# EE769 Intro to ML
# Dimension Reduction

Amit Sethi

Faculty member, IIT Bombay

# Learning objectives

- Write the advantages of dimension reduction

- Write steps for principal component analysis

- Extend PCA to nonlinear methods using the kernel trick

# Dimensionality reduction: what & why?

- Objective:

  , find mapping $f : x_i \rightarrow y_i$, $y_i \in {}^d$, $d<D$

  - Such that, there exists $f^{-1}$ that can (almost) reconstruct

- Advantages:

  - Less redundancy, easier classification

  - Smaller storage, faster search

  - **Unravel meaningful latent variables**

$x \in \mathbb{R}^{100}$

$y \in \mathbb{R}^{5}$

$X \in \mathbb{R}^{N \times D}$

$g : y \rightarrow \hat{x}_i$

$\hat{x}_i \in \mathbb{R}^{D}$

$\| \hat{x}_i - x_i \|_2^2$

3

# Example: Only two latent variables



$$x_i \in \mathbb{R}^{10,000}$$

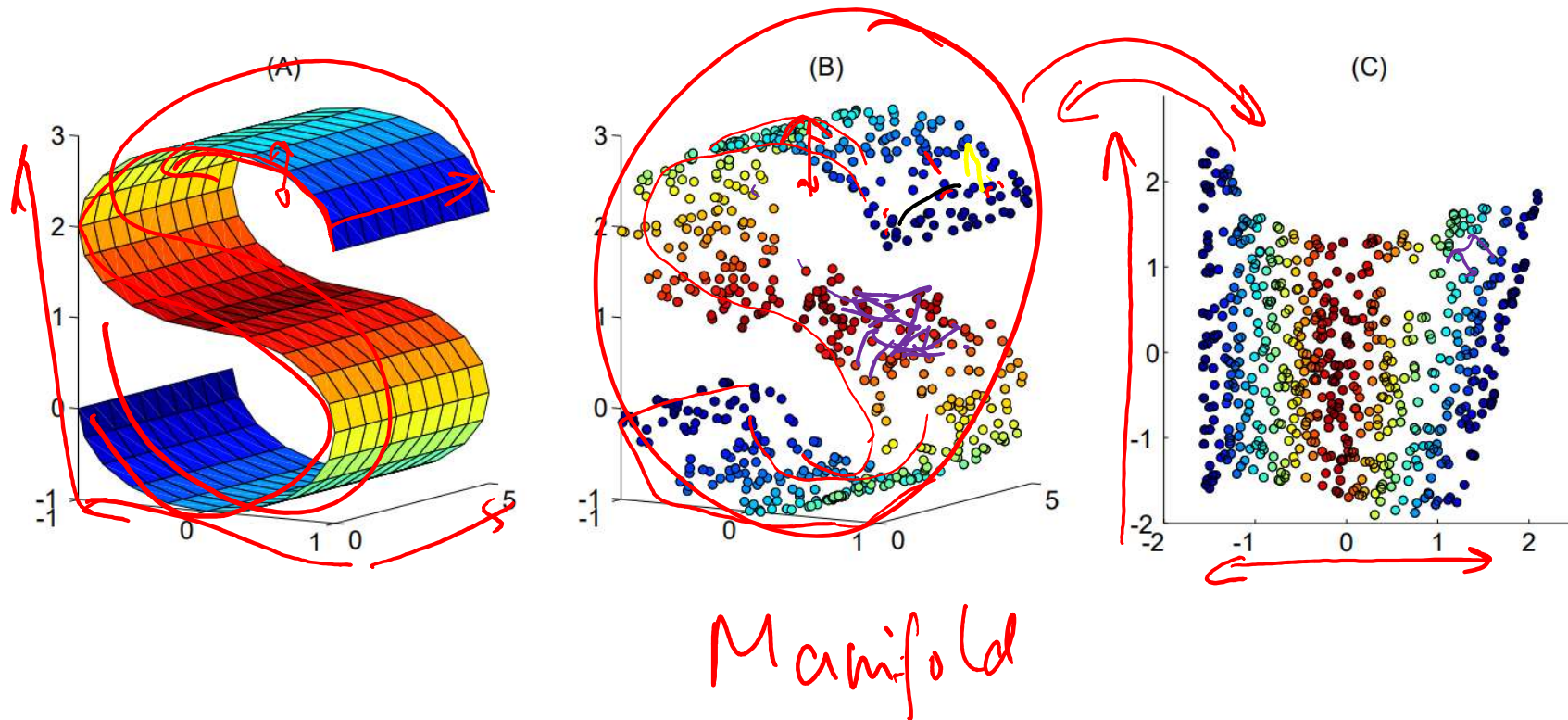$$y \equiv (\text{Size}, \text{Rotation})$$

