# #MAHSA_AMINI

# MACHINE LEARNING

Electrical Summer Workshops (ESW) 2022

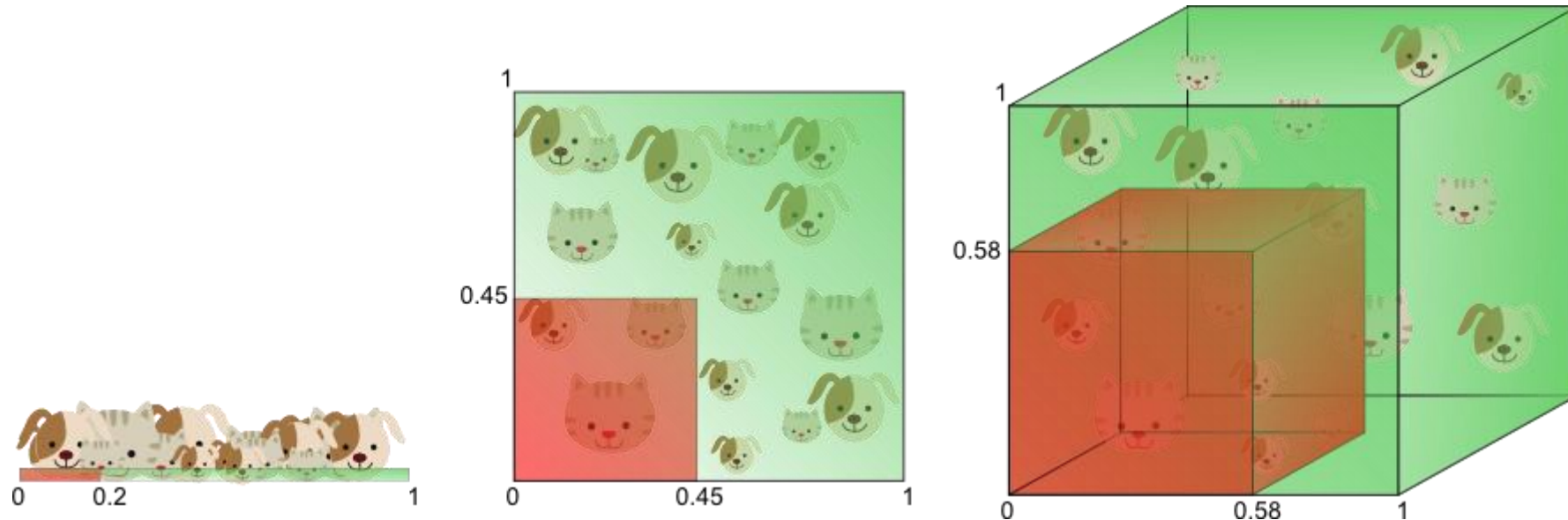Electrical Engineering Department - Sharif University of Technology

