# $t\,SNE$

t-distribution stochastic neighbor embedding

Van der Maaten & Hinton 2008

based on     Hinton & Roweis 2002     SNE

Non-Linear Dimensionality Reduction algo.

(כלליות) Data Science @ סטטיסטיקה בתחום פופולרי

ורבה מאוד שיטה זו. המטרה (בלתי") ווכ'ל אלגוריתם שמקבל נתונים

בממדים גבוהים ומוציא כ׳ים שייצגו (2D) את כ׳ נתונים של הנתונים

ובכך כל נתונים.

שני שלבים: ① בונים מבנה שמייצג קשרים בין הנתונים

② שמים אותם בממד נמוך $\ddot{y}$ כך שמנסים לשמר את המבנה. $D_{KL}$

נתחיל בפרטים של SNE.

Data $X$: $X_i \in \mathbb{R}^p$     High-dim. data

Distances: $\forall_{ij}\ d_{ij} = \|X_i - X_j\|$

נהפוך מרחקים אלו להסתברות שכנות בצורה הבאה:

$$P_{j|i} = \frac{\exp\left(-d_{ij}^{2}/2\sigma_i^{2}\right)}{\sum\limits_{k\neq i} \exp\left(-d_{ik}^{2}/2\sigma_i^{2}\right)}$$

(cont.) ‏מה, $\sigma_i$ ‏ is ‏מס"ר מיצג @ ‏מדביון (cont.)

1. ‏$P_i$ ‏הם ‏ההסתברות ‏שהמסגרת ‏ה"ל ‏ל ‏לשכן ‏j.
(‏המספר ‏מדבל ‏נירטון.)

‏איך ‏נקבע ‏את ‏$\sigma_i$ ? ‏מגדיר ‏ונימי ‏לצדק
‏ולצורך ‏מדבילים ‏ההסתברות @ ‏$P_i$ ‏לפי
‏שימוש ‏ל "‏מספר ‏היעיל ‏השכנים".

$$\underline{\text{Perplexity}} = 2^{H(P_i)}$$

‏מדד ‏ההסתברות @ ‏השכן

$$H(P_i) = -\sum_j P_{j|i} \cdot \log_2 P_{j|i}$$

‏בעצם ‏אנחנו ‏נרצה ‏א-ונפורם / ‏Perplexity ‏בתוך ‏$[5-50] \in$
‏נקבע ‏כזה ‏של $X_i$ ‏את ‏$\sigma_i$ ‏כזה ‏שיהיה ‏$2^{H(P_i)}$ ‏כזה.

small $\sigma_i \to$ "less neighbors" $\to$ more predictable transitions
lower entropy $\to$ lower perplexity

higher $\sigma_i \to \quad \dots \to$ higher perplexity

# Embedding — find low-dimensional representation

$Y: \quad Y_i \in \mathbb{R}^d \quad d \ll p \quad$ כאשר $X$

(t-SNE $\mathscr{l}$ 61)    SNE מתוך פיתוח

(Gauss. Kernel, $\sigma = \frac{1}{\sqrt{2}}$)

בעזרת הסתברויות של $Y$

$$q_{j|i} = \frac{\exp(d_{ij}^2)}{\sum_k \exp(d_{ik}^2)}$$

$$d_{ij}' = \|y_i - y_j\|$$

while __minimizing__ the (KL) divergence of $P_i, Q_i$

score (for SNE):

$$C = \sum_i KL(P_i \| Q_i) = \sum_i \sum_j P_{j|i} \cdot \log_2 \frac{P_{j|i}}{Q_{j|i}} \qquad (*)$$

! Gradient Descent   איך מחשבים את $C$ ? 

$$\frac{\partial C}{\partial y_i} = 2 \sum_j \left( P_{j|i} - q_{j|i} + P_{i|j} - q_{i|j} \right)(y_i - y_j)$$

מתחת למשוואה: המרחק בין שתי הנקודות     כיוון התזוזה

Init. with random points (sampled from Gauss., $r = \frac{1}{\sqrt{2}}$)