$$y \in \mathbb{R}^2$$

# Example 2: Many latent variables, but still less than "pixels"



- What defines the shape of a face?
- Add pose and expression…
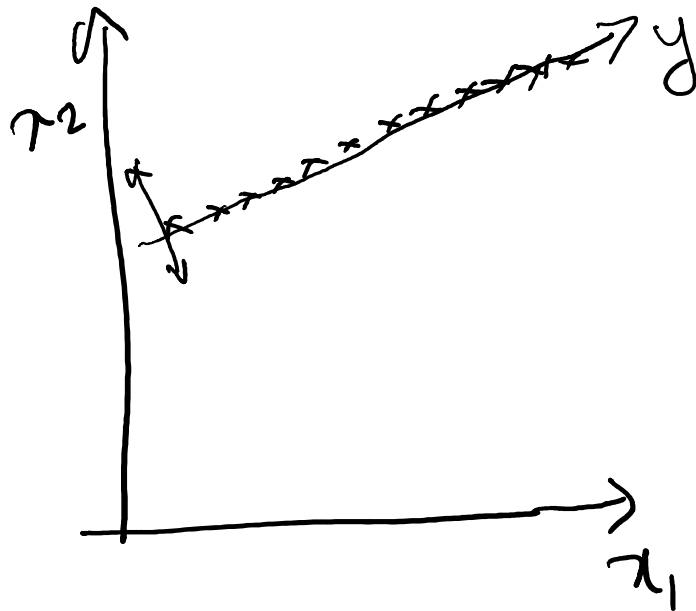- Still a lot less variables than the pixels of a face image

*Image source: Pixabay.com*
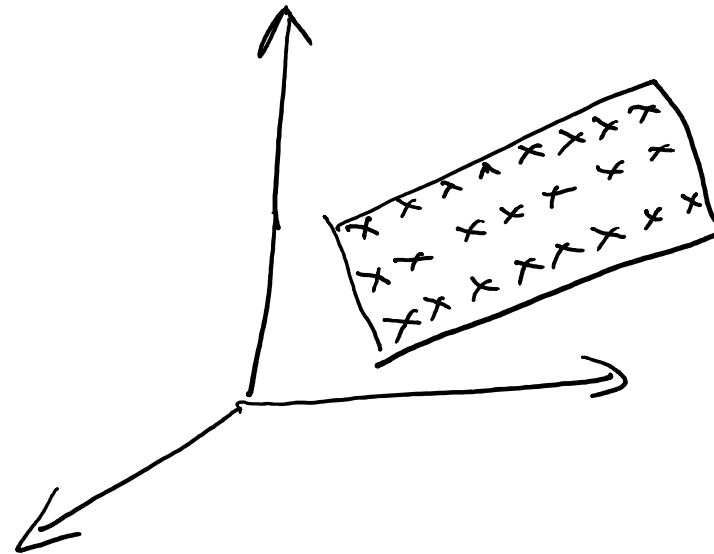
# Visualizing manifolds



