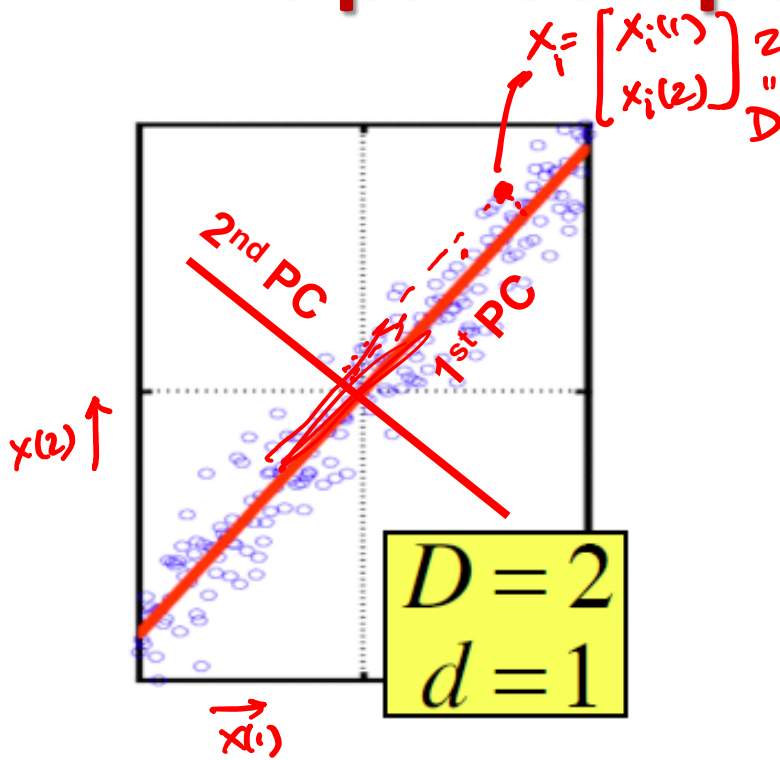


# Principal Component Analysis (PCA)



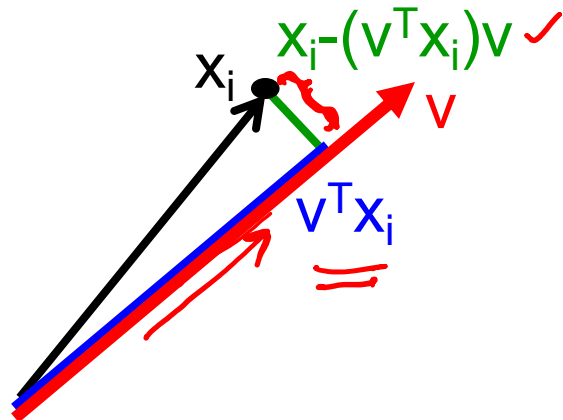
Principal Components (PC) are orthogonal unit norm directions that capture most of the variance in the data

1<sup>st</sup> PC – direction of greatest variability in data

2<sup>nd</sup> PC – Next orthogonal (uncorrelated) direction of greatest variability ✓

(remove all variability in first direction, then find next direction of greatest variability)

And so on ...

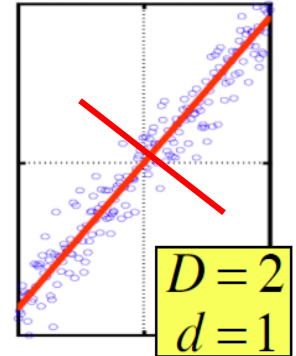


$$D \rightarrow \underline{d} \ll D$$

# Principal Component Analysis (PCA)

Let  $v_1, v_2, \dots, v_d$  denote the principal components

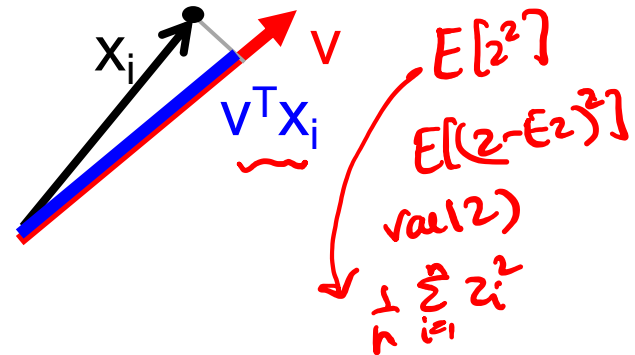
Orthogonal and unit norm  $v_i^T v_j = 0 \quad i \neq j$  ✓ }  
 $v_i^T v_i = 1$  ✓ }



Find vector that maximizes sample variance of projection

$$\frac{1}{n} \sum_{i=1}^n (v^T x_i)^2 = \frac{v^T \overset{d \times d}{\underbrace{XX^T}} v}{n}$$

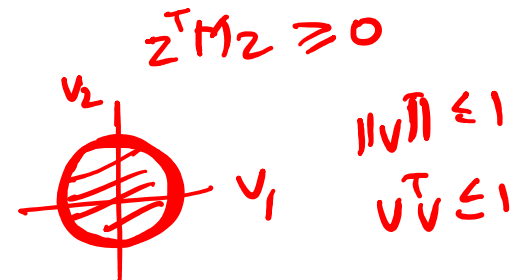
$$\rightarrow \max_v v^T \underbrace{XX^T}_n v \quad \text{s.t.} \quad \underbrace{v^T v}_1 = 1$$



Poll: Convex

$\sum_i v_i^2 = 1$   
 non-convex set

➤ Is this a convex optimization problem?



