# CS182: Introduction to Machine Learning – Dimensionality Reduction

Yujiao Shi

SIST, ShanghaiTech

Spring, 2025

Recall:
$K$-means
Algorithm

- Input: $\mathcal{D} = \{(\boldsymbol{x}^{(n)})\}_{n=1}^{N}, K$

1. Initialize cluster centers $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_K$

How do we set these hyperparameters?

2. While NOT CONVERGED

   a. Assign each data point to the cluster with the nearest cluster center:
   $$z^{(n)} = \underset{k}{\mathrm{argmin}} \left\| \boldsymbol{x}^{(n)} - \boldsymbol{\mu}_k \right\|_2$$

   b. Recompute the cluster centers:
   $$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \,: z^{(n)} = k} \boldsymbol{x}^{(n)}$$
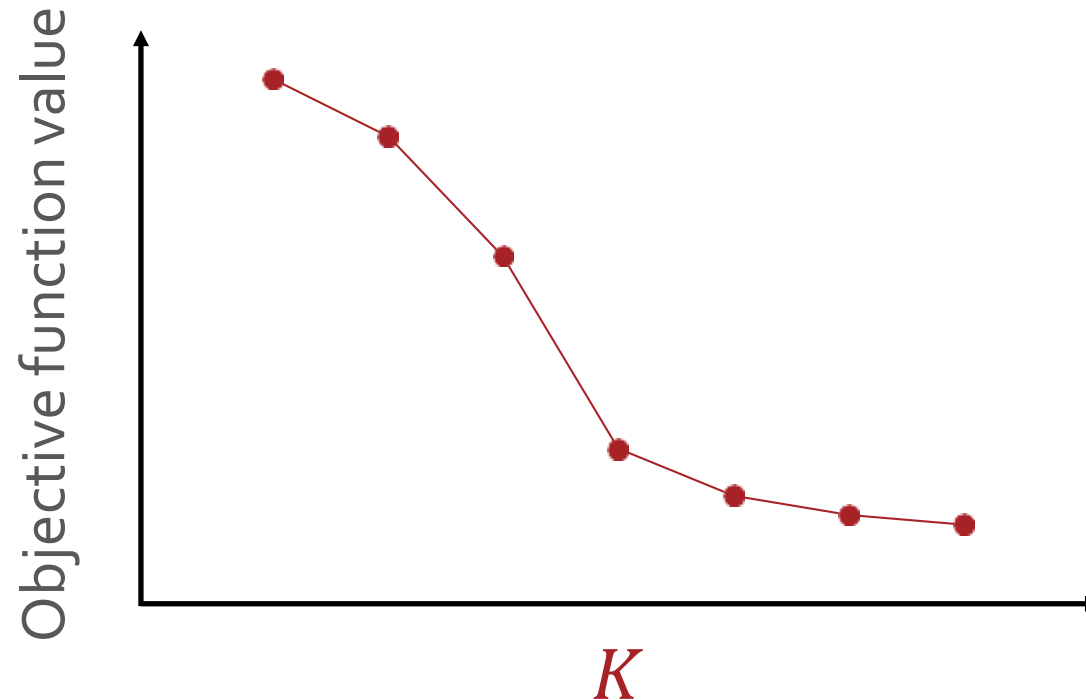
   where $N_k$ is the number of data points in cluster $k$

- Output: cluster centers $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_K$ and cluster assignments $z^{(1)}, \ldots, z^{(N)}$

# Setting $K$

- Idea: choose the value of $K$ that minimizes the objective function



- Better Idea: look for the characteristic "elbow" or largest decrease when going from $K - 1$ to $K$

上 海 科 技 大
ShanghaiTech University

# Initializing $K$-means: Lloyd's Method

- Randomly choose $K$ data points to be the initial cluster centers

上 海 科 技 大
ShanghaiTech University

## Initializing $K$-means: Lloyd's Method

- Randomly choose $K$ data points to be the initial cluster centers

上 海 科 技 大
**ShanghaiTech University**

# Initializing $K$-means: Lloyd's Method

- Randomly choose $K$ data points to be the initial cluster centers

上 海 科 技 大
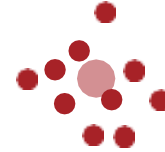ShanghaiTech University

# Initializing $K$-means: Lloyd's Method
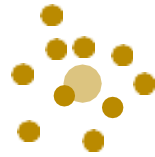
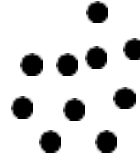- Randomly choose $K$ data points to be the initial cluster centers

上 海 科 技 大
ShanghaiTech University

## Initializing $K$-means: Furthest Point

1. Choose the first cluster center randomly from the data points

2. For each other data point $x$, compute $D(x)$, the distance between $x$ and the closest cluster center to $x$

3. Select the data point with the largest $D(x)$ as the next cluster center

4. Repeat 2 and 3 $K-1$ times

上 海 科 技 大
ShanghaiTech University

## Initializing $K$-means: Furthest Point

1. Choose the first cluster center randomly from the data points

2. For each other data point $x$, compute $D(x)$, the distance between $x$ and the closest cluster center to $x$

3. Select the data point with the largest $D(x)$ as the next cluster center

4. Repeat 2 and 3 $K-1$ times

上 海 科 技 大
ShanghaiTech University

## Initializing $K$-means: Furthest Point

1. Choose the first cluster center randomly from the data points

2. For each other data point $x$, compute $D(x)$, the distance between $x$ and the closest cluster center to $x$

3. Select the data point with the largest $D(x)$ as the next cluster center

4. Repeat 2 and 3 $K - 1$ times

- Works great in the case of well-clustered data!

- Can struggle with outliers…

上海科技大
ShanghaiTech University

# Initializing $K$-means: $K$-means++

1. Choose the first cluster center randomly from the data points

2. For each other data point $x$, compute $D\ (x)$, the distance between $x$ and the closest cluster center to $x$

3. *Sample* the next cluster center with probability proportional to $D(x)^2$

4. Repeat 2 and 3 $K - 1$ times

| $i$ | $D\left(x^{(i)}\right)$ | $D\left(x^{(i)}\right)^2$ | Probability of Being Selected |
|---|---|---|---|
| 1 | 4 | 16 | 16/123 |
| 2 | 7 | 49 | 49/123 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| $N$ | 1 | 1 | 1/123 |
| **Total** | | 123 | $123/123 = 1$ |