Manifold

"An Introduction to Locally Linear Embedding," Lawrence K. Saul
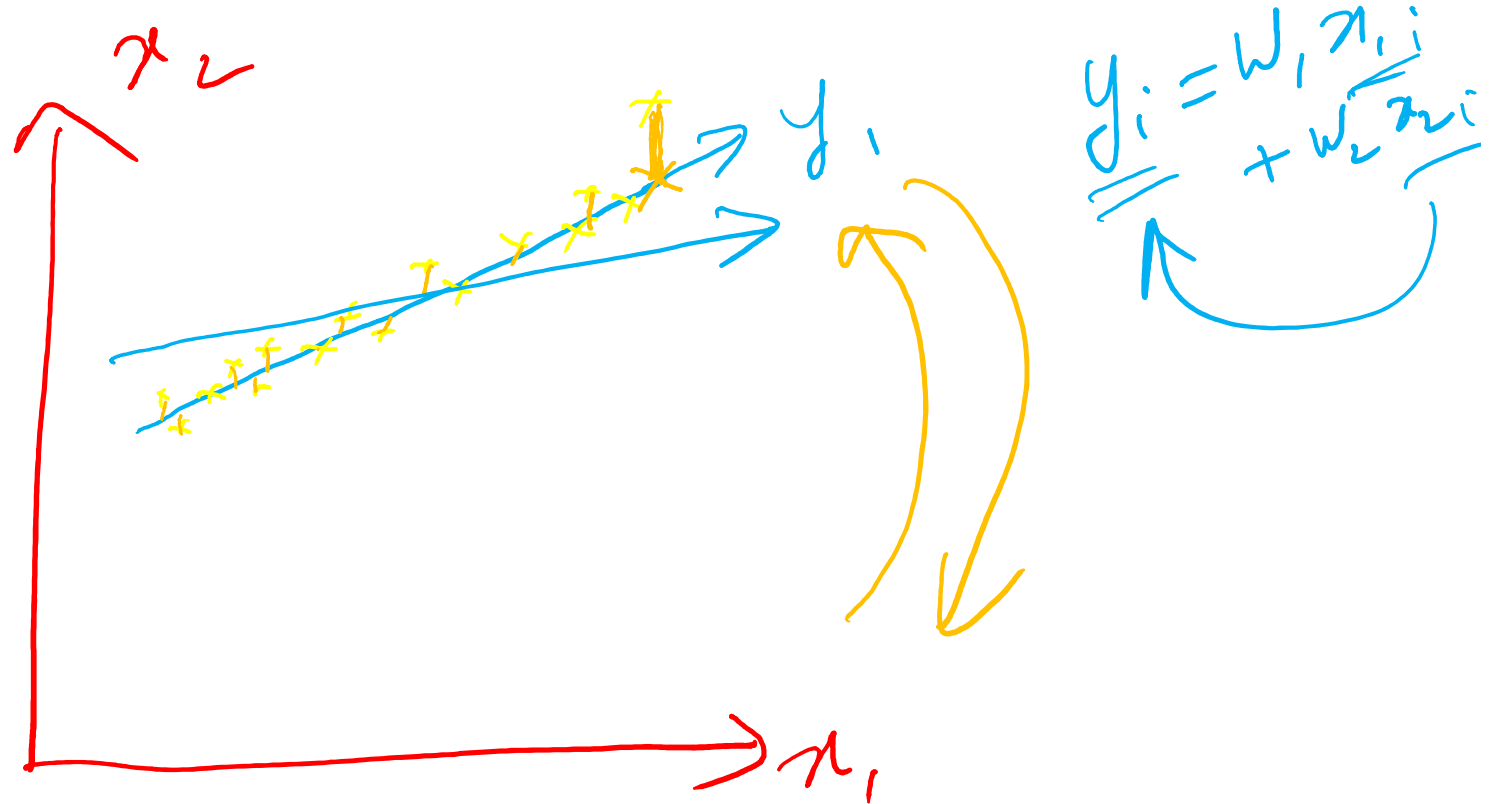
6

# Inherent linear dimension of data

Data along a line on a 2-D plane
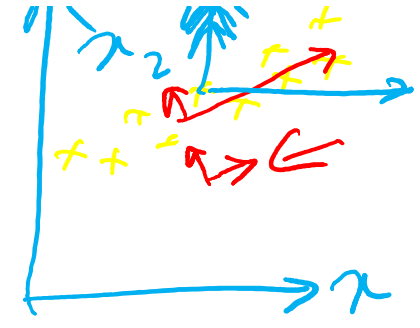
Data along a line or hyperplane in 3-D

# Deviations from the underlying hyperplane



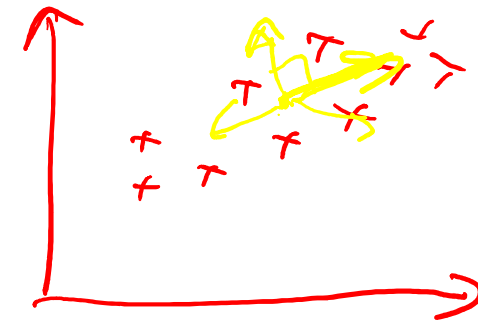$$y_i = w_1 x_{1,i} + w_2 x_{2,i}$$

# Principal component analysis

**Objective:**

- Find a way to select top $d < D$ orthogonal directions that explain the maximum possible variance of the data

