

# PA081: Programování numerických výpočtů

## 8. Analýza hlavních komponent

Aleš Křenek a Jan Fousek

jaro 2013

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

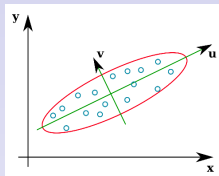
Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

# Úvodní příklad

- ▶ do zemské atmosféry proletěl velký meteorit
- ▶ ještě před dopadem explodoval a rozpadl se na mnoho úlomků
- ▶ na zemi najdeme jen část z nich



- ▶ je třeba určit
  - ▶ jak velký byl průměr shluku úlomků před dopadem
  - ▶ pod jakým úhlem meteorit dopadl
  - ▶ jakým směrem letěl

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ soubor opakovaného měření nějakého jevu apod.
  - ▶ pro každý prvek souboru více měřených hodnot –  $n$  rozměrný prostor
  - ▶ hledáme nezávislé (kolmé) směry zodpovědné za proměnlivost hodnot

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ soubor opakovaného měření nějakého jevu apod.
  - ▶ pro každý prvek souboru více měřených hodnot –  $n$  rozměrný prostor
  - ▶ hledáme nezávislé (kolmé) směry zodpovědné za proměnlivost hodnot
- ▶ v případě dopadu meteoritu
  - ▶ zaznamenáme souřadnice nálezu úlomků
  - ▶ směr největšího rozptylu určuje směr dopadu
  - ▶ rozptyl nejmenší ze všech směrů určuje poloměr shluku
  - ▶ poměr rozptylu ve směru dopadu a v kolmém směru odpovídá úhlu dopadu

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

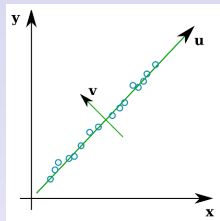
Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

# Analýza hlavních komponent

- ▶ Principal Component Analysis (PCA)
- ▶ hledá nezávislé **hlavní směry** zodpovědné za proměnlivost dat a určuje jejich relativní význam
- ▶ některé směry nemají žádný nebo minimální význam
  - ▶ více veličin než prvků souboru (např. pixely obrázku)
  - ▶ významné závislosti v celém jevu
  - ▶ PCA je dokáže identifikovat - redukce počtu dimenzí



## Úvod

Trocha statistiky

Rozpoznávání tváří

Domácí úkol

Ještě trocha statistiky

Zpracování dat fMRI

Shrnutí

# Trocha statistiky

- ▶ náhodná veličina, soubor vzorků, ...
  - ▶ viz základní kurzy statistiky
- ▶ **střední** (očekávaná) hodnota náhodné veličiny  $X$

$$EX = \bar{x} = \mu_x = \frac{1}{N} \sum_{i=1}^N x_i$$

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ náhodná veličina, soubor vzorků, ...

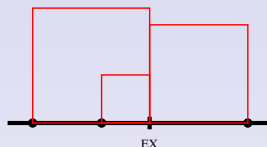
- ▶ viz základní kurzy statistiky

- ▶ **střední** (očekávaná) hodnota náhodné veličiny  $X$

$$EX = \bar{x} = \mu_x = \frac{1}{N} \sum_{i=1}^N x_i$$

- ▶ **rozptyl** (variance)

$$\text{var}(X) = \frac{1}{N-1} \sum_{i=1}^N (x_i - EX)^2$$



- ▶ vystihuje, jak moc jsou prvky souboru vzdáleny od střední hodnoty

Úvod

Trocha  
statistikyRozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistikyZpracování dat  
fMRI

Shrnutí

## ► kovariance

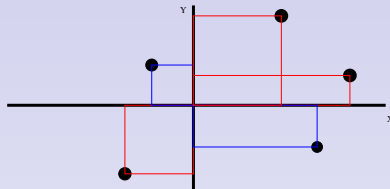
- zobecnění rozptylu na vzájemný vztah dvou veličin

$$\text{cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - EX)(y_i - EY)$$

- malý rozptyl v jedné veličině (rozložení vzorku podél druhé osy) kovarianci zmenšuje bez ohledu na chování druhé veličiny
- současný výskyt  $(x, y)$  a  $(x, -y)$  se ruší
- nejvyšší absolutní hodnota při lineární závislosti  $X = \pm Y$
- triviálně  $\text{cov}(X, X) = \text{var}(X)$  a  $\text{cov}(X, Y) = \text{cov}(Y, X)$



- ▶ jednoduchý příklad graficky



Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

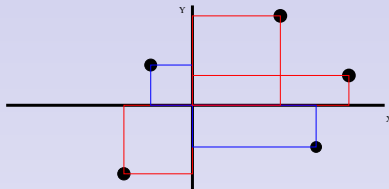
Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ jednoduchý příklad graficky