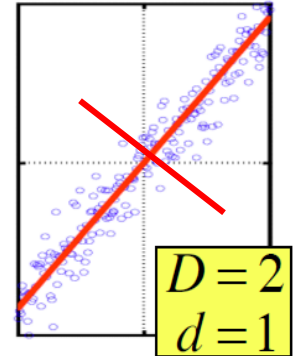
# Principal Component Analysis (PCA)

Let  $v_1, v_2, \dots, v_d$  denote the principal components

Orthogonal and unit norm  $v_i^T v_j = 0 \quad i \neq j$

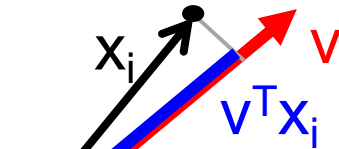
$$v_i^T v_i = 1$$

Find vector that maximizes sample variance of projection



$$\frac{1}{n} \sum_{i=1}^n (v^T x_i)^2 = \frac{v^T X X^T v}{n}$$

$$\max_v \underline{v^T X X^T v} \quad \text{s.t.} \quad \underline{v^T v = 1}$$



Wrap constraints into the objective function

$$\Rightarrow \boxed{(X X^T) v = \lambda v} \quad \begin{matrix} \text{eval}(X X^T) \\ \text{eval}(X X^T) \end{matrix}$$

Lagrangian:  $\max_v v^T X X^T v - \lambda v^T v$

$$\partial / \partial v = 0 \quad \hookrightarrow \quad 2 X X^T v - 2 \lambda v = 0 \quad \Rightarrow \quad (X X^T - \lambda I) v = 0$$

→ Sample var when projecting on v

$$v^T X X^T v = v^T (\lambda v) = \lambda v^T v = \lambda$$

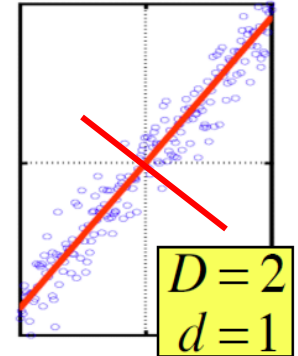
# Principal Component Analysis (PCA)

1) center  $X$

2) Sample cov  $\frac{XX^T}{n}$

$$(XX^T)v = \lambda v$$

Therefore,  $v$  is the eigenvector of sample covariance matrix  $XX^T$



Sample variance of projection  $= v^T XX^T v = \lambda v^T v = \lambda$

Thus, the eigenvalue  $\lambda$  denotes the amount of variability captured along that dimension (aka amount of energy along that dimension).

Eigenvalues  $\lambda_1 > \lambda_2 > \lambda_3 > \dots$

The 1<sup>st</sup> Principal component  $v_1$  is the eigenvector of the sample covariance matrix  $XX^T$  associated with the largest eigenvalue  $\lambda_1$

The 2<sup>nd</sup> Principal component  $v_2$  is the eigenvector of the sample covariance matrix  $XX^T$  associated with the second largest eigenvalue  $\lambda_2$

And so on ...

# Another interpretation

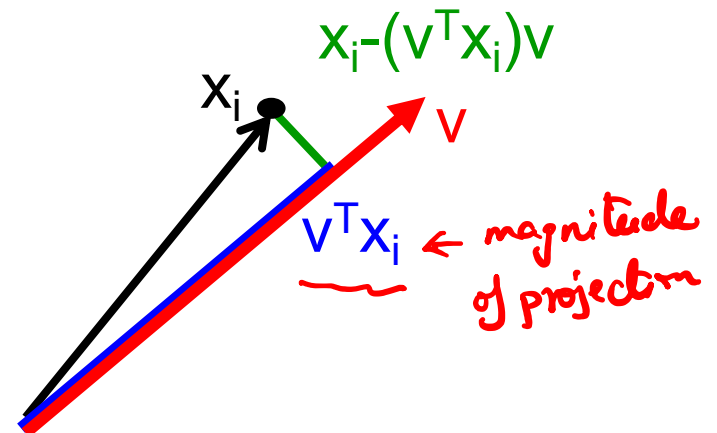
**Maximum Variance Subspace:** PCA finds vectors  $v$  such that projections on to the vectors capture maximum variance in the data

$$\frac{1}{n} \sum_{i=1}^n (v^T x_i)^2 = v^T X X^T v$$

→ **Minimum Reconstruction Error:** PCA finds vectors  $v$  such that projection on to the vectors yields minimum MSE reconstruction

$$\frac{1}{n} \sum_{i=1}^n \|x_i - \underbrace{(v^T x_i)v}_{\tilde{x}_i}\|^2$$

*(Handwritten red annotations: an arrow points to  $x_i$ , another to  $\tilde{x}_i$ )*



argmin  
v

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v})^T (\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v})$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i - 2(\mathbf{v}^T \mathbf{x}_i)^2 + \underbrace{\mathbf{v}^T (\mathbf{v}^T \mathbf{x}_i)^2 \mathbf{v}}_{(\mathbf{v}^T \mathbf{x}_i)^2 \cancel{\mathbf{v}^T \mathbf{v}}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i)^2$$

$$= \arg \min_v \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i)^2$$

$$= \arg \max_v \frac{1}{n} \sum_{i=1}^n \underbrace{(\mathbf{v}^T \mathbf{x}_i)^2}_{\text{projection}}$$

