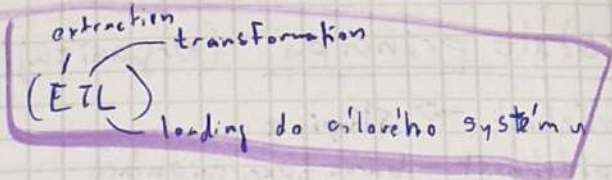


Data science



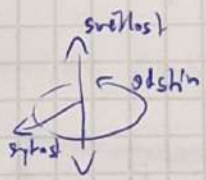
- 1) sečtení a přípravě dat - data wrangling (munging)
- 2) na data se koukáme z různých úhlů, filtrujeme, počítáme
→ trendy, ...
- 3) grafický design, prezentace → efektivně, jasně
- odhalíme strukturu
- 4) interakce s daty - InfoVis, HCI ^{interaction}
- storytelling with data

Exploratorní (EDA) X X

- mezi námi a daty

EXplanatorní vizualizace

- mezi daty a publikem



světlost
obskln + sykost
náštraje, jak něco zvýraznit

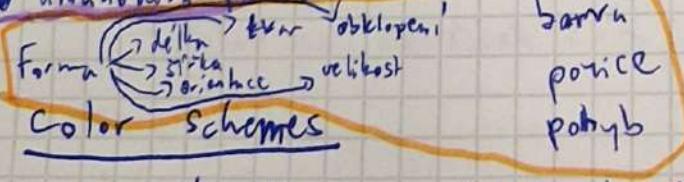
20% , 30% zbyte smysly
Sakadický pohyb očí - nespojitý

Obrazová paměť - předvědomé zpracování

Krátkodobá paměť - vědomé zpracování

Dlouhodobá paměť

barva, velikost, orientace, lokace



Vizuální chunky

= spojení objektů na základě předvědomých atributů

dokážeme si paměť max. 3-4 chunky (v krátkodobé paměti)
- při opakování se dostanou do dlouhodobé paměti

- monochromatic - různé intenzity stejného odstínu
- analogous - vedle sebe
- complementary
- triadic - rovnostranný trojúhelník

→ storytelling je o vztazích těchto

Kvantitativní X X

- např. velikost
→ měřitelné

kategoričké rozdíly

- např. barva, ID, číslo účtu, číslo dveří
- pozorovatelné, ale neměřitelné (nebo spíš je s nimi těžké dělat aritmu)
- nominalní - bez lin. usp.
- ordinální - usp., ale nekvantifikovatelné (např. velikost zvlášť)
- hier. intervalové

Gestalt principy (organizující principy)

- blízkost → skupina
- podobnost - barva, tvar, orientace
- enclosure (ohrada) - co je ohraničené, tvoří jeden celek
- closure (spojení, uzavření) - \square \cup
- continuity
- spojení ⇒ skupina

Efekty kontextu

- víc vnímáme rozdíly než abs. hodn.

Statistiky

- průměr
- medián
- módus

average

$$\text{midrange} = \frac{\text{min} + \text{max}}{2}$$

• stl. odch.

variation

$$\text{spread} = \text{max} - \text{min}$$

↳ ovlivněno extrémny !!

• korelace

-1 → neg. korelace

0 →

1 → poz. korelace

Tabulka X Graf

↓
pokročilejší
čistší
jednodušší
přesnější
hodnoty,
porovnání
je, je
více vizuální
agregace
(Σ, avg, ...)

↓
chceme
ukázat
pattern,
outliery
↑
jako celek,
porovnávání
vůči celku

Kvantitativní hodnoty

- body \Rightarrow scatter plot
- úsečky \Rightarrow $\left\{ \begin{array}{l} \text{spojení bodů} \Rightarrow \text{vzorce} \\ \text{best fits} \Rightarrow \text{trendy} \end{array} \right.$
- bars
- boxes \Rightarrow distribuce - box plot
- plocha \rightarrow pie chart
 \rightarrow bubble chart
- barva \rightarrow bubble plot s různými intenzitami nebo odstíny

Kategorické hodnoty

- pozice
- odstín
- tvar značky
- fill pattern \rightarrow pozice \rightarrow Moiré effect (že to dělá iluzi takového vlnění)
- typ čáry