**Outline:**

- Mean-center the data

- Loop
  - Find direction of maximum variance (that is orthogonal to all previously found directions)
- Until desired percent of variance is captured or number of dimensions

# Complete algorithm

$P\,CA$

$X \in \mathbb{R}^{N \times D}$

- Mean centering: $z_i = x_i - \mu_x$
- Covariance: $N\,C = Z^\top Z$    $C \in \mathbb{R}^{D \times D}$    $C$ is P.S.D.
- Eigen decomposition: $C = U \Lambda U^\top$ ; $\lambda_j u_j = C u_j$ $\longrightarrow$ unit vector
- Selection: $C_d = U \Lambda_d U^\top = U_d \Lambda_d U_d^\top,\ d < D$
  - Drop dimensions: Set lower $D-d$ eigenvalues to zero

$[u_1,\ u_2\ [\cdots u_D]] = U$

- Projection: $Y = Z U_d$

$N \times d$    $N \times D$   $D \times d$    $\lambda_{d+1} \cdots \lambda_D$

$$\begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \cdots \lambda_D & \\ 0 & & & \end{bmatrix} = \Lambda$$

- Reconstruction: $Z_d = Y U_d^\top = Z U_d U_d^\top$    $N \times D$

$\lambda_1 \geq \lambda_2 \geq \cdots \lambda_D \geq 0$

- Reconstruction error: $||Z - Z_d||_2^2 \propto |\Lambda - \Lambda_d|$
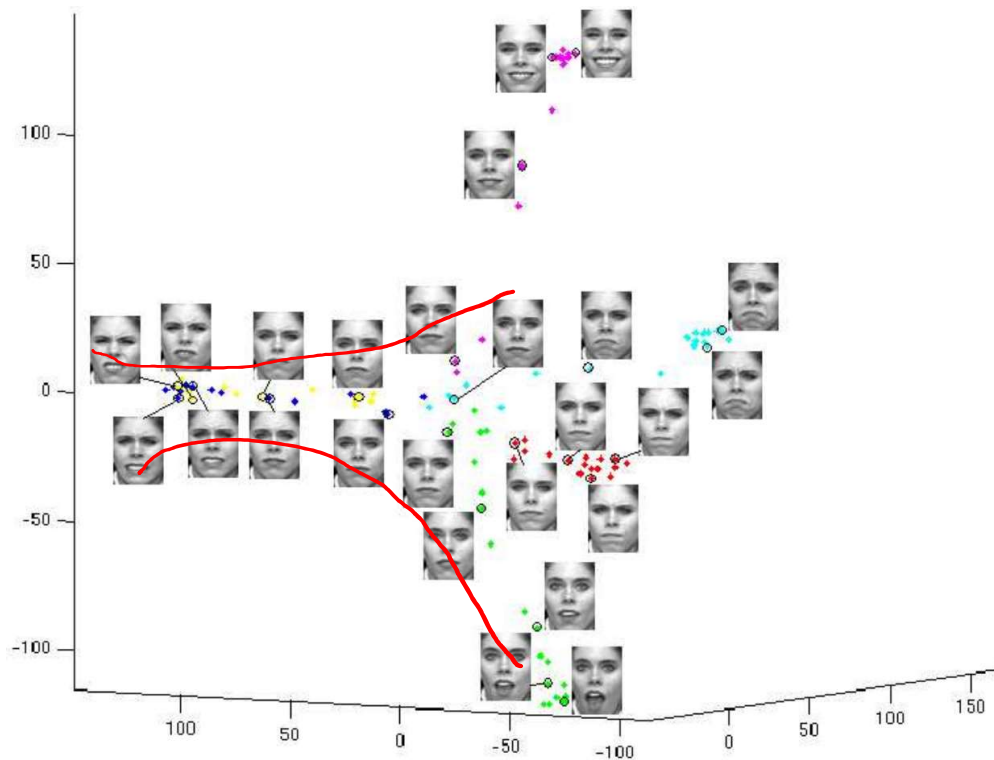
$$= \sum_{i=d+1}^{D} \lambda_i$$