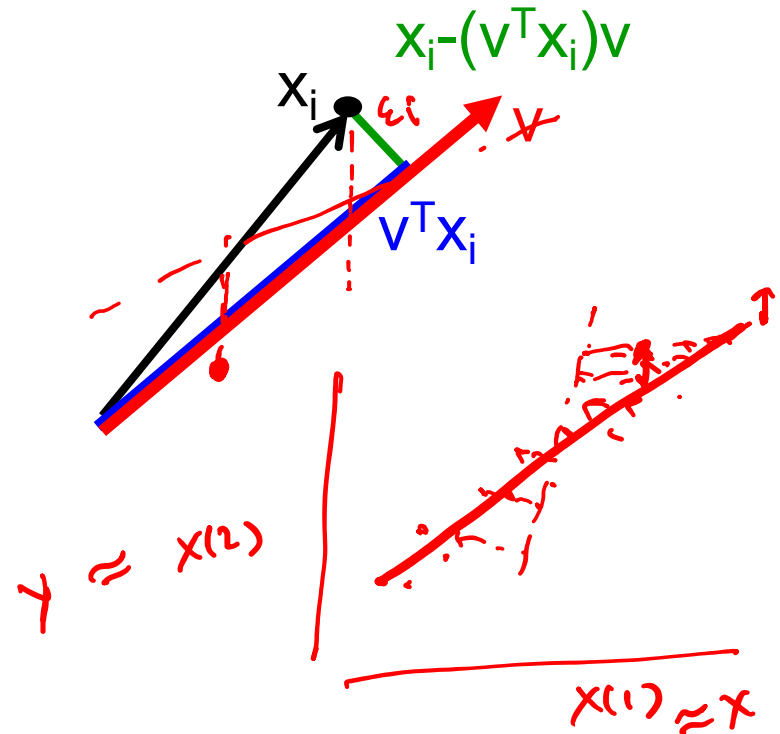


# Another interpretation

**Minimum Reconstruction Error:** PCA finds vectors  $v$  such that projection on to the vectors yields minimum MSE reconstruction

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2$$

$$\frac{1}{n} \sum_{i=1}^n \underbrace{(\gamma_i - \mathbf{v}^T \mathbf{x}_i)^2}_{\epsilon_i}$$



Poll:

- Is the 1<sup>st</sup> PC same as the linear least square fit?

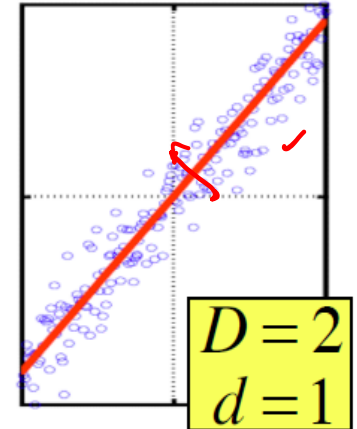
# Dimensionality Reduction using PCA

$$XX^T_{D \times D} \rightarrow v_1 \dots v_D$$

The eigenvalue  $\lambda$  denotes the amount of variability captured along that dimension.

Zero eigenvalues indicate no variability along those directions => data lies exactly on a linear subspace

Only keep data projections onto principal components with non-zero eigenvalues, say  $v_1, \dots, v_d$  where  $d = \text{rank}(XX^T)$



## Original Representation

data point

$$x_i = [x_i^1, x_i^2, \dots, x_i^D]^T$$

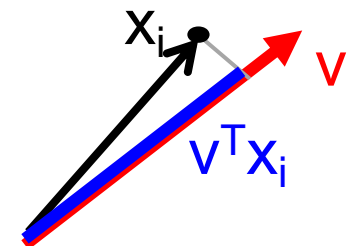
(D-dimensional vector)

## Transformed representation

projections

$$[v_1^T x_i, v_2^T x_i, \dots, v_d^T x_i]$$

(d-dimensional vector)



