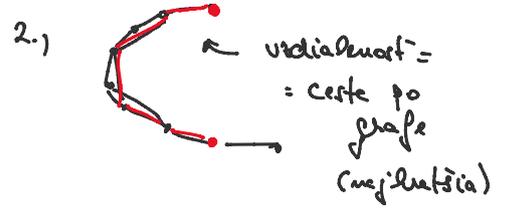
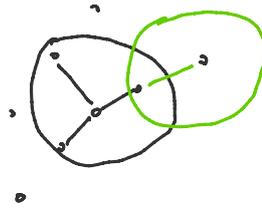


Minule: Isomap
- zachytanie štruktúry

1.) Graf susednosti



3.) RBF na vzdialenostiach grafových

Rko: Swiss roll data
: HWIST

t-SNE

t-distributed stochastic neighbor embedding

- $x_1, \dots, x_n \in \mathbb{R}^d$ - chceme aproximovať pomocou $y_1, \dots, y_n \in \mathbb{R}^2$
→ zachovať vzájomné polohy, štruktúru dát

História: najprv - SNE (Hinton & Roweis, 2002)

potom - t-SNE (van der Maaten & Hinton, 2008)

- vzdialenosť: x_1, \dots, x_n vyjadriť pomocou pp.:

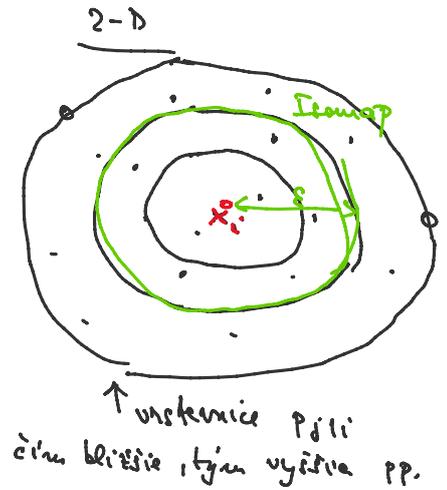
$P_{j|i}$ = pp., že x_i si nadjí bližšie ako x_j ako svojho suseda (v grafe)

→ pravdepodobnostný prístup & Isomap

$$P_{j|i} = \frac{e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}}} \quad \left. \begin{array}{l} \sim \text{proporcionálna hustota } N(x_i, \sigma_i^2 I) \\ \sim \text{normalizácia} \end{array} \right\}$$

(-) Isomap: susedia: \forall body bližšie ako ϵ

• blízkosť x_i a x_j : $p_{ij} = \frac{P_{j|i} + P_{i|j}}{2n}$ $\forall j \neq i$
↑ normalizácia



• vzdialenosť y_i a y_j : pomocou viacrozmerného t-rozdelenia (⇒ "t-" SNE)

$$Q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{n \neq s} \sum_{s \neq n} (1 + \|y_n - y_s\|^2)^{-1}} \quad \left. \begin{array}{l} \sim \text{hustota t-rozdelenia} \\ \sim \text{normalizácia} \end{array} \right\}$$

$$\sum_{n \neq s} \sum (1 + \|y_n - y_s\|^2)^{-1} \} \text{ normalitácia}$$

Odbôcha: $U \sim N(0, \Sigma)$, $V \sim \chi_m^2$, nez. $\Rightarrow Y = \frac{U}{\sqrt{\frac{V}{m}}} + \mu \sim t_m(\mu, \Sigma)$
 pre \bar{x} -reg., $k=2$, $m=1$, ~~že~~ keď Y hustota: \uparrow k -normálne

$$f(y) = \frac{1}{\pi (1 + \|y - \mu\|^2)^{3/2}} \rightarrow \text{pre } \mu = y_i \Rightarrow f(y) = \frac{1}{(\dots \|y - y_i\|^2)^{3/2}}$$

• voľba y_1, \dots, y_n , aby $\begin{matrix} Q \\ \parallel \\ q_{ij} \end{matrix}$ zodpovedalo $\begin{matrix} P \\ \parallel \\ p_{ij} \end{matrix}$

• vzdialenosť pp. rozdelení: Kullbackova - Leiblerova divergencia

$$I(P, Q) = \sum_i \sum_{j \neq i} p_{ij} \ln \left(\frac{p_{ij}}{q_{ij}} \right) \rightarrow \text{ako veľkí sú } P \text{ a } Q \text{ rozdielne}$$

$$I(P, Q) \geq 0 \quad \forall P, Q$$

$$I(P, Q) = 0 \Leftrightarrow P = Q \text{ skoro všade (pre nelo konečné: } \Leftrightarrow Q = P)$$

$$I(Q, P) \neq I(P, Q)$$

• cost function: $C(y) = I(P, Q) \rightarrow \min_y$
 \uparrow \uparrow
 dane' určenie' gypilonmi