- ▶ matice kovariance veličin  $X_1, \dots, X_n$

$$\Sigma = \begin{pmatrix} \text{var}(X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & \text{var}(X_2) & \dots & \text{cov}(X_2, X_n) \\ & & \ddots & \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \dots & \text{var}(X_n) \end{pmatrix}$$

Úvod

Trocha  
statistikyRozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistikyZpracování dat  
fMRI

Shrnutí

- ▶ prvky analyzovaného souboru chápeme jako vektory v  $n$ -rozměrném prostoru
- ▶ hledáme souřadný systém, jehož bázevé vektory jsou hlavní směry daného souboru
- ▶ tedy kovariance mezi různými veličinami jsou nulové
- ▶ matice kovariance je diagonální

$$\Sigma = \mathbf{V}\Lambda\mathbf{V}^T$$

- ▶  $\Sigma$  je symetrická, proto jsou  $\lambda_i$  reálné a  $\mathbf{V}$  ortogonální
- ▶  $\mathbf{V}$  lze volit tak, že  $\lambda_i$  jsou uspořádány od největší

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ prvky souboru jsou normalizované fotografie
  - ▶ stejné rozlišení
  - ▶ obličej zabírá (přibližně) stejné pixely
  - ▶ pozadí je vždy stejné
  - ▶ černá a bílá odpovídá stejným hodnotám
- ▶ proměnné (složky vektoru) jsou jednotlivé pixely
  - ▶ tj. obrázek  $256 \times 256$  znamená 65536 proměnných
- ▶ datová matice  $\mathbf{D}_{N \times n}$  –  $N$  obrázků o  $n$  pixelech

- ▶ vycentrování dat
  - ▶ spočteme „průměrný obrázek“  $\mu$   
(průměr každého sloupce - pixelu)

$$\mu_j = \frac{1}{N} \sum_{k=0}^N d_{kj}$$

- ▶ odečteme od datové matice  $u_{ij} = d_{ij} - \mu_j$
- ▶ matice kovariance

$$\Sigma_{n \times n} = \frac{1}{N-1} \mathbf{U}^T \mathbf{U}$$

- ▶  $\mathbf{U}$  je centrovaná datová matice
  - ▶ jen maticový zápis definice kovariance

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

# Aplikace – rozpoznávání tváří

PA081:  
Programování  
numerických  
výpočtů

A. Křenek,  
J. Fousek

- ▶ výpočet vlastních hodnot  $\Sigma$
- ▶ přímočarý postup
  - ▶ použijeme implementaci pro symetrické matice,
  - ▶ aplikujeme na  $\frac{1}{N-1} \mathbf{U}^T \mathbf{U}$
  - ▶ zbytečně „nafouklá“ data ( $n \gg N$ )

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ výpočet vlastních hodnot  $\Sigma$
- ▶ přímočarý postup
  - ▶ použijeme implementaci pro symetrické matice,
  - ▶ aplikujeme na  $\frac{1}{N-1}\mathbf{U}^T\mathbf{U}$
  - ▶ zbytečně „nafouklá“ data ( $n \gg N$ )

- ▶ obrácený rozklad  $\mathbf{U}\mathbf{U}^T = \mathbf{V}'\Lambda'\mathbf{V}'^T$ 
  - ▶ vynásobení  $\mathbf{V}'$  zprava,  $\mathbf{U}^T$  zleva a ozávkování

$$\mathbf{U}^T\mathbf{U}(\mathbf{U}^T\mathbf{V}') = (\mathbf{U}^T\mathbf{V}')\Lambda'$$

- ▶ tedy  $\mathbf{U}^T\mathbf{V}'$  jsou vlastní vektory  $\Sigma = \mathbf{U}^T\mathbf{U}$
  - ▶  $\Lambda'$  jsou její vlastní hodnoty
  - ▶ ostatní vlastní hodnoty jsou nulové
- ▶ počítali jsme rozklad podstatně menší matice  $N \times N$

- ▶ vlastní vektory  $\Sigma$  jsou „typické tváře“
  - ▶ jsou vzájemně nezávislé
  - ▶ charakterizují nějaké rysy tváří
- ▶ ne všechny jsou stejně významné
  - ▶ přispívají k variabilitě vstupu úměrně velikosti vlastních hodnot
  - ▶ zvolíme (experimentálně) prahovou hodnotu  $0 < \sigma < 1$  (např. 0.95)
  - ▶ uvažujeme pouze  $\lambda_1, \dots, \lambda_m$  tak, že

$$\sum_{i=1}^m \lambda_i = \sigma \sum_{i=1}^N \lambda_i$$

- ▶ ostatní  $\lambda_{m+1}, \dots, \lambda_N$  zanedbáme
- ▶ typicky pro rozumně velký počet  $N$  platí  $m \ll N$