2 Distribuce

box ploty \rightarrow

- \rightarrow median,
- min, max,
- 1. a 3. kvartil

histogramy \rightarrow

\rightarrow chlívěčky

density histogramy \rightarrow

\rightarrow hodiny procenta
z celku

kernel density \rightarrow

Function (KDF)

\rightarrow místo chlívěček už jednotlivé hodnoty, on kde se mají a poselují se

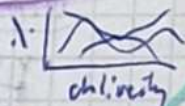
\rightarrow kum. distrib. funkce
 \rightarrow \neq hodnota

QQ plot

\rightarrow kvantily proti sobě

Ještě víc distribucí

• freq. polygon



chlívěčky

• ridgeline plot

- prostě víc těch KDF



• stacked density chart

• candlestick chart

• violin plot

- KDF \cdot 90° \rightarrow zrcadleno



mnoho barev s histogramem \rightarrow

Nejistota



Parallel coords



Polar chart



Circle tree



Lhavi

Lie Factor = $\frac{\text{velikost efektu ve vizualizaci}}{\text{velikost efektu v datech}}$

< 0.95 = 0 > 1.05 => FUJ

Sunburst

-> hierarchie - listy na obvodu



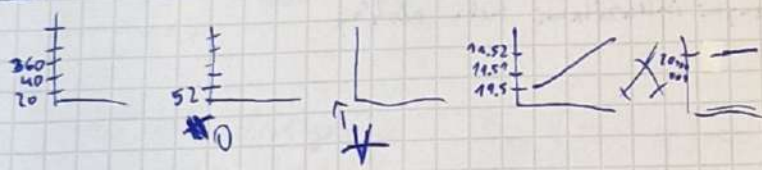
Circle packing



Sankey diagram

=> tok dat

Manipulace s Y



Parallel sets



Waterfall chart

- bar zprava
- bar vlevo
- predchozi stav
=> zmeny

Gantt chart

-> "rozk" (tasks)

Sprinkles

-> data words
-> "text a minigraf"

Slopegraph

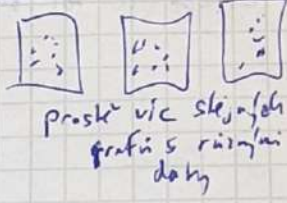


Taylor cloud

Jointplot



Small multiple



proste vic stejnych grafu s ruznymi daty

Stem-and-leaf plot

"histogram" ale v ch. urovn. c. ch
(je sou ty hodnoty)
-> tabulka

Arc diagram



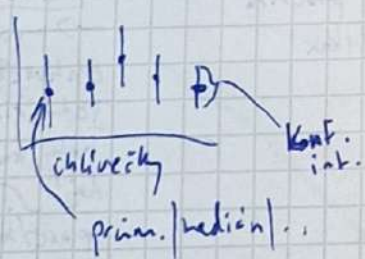
-> 3D: arc map
- takove chobotnicky

-> kdyz spojime konce do kruhu: radial chart



muze byt i cast hierarchical edge bundling

Binned scatterplot



ch. urovn. prim. median...
kont. int.