Instructors: *Alireza Gargoori Motlagh – Amir Mirrashid – Ali Nourian*

# DIMENSIONALITY REDUCTION

# Motivation

- High dimensional data problems:
  - High computation cost and runtime of training
  - Increase the chance of overfitting
  - Correlated variables
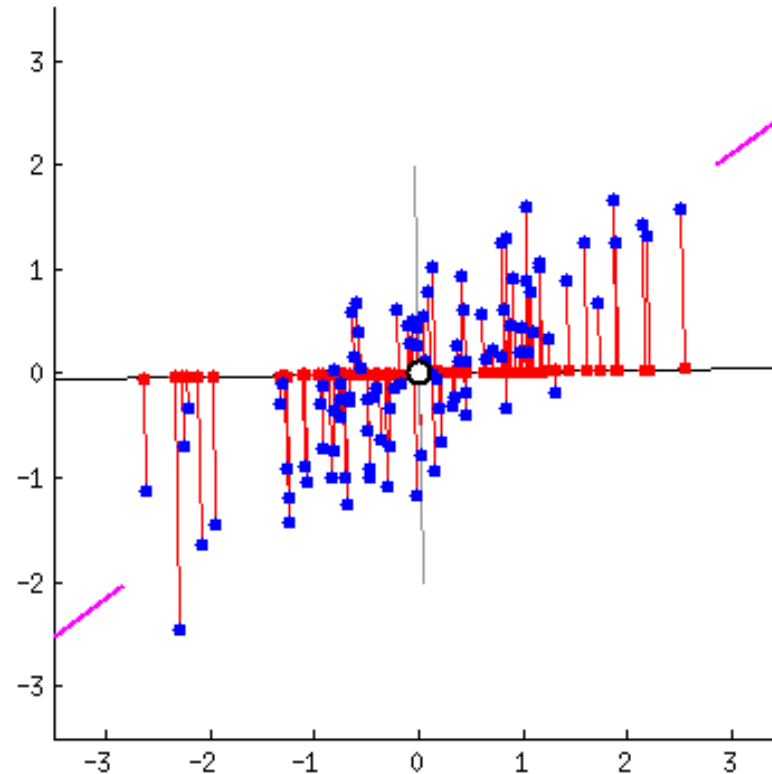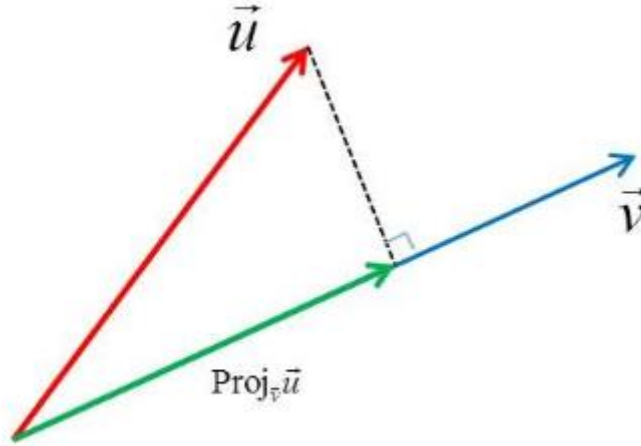  - Curse of dimensionality!

# *PCA*

## (PRINCIPAL COMPONENT ANALYSIS)

# Principal Components

- Goal: Trying to find the directions with the most variance (after projection).

- The direction with the most variance is called *first principal component (PC1)*, the direction with the second highest variance and *orthogonal* to PC1 is called *second principal component (PC2)* and so on.

