

PCA Wrap-Up

Projection perspective

$$\tilde{x}_i = Bz_i \quad B \text{ again orthonormal}$$

Want to minimize reconstruction error

$$J_M \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$$

What are the optimal coordinates z_i for x_i w.r.t. B ?

$$\frac{\partial J_M}{\partial z_{ji}} = \frac{\partial J_M}{\partial \tilde{x}_i} \frac{\partial \tilde{x}_i}{\partial z_{ji}}$$

$$\frac{\partial J_M}{\partial \tilde{x}_i} = -\frac{2}{N} (x_i - \tilde{x}_i)^T \in \mathbb{R}^{1 \times D}$$

$$\frac{\partial \tilde{x}_i}{\partial z_{ji}} = \frac{\partial}{\partial z_{ji}} \left(\sum_{m=1}^M z_{mi} b_m \right) = b_j$$

So

$$\frac{\partial J_M}{\partial z_{ji}} = -\frac{2}{N} (x_i - \tilde{x}_i)^T b_j$$

$$= -\frac{2}{N} \left(x_i - \sum_{m=1}^M z_{mi} b_m \right)^T b_j \quad (b_m^T \cdot b_j = \delta_{mj} \text{ if } m \neq j)$$

$$= -\frac{2}{N} (x_i^T b_j - z_{ji} b_j^T b_j) = -\frac{2}{N} (x_i^T b_j - z_{ji})$$

$$\text{Set to } \delta \rightarrow z_{ji} = b_j^T x_i$$

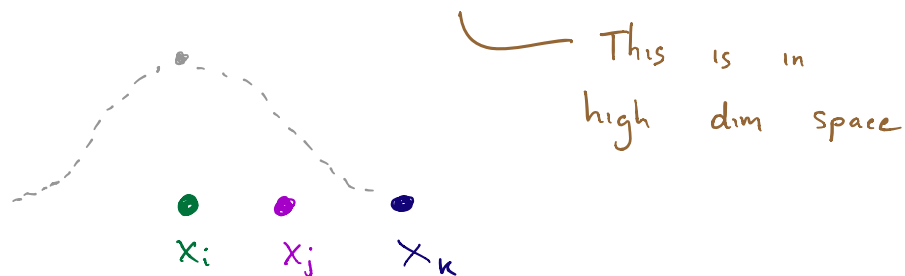
A similar argument for choice of B (basis) can be made, yielding again the M largest Eigenvectors (see reading!)

Stochastic Neighbor Embedding (SNE)

Aim: $X \xrightarrow{\text{(very) low dim}} Y$
(very) high dim

Define a conditional probability that encodes similarity

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



Similarly in the map y

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

Ideally: $P_{j|i} \approx q_{j|i} \quad \forall i, j$

Formalize this with KL-Divergence (KL)

$$KL(q||p) \doteq \sum_x q(x) \log \frac{q(x)}{p(x)}$$

"How different is the distribution q from p ?"

Properties:

(1) $KL(q||p) \geq 0$

(2) if $KL(q||p) = 0 \rightarrow q = p$

(3) $KL(q||p) \neq KL(p||q)$

Cost function:

$$C = \sum_i KL(Q_i || P_i) = \sum_i \sum_j P_{j|i} \log \frac{P_{j|i}}{q_{j|i}}$$

Same, but
in M .

Conditional distribution
over all other points j
given i in D

to place points, find y to minimize C

Using gradient descent $\nabla_y C$.

Symmetric SNE

Instead of **Conditionals** $P_{i|j}, q_{i|j}$ define
Joint distributions P_{ij}, q_{ij}

$$q_{ij} = \frac{\exp\{-\|y_i - y_j\|^2\}}{\sum_{k \neq l} \exp\{-\|y_k - y_l\|^2\}}$$

Normalizes dist
by all pairs

$$\rightarrow q_{ij} = q_{ji}$$

For high-dim space, outliers pose a challenge
because **denominator** will be large $\rightarrow P_{ij}$ small
 $\rightarrow y_{ij}$ unimportant.