Tree map



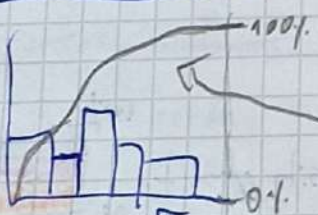
=> hierarchie / cast celku

Heatmap

-> pomeri oblastu kvantitativni atributu
- korelace

- Easte s dendrogramem
-> sklucuje podobne objekty

Pareto chart



- proste neco - treba bar chart + kum. % nebo prubehny soucet
- treba line chart

- correlogram = kdyz udelame heatmapu korelaci vic atributu

Hive plots



BioFabric



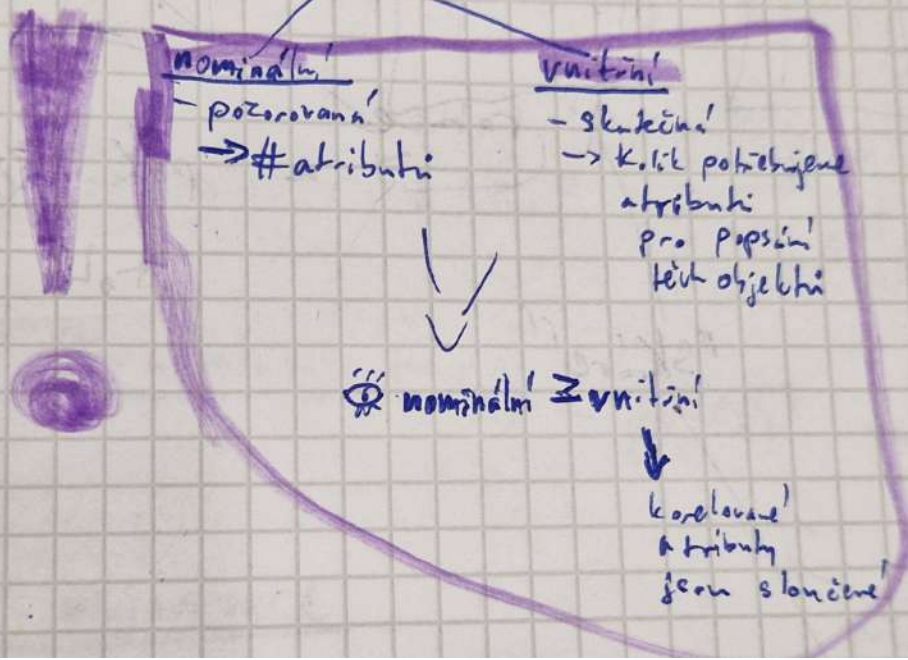
vrchol

Bag plot



3x naf.
konu: dal v horn 3x naf.

Redukce dimenze



PCA

přít. uí vycentrované

$$\text{var}(X) = \frac{\sum_{i=1}^n (X_i - E[X])^2}{n} = \frac{\sum_{i=1}^n X_i^2}{n} = E[(X - E[X])^2] = E[X^2]$$

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - E[X])(Y_i - E[Y])}{n} \rightsquigarrow \text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{cov}(X, X)\text{cov}(Y, Y)}}$$

kov. matice $\Sigma = XX^T$
 pro X je to XX^T

1. PC: $Y_1 = a_1^T X$ t.j. $\text{var}(Y_1)$ je max. při $a_1^T a_1 = 1$

$$\text{var}(Y_1) = \text{var}(a_1^T X) = E[(a_1^T X)(a_1^T X)^T] = E[a_1^T X X^T a_1] = a_1^T E[\Sigma] a_1 = a_1^T \Sigma a_1$$

Lagrangeovy multiplikatory - max. $f(x, y)$ při $g(x, y) = c$

$$\Delta(x, y, \lambda) = f(x, y) - \lambda(g(x, y) - c) \leftarrow \text{chceme max.}$$

$$\hookrightarrow \text{tedy } \Delta(a_1, \lambda) = a_1^T \Sigma a_1 - \lambda(a_1^T a_1 - 1)$$

$$\text{chceme: } \frac{\partial \Delta(a_1, \lambda)}{\partial a_1} = 0$$

$$2\Sigma a_1 - 2\lambda a_1$$

$\Sigma a_1 = \lambda a_1 \Rightarrow a_1$ je vl. vektor Σ ,
 λ je k němu přísl. vl. číslo

$$\text{var}(Y_1) = a_1^T \Sigma a_1 = \lambda \underbrace{a_1^T a_1}_1 = \lambda \Rightarrow \lambda = \lambda_1$$

2. PC: $Y_2 = a_2^T X$, s mx. $\text{Var}(Y_2)$ pri $a_2^T a_2 = 1$
 $\wedge \text{cov}(Y_2, Y_1) = 0$

$$\Lambda(a_2, \lambda, \kappa) = a_2^T \Sigma a_2 - \lambda(a_2^T a_2 - 1) - \kappa(a_2^T a_1)$$

$$\frac{d}{da_2} = 0$$

$$0 = a_2^T \Sigma a_1 = a_2^T \lambda_1 a_1 = \lambda_1 a_2^T a_1 \rightarrow a_2^T a_1 = 0$$