- ▶ každý obrázek  $\mathbf{p}$  lze reprezentovat souřadnicemi v prostoru PCA

$$\mathbf{p}' = \mathbf{V}^T (\mathbf{p} - \boldsymbol{\mu})$$

- ▶ přitom posledních  $n - m$  komponent  $\mathbf{p}'$  lze zanedbat
  - ▶ lze použít k výrazné kompresi
- ▶ rozpoznání neznámé tváře
  - ▶ spočteme  $\mathbf{p}'$  pro všechny vstupní obrázky
  - ▶ spočteme  $\mathbf{p}'$  pro neznámý vstup
  - ▶ hledání nejbližšího prvku provedeme pouze v  $m$  dimenzích

# Domácí úkol

- ▶ implementujte popsanou metodu rozpoznávání tváří
- ▶ použijte fotografie kamarádů/spolužáků z IS MU
  - ▶ jsou přijatelně uniformní
  - ▶ převed'te na černobílou, snažte se dosáhnout stejnoměrného jasu
- ▶ vyzkoušejte účinnost na jiných fotografiích stejných lidí
- ▶ zjistěte, jaké rysy tváře hlavní komponenty postihují
  - ▶ zvolíme libovolné hodnoty prvních  $m$  souřadnic,  $\mathbf{p}'$  doplníme nulami
  - ▶ provedeme zpětnou transformaci na obrázek

$$\mathbf{p} = \mathbf{V}\mathbf{p}' + \boldsymbol{\mu}$$

- ▶ pracujte i týmově, výsledky můžete předvést na některé z příštích přednášek

# Ještě trocha statistiky

- ▶ praktické nevýhody kovariance
  - ▶ závislost na absolutní hodnotě veličin
  - ▶ neintuitivní jednotky ( $m^2$  v úvodním příkladu, ale např. kgm, kdybychom uvažovali také hmotnost úlomku)

# Ještě trocha statistiky

- ▶ praktické nevýhody kovariance
  - ▶ závislost na absolutní hodnotě veličin
  - ▶ neintuitivní jednotky ( $m^2$  v úvodním příkladu, ale např. kgm, kdybychom uvažovali také hmotnost úlomku)
- ▶ zavádíme pojem **korelace**

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$$

- ▶ existuje více podob, toto je nejběžnější tzv. „Pearsonova“
  - ▶ nadále vypovídá o vztahu veličin
  - ▶ je normovaná do intervalu  $[-1, 1]$  a bez jednotek
  - ▶  $\pm 1$  - plná lineární závislost, 0 - zcela nezávislé
- ▶ PCA provádíme obdobně na **korelační matici**
  - ▶ numericky potenciálně lepší vlastnosti kvůli rovnoměrnému škálování
  - ▶ zpravidla stejně převáží nepoměr vlastních hodnot :-)

# Zpracování dat fMRI

- ▶ funkční magnetická rezonance
- ▶ 3D diagnostická metoda funkce mozku
  - ▶ založena na NMR – rezonance atomových jader v magnetickém poli
  - ▶ v různém chemickém kontextu různé spektrum
  - ▶ specificky fMRI detekuje množství okysličené krve
- ▶ výstupní data
  - ▶ 3D struktura voxelů, typicky desítky až stovky v jednom rozměru
  - ▶ v každém voxelu jedna skalární hodnota
  - ▶ série desítek až stovek takových snímků
- ▶ použití PCA
  - ▶ detekce vzorů chování bez apriorního zavedení nějakého modelu
  - ▶ upozorní na artefakty v datech, které by jinak zkreslovaly další analýzu
  - ▶ může napovědět, jaký funkční model zvolit

- ▶ zpracovávaná data
  - ▶ série desítek až stovek 3D snímků
  - ▶ každý desítky až stovky voxelů ve všech dimenzích
  - ▶ zpravidla ruční výběr zpracovávané oblasti

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ zpracovávaná data
  - ▶ série desítek až stovek 3D snímků
  - ▶ každý desítky až stovky voxelů ve všech dimenzích
  - ▶ zpravidla ruční výběr zpracovávané oblasti
- ▶ prostorová korelace
  - ▶ matice korelace z  $U^T U$  (odpovídá rozpoznávání tváří)
  - ▶ voxely jsou nezávislé náhodné veličiny
  - ▶ snímky chápeme jako jednotlivá pozorování
  - ▶ vysoká korelace – voxely se v čase chovají na sobě závisle
  - ▶ významnější vlastní vektory – reprezentativní snímky, z nichž dokážeme ostatní poskládat

