

## MORE ON DIM. RED.

① PCA: Eigenvalue - Eigenvector  
Covariance matrix

- Rotation
- Linear transformation

② MDS  $\rightarrow$  Reproduce the distances from the embedding space as euclidean distances in the projected space

Ⓐ Classical MDS:

$$S = \| \Theta - \Delta \|^2 \quad \swarrow \text{diagonal}$$

$$\Delta \rightarrow G \rightarrow G = V \Lambda V^T$$

$$G = Y Y^T$$

$$\downarrow$$

$$Y = V \Lambda^{1/2}$$

If we use Euclidean distances also in the embedding space

Classical MDS is equiv. to PCA

① Metric MDS:

$$S = \sum_{ij} \left( \frac{\Delta_{ij} - \theta_{ij}}{\Delta_{ij}} \right)^2$$

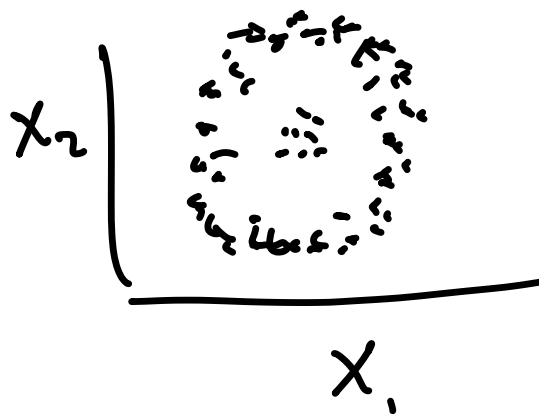
Minimize  $S$  with GP

② Non-metric MDS

$$S = \sum_{ij} \left( \frac{g(\Delta_{ij}) - f(\theta_{ij})}{g(\Delta_{ij})} \right)^2$$

Kernel PCA:

The Kernel Trick



$$(x_1, x_2) \xrightarrow{\phi} (x_1, x_2, \underline{x_1^2 + x_2^2})$$



$\phi$ 

$$K_{ij} = K(x^i, x^j) = \langle \phi^i, \phi^j \rangle$$

$$x^i \rightarrow \phi^i$$

$$K_{ij} = \langle \phi^i, \phi^j \rangle$$

$$K(x^i, x^j) = x^i x^{jT} \quad \text{linear Kernel}$$

$$K(x^i, x^j) = (x^i \cdot x^{jT})^\delta \quad \vee, \quad \text{Polynomial Kernel}$$

$$(x^i x^{jT} + 1)^\delta \quad \text{V}_2$$

$$\delta = 2 \quad D = 2$$

$$K_{i\ell} = \left( 1 + (x_1^i, x_2^i) \cdot (x_1^\ell, x_2^\ell) \right)^2 =$$

$$1 + (x_1^i, x_1^\ell)^2 + (x_2^i, x_2^\ell)^2 + 2(x_1^i x_1^\ell) + 2(x_2^i x_2^\ell) + 2x_1^i x_1^\ell x_2^i x_2^\ell =$$

$$\rightarrow \left( 1, \sqrt{2} \underline{x_1^i}, \sqrt{2} \underline{x_2^i}, \underline{(x_1^i)^2}, \sqrt{2} (\underline{x_1^i} \underline{x_2^i}), \underline{(x_2^i)^2} \right).$$

$$\left( 1, \sqrt{2}, x_1^\ell, \sqrt{2} x_2^\ell, (x_1^\ell)^2, \sqrt{2} (x_1^\ell x_2^\ell), (x_2^\ell)^2 \right)$$

$$= \underline{\phi^i} \cdot \phi^\ell \rightarrow D(\phi) = \frac{(\delta \cdot D)^2 - (\delta D)}{2}$$