$$2 \Sigma a_2 - 2 \lambda a_2 - \kappa a_1 = 0$$

$$\left| \cdot \frac{a_1^T}{2} \right.$$

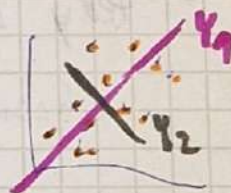
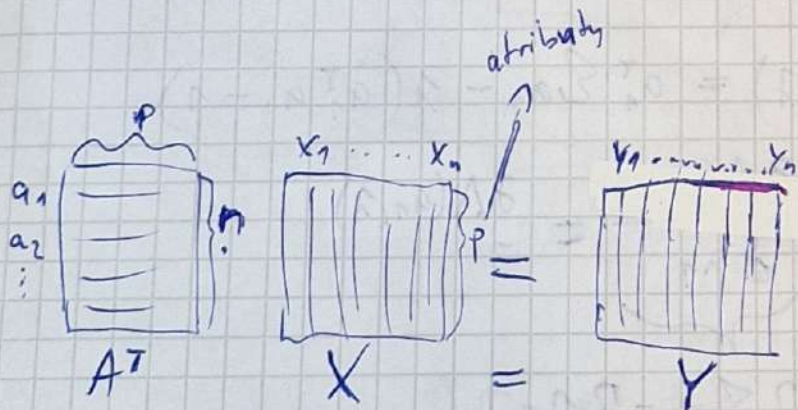
$$\underbrace{a_1^T \Sigma a_2}_0 - \underbrace{\lambda a_1^T a_2}_0 - \underbrace{\frac{\kappa}{2} a_1^T a_1}_1 = 0$$

$$\Rightarrow \kappa = 0$$

$$\Rightarrow \Sigma a_2 = \lambda a_2 \rightarrow \lambda := \lambda_2$$

k. PC: $Y_k = a_k^T X$ s mx. $\text{Var}(Y_k)$ pri $a_k^T a_k = 1$

$\wedge \forall l < k: \text{cov}(Y_k, Y_l) = 0$



loadings

↳ loadings plot pro. PCi:

score plot

λ_i

kolik variance nose i-th PC

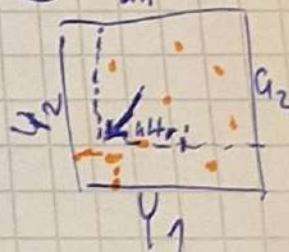
$$\sum_{j=1}^p \lambda_j$$

screen plot

variance

PC

biplot



MDS \rightarrow vzdálenosti \rightsquigarrow souřadnice

Klasické - vzdálenosti jsou Euklidovské

$$B := X^T X$$

matice vzdáleností D , souřadnice (neznáme) X \rightarrow předp. centrovane
ve sloupcích

$$D_{ij}^2 = (x_i - x_j)^T (x_i - x_j) = x_i^T x_i - x_j^T x_i - x_i^T x_j + x_j^T x_j =$$

$$= \langle x_i, x_i \rangle + \langle x_j, x_j \rangle - 2 \underbrace{x_i^T x_j}_{B_{ij}}$$

$$m := (\langle x_1, x_1 \rangle, \dots, \langle x_n, x_n \rangle)$$

$$D^2 = m \mathbb{1}_n^T + \mathbb{1}_n m^T - 2B$$

$$C := I_n - \frac{1}{n} \mathbb{1}_n \mathbb{1}_n^T = \begin{pmatrix} 1 & 0 \\ 0 & \ddots & 0 \\ 0 & \dots & 1 \end{pmatrix} - \begin{pmatrix} \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \\ \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} \end{pmatrix} \dots \text{kontrovací matice}$$

$$\text{👁} - \frac{1}{2} C D^2 C = -\frac{1}{2} C (m \mathbb{1}_n^T + \mathbb{1}_n m^T - 2B) C$$

$$= -\frac{1}{2} C m \mathbb{1}_n^T C - \frac{1}{2} C \mathbb{1}_n m^T C + C B C$$

$$\text{👁} \mathbb{1}_n^T C = \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \dots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \dots & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \dots & 1 - \frac{1}{n} \end{pmatrix} = \begin{pmatrix} 0 & 0 & \dots & 0 \end{pmatrix} = \mathbb{0}_n$$