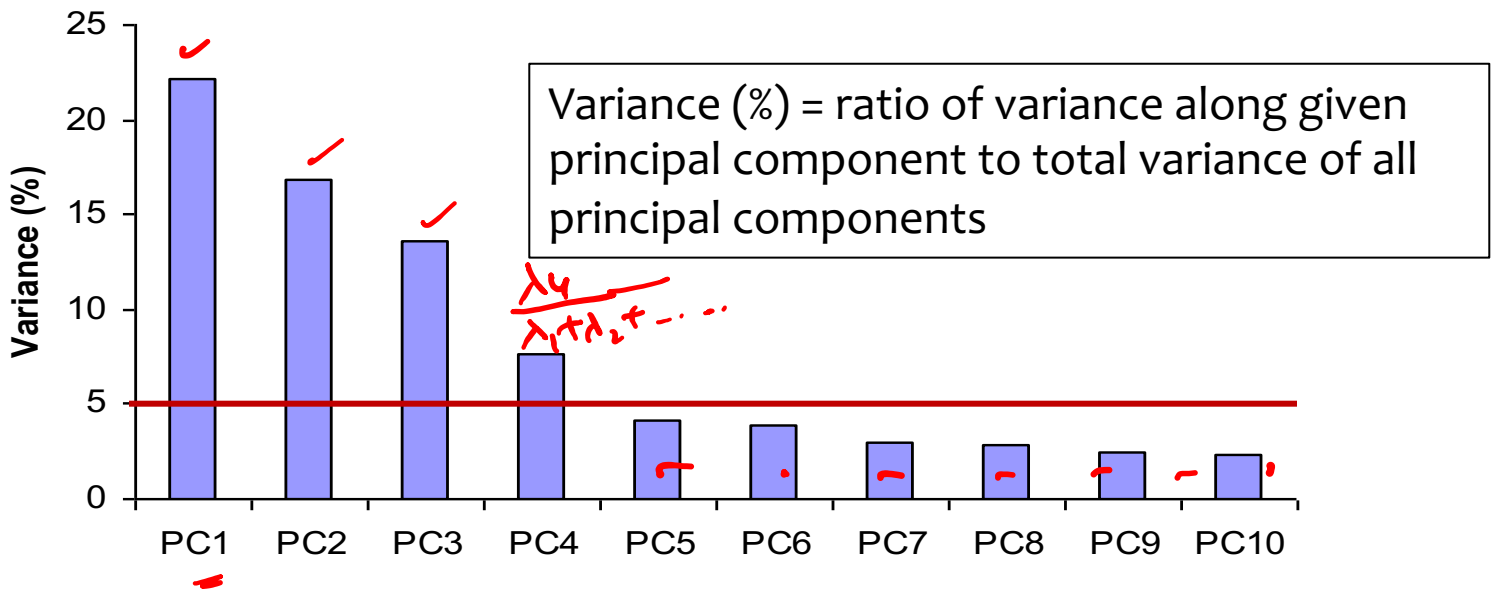


# Dimensionality Reduction using PCA

In high-dimensional problem, data usually lies near a linear subspace, as noise introduces small variability

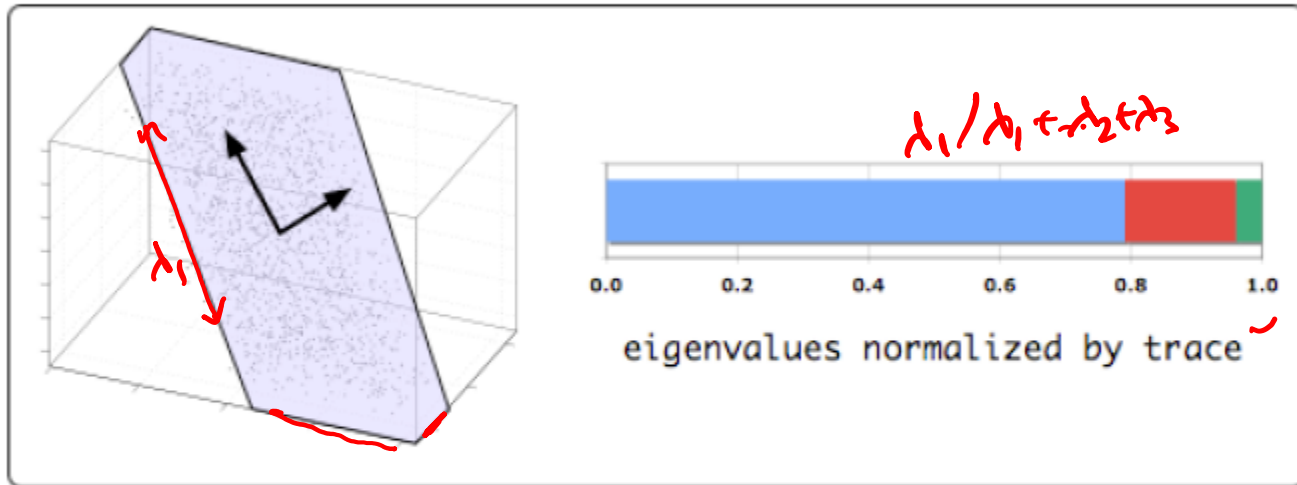
Only keep data projections onto principal components with large eigenvalues

Can *ignore* the components of lesser significance.



You might *lose some information*, but if the eigenvalues are small, you don't lose much

# Example of PCA



**Eigenvectors and eigenvalues of covariance matrix for  $n=1600$  inputs in  $d=3$  dimensions.**

$$D = 256 \times 256$$

$V_{D \times 1}$

## Example: faces



**Eigenfaces**

from 7562  
images:

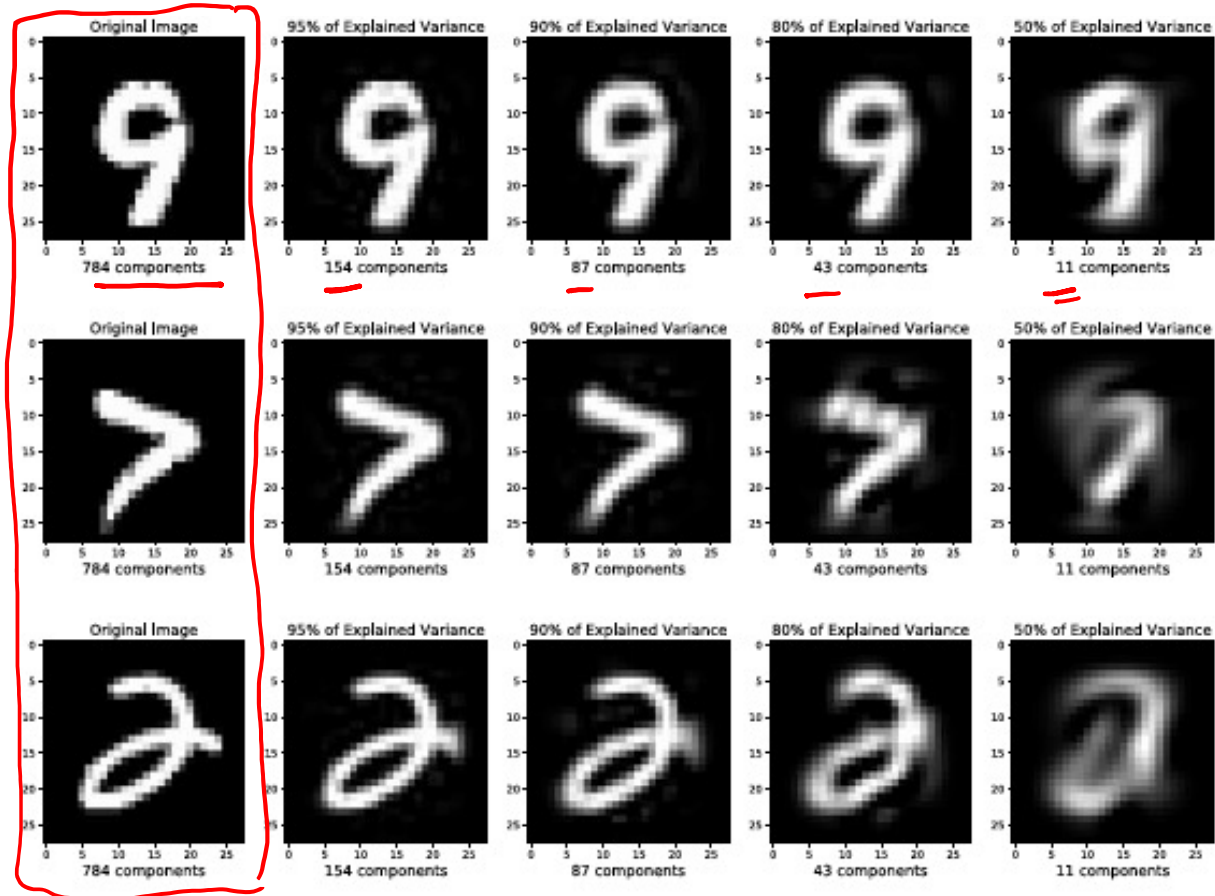
top left image  
is linear  
combination  
of rest.