# Gaussian kernel

$$k_{i,e} = \exp\left(-\frac{\|x^i - x^e\|^2}{2\sigma^2}\right)$$

$$D=2$$

$$k_{i,e} = \exp\left(-\left(x_1^i - x_1^e\right)^2 - \left(x_2^i - x_2^e\right)^2\right) =$$

$$\exp\left(-\left(x_1^i\right)^2 + 2x_1^i x_1^e - \left(x_1^e\right)^2 - \left(x_2^i\right)^2 + 2x_2^i x_2^e - \left(x_2^e\right)^2\right) = \underbrace{\exp\left(-\|x^i\|^2\right)}_{\phi_i} \cdot \underbrace{\exp\left(-\|x^e\|^2\right)}_{\phi_e}$$

$$\exp\left(2x^i \cdot x^e\right) = \underbrace{\exp\left(-\|x^i\|^2\right)}_{\phi_i} \cdot \underbrace{\exp\left(-\|x^e\|^2\right)}_{\phi_e}$$

$$\cdot \sum_{n=0}^{\infty} \left( \frac{(2x^i \cdot x^e)^n}{n!} \right)$$

$$\underbrace{\phi_i \cdot \phi_e}_{\phi_i \cdot \phi_e}$$

$$D(\phi) = \infty$$

$$\phi_i = \sum_{n=0, \dots, \infty} \frac{(x^i)^n}{\sqrt{n!}} \dots$$

$$\phi^e =$$

Linear kernel  $\rightarrow$  kernel PCA  $\sim$  PCA

Non-linear kernel  $\rightarrow$  Non-linear Dim. Red.

① Mercer theorem  $K \rightarrow$  positive semi def.  
sym.

$$k(x^i, x^l) = \langle \phi^i, \phi^l \rangle$$

$$\tilde{\phi}^i = \phi^i - \frac{1}{N} \sum_k \phi^k$$

$$\tilde{k} \langle \phi^i, \phi^l \rangle = \langle \phi^i, \phi^l \rangle - \frac{1}{N} \sum_k \langle \phi^i, \phi^k \rangle - \frac{1}{N} \sum_k \langle \phi^l, \phi^k \rangle + \frac{1}{N^2} \sum_{k,m} \langle \phi^k, \phi^m \rangle$$

Kernel PCA in a nutshell

① Pick a kernel

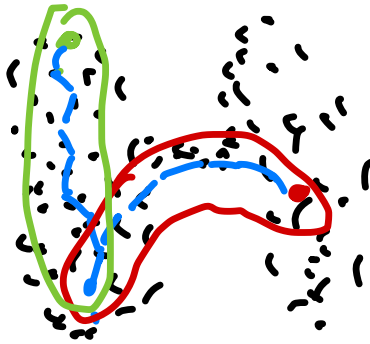
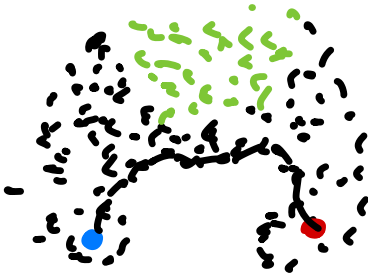
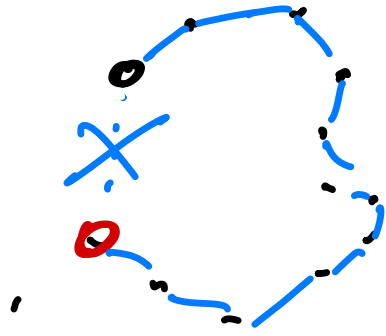
② Build the kernel matrix  $K$

③ Recover  $G$  in the feature space by double centering  
 $\tilde{K}$

④  $\tilde{K} \alpha_i = \lambda_i \alpha_i$

⑤  $y = \Lambda^{1/2} \alpha$

# \* Diffusion Map



Geodesic distance  $\rightarrow$  Diffusion distance  
(ISOMAP)

$$k_{ij} = \exp \left( \frac{-\|x^i - x^j\|^2}{2\sigma^2} \right)$$

$\sigma$  max. distance  
with a single  
jump

$$\tilde{k}_{ie} = \frac{k_{ie}}{\sqrt{\sum_k k_{ik} \cdot \sum_m k_{em}}}$$

$$p_{ie} = \frac{\tilde{k}_{ie}}{\sum_n \tilde{k}_{in}}$$

$$\underline{P} \underline{v} = \lambda \underline{v}$$

we choose as new  
coordinates  
the  $d$  eigenvectors  
corresponding to  
the highest  $\lambda$