$$\underbrace{\frac{1 - \frac{1}{n}}{n-1}}_{\substack{\text{---} \\ (n-1)\text{-krát}}} \underbrace{- \frac{1}{n} - \frac{1}{n} \dots}_{\substack{\text{---} \\ (n-1)\text{-krát}}} = \frac{(n-1) - (n-1)}{n} = 0$$

$$= C B C = B \quad (\text{předpokládali jsme už vycentrovane})$$

$$X^T X = B = Q \Lambda Q^T = \underbrace{(Q \Lambda^{\frac{1}{2}})}_{X^T} \underbrace{(\Lambda^{\frac{1}{2}} Q^T)}_X$$

$$GOF = \frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^n \lambda_i}$$

a z nich jen prvních m

- \Rightarrow Alg:
- 1) $D \rightarrow D^2$
 - 2) $B := -\frac{1}{2} C D^2 C$ \uparrow (vezmeme jen nenulové)
 - 3) spektrální rozklad $B = Q \Lambda Q^T$

(ne)metrické - máme δ_{ij} ... dissimilarities # bodů \nearrow dim. prostoru
 když už máme nějakou konfiguraci X - matici $n \times m$

$$D_{ij}(X) := \sqrt{\sum_{a=1}^m (x_{ia} - x_{ja})^2}$$

\hookrightarrow chyba je $\sigma^2(X) := \sum_{i < j} (D_{ij}(X) - \delta_{ij})^2$

lepe: raw stress $\sigma_r^2(X) := \sum_{i < j} w_{ij} (D_{ij}(X) - \delta_{ij})^2$
 \downarrow
 váhy - kvůli chybějícím datům

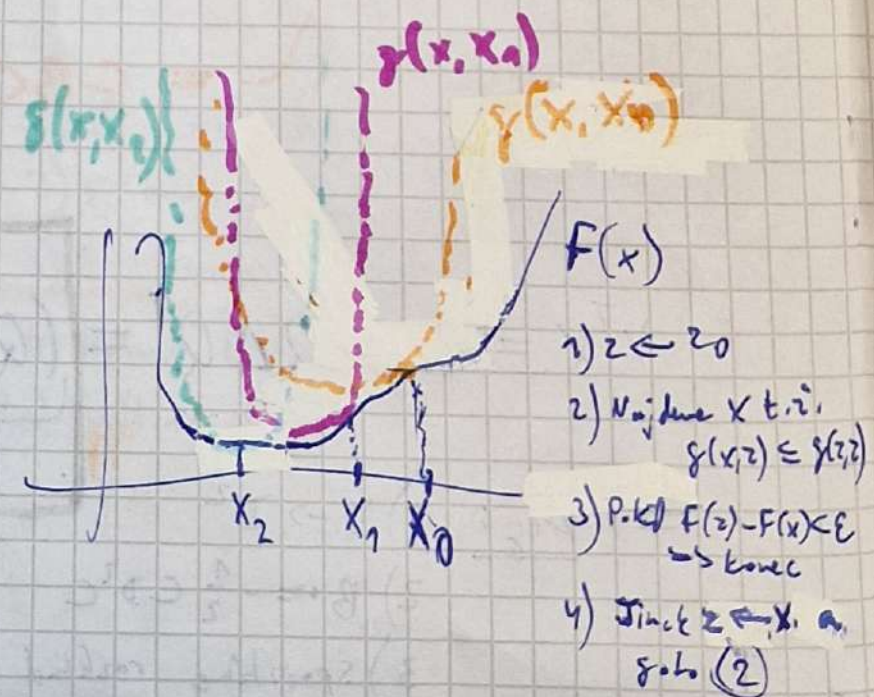
normalizovaný stress $\sigma_n^2(X) := \frac{\sigma_r^2(X)}{\sum_{i < j} w_{ij} \delta_{ij}^2}$

Kruskalov stress $\sigma_\gamma(X) := \sqrt{\sigma_n^2(X)}$
 0 ... super
 > 0.2 ... BLE