Sirovich & Kirby (1987)

Turk & Pentland (1991)

# Example: MNIST digits

- 28x28 images = 784 PCA vectors
- Project to K dimensional space and then project back up

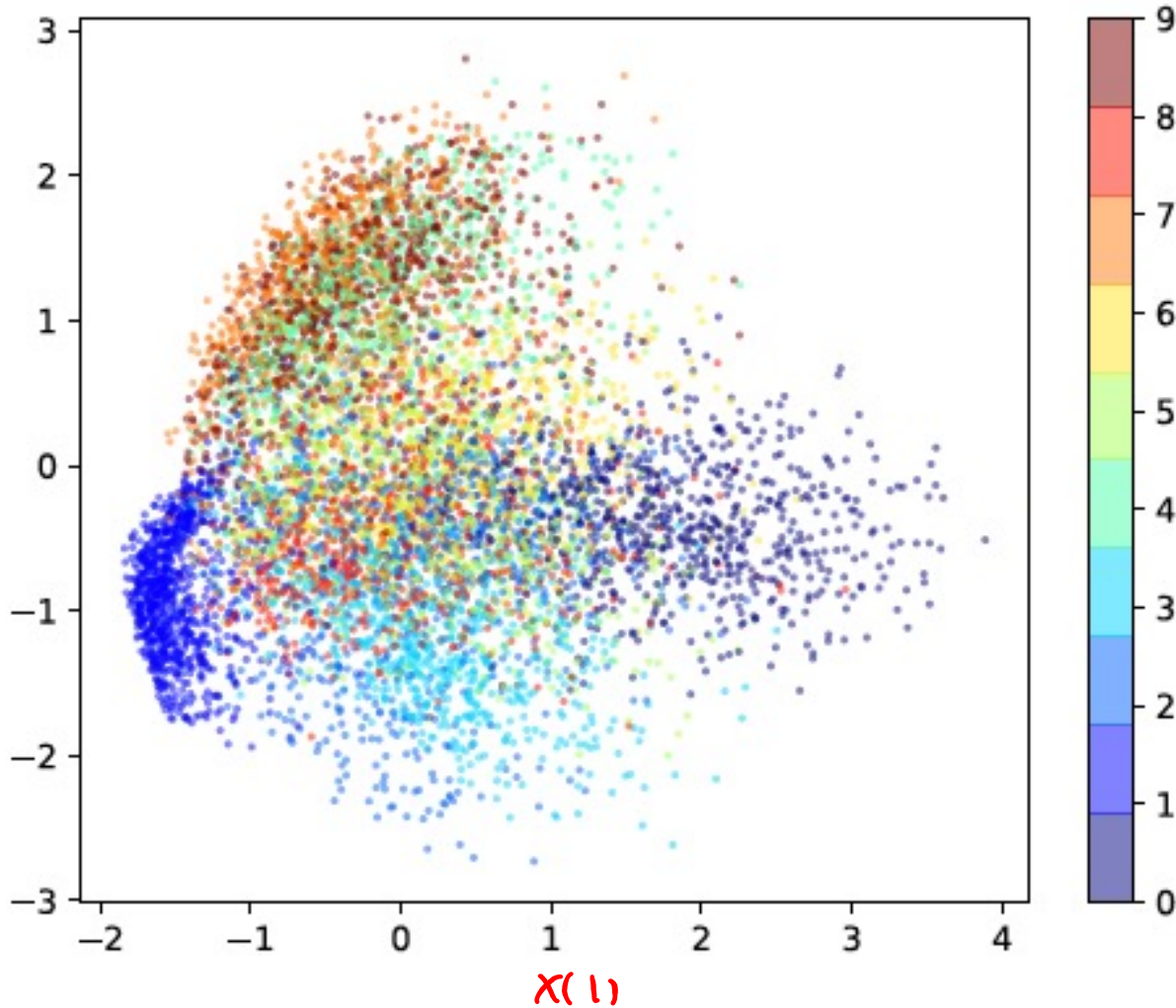


# Projecting MNIST digits

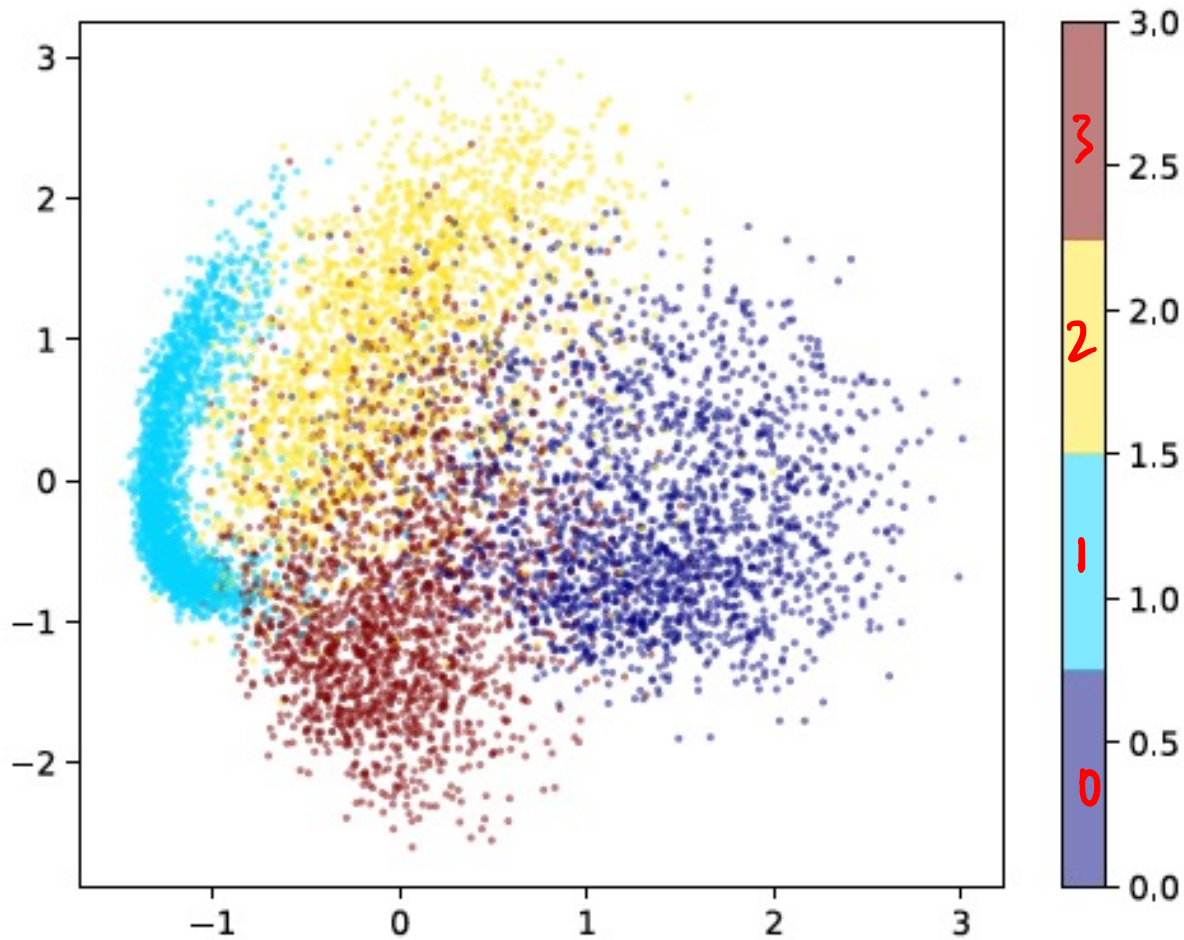
$D=28 \times 28$

$d=2$

$x(2)$



# Projecting MNIST digits



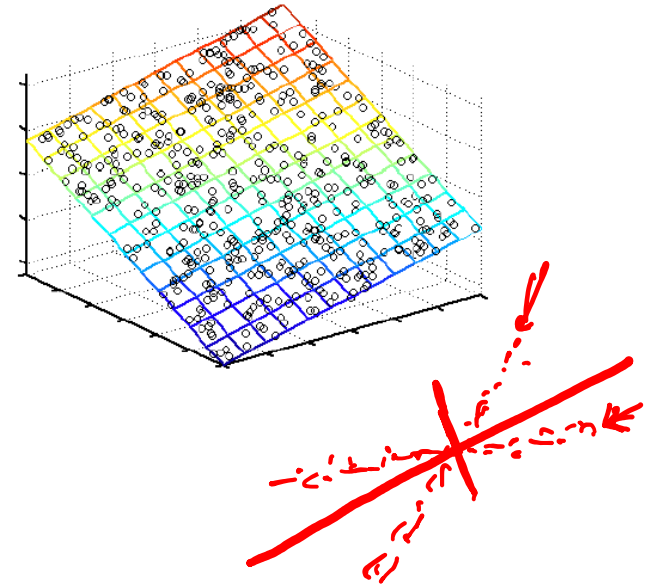
# Unsupervised Dimensionality Reduction

- Linear

**Principal Component Analysis (PCA)**

Factor Analysis

Independent Component Analysis (ICA)





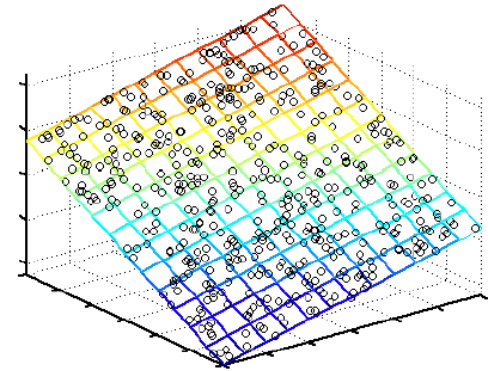
# Unsupervised Dimensionality Reduction

- Linear

  - Principal Component Analysis (PCA)**

  - Factor Analysis

  - Independent Component Analysis (ICA)