Sketch Map : Non metric MDS

$$S = \sum_{ie} \frac{f(\theta_{ie}) - g(\Delta_{ie})}{g(\Delta_{ie})}$$

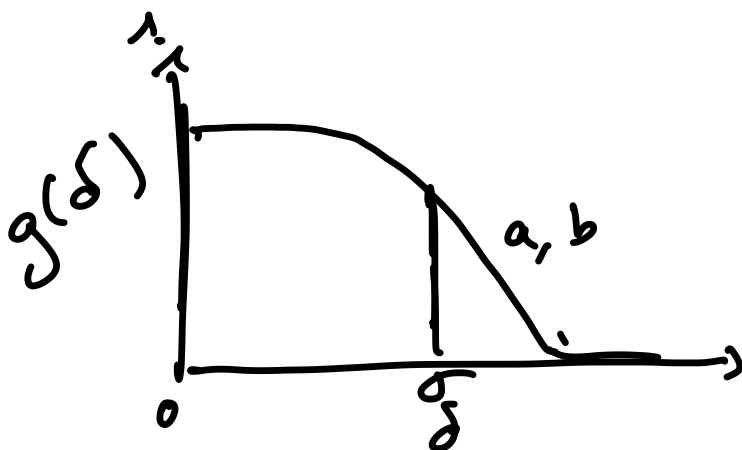
$$S = \frac{\sum_{ie} \omega_i \omega_e [g(\Delta_{ie}) - g(\theta_{ie})]}{\sum_{ie} \omega_i \omega_e}$$

$$g(\delta) = 1 - \left[ (2^{a/b} - 1) \left( \frac{\delta}{\sigma} \right)^a \right]^{-b/a}$$

$\omega_i \Rightarrow$  relative importance of point  $i$ ;

$\sigma =$  distance scale that should be reproduced

$a, b \rightarrow$  rate of changing



Reproducing distances of the embedding space into the projected space

↓ probabilistic view

| Probabilities of being a neighbor

t - SNE

t - (distributed) Stochastic Neighbor Embedding

① maximizes the similarities between the probabilities distributions of being neighbors in the original and projected spaces

Probability that point  $l$  is a neighbor of the point  $i$ :

$$p_{li} = \frac{\exp\left(-\frac{1}{2}\left(\frac{\Delta_{il}}{\sigma_i}\right)^2\right)}{\sum_{k \neq i} \exp\left(-\frac{1}{2}\left(\frac{\Delta_{ik}}{\sigma_i}\right)^2\right)}$$



$$S_{\text{sym}} = P_{ie} = \frac{P_{ie} + P_{ei}}{2N} \quad \left\{ \begin{array}{l} P_{ie} = P_{ei} \\ \sum_{i=0} P_{ie} = 1 \end{array} \right.$$

$\downarrow$   
 total number of data points

$\sigma_i$

$$\text{Perplexity} = 2^{\frac{-\sum_e P_{ei} \log_2 P_{ei}}{e}}$$

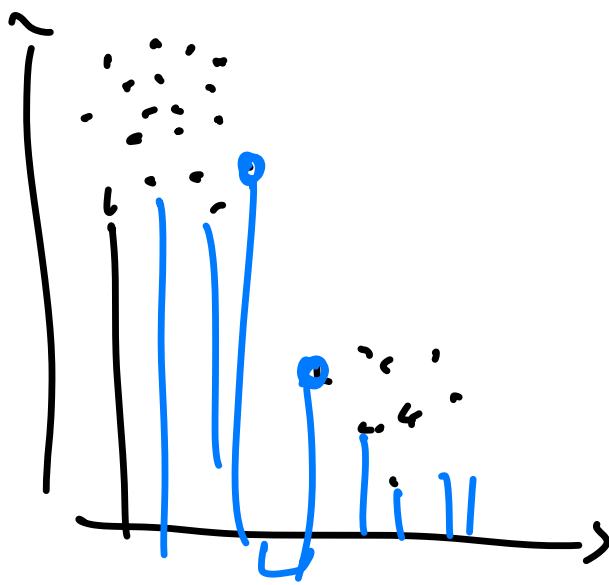
$\sim$  Preserve the distances up to the perp NN