Use __momentum__

מומנטום

$$Y^{(t)} = Y^{(t-1)} + \eta \frac{\partial C}{\partial Y} + \alpha(t) \cdot \left[ Y^{(t-1)} - Y^{(t-2)} \right]$$

learning rate $\rightarrow \eta$    gradient

# Comment 1 - Why KL-divergence?

Kullback-Leibler (Relative Entropy)

In Inf. theory: expected #bits (per sample)
when compressing samples from P using
code for Q (compared to a code for P)

In Bayesian Inference: Info: gained when moving from
a prior dist Q to a prior P (for data from P).

$$D_{KL}(P \parallel Q) = -\sum_j P_j \log Q_j + \sum_j P_j \log P_j = H(P,Q) - H(P)$$

cross entropy $\nearrow$   entropy $\nearrow$

Symmetric version:   Jensen-Shannon Divergence

$$D_{JS}(P \parallel Q) = D_{JS}(Q \parallel P) = \tfrac{1}{2}D(P \parallel M) + \tfrac{1}{2}D(Q \parallel M)$$

$$M = \frac{P+Q}{2}$$

.Q ·| P ·√ M ‏פֿשודֿ‎ Ⓔ ‏ני﬘‎ ‏פֿשﬢ‎ ‏אﬥⅠⅠG‎

# Comment 2 - So Why KL? (for SNE)

large $P_{j|i}$ are important. w/ small $q_{j|i}$ → BIG penalty

Small $P_{j|i}$, even if large $q_{j|i}$ yield small penalty

# Moving to t-SNE

Solving two problems    ①   <u>Outliers</u>
                 ②   <u>Crowding</u>
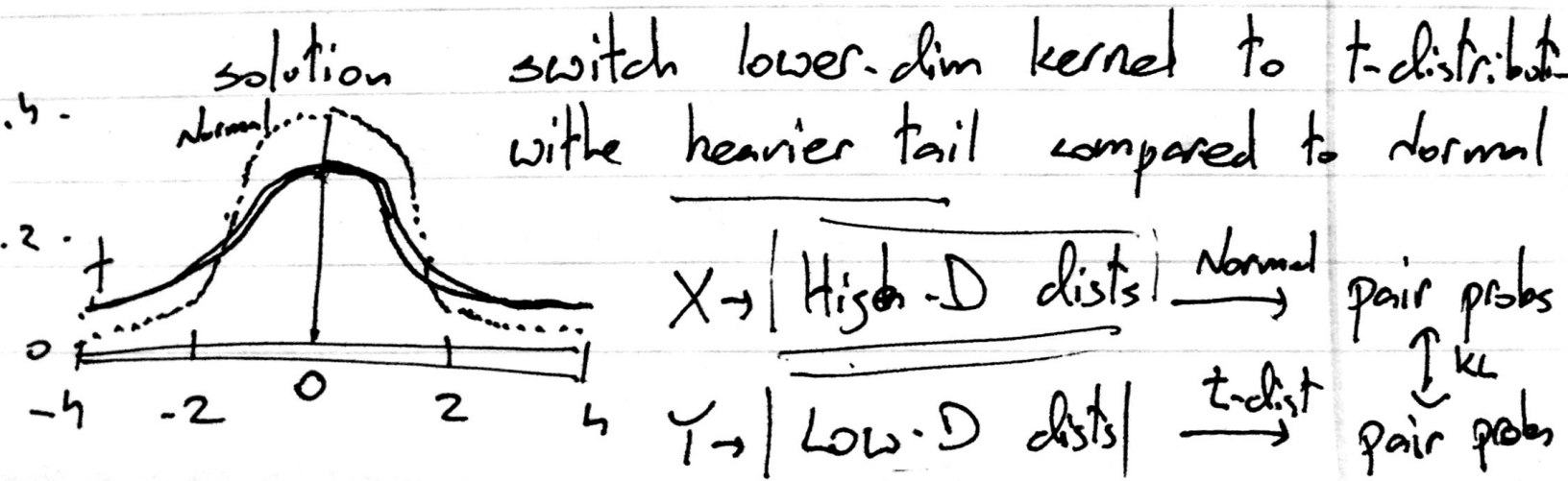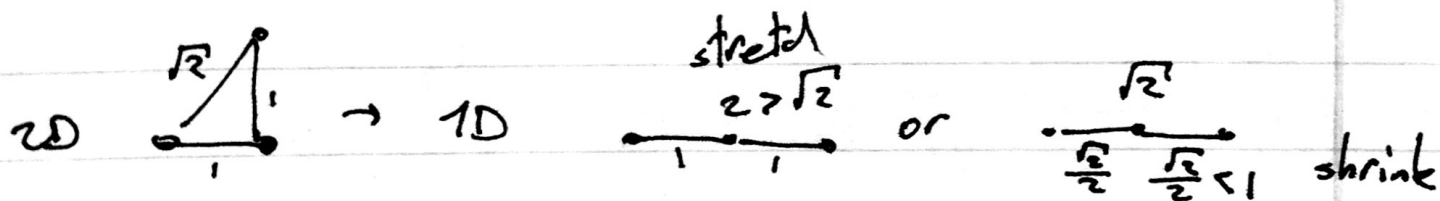
①   $x_i$ outlier → add noise to $C$    (far from all)