# Projection Along a Direction



$$proj_v u = \left( \frac{u.v}{\|v\|^2} \right) v$$
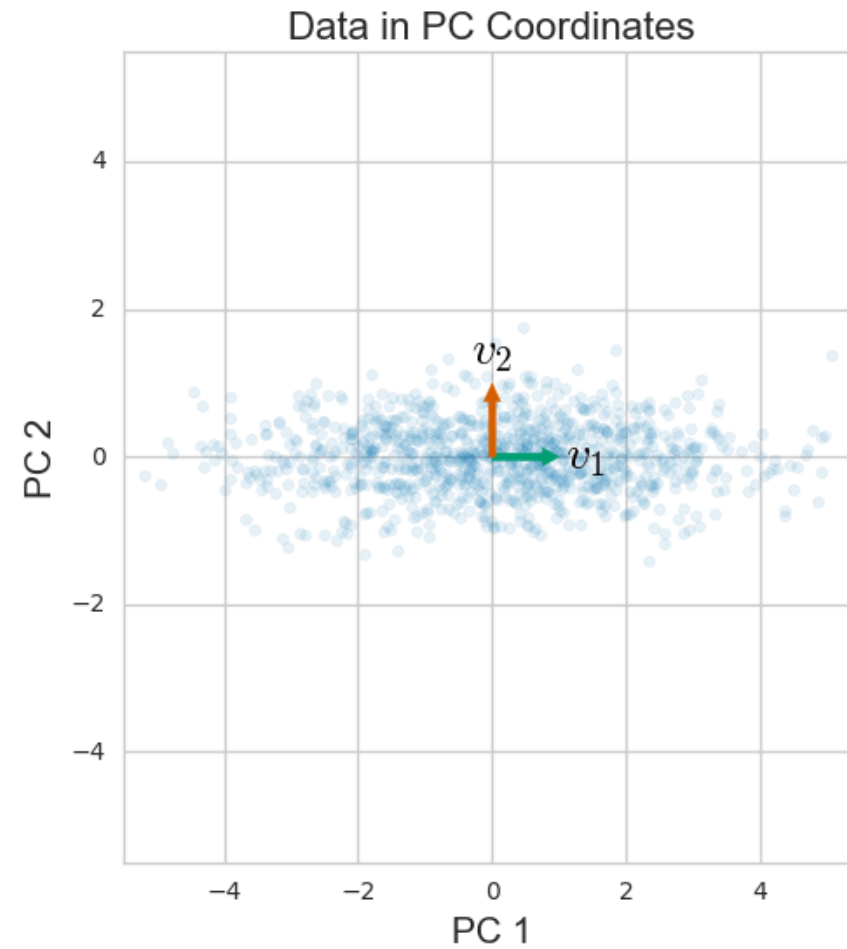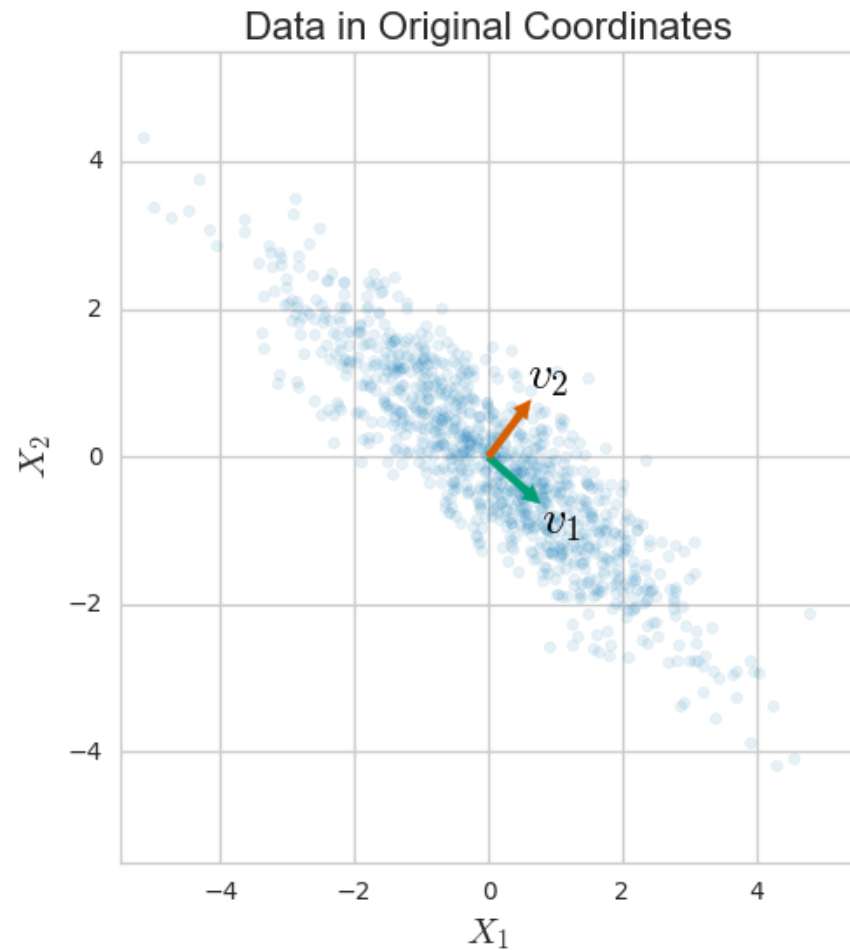
Since we care about the direction itself, we can consider norm of $v$ to be one:
$$\|v\|_2 = v^T v = 1$$

→ So the principal component score is:
$$PC\ score = \langle u, v \rangle = u^T v$$

# PC Transform

# How to find PC1, PC2, ... ?

The projected samples are: $\boldsymbol{v}^T \boldsymbol{x}^{(i)}$ with the mean: $\boldsymbol{v}^T \overline{\boldsymbol{x}}$

Considering all the data matrix $\boldsymbol{X}$:

$$var(\boldsymbol{Xv}) = cov(\boldsymbol{Xv}, \boldsymbol{Xv}) = \mathbb{E}\left[(\boldsymbol{Xv} - \mathbb{E}(\boldsymbol{Xv}))(\boldsymbol{Xv} - \mathbb{E}(\boldsymbol{Xv}))^T\right] =$$

$$\boldsymbol{v}^T \mathbb{E}\left[(\boldsymbol{X} - \mathbb{E}(\boldsymbol{X}))(\boldsymbol{X} - \mathbb{E}(\boldsymbol{X}))^T\right)\boldsymbol{v} = \boldsymbol{v}^T \boldsymbol{C} \boldsymbol{v}$$

where $\boldsymbol{C} = cov(\boldsymbol{X})$ is a symmetric matrix as follows:

$$\boldsymbol{C} = \begin{bmatrix} var(\boldsymbol{X_1}) & \cdots & cov(\boldsymbol{X_1}, \boldsymbol{X_p}) \\ \vdots & \ddots & \vdots \\ cov(\boldsymbol{X_p}, \boldsymbol{X_1}) & \cdots & var(\boldsymbol{X_p}) \end{bmatrix}$$