$P \rightarrow$  prob. dist. in the embedding space

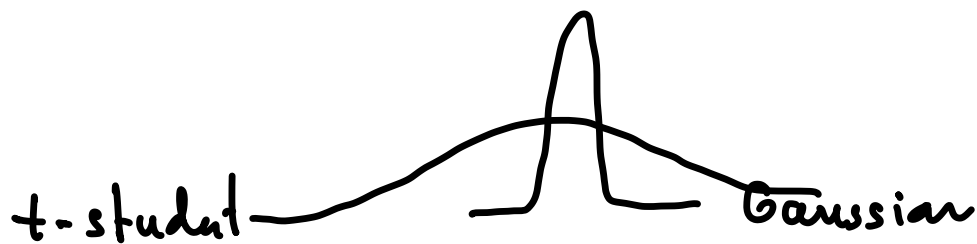
$Q \rightarrow$  prob. dist. in the proj. space

minimize the KL divergence between P & Q

$$KL(P||Q) = \sum_{e \neq i} P_{ie} \log \frac{P_{ie}}{Q_{ie}}$$



Over crowding



$$q_{ie} = \frac{(1 + \theta_{ie}^2)^{-1}}{\sum_{k \neq i} (1 + \theta_{ik}^2)^{-1}}$$

$\theta$  = Euclidean distance in  $\mathcal{X}$

$$KL(P \parallel Q) = \sum_{i \neq j} p_{ie} \lg \frac{p_{ie}}{q_{ie}(\gamma)}$$

Gradient Descent  $\rightarrow$  obtain

$$\gamma = \min_{\gamma} KL(P \parallel Q)$$

Problems:

- ① Projection is different at each run
- ②  $d > 2$  convergence problems
- ③ Loss the original structure

UMAP Uniform Manifold  
Approximation & Projection

variation t-SNE

①  $p_{i|e} = \exp\left(\frac{-\Delta_{ie} - p}{\sigma_i}\right)$

① we do not normalize

②  $p \rightarrow \text{min-dist}$   
if  $\Delta_{ie} < p$   $y_i \sim y_e$

②  $\sigma_i$

$k = 2^{\sum p_{i|e}} \rightarrow \sigma_i$   
↓  
number of neighbors  
of interest

③ Symmetrizing scheme

$$t\text{-SNE} = p_{ie} = \frac{p_{ie} + p_{ei}}{2N}$$

$$\text{UMAP} = p_{ie} = p_{ie} + p_{ei} - p_{ie} \cdot p_{ei}$$

$$\textcircled{4} \quad q_{ie} = \left( 1 + a \Delta_{ie}^{2b} \right)^{-1}$$

$a, b \rightarrow \text{fit}$

$$\left( 1 + a (\Delta_{ie})^{2b} \right)^{-1} = \begin{cases} 1 & \text{if } \Delta_{ie} \leq \rho \\ e^{-\Delta_{ie} - \rho} & \text{if } \Delta_{ie} > \rho \end{cases}$$

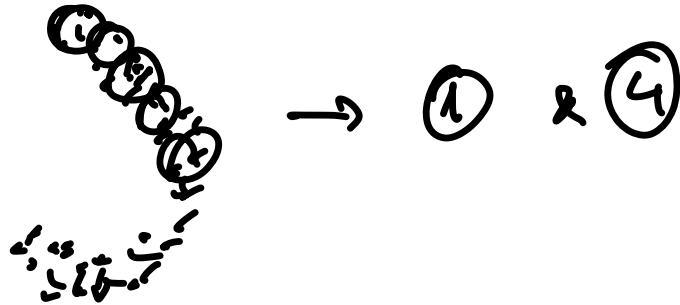
⑤ Instead of KL divergence  
binary cross entropy

$$L = \sum_{i,e} \left[ \underbrace{p_{ie} \log \frac{p_{ie}}{q_{ie}}}_{\text{KL}} + (1 - p_{ie}) \log \frac{(1 - p_{ie})}{(1 - q_{ie})} \right]$$

⑥ Init  $\gamma$  from graph Laplacian

Why UMAP performs better than t-SNE?

① 1 hyperparameter



② Faster & Not changes from run to run  
①, ②, ③ ⑥

③ Preserves better the global structure  
⑤