# Example: Eigenfaces



↗ P.C. of face images

Turk, Matthew, and Alex Pentland. "Eigenfaces for recognition." *Journal of cognitive neuroscience* 3.1 (1991): 71-86.
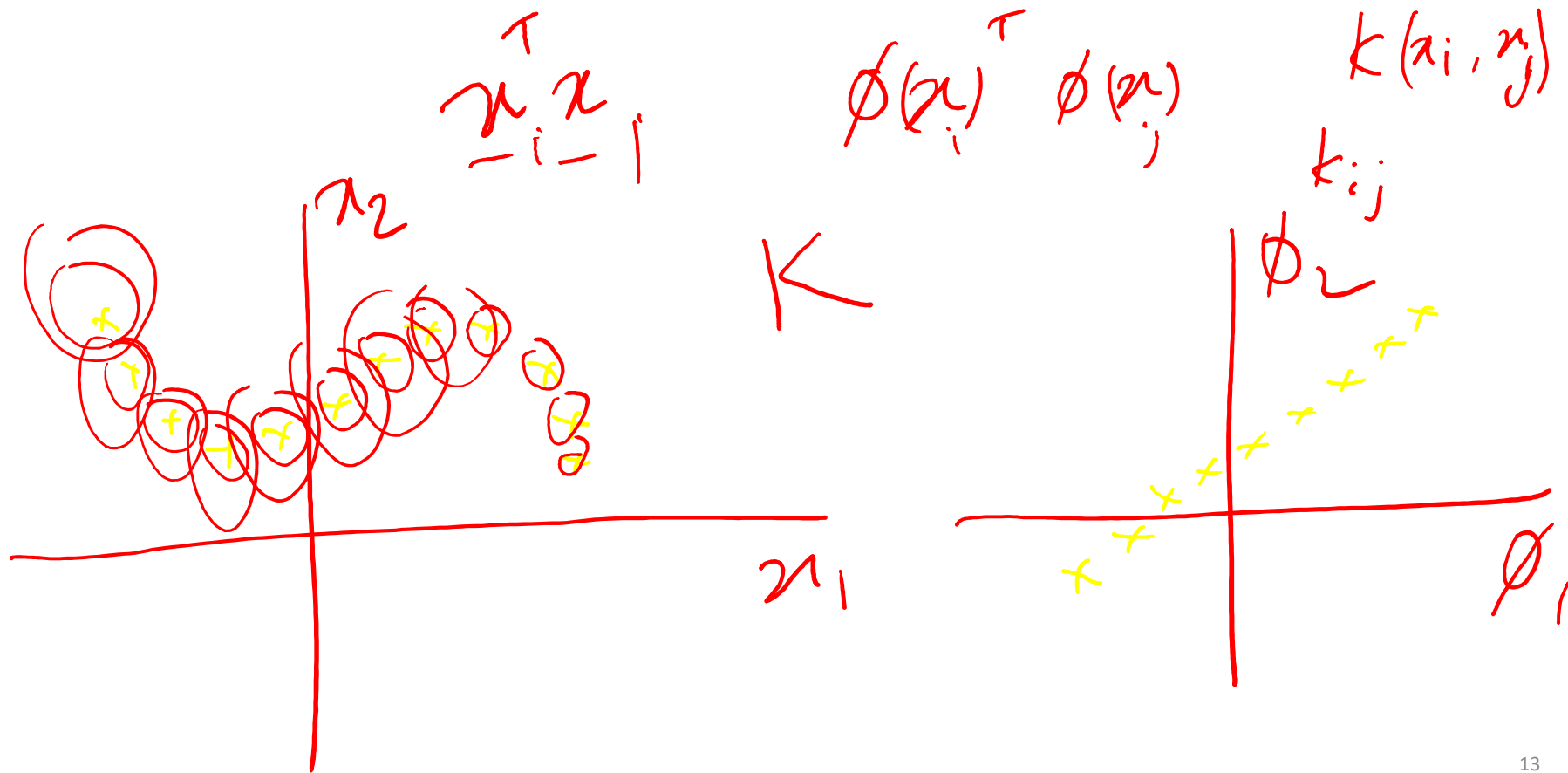
# Data often lies around a nonlinear manifold



$$x_2 = \sin(x_1)$$

$$y = f(x_1)$$

Shan, Caifeng, Shaogang Gong, and Peter W. McOwan. "Appearance manifold of facial expression." *International Workshop on Human-Computer Interaction*. Springer, Berlin, Heidelberg, 2005.