So our problem would be:

$$\max_{\boldsymbol{v}} \quad \boldsymbol{v}^T \boldsymbol{C} \boldsymbol{v}$$
$$s.t. \ \|\boldsymbol{v}\|_2 = 1$$

# Finding PCs

$$\max_{v} \quad v^T C v$$
$$s.t. \ ||v||_2 = 1$$

using lagrange multipliers:

$$\mathcal{L}(v, \lambda) = v^T C v + \lambda(1 - v^T v)$$

$$\nabla \mathcal{L}_v = \frac{\partial \mathcal{L}(v, \lambda)}{\partial v} = 0 \ \rightarrow 2Cv - 2\lambda v = 0 \rightarrow Cv = \lambda v$$

Since $Cv = \lambda v$, $v$ is the unit eigenvector of covariance matrix with the eigenvalue $\lambda$.

$$\rightarrow var(Xv) = v^T C v = v^T \lambda v = \lambda v^T v = \lambda$$

$\Rightarrow$ the first k-components would be the first k *unit eigenvectors* corresponding to the biggest k *eigenvalues* of the covariance matrix of **X**.

# Proportion of Variance Explained (PVE)

Explained variance can be represented as a **function of ratio of related eigenvalue and sum of eigenvalues of all eigenvectors.**
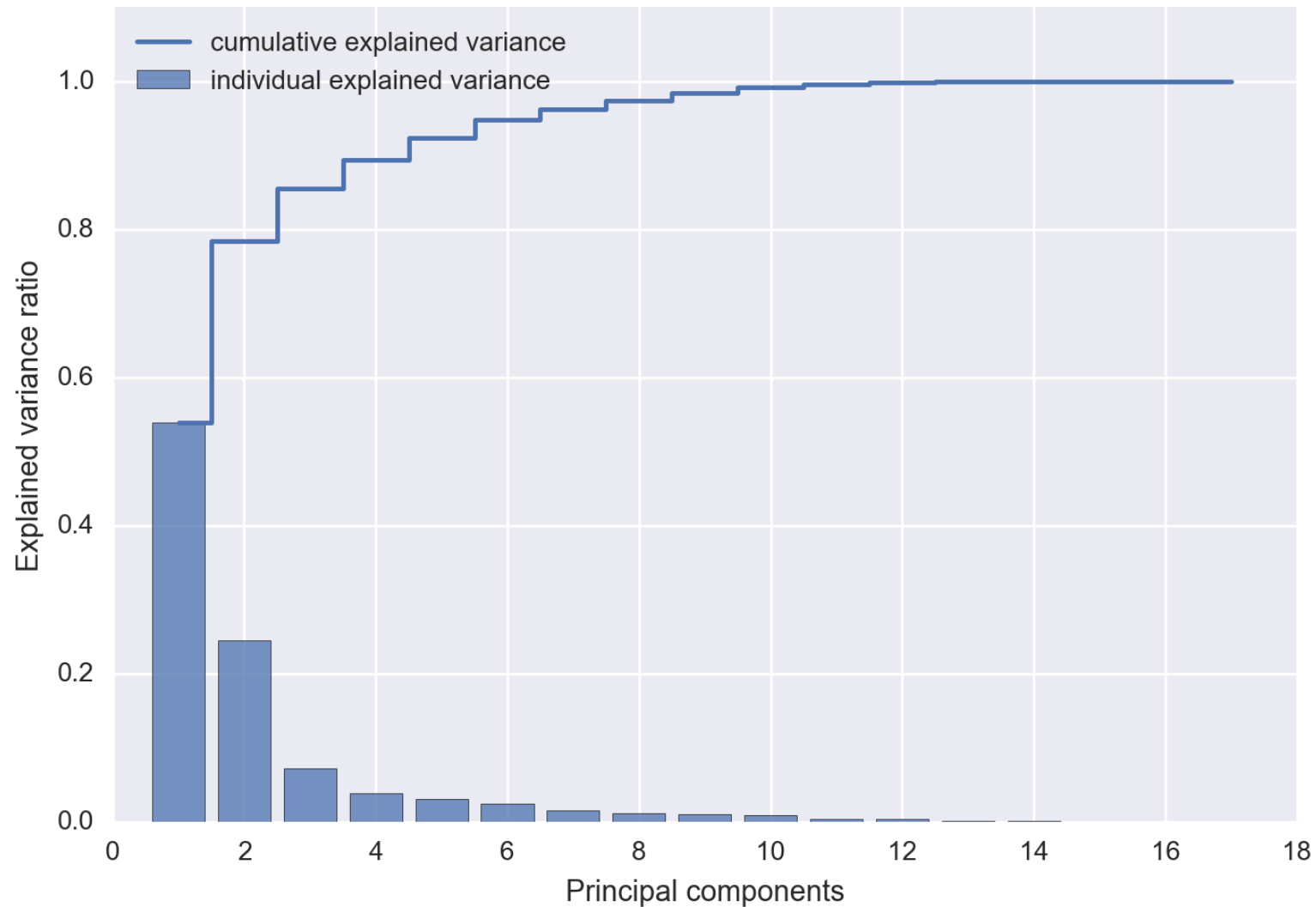
The total variance of $X$ can be explained using all the eigenvectors of covariance matrix:

