa naučil identitu

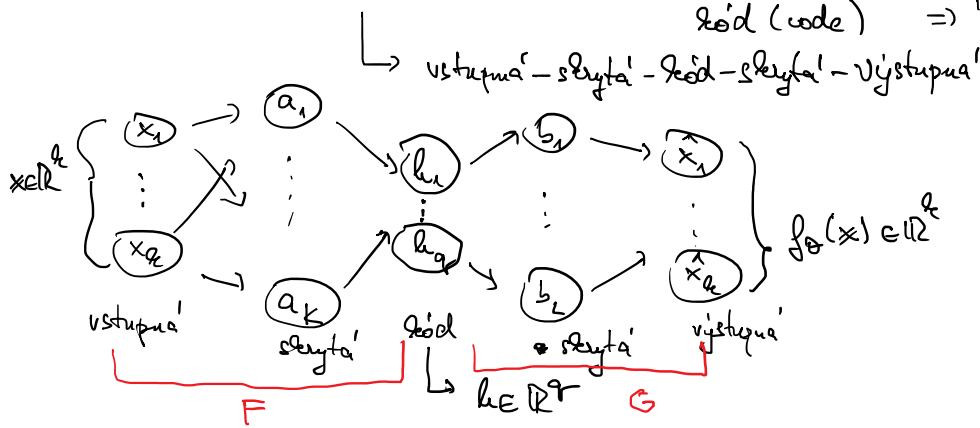
trick: v skrytej vrstve menej uzlov



• všeob.: skrytá vrstva:  $q_1 < q_2$  uzlov

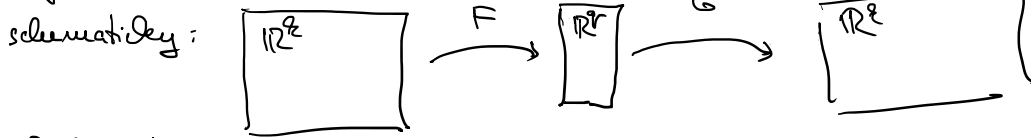
(3)  $\rightarrow (x_3)$

- všeob. : skrytá vrstva :  $q < p$  uzlov
- ak  $\varphi$  - lineárna (identita)  $\Rightarrow$  : autoencoder  $\Leftrightarrow$  PCA
  - $\uparrow$
  - výstupy  $\hat{x}_1, \dots, \hat{x}_n \in \mathbb{R}^q$  ležia v priestore princípu  $q$  l. vektorov  $S$
- ak iba 1 skrytá vrstva : v teórii by sme mohli dostať PCA
- $\Rightarrow$  autoencoder - nelineárna  $\varphi$ 
  - typický aspekt 5 vrstiev :  $\underbrace{\text{ta}}_{\text{kód (code)}}$  v strede ;  $q < p$  uzlov

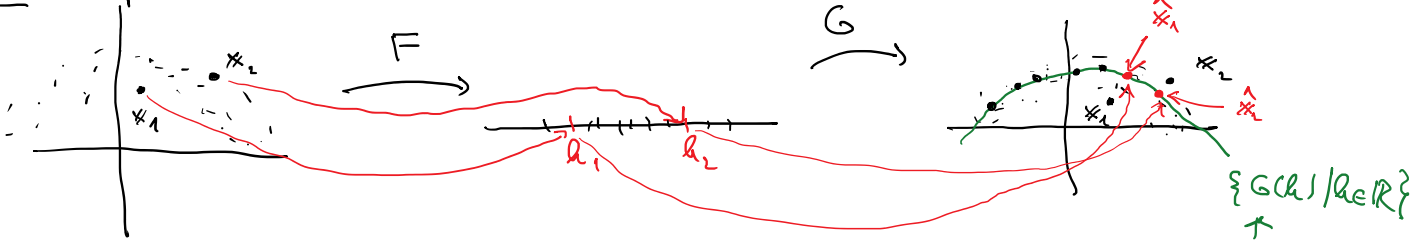


1.)  $h = F(x)$   $F: \mathbb{R}^p \rightarrow \mathbb{R}^q$   
 2.)  $\hat{x} = G(h)$   $G: \mathbb{R}^q \rightarrow \mathbb{R}^p$

$\left. \begin{matrix} 1.) \\ 2.) \end{matrix} \right\} f_{\hat{x}}(x) = G(F(x))$



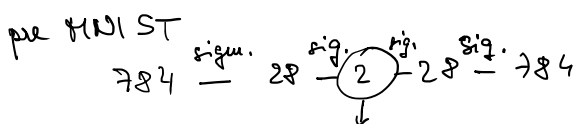
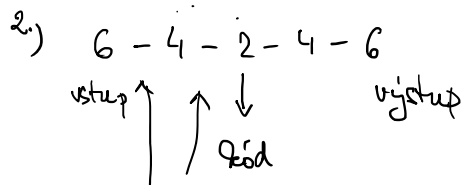
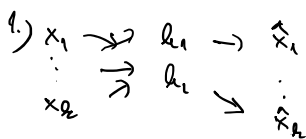
Pr.  $p=2, q=1$



- $q$  - rozmerová reprezentácia : výstupy kódu :  $(h_1, \dots, h_n) \in \mathbb{R}^q$
- autoencoder : "nelineárna PCA" - niekedy sa táž aj využije použitie : vyčistenie
  - : zabránenie pretrénovaniu
  - : odčumenie / vyčistenie dát, obrázkov

•  $\exists$  iné varianty

Príklad :



$\varphi(t) = \text{"sigmoid"} : \mathbb{R} \rightarrow [0, 1]$

↓  
Q'd