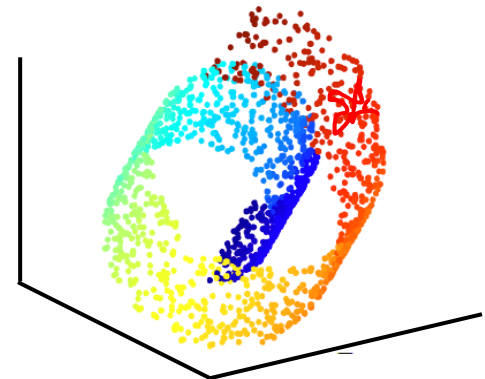
- Nonlinear

  - Kernel PCA** ←

  - Laplacian Eigenmaps, ISOMAP, LLE ←

  - Autoencoders ←

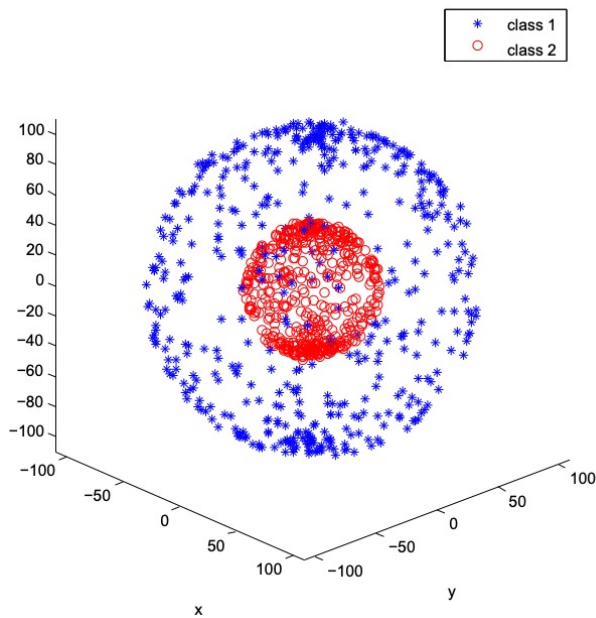
$$\begin{matrix} \times & 0 & & & 0 & \times \\ & 0 & & & 0 & \\ & 0 & & & 0 & \\ & 0 & & & 0 & \\ & 0 & & & 0 & \end{matrix} \rightarrow (\tilde{x} - x)^2$$



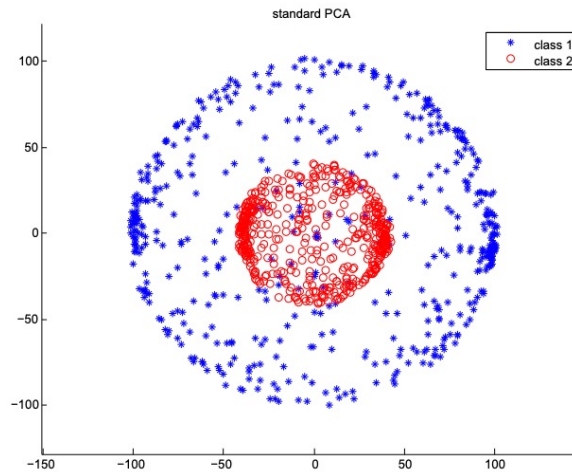


# Kernel PCA

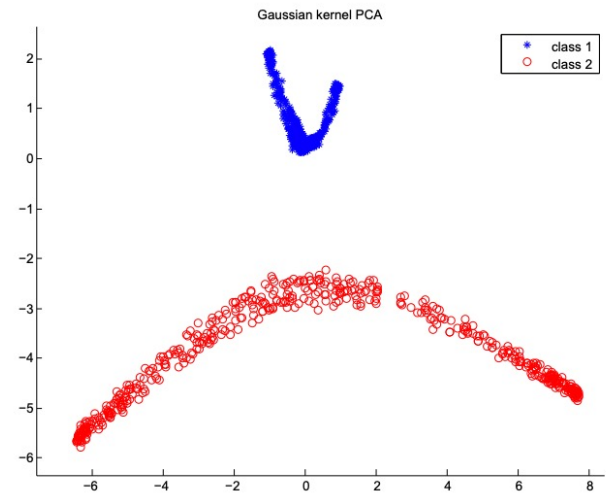
**Latent features:** linear in  $\phi(x)$  where  $\phi(x) \cdot \phi(x') = K(x, x')$  that capture maximum variance or minimum reconstruction error



Original data points



PCA



Kernel PCA  
Gaussian/RBF kernel

# Kernel PCA

**Latent features:** linear in  $\phi(x)$  where  $\phi(x) \cdot \phi(x') = K(x, x')$  that capture maximum variance or minimum reconstruction error

PCA:

Top d eigenvectors (each D dimensional) of sample covariance  $XX^T$

Low d-dimensional embedding of a point:  $[v_1^T x_i, v_2^T x_i, \dots, v_d^T x_i]$

$D \times D$

$d \times 1$

Kernel PCA:

