

Přednáška 24 - SVD, PCA & data analysis

Motivace: máme dataset → uložený obrázek / zdravotní data / data preference hudby/filmů/... / data z burzy/...

Problém → máme hudné dat & nevíme jak je analyzovat & & nelze je prakticky dlaně sklopit

- potřebujeme
 - kompresi dat, která zachová důležité rysy
 - analýzu dat → co z nich lze vycílit

→ ve spoustě případů máme dataset v tabulce:

12	12	203	12	12
12	12	203	12	12
203	203	203	203	203
12	12	203	12	12
12	12	203	12	12

odpovídá černobílému obrázku → čísla na škále

~ 0 až 255 odpovídají odstínům šedi ($0 \sim$ černá) ($255 \sim$ bílá)

=> tady máme „světlý kříž na tmavém poli“

	věk	váha	krevní tlak
pacient 1			
pacient 2			
pacient 3			
:			

	akcie A	akcie B	akcie C
simulace I			
simulace II			
simulace III			
:			

	věk	bydliště	příjem
uživatel 1			
uživatel 2			
uživatel 3			
:			

=> máme data v matici a chceme

- komprimovat velikost matice
- „analyzovat, co nám říká“

Hlavní nástroj: singulární rozklad (SVD) pro $A \in \mathbb{R}^{m \times n}$.

Linegebra 2: $\forall A \in \mathbb{R}^{m \times n} \exists U \in \mathbb{R}^{m \times m}, V \in \mathbb{R}^{n \times n}, \Sigma \in \mathbb{R}^{m \times n} :$

- U, V jsou unitární
- $\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \dots \\ & & \sigma_n \end{bmatrix}$ a $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$
- $A = U \Sigma V^T$

Komprese dat

$$A = U \Sigma V^T \quad \& \quad U = \begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_m \end{bmatrix}, \quad V = \begin{bmatrix} \vec{v}_1 & \dots & \vec{v}_n \end{bmatrix} \Rightarrow$$

$$\Rightarrow A = \sigma_1 \cdot \vec{u}_1 \vec{v}_1^T + \sigma_2 \cdot \vec{u}_2 \vec{v}_2^T + \dots + \sigma_n \cdot \vec{u}_n \vec{v}_n^T \quad (\text{tr. dyadický rozvoj } A)$$

kde

- σ_i jsou nezáporné a nerostoucí

- $\|\vec{u}_i \vec{v}_i^T\| = 1$... unitární matice

$$\Rightarrow \text{idea approximace } A: \quad A \approx \underbrace{\sigma_1 \vec{u}_1 \vec{v}_1^T + \dots + \sigma_r \vec{u}_r \vec{v}_r^T}_{A_r} =: A_r$$

memory cost $A \sim m \cdot n$

memory cost $A_r \sim r \cdot (m+n)$... (v praxi $\tilde{v}_k := \sigma_k \vec{v}_k$ můžeme ušetřit $r \cdot (m+n)$)

$$A_r = U_{:,1:r} \sum_{1:r,1:r} (V_{:,1:r})^T$$

Věta (Eckhart-Young-Firsky)

Platíme $A \in \mathbb{R}^{m \times n}$, $r \leq n \leq m$. Nechť $A = U \Sigma V^T$ je singulární rozklad.

Pak nejpřesnější approximace A matice hodnoty r je právě matice A_r dáná prvnimi r členy dyadického rozvoje A .

Můžeme zapsat jako $\forall X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r}$:

$$\|A - XY^T\| \geq \|A - A_r\|$$

$$\text{kde } A_r = U_{:,1:r} \sum_{1:r,1:r} (V_{:,1:r})^T = \sigma_1 \vec{u}_1 \vec{v}_1^T + \dots + \sigma_r \vec{u}_r \vec{v}_r^T.$$

matice hodnoty r jsou právě matice formy $X \cdot Y^T$ pro nejaké $X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r} \rightarrow$ tedy právě matice s memory cost $r \cdot (m+n)$

$$\text{Navíc: } \|A - A_r\|_2 = \sigma_{r+1} \quad \& \quad \|A - A_r\|_F = \sqrt{\sigma_{r+1}^2 + \sigma_{r+2}^2 + \dots + \sigma_{\min(m,n)}^2}$$

\Rightarrow spectru singulárního rozkladu (\equiv singular value decomposition = SVD)

jsme schopni majít nejpřesnější komprezi dat v

Eukleidovské normě z $m \cdot n$ dat na $r \cdot (m+n)$ dat.

Python Demo: image compression

Analýza dat

data = řádky matice = body $\vec{a}_i \in \mathbb{R}^n$ & máme jich $m \rightarrow i=1, \dots, m$

motivace pro $n=3$: • když budou $\vec{a}_1, \dots, \vec{a}_m$ všechny ležet na jedné přímce

\Rightarrow směrový vektor \vec{s} té přímky mi o těch datech hodně odpovídá

$$\Rightarrow \vec{a}_i = \alpha_i \cdot \vec{s} \Rightarrow \begin{bmatrix} -\vec{a}_1- \\ \vdots \\ -\vec{a}_m- \end{bmatrix} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix} \cdot \vec{s}^\top \quad \dots \text{rank } 1$$

\rightsquigarrow slovy: všechna data mají stejný „poměr“ mezi hodnotami „features“
 $(=\text{ratio})$ \Rightarrow liší se jen škalováním

nebo-li rozptyl v datech lze „vysvětlit“ pozorováním
 $(=\text{variance})$ pouze tohoto škalování

• když $\vec{a}_1, \dots, \vec{a}_m$ leží „okolo přímky“ \Rightarrow všechny rovnosti výše lze nahradit „ \approx “ a dostaneme

$$\Rightarrow \begin{bmatrix} -\vec{a}_1- \\ \vdots \\ -\vec{a}_m- \end{bmatrix} \approx \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix} \cdot \vec{s}^\top$$

\rightsquigarrow slovy: data mají podobný „poměr“ mezi hodnotami „features“ a tedy rozptyl / rozdíly v datech odpovídají pouze škalováním tohoto poměru

Pozorování 1: Pokud by data ležela přibližně v rovině dané vektory \vec{s}_1, \vec{s}_2

pak máme $\begin{bmatrix} -\vec{a}_1- \\ \vdots \\ -\vec{a}_m- \end{bmatrix} \approx \begin{bmatrix} \alpha_1 & \alpha_2 \\ \vdots & \vdots \\ \alpha_m & \alpha_m \end{bmatrix} \cdot \begin{bmatrix} -\vec{s}_1- \\ -\vec{s}_2- \end{bmatrix}$... rank-2 ... a opět platí, že

rozptyl v datech lze „přibližně vysvětlit“ pouze škalováním poměru hodnot prvků vektorů \vec{s}_1, \vec{s}_2 , tj. poměru konkrétních hodnot našich měřených „features“

Pozorování 2: V \mathbb{R}^n funguje stejná idea, jen „nejde vizualizovat“

Pozorování 3: Analyzovat data = majit co nejménší počet $\vec{s}_1, \dots, \vec{s}_r$, které mají fakt „vysvětlit“ ty data. Vektorům $\vec{s}_1, \dots, \vec{s}_r$ se říká „principal components (of the dataset)“.

Jejich malezení odpovídá malezení nejlepší rank-r approximaci matice A \Rightarrow stačí spočítat SVD! (tj. PCA = ^{principal component analysis})

Python Demo: Iris dataset