$$tot_{var}(\boldsymbol{X}) = tr\big(cov(\boldsymbol{X})\big) = \sum_{i=1}^{p} \lambda_i$$

So the explained variance of the k-th eigenvector is:

$$PVE(k) = \frac{\lambda_k}{\sum_{i=1}^{p} \lambda_i}$$

# Choosing the new dimension using PVE

# Transform into PC coordinates

- $C$ with eigenvalue decomposition:

$$C = VDV^T = [v_1 \quad v_2 \quad \cdots \quad v_p] \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_p \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_p^T \end{bmatrix}$$

- After choosing the k (with PVE or any other method), the rotation matrix $W$ with the first k eigenvectors would be:

$$W = [v_1 \quad v_2 \quad \cdots \quad v_k]$$

- Samples are transformed from the original feature space into new PC coordinates as:

$$x_{new}^{(i)} = y^{(i)} = W^T x^{(i)}$$

- Considering all the samples, we can write:

$$X_{new} = W^T X$$

# Reconstruction using PCA

$$PCA\ reconstuction = PC\ scores \times eigenvectors^T + mean$$

$$\hat{x}^{(i)} = \sum_{j=1}^{k} y_j^{(i)} v_j + \overline{x}$$
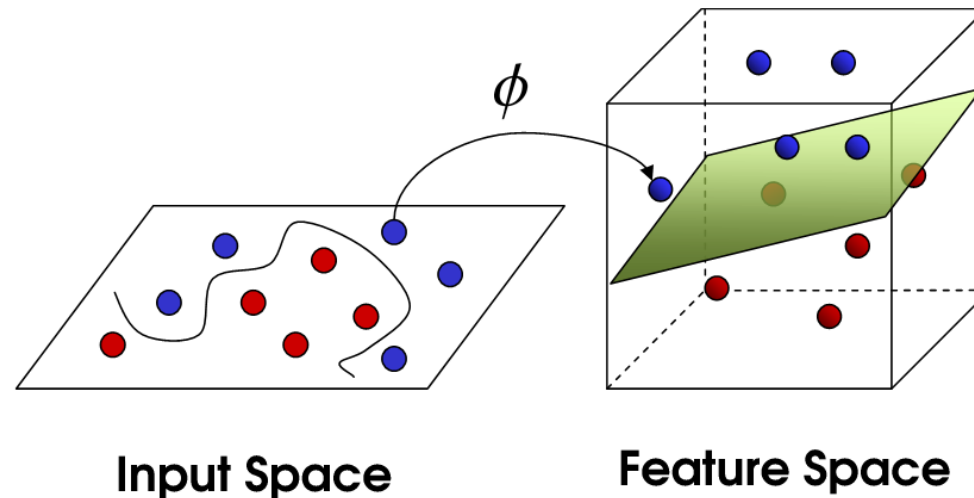


Original faces

Recovered faces

# KERNEL PCA

## (K-PCA)

# Making PCA non-Linear

- Suppose that instead of using the samples themselves, we wanted to go to some different feature space $\phi(x^{(i)}) \in \mathbb{R}^p$

- In the new feature space, then we can do PCA.
- This result will be non-linear in the original data space!