上 海 科 技 大
ShanghaiTech University

# Initializing $K$-means: $K$-means++

1. Choose the first cluster center randomly from the data points

2. For each other data point $x$, compute $D(x)$, the distance between $x$ and the closest cluster center to $x$

3. *Sample* the next cluster center with probability proportional to $D(x)^2$
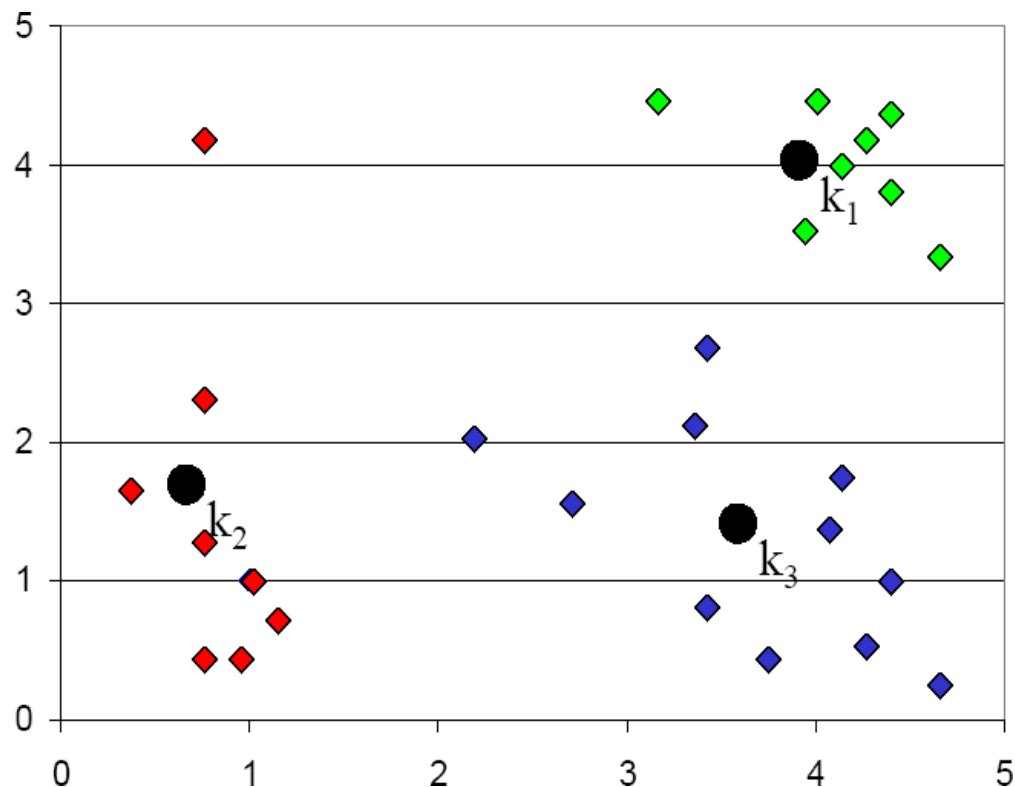
4. Repeat 2 and 3 $K-1$ times

- $K$-means++ achieves a $O(\log K)$ approximation to the optimal clustering in expectation!

- All initialization methods can benefit from multiple random restarts

# Problem with K-means

|  | Cluster 1 | Cluster 2 | Cluster 3 |
|---|---|---|---|
| Individual 1 | 1 | 0 | 0 |
| Individual 2 | 0 | 1 | 0 |
| Individual 3 | 0 | 1 | 0 |
| Individual 4 | 1 | 0 | 0 |
| Individual 5 | … | … | … |
| Individual 6 | … | … | … |
| Individual 7 | … | … | … |
| Individual 8 | … | … | … |
| Individual 9 | … | … | … |
| Individual 10 | … | … | … |

Hard assignment of samples into three clusters

# Probabilistic Sof-Clustering of Samples into Three Clusters

| Probability of | Cluster 1 | Cluster 2 | Cluster 3 | Sum |
|---|---|---|---|---|
| Individual 1 | 0.1 | 0.4 | 0.5 | 1 |
| Individual 2 | 0.8 | 0.1 | 0.1 | 1 |
| Individual 3 | 0.7 | 0.2 | 0.1 | 1 |
| Individual 4 | 0.10 | 0.05 | 0.85 | 1 |
| Individual 5 | … | … | … | 1 |
| Individual 6 | … | … | … | 1 |
| Individual 7 | … | … | … | 1 |
| Individual 8 | … | … | … | 1 |
| Individual 9 | … | … | … | 1 |
| Individual 10 | … | … | … | 1 |

• Each sample can be assigned to more than one clusters with a certain probability.
• For each sample, the probabilities for all clusters should sum to 1. (i.e., each row should sum to 1.)
• Each cluster is explained by a cluster center variable (i.e., cluster mean)

# Mixture Model

- A density model $\mathbf{p(x)}$ may be multi-modal.



Multi--model: how do we model this?

Unimodal -Gaussian

# Mixture Model

- We may be able to model it as a mixture of uni-modal distributions (e.g., Gaussians).
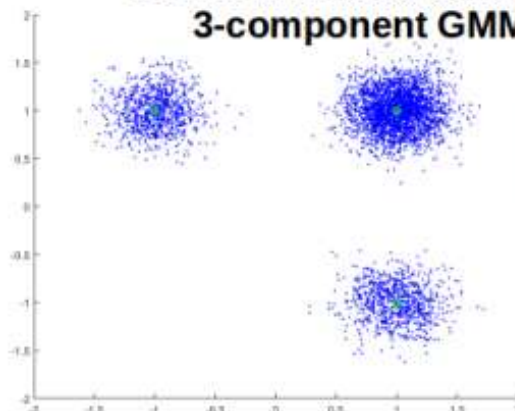- Each mode may correspond to a different sub-population (e.g., male and female).

上 海 科 技 大

**ShanghaiTech University**

Recall...

- The generative story for each $\boldsymbol{x}_n$, $n = 1, 2, \ldots, N$

  - First choose one of the $K$ mixture components as

    $$\boldsymbol{z}_n \sim \text{Multinomial}(\boldsymbol{z}_n | \pi) \qquad \text{(from the prior } p(\boldsymbol{z}) \text{ over } \boldsymbol{z})$$

  - Suppose $\boldsymbol{z}_n = k$. Now generate $\boldsymbol{x}_n$ from the $k$-th Gaussian as

    $$\boldsymbol{x}_n \sim \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \qquad \text{(from the data distr. } p(\boldsymbol{x}|\boldsymbol{z}))$$

**Some simulated data from a 3-component GMM**

**Directed Graphical Model for a K-component GMM**

Note: Arrow-heads point towards the dependent nodes in a directed graphical model

Shaded nodes: Observed

White nodes: Unknowns

*Recall...*

- Consider the 'incomplete" data log likelihood

$$\log p(\mathbf{X}|\Theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some dist.)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} \quad \text{(concave } f, \text{ Jensen's Ineq.: } f(\sum \lambda_i x_i) \geq \sum \lambda_i f(x_i))$$

$$\log p(\mathbf{X}|\Theta) \geq \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \Theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.}$$

- If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta) p(\mathbf{X}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}|\Theta)$$

$$= \log p(\mathbf{X}|\Theta) \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) = \log p(\mathbf{X}|\Theta)$$

- Thus for $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \Theta)$, we have

$$\log p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta) \log p(\mathbf{X}, \mathbf{Z}|\Theta) + \text{const.} = \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] + \text{const.}$$

- Thus $\log p(\mathbf{X}|\Theta)$ is tightly lower-bounded by $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)]$ which EM maximizes

Recall...

Initialize the parameters: $\Theta^{old}$. Then alternate between these steps:

- **E (Expectation) step:**
  - Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ over latent variables $\mathbf{Z}$ using $\Theta^{old}$
  - Compute the expected complete data log-likelihood w.r.t. *this* posterior

$$\mathcal{Q}(\Theta, \Theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \Theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

- **M (Maximization) step:**
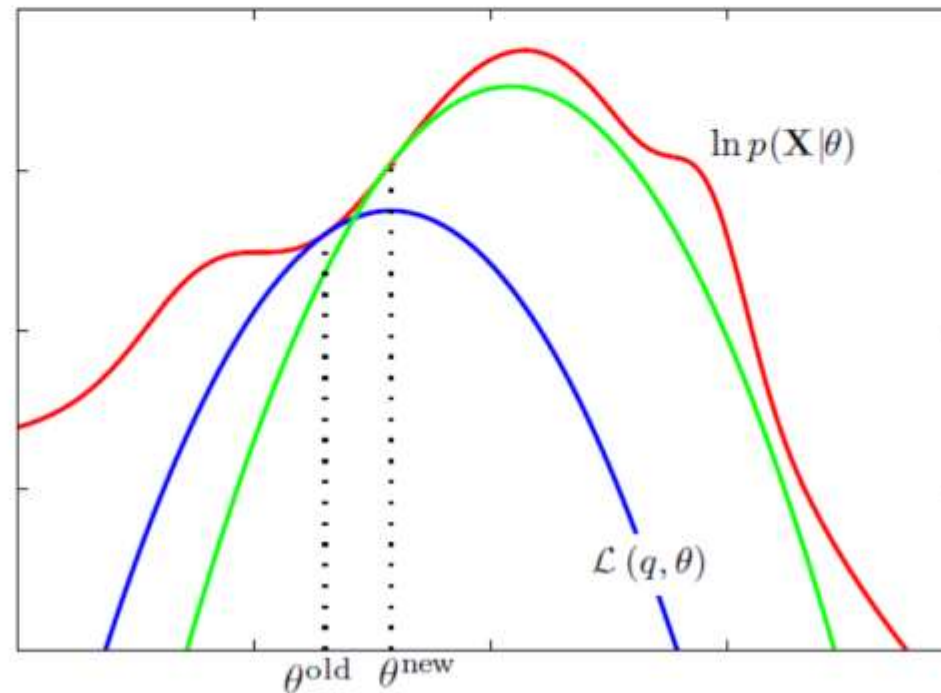  - Maximize the expected complete data log-likelihood w.r.t. $\Theta$