- Alg.:
- 1) X náhodně
 - 2) Opakujeme, dokud se nám to nelíbí:
 - 3) Spočítáme stress X
 - 4) Upravíme X , abychom zmenšili stress

iterativní majorizace

$f(x) \rightsquigarrow g(x, z)$
 majorizační funkce,
 dotýká se f
 v bodě z
 $(f(z) = g(z, z))$
 $f(x) \leq g(x, z)$



- 1) $z \leftarrow z_0$
- 2) Najdeme x t. i. $g(x, z) \leq f(x)$
- 3) P. k. $f(z) - f(x) < \epsilon \rightarrow$ konec
- 4) Jinak $z \leftarrow x$ a. $g(x, z)$

$$Q(X) = \sum_{i < j} w_{ij} (D_{ij}(X) - \delta_{ij})^2 =$$

$$= \underbrace{\left(\sum_{i < j} w_{ij} \delta_{ij}^2 \right)}_{\eta^2_{\delta}} + \underbrace{\left(\sum_{i < j} w_{ij} D_{ij}^2(X) \right)}_{\eta^2(X)} - 2 \cdot \underbrace{\left(\sum_{i < j} w_{ij} \delta_{ij} D_{ij}(X) \right)}_{\rho(X)} =: \mathcal{J}(X, Z)$$

Kvadr. majorizační funkce

SMACOF
Function convex

$$\text{trace}(X^T V X)$$

$$w_{ij} A_{ij} \rightarrow \begin{pmatrix} 1 & -1 \\ -1 & 1 \\ & & \ddots & \ddots \\ & & & 1 & -1 \\ & & & & & 1 \end{pmatrix}$$

$$\text{trace}(X^T B(Z) Z)$$

$$\frac{w_{ij} \delta_{ij}}{D_{ij}(Z)}$$

$\delta_{ij}(p_{ij})$

Máme ale být potřeba transformace těch "vzdáleností" (proximities) \rightarrow disparities

\hat{D}_{ij}

metrické: $\hat{D}_{ij} = f(p_{ij})$

ne-metrické:

ordinální

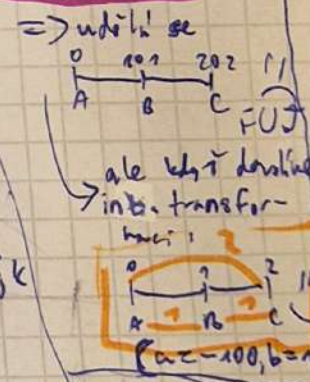
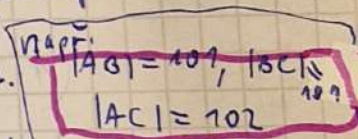
$p_{ij} < p_{jk} \rightarrow \hat{D}_{ij} < \hat{D}_{jk}$

nominální

$p_{ij} = p_{jk} \Rightarrow \hat{D}_{ij} = \hat{D}_{jk}$

vstupních dat

- b · p_{ij} ... poměr.
- a + b · p_{ij} ... int.
- a + b · log(p_{ij}) ... log.



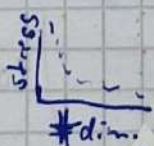
disparity (n x (n-1)/2)

$$Q_r(\hat{D}, X) = \sum_{i < j} w_{ij} (D_{ij}(X) - \hat{D}_{ij})^2 = \eta^2_{\delta} + \eta^2(X) - 2\rho(\hat{D}, X)$$

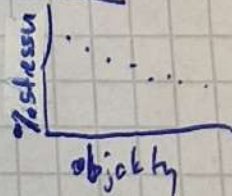
Configuration plot

- prostě ten výsledek - umístění bodů na ty souřadnice

Scree plot



Stress plot



Sheppard plot

- u ne-metrického MDS



Bubble plot

= configuration plot + ty body mají velikost podle toho, jaký přísl. stress

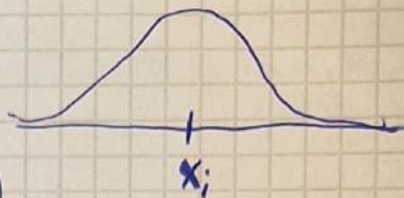
musíme opt. minimalizovat $\hat{D}; X$

treba štírdavě

t-SNE → pro vizualizace (ne redukci dim.)

