# Dimensionality Reduction

* ## Motivation:

  - Complexity of most estimators' depends on the number of inputs.

  - Affects time and space complexity.

  - Reducing the dimensionality of the input lowers the complexity.

  Goal: Porune part of the feature space that contribute little to the solution.

* ## Arguments for Dimensionality Reduction for Classifications

  - Reduces complexity of the classifier.

  - Irrelevant dimensions add to the variance.

  - If data is explained with fewer features, we can get a better idea about the underlying process.

  - What is relevant for computing a solution to our problem?

**\* Feature Selection vs Feature Extraction**

(Feature Selection) → We try to find the K out of D dimensions that contain most of the information. We discard the other (D-K) dimensions.

(Feature Extraction) → We try to find a new set of K dimensions that are combindies of the D dimensions and yield most of the information.

**\* Popular Dimensionality Reduction Technique for Feature Extraction**

- **PCA** : Principal Component Analysis
- **Fisher-LDA** : Fisher's Linear Discriminant Analysis.
- **LLE:** Locally Linear Embedding

# * Principal Component Analysis
## (PCA)

### # Idea of PCA

→ Find the mapping from the original $D$ to $K$ dimensional space with minimum loss of information $(K < D)$.

→ PCA analysis the spread of the data and tries to find new dimensions that cover the spread best.

### # Data Matrix

→ Matrix of $N$ vectors with $D$ dimensions

$$\underset{N \times D}{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \longrightarrow \begin{array}{l} D\text{-dimensionals} \\ \text{Vectors.} \end{array}$$

→ Goal: Find $K$ dimensions that represents the data as good as possible.

$$\underset{N \times D}{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \Rightarrow \underset{N \times K}{W} = \begin{bmatrix} w_1^T \\ \vdots \\ w_N^T \end{bmatrix}$$

* <u>Mean Reduced Features</u>

$$\mu_x = \frac{1}{N} \sum_{m=1}^{N} x_m = \frac{1}{N} X^T 1_N$$

$\Rightarrow$ Mean reduced features

$$\bar{x}_m = x_m - \mu_x$$

$\rightarrow$ In matrix form

$$\bar{X} = X - 1_N \mu_x^T$$

* <u>Covariance Matrix of Mean-Reduced Features</u>

$$\Sigma_{xx} = \frac{1}{N-1} \bar{X}^T \bar{X}$$

$\Rightarrow$ We describe our data using first and second central moment in the D-dimensional space. (i.e. $\mu, \Sigma_{xx}$)

$$\downarrow$$
$$D\text{-dimensional}$$

<u>Question</u>: What are the best $K \ll D$ dimensions to approximate the data?

**★ Eigenvector and Eigenvalue**

→ The Eigenvector $u_1$ Corresponding to the largest the larger Eigenvalue of $E_{xx}$ is the direction of maximum spread.

→ Make this Eigenvector $u_1$ the first principal component.

→ All other Eigenvectors are orthogonal to $u_1$.

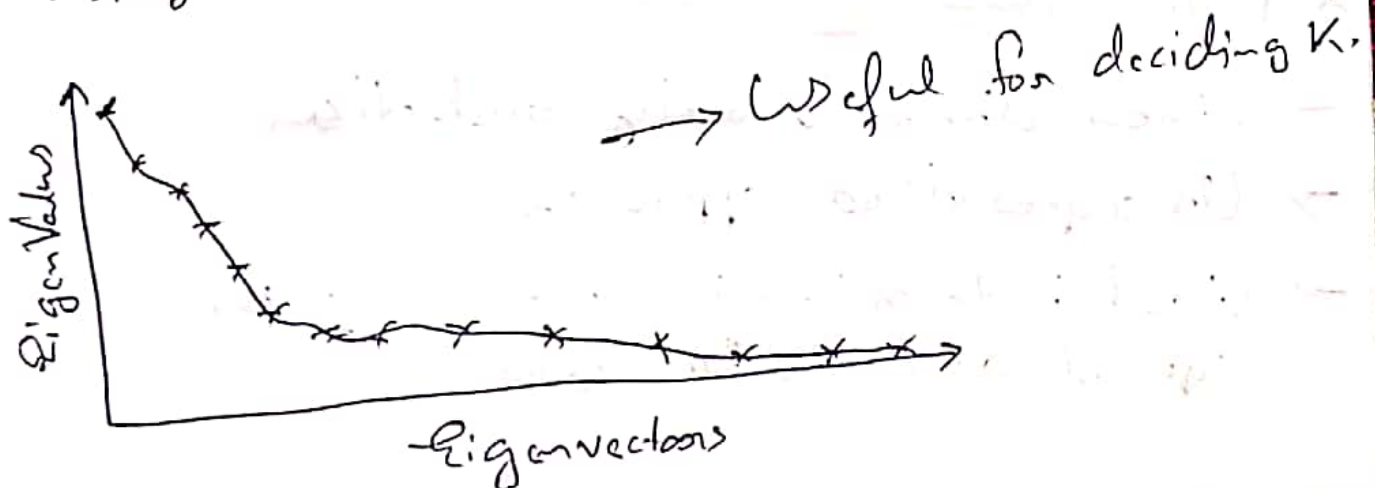→ Repeat the process $k$ times for the remaining Eigenvalues/Vectors.