$$\Theta^{new} = \arg\max_{\Theta} \mathcal{Q}(\Theta, \Theta^{old}) \quad \text{(if doing MLE)}$$

$$\Theta^{new} = \arg\max_{\Theta}\{\mathcal{Q}(\Theta, \Theta^{old}) + \log p(\Theta)\} \quad \text{(if doing MAP)}$$

- If the log-likelihood or the parameter values not converged then set $\Theta^{old} = \Theta^{new}$ and go to the E step.

上 海 科 技 大

ShanghaiTech University

Recall...

- E-step: Update of $q$ makes the $\mathcal{L}(q, \Theta)$ curve touch the $\log p(\mathbf{X}|\Theta)$ curve
- M-step gives the maxima $\Theta^{new}$ of $\mathcal{L}(q, \Theta)$
- Next E-step readjusts $\mathcal{L}(q, \Theta)$ curve (green) to meet $\log p(\mathbf{X}|\Theta)$ curve again
- This continues until a local maxima of $\log p(\mathbf{X}|\Theta)$ is reached

# Learning GMM by EM

1. Initialize the Parameters $\Theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^{K}$ randomly, or using K-means

2. Iterate until convergence (e.g., when $\log p(X|\Theta)$ ceases to increase

   a) E-step:

$$\gamma_{ik} = p(z_i = k|x_i, \Theta^{old}) = \frac{\pi_k^{old} \mathcal{N}\left(x_i\middle|\mu_k^{old}, \Sigma_k^{old}\right)}{\sum_{j=1}^{K} \pi_j^{old} \mathcal{N}\left(x_i\middle|\mu_j^{old}, \Sigma_j^{old}\right)}$$

   b) M-step:

$$\pi_k^{new} = \frac{1}{N}\sum_{i=1}^{N} \gamma_{ik}$$

$$\mu_k^{new} = \frac{\sum_{i=1}^{N} \gamma_{ik} x_i}{\sum_{i=1}^{N} \gamma_{ik}}$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^{N} \gamma_{ik}(x_i - \mu_k^{new})(x_i - \mu_k^{new})^\mathsf{T}}{\sum_{i=1}^{N} \gamma_{ik}}$$

# DIMENSIONALITY REDUCTION

## Examples of high dimensional data:

– High resolution images (millions of pixels)

## Examples of high dimensional data:

– Multilingual News Stories
(vocabulary of hundreds of thousands of words)

## Examples of high dimensional data:
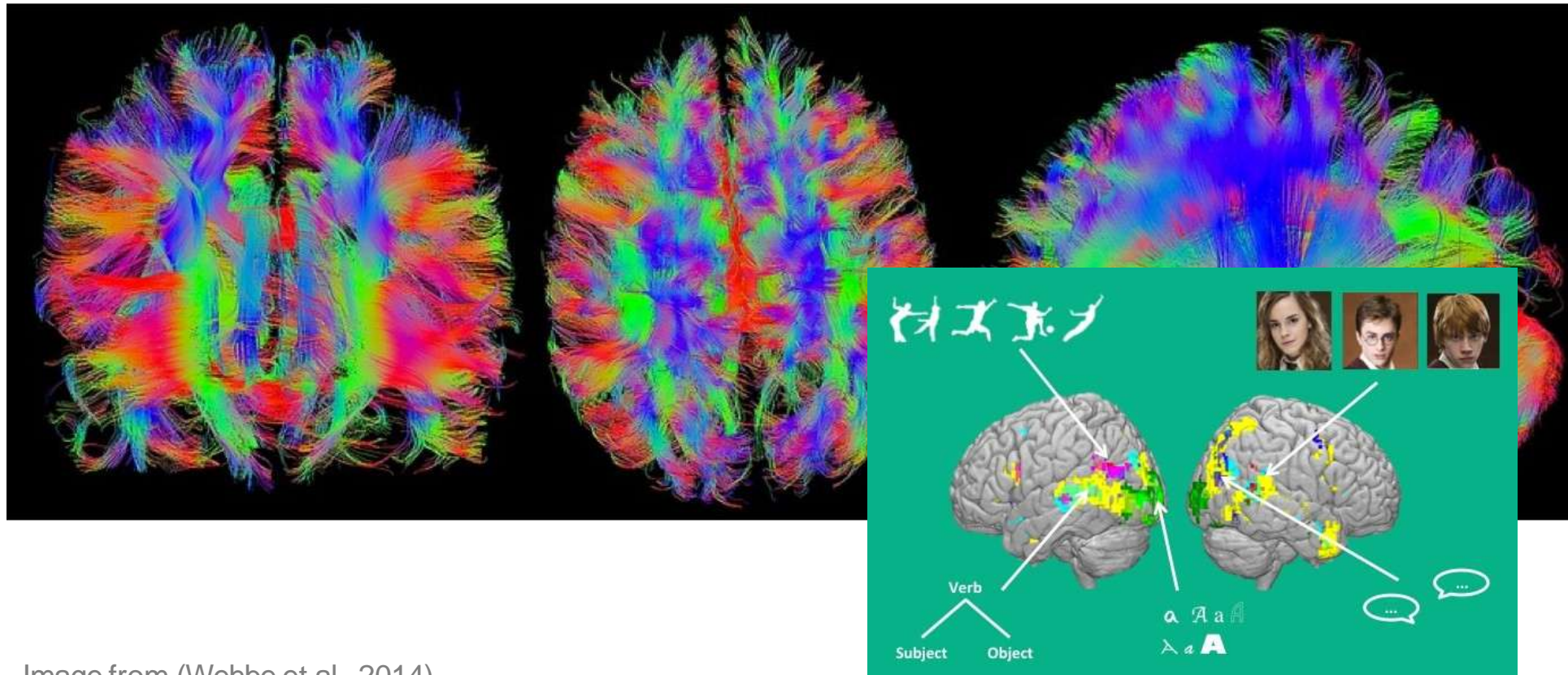– Brain Imaging Data (100s of MBs per scan)



Image from (Wehbe et al., 2014)