- Similar to SVM kernels idea!



**Input Space**          **Feature Space**

# Kernels & Inner Products

- Suppose the zero-mean samples, then covariance matrix would be:
$$C = X^T X$$

- Also, the inner product matrix is:

$$K = \begin{bmatrix} x^{(1)^T} x^{(1)} & \cdots & x^{(1)^T} x^{(N)} \\ \vdots & \ddots & \vdots \\ x^{(N)^T} x^{(1)} & \cdots & x^{(N)^T} x^{(N)} \end{bmatrix} = XX^T$$

where $K_{i,j}$ is the inner product of $i^{th}$ & $j^{th}$ samples.

This can be generalized into other kernels, such as RBF kernel as follows:

$$K = \begin{bmatrix} \langle \phi(x^{(1)}), \phi(x^{(1)}) \rangle & \cdots & \langle \phi(x^{(1)}), \phi(x^{(N)}) \rangle \\ \vdots & \ddots & \vdots \\ \langle \phi(x^{(N)}), \phi(x^{(1)}) \rangle & \cdots & \langle \phi(x^{(N)}), \phi(x^{(N)}) \rangle \end{bmatrix}$$

where $K_{i,j} = K(x^{(i)}, x^{(j)})$

# Covariance Matrix & Inner Products

- Suppose $\boldsymbol{u}$ is an eigenvector of $\boldsymbol{K}$:

$$\boldsymbol{Ku} = \mu\boldsymbol{u} \rightarrow \boldsymbol{XX^T u} = \mu\boldsymbol{u}$$

then:

$$\boldsymbol{X^T XX^T u} = \mu\boldsymbol{X^T u} \rightarrow \boldsymbol{CX^T u} = \mu\boldsymbol{X^T u}$$

- We only need inner products!

To find the PC coefficients for a sample $x$:

$$y = \boldsymbol{x^T}(\boldsymbol{X^T u}) = (\boldsymbol{Xx})^T\boldsymbol{u} = [\boldsymbol{x^T x^{(1)}} \quad \boldsymbol{x^T x^{(2)}} \quad \dots \quad \boldsymbol{x^T x^{(N)}}]\boldsymbol{u}$$

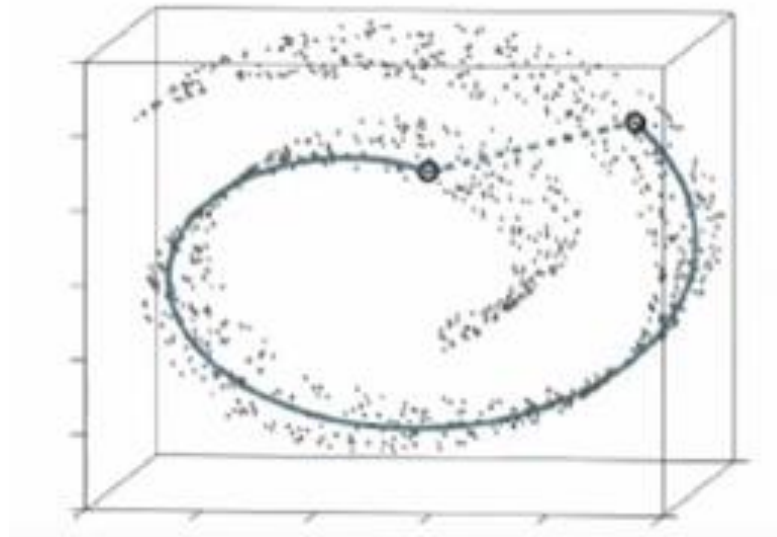- If K-PCA, we always have to use normalized kernel matrix!

# *T-SNE*

## (T-DISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING)

# t-SNE: Motivation

t-Distributed Stochastic Neighbor Embedding (t-SNE) *(L.J.P. van der Maaten and G.E. Hinton, 2008)* is a technique for dimensionality reduction that is particularly well suited for the visualization of high-dimensional datasets.

The t-SNE algorithm calculates a similarity measure between pairs of instances in the high dimensional space and in the low dimensional space. It then tries to optimize these two similarity measures using a cost function (KL divergence).

# PCA vs t-SNE