nejdřív SNE:

$$P_{j|i} := \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$



Pravděpodobnost, že si x_i vezme x_j za svého souseda

$$q_{i|j} := \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

($\sigma := \frac{1}{\sqrt{2}}$)

entropie distribuce p: $H(p) := -\sum_{i=1}^N P(x_i) \log p(x_i)$

↳ KL divergence

→ "vyjádření, jak moc se dvě distribuce liší"

$$D_{KL}(p \parallel q) := \sum_{i=1}^N P(x_i) (\log p(x_i) - \log q(x_i)) = \sum_{i=1}^N P(x_i) \log \frac{P(x_i)}{q(x_i)}$$

↳ očekávané množství informace, kterou ztratíme, když použijeme q místo p

cost function: $C := \sum_a KL(p_a \parallel q_a) = \sum_i \sum_j P_{j|i} \log \frac{P_{j|i}}{q_{j|i}}$

↳ pomocí gradient descent minimalizujeme "pružiny"

ALE někde to může být málo škála, a jinde zase ne

⇒ u každého bodu nastavíme tu σ_i podle perplexity

$Perp(P_i) = 2^{H(P_i)}$
dáno

bin. vyhl. nejbone s σ_i , aby to vyšlo

kolik přibližně má každý bod s n sousedů (~ 5-50) **PARAMETR**

SNE → t-SNE

- zesymetrizováno → rychlejší konv. + méně problémů s outliery

$$q_{ij} = q_{ji} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

ne jen vzdálenosti do i , ale všechny

- t-distrib. pro q ⇒ méně crowdingu + lépe se posítá

$$P_{ij} = P_{ji} = \frac{P_{j|i} + P_{i|j}}{2n}$$

(aby nebyl problém s outliery)

UMAP

k -simplex $\equiv k$ -dim. Δ



simplexový komplex \equiv skupina simplexů, které se dohledují na nižší-dim. simplexu

Čechův

- bod \rightarrow 0-simplex
- $k+1$ prokládajících se koulí \rightarrow k -simplex

Vietoris-Ripsův

- jednoduší -
- z těch $k+1$ koulí se nemusí prohnat všechny, ale stačí po párech

Pro ten poloměr koulí bude mít každý bod svou metrickou, aby ta koule zasahovala ke **k**-tému nejbližšímu sousedovi (k -NN)

PARAMETR

\rightarrow ale NE - uděláme, že 1. NN tam bude patřit a ti ostatní postupně se snižující se pravděpod. + protože NN-vztah není symetrický, tak zesymetrizujeme (např. $P(A) \times P(B) := P(A) + P(B) - P(A) \cdot P(B)$)

plně určen θ a 1-simplex \Rightarrow dobře reprezentovatelný jako graf
 křivka \rightarrow spojená do jednoho simplexu

graf s hranami: ohodnocenými pravděpodobnostmi.

Alg: 1) V bod si spočítáme k -NN, vzdálenost k nejbližšímu sousedovi ρ_i

a σ_i tak aby $\sum_{j=1}^k \exp\left(\frac{-\max(0, d(x_i, x_j) - \rho_i)}{\sigma_i}\right) = \log_2 k$ aby to do 1. NN bylo $e^0 = 1$

2) Sestrojíme graf $H = (X, E, w)$: $E \dots$ mezi vrcholem a jeho k -NN

$$w_h(x_i, x_j) := \exp\left(\frac{-\max(0, d(x_i, x_j) - \rho_i)}{\sigma_i}\right)$$

3) Body rozmístíme do toho nižko-dim. prostoru (nějak přibližně)

\Rightarrow na základě Euklid. vzdal. spočítáme (váhy hran) (do **PARAMETR** váha 1, pak klesá)

4) Minimalizujeme cross-entropii \rightarrow SGD + negative sampling

$$\sum_{e \in E} w_h(e) \log\left(\frac{w_h(e)}{w_e(e)}\right) + (1 - w_h(e)) \log\left(\frac{1 - w_h(e)}{1 - w_e(e)}\right)$$

pritažlivá síla pro body blízko ve vyš. dim. prost.

odpudivá síla

4 zase
47 až
42 až