Image from  https://pixabay.com/en/brain-mrt-magnetic-resonance-imaging-1728449/
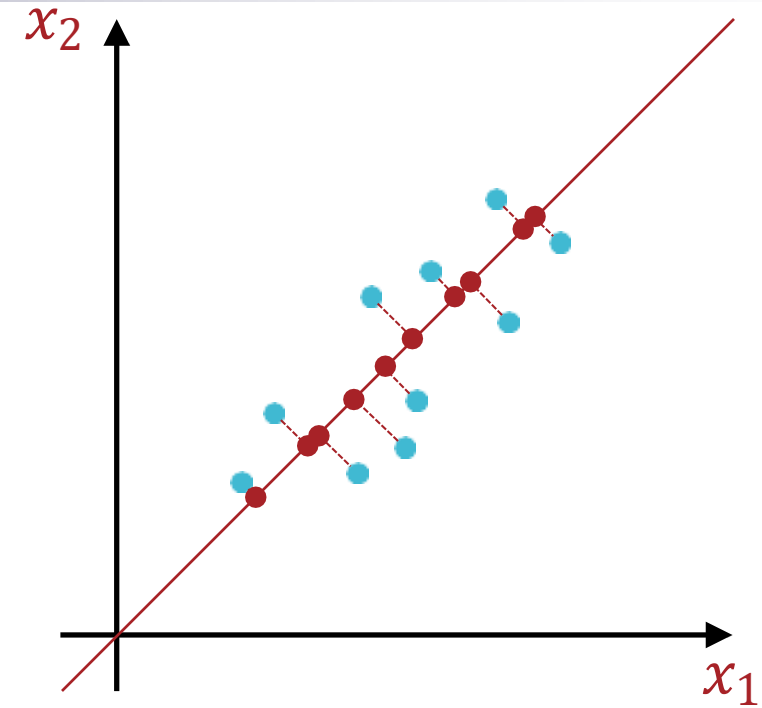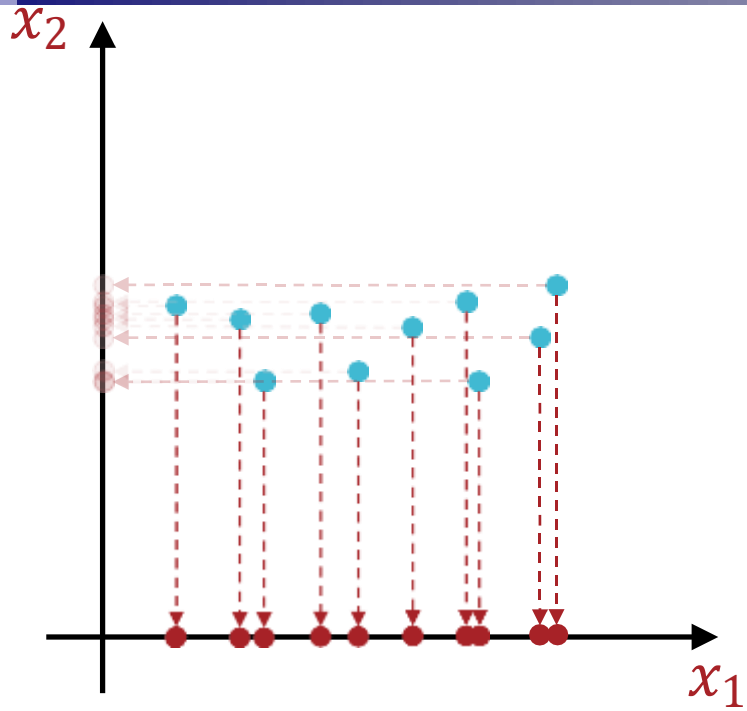
**Dimensionality Reduction Algorithms:**

Powerful (often unsupervised) learning techniques for extracting hidden (potentially lower dimensional) structure from high dimensional datasets.
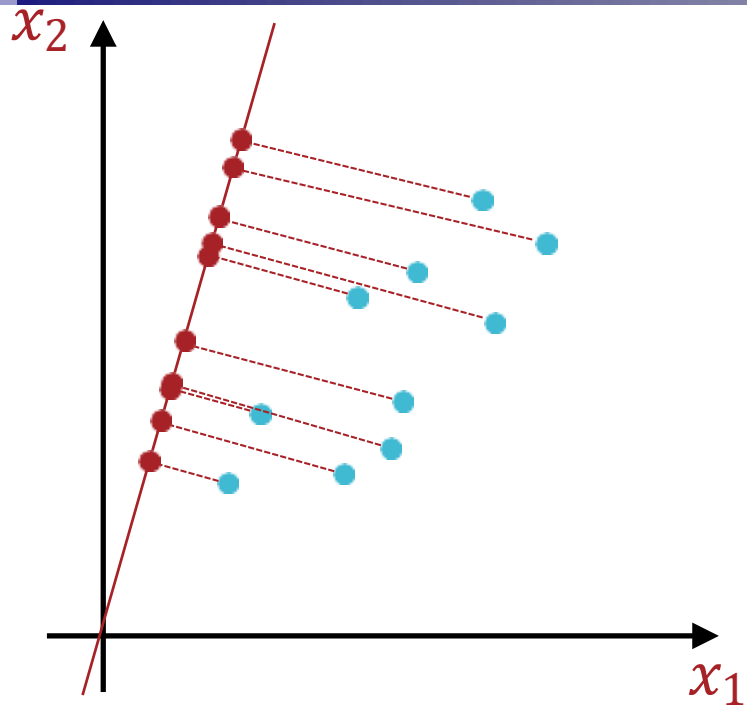
**Examples:**

PCA, Kernel PCA, ICA, CCA, t-SNE, Autoencoders, VAEs
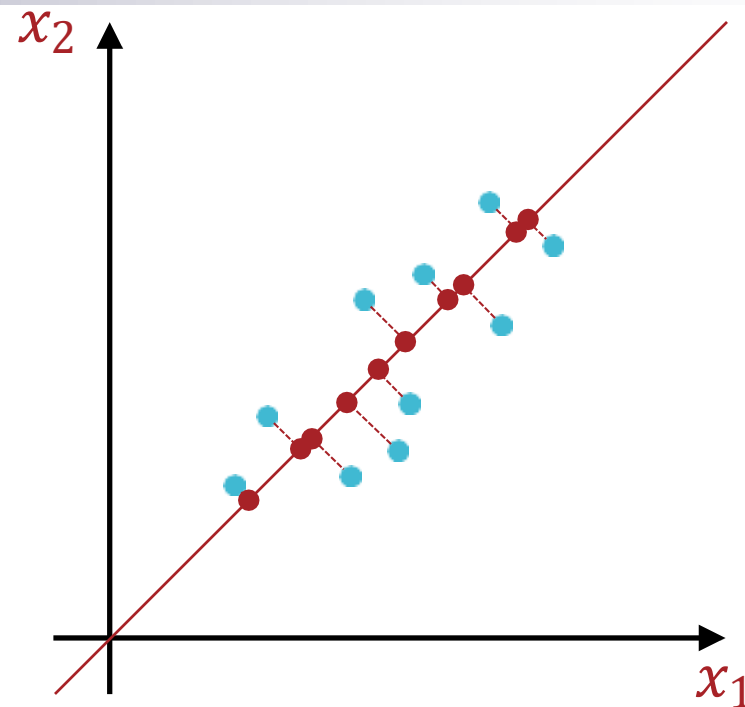
**Useful for**:
- Visualization
- More efficient use of resources (e.g., time, memory, communication)
- Statistical: fewer dimensions → better generalization
- Noise removal (improving data quality)

Feature Elimination ∈ Dimensionality Reduction

Option A                    Option B (TOXIC)                    Option C

# Which projection do you prefer (Q1) and why (Q2)?

## Background: Sample Variance and Covariance

- Given a collection of $N$ 1-dimensional samples $\left[x^{(1)}, x^{(2)}, \ldots, x^{(N)}\right]$ from some random variable, the **sample variance** is

$$\hat{\sigma}^2 = \frac{1}{N}\sum_{i=1}^{N}\left(x^{(i)} - \hat{\mu}\right)^2 = \frac{1}{N}\sum_{i=1}^{N}\left(x^{(i)} - \frac{1}{N}\sum_{n=1}^{N}x^{(n)}\right)^2$$