# Kernel PCA introduces nonlinearity by mapping data to a feature space

$$\underline{x}_i^T \underline{x}_j$$

$$\phi(x_i)^T \phi(x_j)$$

$$k(x_i, x_j)$$

$$K$$

$$k_{ij}$$

# Kernel PCA avoids calculating features directly by using the kernel trick

- Recall PCA:            $C\,u_i = \lambda_i u_i\,;\; N\,C = Z^{\mathsf{T}} Z$

- Substituting features:    $\sum_n \varphi(z_n)\,\varphi(z_n)^{\mathsf{T}}\,u_i = N\,\lambda_i u_i$

- Substitute:           $u_i = \sum_n a_{in}\,\phi(z_n)$

- This means:           $\sum_n k_{l,n} \sum_m a_{im}\,k_{n,m} = N\,\lambda_i \sum_n a_{in}\,k_{l,n}$

- In matrix form:       $K^2\,a_i = \lambda_i\,N\,K\,a_i$

- Solve eigenvalue problem:    $K\,a_i = \lambda_i\,N\,a_i$

- There is a mean-centering step:
    - Ref: https://www.ics.uci.edu/~welling/classnotes/papers_class/Kernel-PCA.pdf

Schölkopf, Bernhard (1998). "Nonlinear Component Analysis as a Kernel Eigenvalue Problem". *Neural Computation*. **10** (5): 1299–1319.
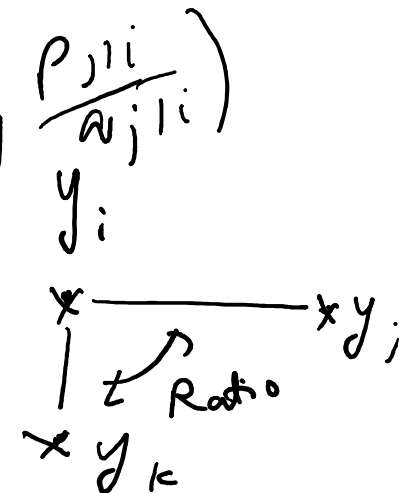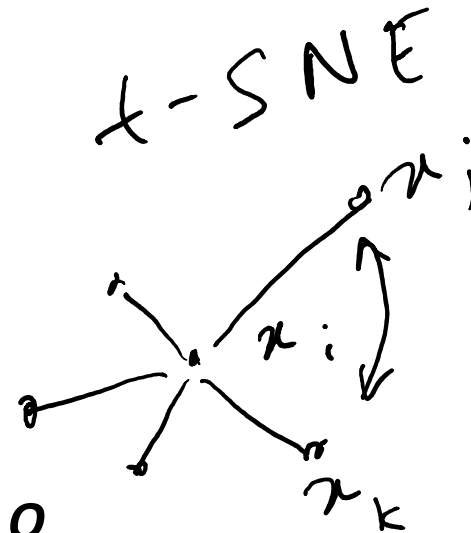
# Stochastic Neighbor Embedding

$t$-SNE

- $p_{j|i} = softmax(-||\textbf{x}_i - \textbf{x}_j||^2 / 2\sigma_i^2)$
- $q_{j|i} = softmax(-||\textbf{y}_i - \textbf{y}_j||^2)$

- Objective: Reduce KL divergence between $\textbf{P}$ and $\textbf{Q}$

$$\frac{e^{-||x_i - x_j||/2\sigma_i^2}}{\sum_{k \neq i} e^{-||x_i - x_k||/2\sigma_i}}$$

$$\frac{(1+||y_i - y_j||)^{-1}}{\sum_k \cdots}$$

$N \times N \quad N \times N$

$\left( \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}} \right)$

$y_i$

"Stochastic Neighbor Embedding," Geoffrey Hinton and Sam Roweis

# Auto-encoder bottlenecks for dimension reduction



$MSE$

$$\sum_i \|\hat{x}_i - x_i\|_2^2$$

Encoder → Contraction

Decoder → expansion