\rightarrow nekonnečné úloha \Rightarrow globálna numerická optimalizácia

- upravená gradientná metóda (pomocou $\frac{\partial C}{\partial y_i}$)

- v poč. iteráciách: pridať nové q gypilonom náhodným spôsobom - aby sme neuviazli v lok. optime

- poč. riešenie: $y_1, \dots, y_n \sim N(0, \sigma^2 I)$ pre malé $\sigma > 0$

- trik: early exaggeration - skoré zväčšenie

\hookrightarrow prvých niekoľko iterácií: $p_{ij} \mapsto c \cdot p_{ij}$ pre nejaké $c > 1$ (napr. $c = 12$)

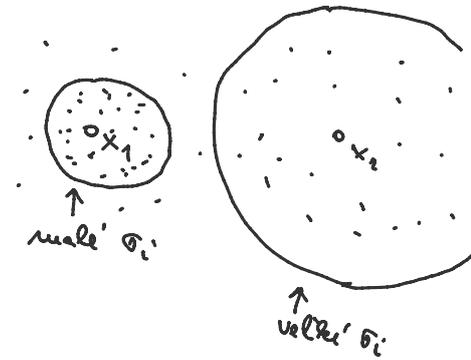
$\Rightarrow Q$ nedokáže zachytiť $P \rightarrow$ optimalizácia sa porieši zachytiť hlavne veľké p_{ij} pomocou veľkých q_{ij} (ostatné $q_{ij} \approx 0$)

\Rightarrow vytvárajú zhluky - kuste'

- celkom dobre od seba oddelené

Vol'ba σ_i :

- čím väčšie σ_i , do tým väčšej diaľky sa x_i posera'
- ak sú body bližšie sebe, stačí mali' σ_i ,
ak sú niekde, tak treba veľké' - aby mali porovnateľnú veľkú "susedov"



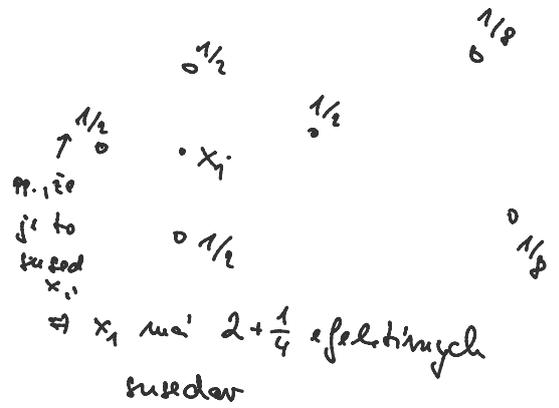
→ $\text{Pexp}(P_i) = 2^{H(P_i)}$ $H(P_i) = - \sum_{j \neq i} p_{ji} \cdot \log_2(p_{ji})$: entropia pravdepodobnosti P_i
 (p_{1i}, \dots, p_{ni})

↓
 perplexity ↓ rozdelenie p_{ji}

: čím väčšie, tým menej rôznorode' je rozdelenie

⇒ $\text{Pexp}(P_i)$: vyjadruje efektívny # susedov predane' x_i .

: čím väčšie σ_i , tým väčšie P_i



→ σ_i nastavíme, aby všetky x_i mali rovnako veľkú "susedov", teda aby $\text{Pexp}(P_i)$ bolo rovnaké H_i (nepr. = 10)

t-SNE: veľmi dobre rozlišuje zhluky

: môže zachytiť zhluky aj ak v dátach žiadne nie sú

: vol'ba t-rozdelenia - v mnohorozmere je "veľká priestora": body sú vo všeob. ďaleko od seba → ako ich posunúť do 2-rozmeru?

- crowding problem

→ v 2-D (q-D) volíme t-rozd. s ľahkými chvostami a tak sa trváime, že body sú od seba ďalej (aby sme udeľili možnosť vzdialenejším body)

Reo: MNIST

→ tsne : veľmi pomalá

: umožňuje prebrúne vyhodnotenie

: medzivremu robíte, lebo štandardizáciu som už spravil pomocou PCA

→ Rtsne : C++ foel oel van der haatene volenig' cet RLo
: tief spreu' PCA unetri (default)