- Given a collection of $N$ $D$-dimensional samples $\left[\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots, \boldsymbol{x}^{(N)}\right]$ from some random variable, the **sample covariance** between dimension $j$ and $k$ is

$$\Sigma_{jk} = \frac{1}{N}\sum_{i=1}^{N}\left(x_j^{(i)} - \hat{\mu}_j\right)\left(x_k^{(i)} - \hat{\mu}_k\right) \text{ where } \hat{\mu}_d = \frac{1}{N}\sum_{n=1}^{N}x_d^{(n)}$$

# Background: Sample Variance and Covariance

- Given a collection of $N$ 1-dimensional samples $\left[x^{(1)}, x^{(2)}, \dots, x^{(N)}\right]$ from some random variable, the **sample variance** is

$$\hat{\sigma}^2 = \frac{1}{N}\sum_{i=1}^{N}\left(x^{(i)} - \hat{\mu}\right)^2 = \frac{1}{N}\sum_{i=1}^{N}\left(x^{(i)} - \frac{1}{N}\sum_{n=1}^{N}x^{(n)}\right)^2$$

- Given a collection of $N$ $D$-dimensional samples $\left[\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \dots, \boldsymbol{x}^{(N)}\right]$ from some random variable, the **sample covariance matrix** is

$$\Sigma = \frac{1}{N}X^T X \quad \text{where} \quad X = \begin{bmatrix} (\boldsymbol{x}^{(1)} - \boldsymbol{\mu})^T \\ (\boldsymbol{x}^{(2)} - \boldsymbol{\mu})^T \\ \vdots \\ (\boldsymbol{x}^{(N)} - \boldsymbol{\mu})^T \end{bmatrix}$$

上 海 科 技 大 学
ShanghaiTech University

# Centering the Data

- To be consistent, we will constrain principal components to be *orthogonal unit vectors* that begin at the origin

- Preprocess data to be centered around the origin:

1. $\boldsymbol{\mu} = \dfrac{1}{N}\sum\limits_{n=1}^{N} \boldsymbol{x}^{(n)}$

2. $\widetilde{\boldsymbol{x}}^{(n)} = \boldsymbol{x}^{(n)} - \boldsymbol{\mu} \; \forall \, n$

3. $X = \begin{bmatrix} \widetilde{\boldsymbol{x}}^{(1)^T} \\ \widetilde{\boldsymbol{x}}^{(2)^T} \\ \vdots \\ \widetilde{\boldsymbol{x}}^{(N)^T} \end{bmatrix}$

## Reconstruction Error

- The projection of $\widetilde{x}^{(n)}$ onto a vector $v$ is

$$z^{(n)} = \left( \frac{v^T \widetilde{x}^{(n)}}{\|v\|_2} \right) \frac{v}{\|v\|_2}$$

Length of projection     Direction of projection

# Reconstruction Error

- The projection of $\widetilde{x}^{(n)}$ onto a unit vector $v$ is

$$z^{(n)} = \left(v^T \widetilde{x}^{(n)}\right)v$$

$$\hat{v} = \operatorname*{argmin}_{v:\|v\|_2^2=1} \sum_{n=1}^{N} \left\| \widetilde{x}^{(n)} - \left(v^T\widetilde{x}^{(n)}\right)v \right\|_2^2$$

$$\left\| \widetilde{x}^{(n)} - \left(v^T\widetilde{x}^{(n)}\right)v \right\|_2^2$$

$$= \widetilde{x}^{(n)^T}\widetilde{x}^{(n)} - 2\left(v^T\widetilde{x}^{(n)}\right)v^T\widetilde{x}^{(n)} + \left(v^T\widetilde{x}^{(n)}\right)\left(v^T\widetilde{x}^{(n)}\right)v^Tv$$

$$= \widetilde{x}^{(n)^T}\widetilde{x}^{(n)} - \left(v^T\widetilde{x}^{(n)}\right)v^T\widetilde{x}^{(n)}$$

$$= \left\| \widetilde{x}^{(n)} \right\|_2^2 - \left(v^T\widetilde{x}^{(n)}\right)^2$$

# Minimizing the Reconstruction Error ⇕ Maximizing the Variance

$$\widehat{v} = \operatorname*{argmin}_{v:\|v\|_2^2=1} \sum_{n=1}^{N} \left\| \widetilde{x}^{(n)} - \left(v^T \widetilde{x}^{(n)}\right)v \right\|_2^2$$

$$= \operatorname*{argmin}_{v:\|v\|_2^2=1} \sum_{n=1}^{N} \left\| \widetilde{x}^{(n)} \right\|_2^2 - \left(v^T \widetilde{x}^{(n)}\right)^2$$

$$= \operatorname*{argmax}_{v:\|v\|_2^2=1} \sum_{n=1}^{N} \left(v^T \widetilde{x}^{(n)}\right)^2 \longleftarrow$$

Variance of projections ($\widetilde{x}^{(n)}$ are centered)

$$= \operatorname*{argmax}_{v:\|v\|_2^2=1} v^T \left( \sum_{n=1}^{N} \widetilde{x}^{(n)} \widetilde{x}^{(n)T} \right) v$$

$$= \operatorname*{argmax}_{v:\|v\|_2^2=1} v^T (X^T X) v$$

34

# Maximizing the Variance

$$\hat{v} = \operatorname*{argmax}_{v:\|v\|_2^2=1} v^T(X^TX)v$$

$$\mathcal{L}(v,\lambda) = v^T(X^TX)v - \lambda(\|v\|_2^2 - 1)$$
$$= v^T(X^TX)v - \lambda(v^Tv - 1)$$

$$\frac{\partial\mathcal{L}}{\partial v} = 2(X^TX)v - 2\lambda v$$

$$\rightarrow 2(X^TX)\hat{v} - 2\lambda\hat{v} = 0 \rightarrow (X^TX)\hat{v} = \lambda\hat{v}$$
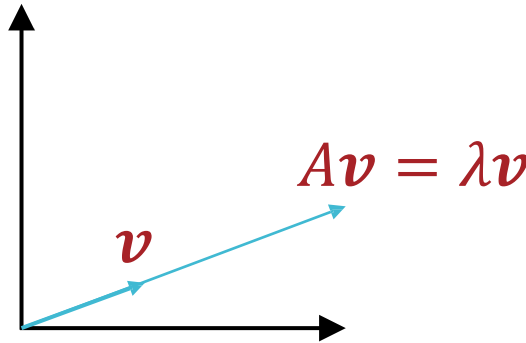
- $\hat{v}$ is an eigenvector of $X^TX$ and $\lambda$ is the corresponding eigenvalue!

## Background: Eigenvectors & Eigenvalues

- Given a square matrix $A \in \mathbb{R}^{N \times N}$, a vector $\boldsymbol{v} \in \mathbb{R}^{N \times 1}$ is an

  **eigenvector** of $A$ iff there exists some scalar $\lambda$ such that

$$A\boldsymbol{v} = \lambda\boldsymbol{v}$$



Intuition: $A$ scales or stretches $\boldsymbol{v}$ but does not rotate it

- Key property: the eigenvectors of <u>symmetric</u> matrices (e.g., the covariance matrix of a data set) are orthogonal!

# Maximizing the Variance

$$\hat{\boldsymbol{v}} = \underset{\boldsymbol{v}:\|\boldsymbol{v}\|_2^2=1}{\arg\max} \; \boldsymbol{v}^T (X^T X) \boldsymbol{v}$$

$$\mathcal{L}(\boldsymbol{v}, \lambda) = \boldsymbol{v}^T (X^T X) \boldsymbol{v} - \lambda(\|\boldsymbol{v}\|_2^2 - 1)$$

$$= \boldsymbol{v}^T (X^T X) \boldsymbol{v} - \lambda(\boldsymbol{v}^T \boldsymbol{v} - 1)$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = 2(X^T X)\boldsymbol{v} - 2\lambda \boldsymbol{v}$$

$$\rightarrow 2(X^T X)\hat{\boldsymbol{v}} - 2\lambda\hat{\boldsymbol{v}} = 0 \rightarrow (X^T X)\hat{\boldsymbol{v}} = \lambda\hat{\boldsymbol{v}}$$

- $\hat{\boldsymbol{v}}$ is an eigenvector of $X^T X$ and $\lambda$ is the corresponding eigenvalue!
- But which one?