**★ Eigenvalue Decomposition**

⇒ Eigenvalue decomposition yields:

$$\Sigma_{xx} = R S^2 R^T$$

rotation matrix

diagonal matrix

⇒ Eigenvalues/Vectors are computed via SVD - Single value decomposition (or special variants).

→ Useful for deciding $K$.



Eigenvalues (y-axis) vs Eigenvectors (x-axis)

**\* <u>Mapping to the Reduced Space</u>**

$$x \rightarrow [u_1^T(x-\mu), \cdots - u_k^T(x-\mu)]$$

$$\downarrow$$

$$= [\omega_1, \cdots \omega_k]$$

$\searrow$ K dimensional data points

$\downarrow$

D dimensional data point

**\* <u>Mapping to the Original Space</u>**

$\Rightarrow$ We can also map from the reduced K-dimensional space to the original one.

$$\hat{x} = \mu + \sum_{i=1}^{K} \omega_i u_i$$

$\Rightarrow$ which yields the squared reconstruction error:

$$e(x) = \|\hat{x} - x\|^2$$

**\* <u>PCA Summary</u>**

$\rightarrow$ Linear dimensionality reduction

$\rightarrow$ Unsupervised approach

$\rightarrow$ Goal is to minimize the sum of the squared reconstruction error

→ Principal Components are the directions of maximum spread.

→ Computed via Eigenvalues/vectors of the Covariance matrix of the training data Points.

——————✗———————✗————————

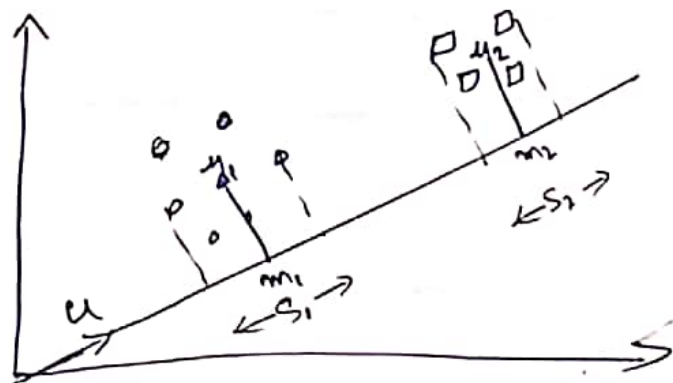## Fisher's Linear Discriminant Analysis

Limitation of PCA

↳ The direction of maximum variance is not always good for classification.

⇒ Fisher's LDA tries to find the best direction to separate classes.

★ Idea of Fisher-LDA

⇒ Find the $u$ that
- maximizes $\|m_1 - m_2\|$
- minimizes $S_1$ and $S_2$



★ LDA for 2 classes and $k=1$

⇒ Compute the means of the classes

$$m_1 = \frac{\sum_t u^T x^t c^t}{\sum_t c^t} \qquad m_2 = \frac{\sum_t u^T x^t (1-c^t)}{\sum_t (1-c^t)}$$

$$c^t = 1 \rightarrow m_1 \qquad c^t = 0 \rightarrow m_2$$

$\Rightarrow$ and the scatter of samples in the 1-dim space $(s_i^2 = \sigma_i^2 (\Sigma_t c^t - 1))$

$$S_1^2 = \sum_t (u^T x^t - m_1)^2 c^t$$

$$S_2^2 = \sum_t (u^T x^t - m_2)^2 (1 - c^t)$$

* ## Objective Function

$$J(u) = \frac{(m_1 - m_2)^2}{S_1^2 + S_2^2}$$

$\searrow$ maximize for LDA

$$(m_1 - m_2)^2 = (u^T \mu_1 - u^T \mu_2)^2$$

$$= u^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T u$$

$$= u^T S_B u$$

$\searrow$ between-class scatter matrix

$$S_1^2 = \sum_t (u^T x^t - m_1)^2 c^t$$

$$= u^T \sum (x^t - \mu_1)(x^t - \mu_2)^T c^t \, u$$

$$= u^T S_1 u$$

$\searrow$ within-class scatter matrix of class $i$.

$\Rightarrow$ Similarly for $S_2^2$.

$$S_1^2 + S_2^2 = u^T S_1 u' + u^T S_2 u$$
$$= u^T (S_1 + S_2) u$$
$$= u^T (S_u) u$$

↳ total within-class Scatter matrix

$$J(u) = \frac{u^T S_B u}{u^T S_u u}$$

$$u^* = \underset{u}{argmax} \ J(u) = S_u^{-1} (\mu_2 - \mu_1)$$

→ Fisher-LDA is the optimal solution if both features are normally distributed.

→ Also applicable for non-normally distributed features.

→ Can be easily generalized to N classes and k>1.

→ Linear dimensionality reduction.

* <u>Fischer's LDA vs PCA</u>

▫ PCA minimizes the reconstruction error

▪ PCA is the standard choice for unsupervised Problem (no labels)

* Fisher-LDA exploite class labels to find a subspace so that separates the classes as good as possible.

---

* ## Locally Linear Embedding (LLE)

* Technique for un supervised non-linear dimensionality reduction

* Compute for each input data point a coordinate on a low-dimensional manifold.

* ## LLE Key Steps

1. Step: Find neighbors for each input point.

2. Step: Compute for each input point a weight vector that best recovers the point itself from its neighbors.

3. Step: For each point, find latent coordinates such that the same weights can be used for re-construction.