Update: $C = KL(P \| Q) = \sum_{ij} P_{ij} \log P_{ij} / Q_{ij}$

Hence: move from $N$ dist. $P_{j|i}$ to one $P_{ij}$ (pairwise)

$$P_{ij} = \frac{P_{j|i} + P_{i|j}}{2 \cdot N}$$

         → large $P_{ij}$ (= close points) get stronger

②   High-to-low dimen. embeddings expon. reduce the "space" for neighboring points. Causing <u>under-est.</u> for ~~small~~ close points, and <u>over-est.</u> of dist for dist.

$2D$  → $1D$    stretch $\underset{1 \quad 1}{\underline{\bullet \;\; 2 > \sqrt{2} \;\; \bullet}}$   or   $\underset{\frac{\sqrt{2}}{2} \; \frac{\sqrt{2}}{2} < 1}{\underline{\bullet \;\; \sqrt{2} \;\; \bullet}}$   shrink

solution    switch lower-dim kernel to t-distribut.
with heavier tail compared to normal



$X \rightarrow \boxed{\text{High-D dists}} \xrightarrow{\text{Normal}}$ pair probs
$\qquad\qquad\qquad\qquad\qquad\qquad \updownarrow KL$
$Y \rightarrow \boxed{\text{Low-D dists}} \xrightarrow{\text{t-dist}}$ pair probs

Students t-distribution

$$\text{PDF}: \quad f(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\,\Gamma\left(\frac{\nu}{2}\right)} \cdot \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}, \quad \text{w/ param } \nu$$

$$\Gamma(x) = (x-1)! \quad, \quad \Gamma(1) = 1 \quad, \quad \Gamma\left(\tfrac{1}{2}\right) = \sqrt{\pi}.$$

$$\text{w/ } \nu=1 \quad f(t) = \frac{\overset{=1}{\Gamma(1)}}{\underset{\text{``}\sqrt{\pi}}{\sqrt{\pi}\cdot\Gamma\left(\tfrac{1}{2}\right)}}\left(1+t^2\right)^{-1} = \frac{\left(1+t^2\right)^{-1}}{\pi}$$

in t-SNE $Q$ is defined with a t-dist kernel

$$q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k\neq l}\left(1 + \|y_k - y_l\|^2\right)^{-1}}$$

Gradient Descend ןיבנהל ה.3ט.ש'כ'ו10 ejvei ובוא

$$\frac{\partial C}{\partial y_i} = 4\sum_j \underbrace{\left(P_{ij} - q_{ij}\right)}_{\substack{\text{nמב קוلın} \\ \text{HiD-LowD}}} \underbrace{\left(y_i - y_j\right)}_{\substack{ןוויכ \\ ро}} \underbrace{\left(1 + \|y_i - y_j\|^2\right)^{-1}}_{\text{"הלולה ברוצ' ה"}}$$

ול'וח ונוחמ קוחר
.נווחה ליב ?י

[ for speed-up , use Barnes - Hut simulation
where neighboring points are grouped to compute forces ]