# Maximizing the Variance

$$\widehat{\boldsymbol{v}} = \underset{\boldsymbol{v}:\|\boldsymbol{v}\|_2^2=1}{\arg\max}\ \boldsymbol{v}^T(X^TX)\boldsymbol{v}$$

$$(X^TX)\widehat{\boldsymbol{v}} = \lambda\widehat{\boldsymbol{v}} \ \rightarrow\ \widehat{\boldsymbol{v}}^T(X^TX)\widehat{\boldsymbol{v}} = \lambda\widehat{\boldsymbol{v}}^T\widehat{\boldsymbol{v}} = \lambda$$

- The first principal component is the eigenvector $\widehat{\boldsymbol{v}}_1$ that corresponds to the largest eigenvalue $\lambda_1$
- The second principal component is the eigenvector $\widehat{\boldsymbol{v}}_2$ that corresponds to the second largest eigenvalue $\lambda_2$
  - $\widehat{\boldsymbol{v}}_1$ and $\widehat{\boldsymbol{v}}_2$ are orthogonal
- Etc …
- $\lambda_i$ is a measure of how much variance falls along $\widehat{\boldsymbol{v}}_i$

# Principal Components: Example

How can we efficiently find principal components (eigenvectors)?

# Singular Value Decomposition (SVD) for PCA

- Every real-valued matrix $X \in \mathbb{R}^{N \times D}$ can be expressed as

$$X = USV^T$$

where:

1. $U \in \mathbb{R}^{N \times N}$ - columns of $U$ are eigenvectors of $XX^T$

2. $V \in \mathbb{R}^{D \times D}$ - columns of $V$ are eigenvectors of $X^TX$

3. $S \in \mathbb{R}^{N \times D}$ - diagonal matrix whose entries are the eigenvalues of $X$ $\rightarrow$ squared entries are the eigenvalues of $XX^T$ and $X^TX$

# PCA Algorithm

- Input: $\mathcal{D} = \{(\boldsymbol{x}^{(n)})\}_{n=1}^{N}, \rho$

1. Center the data

2. Use SVD to compute the eigenvalues and eigenvectors of $X^T X$

3. Collect the top $\rho$ eigenvectors (corresponding to the $\rho$ largest eigenvalues), $V_\rho \in \mathbb{R}^{D \times \rho}$

4. Project the data into the space defined by $V_\rho$, $Z = X V_\rho$

- Output: $Z$, the transformed (potentially lower-dimensional) data

# How many PCs should we use?

- Input: $\mathcal{D} = \{(\boldsymbol{x}^{(n)})\}_{n=1}^{N}, \rho$

1. Center the data

2. Use SVD to compute the eigenvalues and eigenvectors of $X^T X$

3. Collect the top $\rho$ eigenvectors (corresponding to the $\rho$ largest eigenvalues), $V_\rho \in \mathbb{R}^{D \times \rho}$

4. Project the data into the space defined by $V_\rho$, $Z = X V_\rho$

- Output: $Z$, the transformed (potentially lower-dimensional) data

# Choosing the number of PCs

- Define a percentage of explained variance for the $i$th PC:

$$\lambda_i \Big/ \sum \lambda_j$$

- Select all PCs above some threshold of explained variance, e.g., 5%

- Keep selecting PCs until the total explained variance exceeds some threshold, e.g., 90%
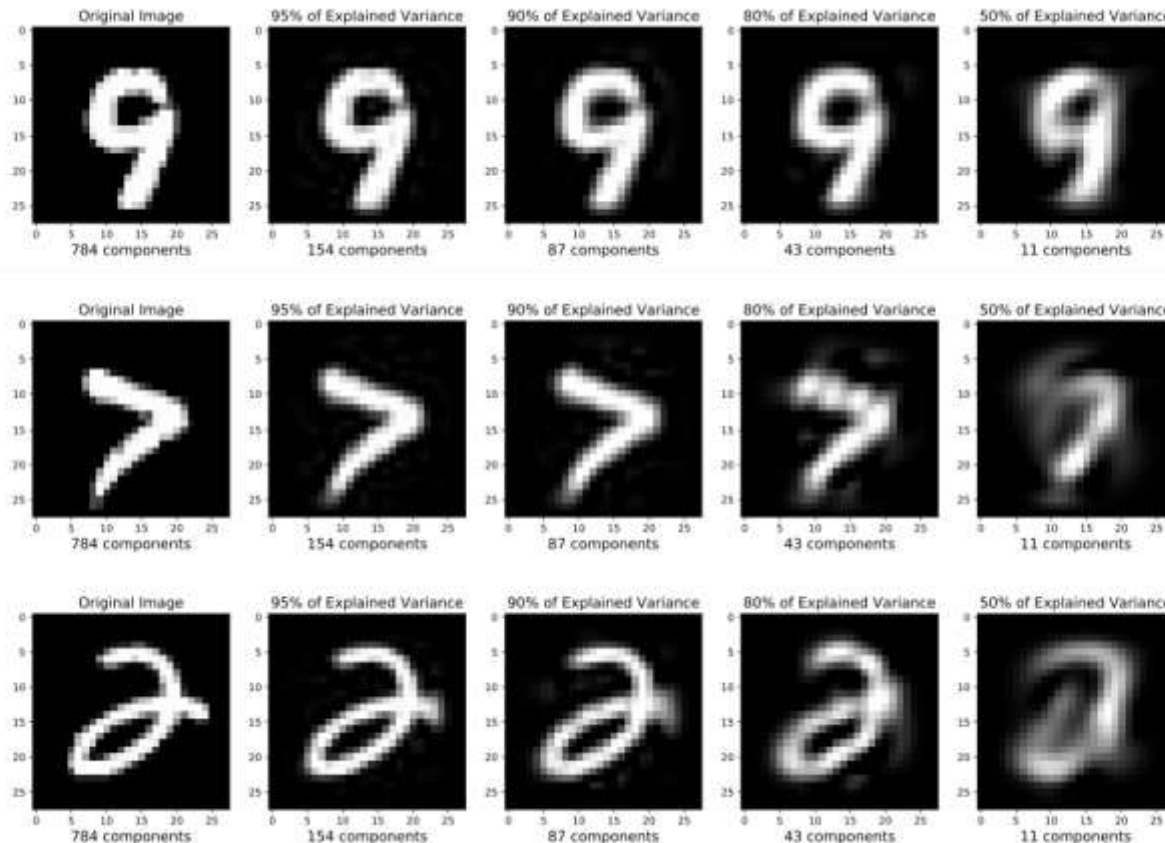
- Evaluate on some downstream metric

# PCA EXAMPLES

# Projecting MNIST digits

**Task Setting:**
1. Take each 28x28 image of a digit (i.e. a vector $\mathbf{x}^{(i)}$ of length 784) and project it down to K components (i.e. a vector $\mathbf{u}^{(i)}$)
2. Report percent of variance explained for K components
3. Then project back up to 28x28 image (i.e. a vector $\tilde{\mathbf{x}}^{(i)}$ of length 784) to visualize how much information was preserved