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ časová korelace
  - ▶ (podstatně méně intuitivní)
  - ▶ matice korelace z  $UU^T$
  - ▶ body v čase jsou nezávislé náhodné veličiny
  - ▶ voxely jsou jednotlivá místa jejich pozorování
  - ▶ poloha bodu v prostoru vypovídá o jeho chování v čase
  - ▶ významné vlastní vektory – reprezentativní vzory chování voxelů



- ▶ série 120 snímků mozku
- ▶ nejvýznamnější komponenty časové korelace

Úvod

Trocha  
statistiky

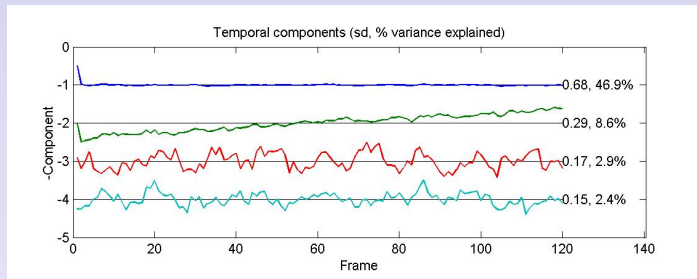
Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

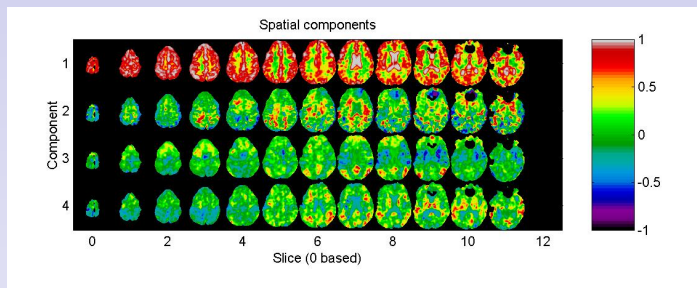
Shrnutí



- ▶ první zachycuje zjevně artefakt na začátku snímání

# Zpracování dat fMRI

- ▶ stejný efekt je patrný na prostorových hlavních komponentách



- ▶ červené oblasti označují vzájemně závislé voxely
  - ▶ sklon chovat se v čase stejně
  - ▶ představte si na 3 voxelech, rozmístění 120 bodů odpovídajících snímkům v prostoru, a vektorech hlavních komponent
- ▶ 1. komponenta: artefakt měření je v krajích a ve svíslé ose

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

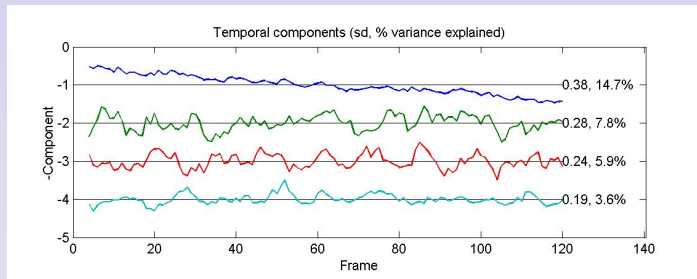
Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ ze souboru odstraníme prvních několik snímků



- ▶ první komponenta ukazuje už jasný lineární trend
- ▶ další jsou víceméně pravidelné oscilace

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

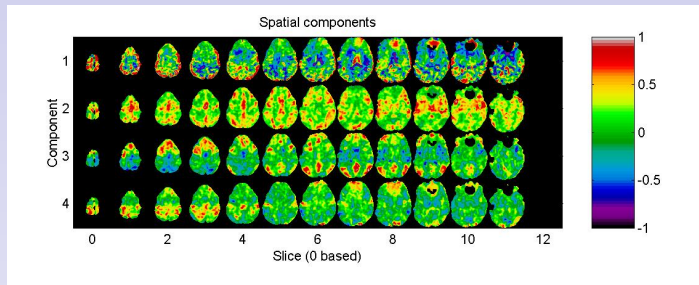
Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ hlavní komponenty prostorových korelací



- ▶ červená (resp. modrá) znamená voxely, které definují danou komponentu
- ▶ variabilita v čase v konkrétním hlavním směru je dána těmito voxely
- ▶ směry jsou vzájemně nezávislé
  - ▶ každý voxel je výrazně modrý nebo červený jen v jednom řádku

Úvod

Trocha  
statistiky

Rozpoznávání  
tváří

Domácí úkol

Ještě trocha  
statistiky

Zpracování dat  
fMRI

Shrnutí

- ▶ analýza hlavních komponent
- ▶ nalezení vzájemně nezávislých směrů variability vzorku
- ▶ identifikace jejich významu - snížení počtu dimenzí
- ▶ aplikace ve zpracování obrazu
- ▶ 2D - rozpoznávání tváří
- ▶ 3D - fMRI scany mozku