Instead:

$$P_{ij} \stackrel{\text{def}}{=} \frac{P_{j|i} + P_{i|j}}{2n}$$

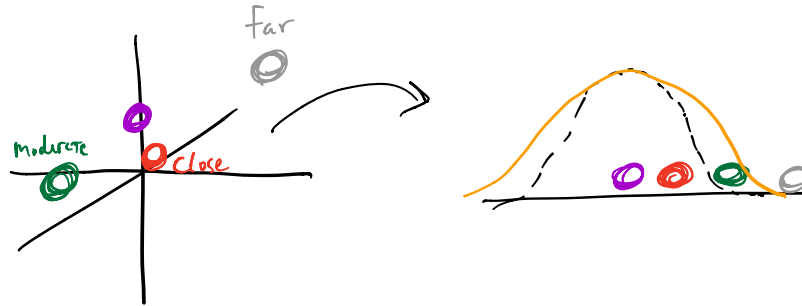
$$\rightarrow \sum_j P_{ij} \geq \frac{1}{2n} \quad \forall x_i$$

Yields a nicer ∇_y

$$\nabla_{y_i} C_{\text{symmetric}} = 4 \sum_j (P_{ij} - q_{ij})(y_i - y_j)$$

The Crowding problem

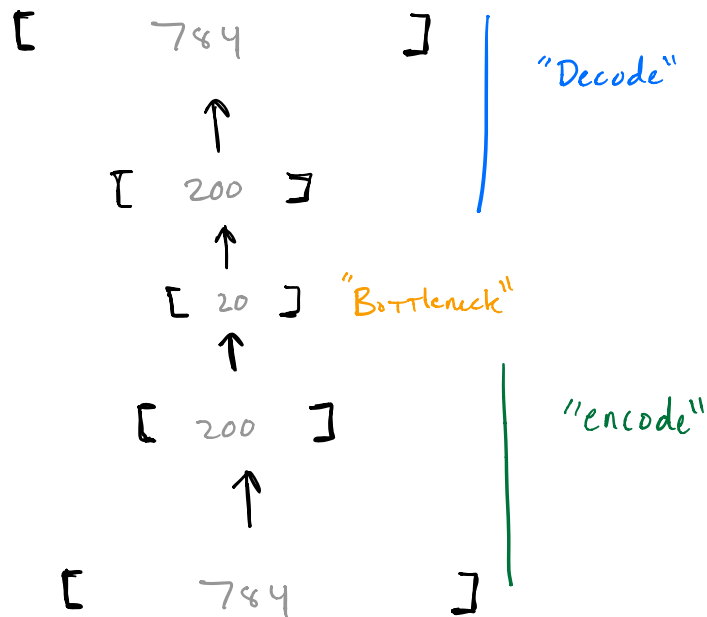
Not enough space in lower dims



In t-SNE we model the joint probability q_{ij} using a Student T-distribution which has heavier tails \rightarrow moderate distance in $X \rightarrow$ big distance in Y is ok -- does not force "moderate" distances in X to yield small distances in Y

Auto-encoders

Design a network that consumes x , and
Then re-constructs it.



Simplest version

$$m \begin{bmatrix} x \\ \vdots \\ \vdots \end{bmatrix} \xrightarrow{W} \begin{bmatrix} \vdots \\ d \\ \vdots \end{bmatrix} \xrightarrow{V} \begin{bmatrix} \tilde{x} \\ \vdots \\ \vdots \end{bmatrix} \quad \mathcal{L}(x, \tilde{x}) = \|x - \tilde{x}\|^2$$

Note: This is just linear dim reduction, and
should look familiar

$$\underset{m}{W} x \rightarrow \underset{d}{z}; \quad \underset{m}{\tilde{x}} = \underset{d}{V} z$$

$\begin{matrix} / & | & | & & | \\ dxm & mxl & dxl & & mxd \end{matrix}$

More next time!