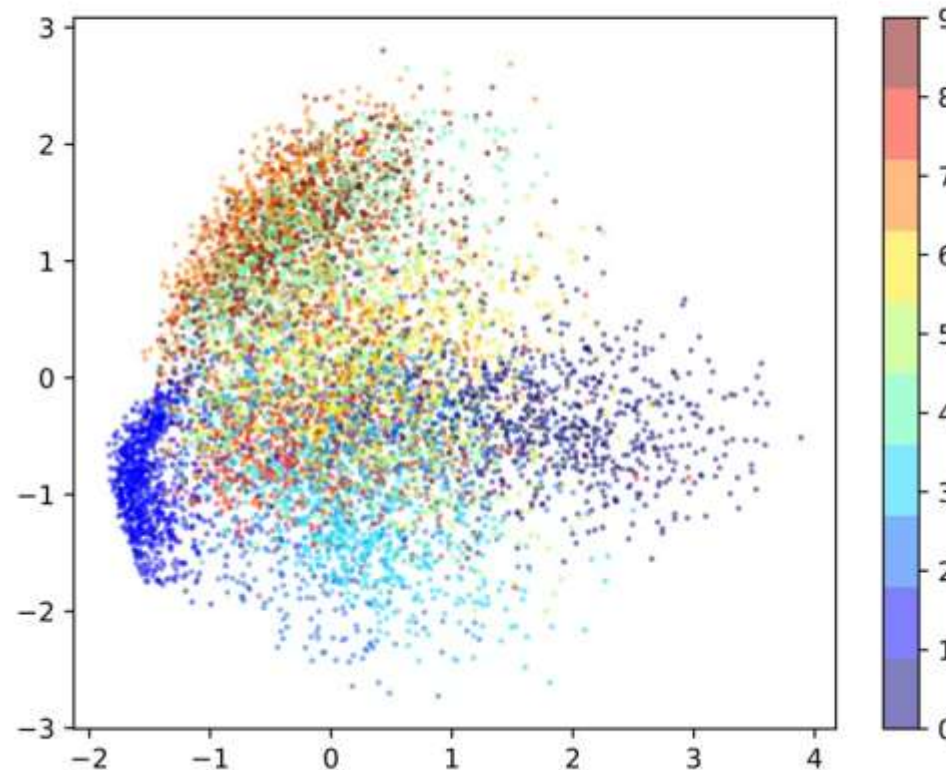**Takeaway**: Using fewer principal components K leads to higher reconstruction error.

But even a small number (say 43) still preserves a lot of information about the original image.

# Projecting MNIST digits

**Task Setting:**
1. Take each 28x28 image of a digit (i.e. a vector $\mathbf{x}^{(i)}$ of length 784) and project it down to K=2 components (i.e. a vector $\mathbf{u}^{(i)}$)
2. Plot the 2 dimensional points $\mathbf{u}^{(i)}$ and label with the (unknown to PCA) label $y^{(i)}$ as the color
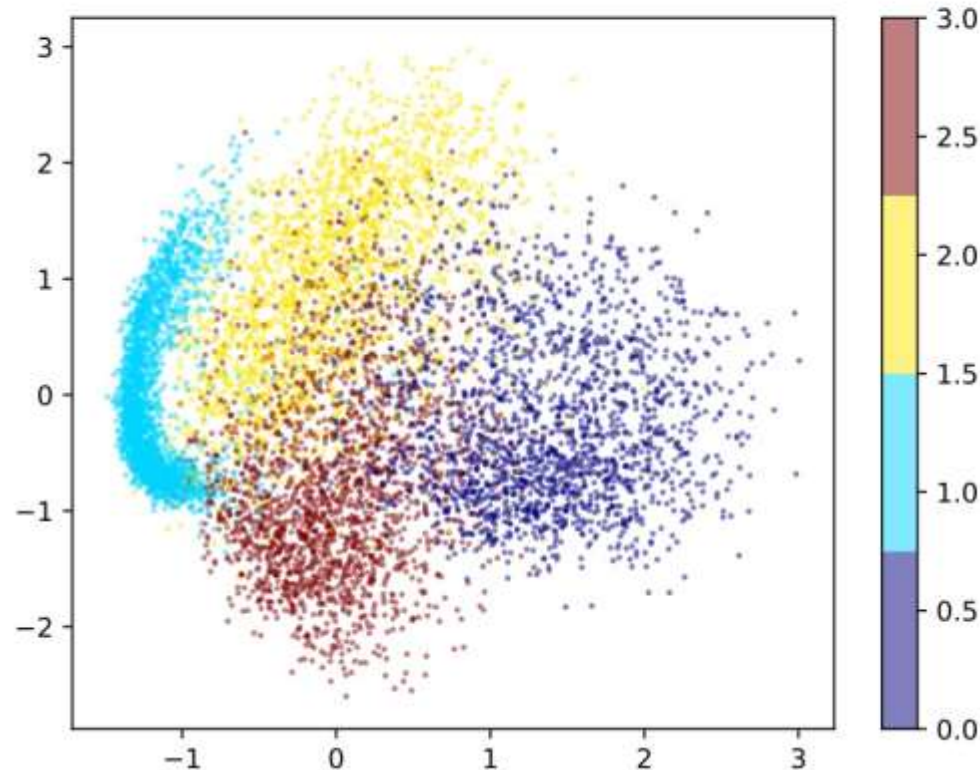3. Here we look at all ten digits 0 - 9



**Takeaway**: Even with a tiny number of principal components K=2, PCA learns a representation that captures the *latent* information about the type of digit

# Projecting MNIST digits

**Task Setting:**

1. Take each 28x28 image of a digit (i.e. a vector $\mathbf{x}^{(i)}$ of length 784) and project it down to K=2 components (i.e. a vector $\mathbf{u}^{(i)}$)

2. Plot the 2 dimensional points $\mathbf{u}^{(i)}$ and label with the (unknown to PCA) label $y^{(i)}$ as the color

3. Here we look at just four digits 0, 1, 2, 3



**Takeaway**: Even with a tiny number of principal components K=2, PCA learns a representation that captures the *latent* information about the type of digit
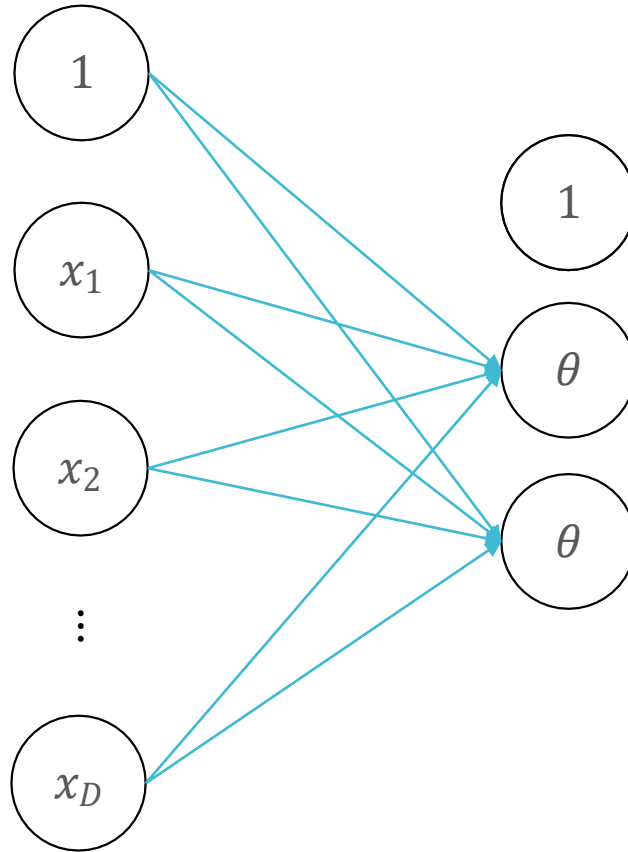
**Dimensionality Reduction / PCA**
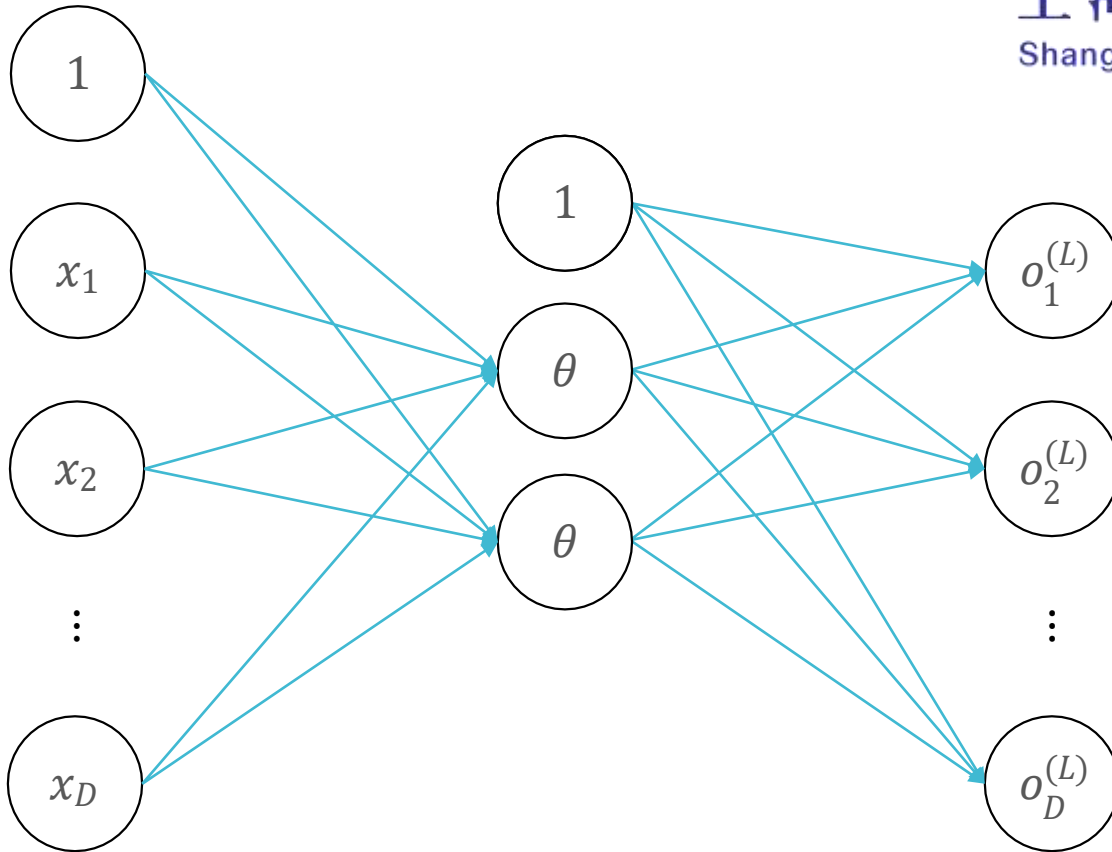
*You should be able to…*

1.  Define the sample mean, sample variance, and sample covariance of a vector-valued dataset
2.  Identify examples of high dimensional data and common use cases for dimensionality reduction
3.  Draw the principal components of a given toy dataset
4.  Establish the equivalence of minimization of reconstruction error with maximization of variance
5.  Given a set of principal components, project from high to low dimensional space and do the reverse to produce a reconstruction
6.  Explain the connection between PCA, eigenvectors, eigenvalues, and covariance matrix
7.  Use common methods in linear algebra to obtain the principal components

# Autoencoders



Insight: neural networks implicitly learn low-dimensional representations of inputs in hidden layers
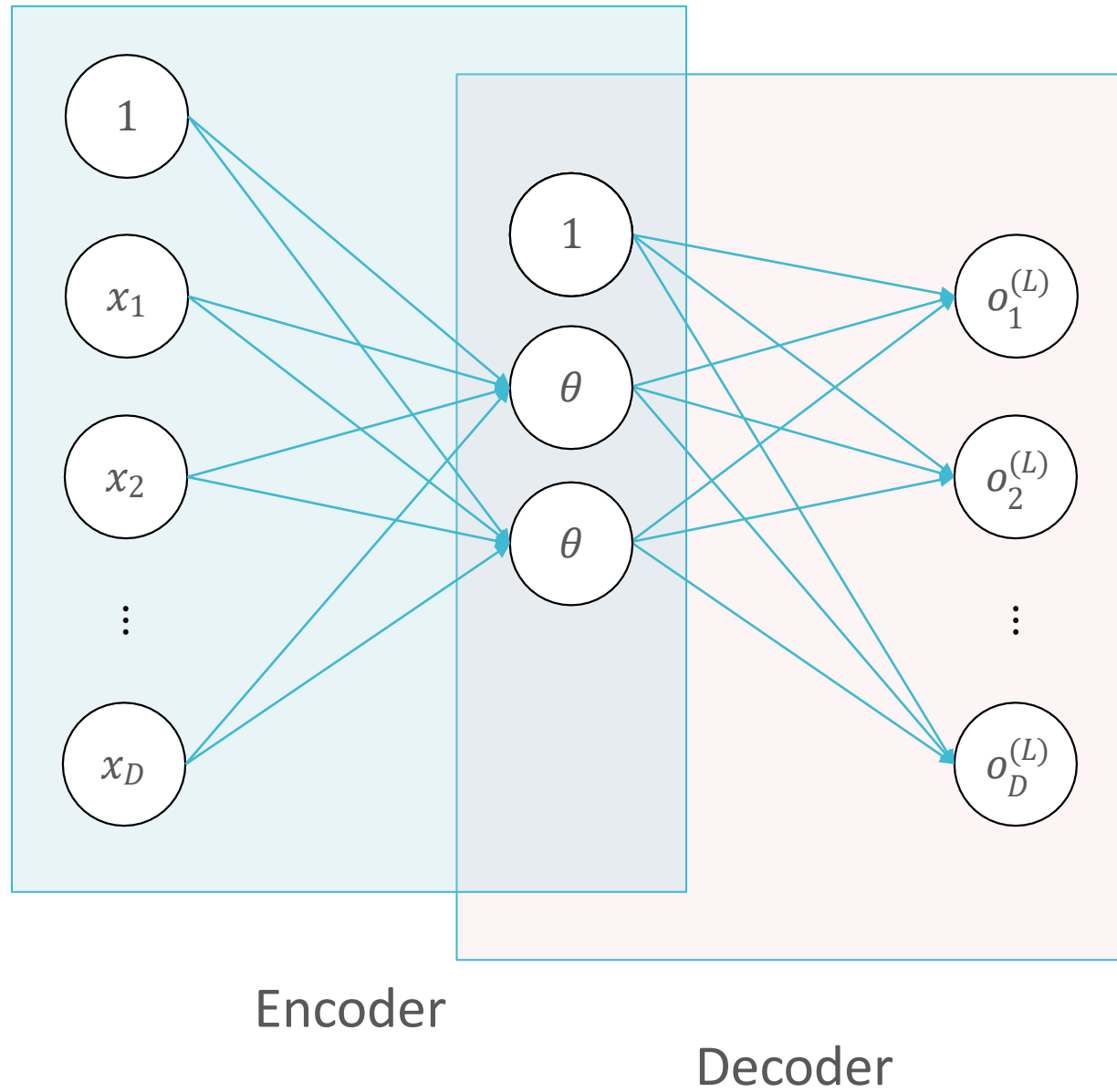
# Autoencoders



- Learn the weights by minimizing the reconstruction loss:

$$e(\boldsymbol{x}) = \left\| \boldsymbol{x} - \boldsymbol{o}^{(L)} \right\|_2^2$$

# Autoencoders



Encoder

Decoder

# Deep Autoencoders