Top d eigenvectors (each n dimensional) of kernel matrix  $K(X, X) = X^T X$

$\phi(x)^T \phi(x)$

$n \times n$

Low d-dimensional embedding of a point:  $[v_1(i), v_2(i), \dots, v_d(i)]$

Eigenvectors are not PCs but projections of data points

# PCA Summary

- PCA finds latent features linear in original features  $x$  that capture
  - Maximum variance amongst all linear features
  - Minimum reconstruction error when recovering points from PC projections
- Non-convex problem with simple solution:
  - PCs = eigenvectors of sample covariance matrix
  - Lower ( $d < D$ ) dimensional embedding of data point = projection of data point onto  $d$  PCs
- Kernel PCA: latent features linear in  $\phi(x)$  where  $\phi(x)$ .  $\phi(x') = K(x, x')$  that capture maximum variance or minimum reconstruction error
  - Directly get projections of data points

# Clustering

Aarti Singh

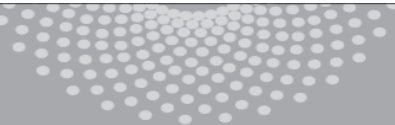
Machine Learning 10-701

Apr 24, 2023

Some slides courtesy of Eric Xing, Carlos Guestrin



**MACHINE LEARNING** DEPARTMENT

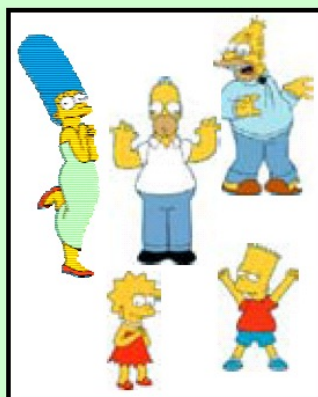


**Carnegie Mellon.**  
School of Computer Science

# What is clustering?

- Clustering: the process of grouping a set of objects into classes of similar objects
  - high intra-class similarity
  - low inter-class similarity
  - It is the most common form of **unsupervised learning**

## Clustering is subjective



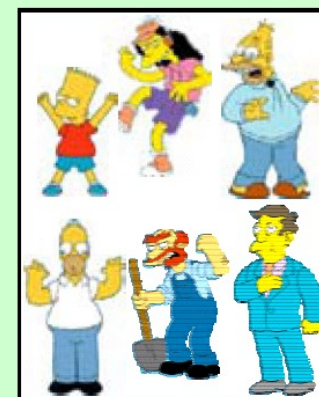
Simpson's Family



School Employees



Females



Males

# What is Similarity?

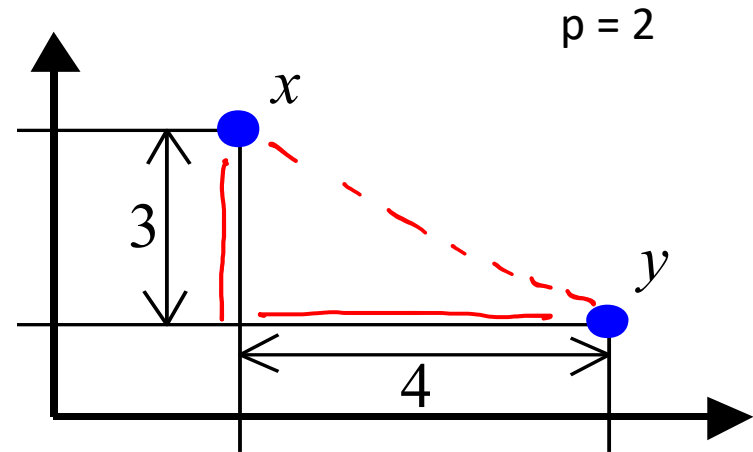


Hard to  
define! But *we*  
*know it when*  
*we see it*

- The real meaning of similarity is a philosophical question. We will take a more pragmatic approach - think in terms of a distance (rather than similarity) between vectors or correlations between random variables.

# Distance metrics

$$x = (x_1, x_2, \dots, x_p)$$
$$y = (y_1, y_2, \dots, y_p)$$



Euclidean distance

$$d(x, y) = \sqrt{\sum_{i=1}^p |x_i - y_i|^2}$$

5 ✓

Manhattan distance

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

7

Sup-distance

$$d(x, y) = \max_{1 \leq i \leq p} |x_i - y_i|$$

4 ✓

# Correlation coefficient

$$x = (x_1, x_2, \dots, x_p)$$

$$y = (y_1, y_2, \dots, y_p)$$

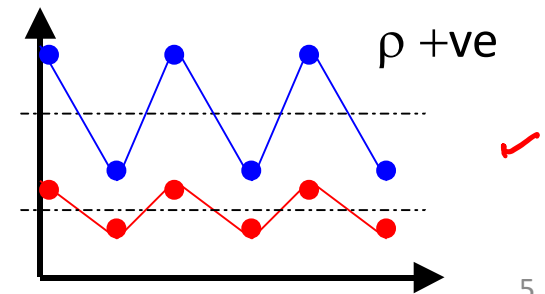
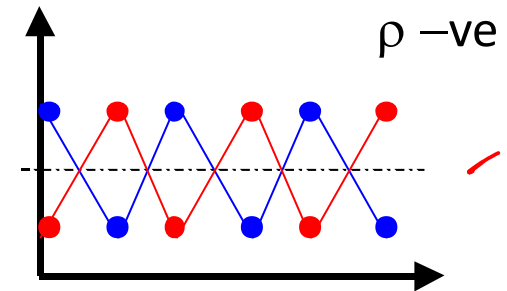
$I(x, y)$

Random vectors (e.g. expression levels of two genes under various drugs)

Pearson correlation coefficient

$$\rho(x, y) = \frac{\sum_{i=1}^p (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^p (x_i - \bar{x})^2 \times \sum_{i=1}^p (y_i - \bar{y})^2}}$$

where  $\bar{x} = \frac{1}{p} \sum_{i=1}^p x_i$  and  $\bar{y} = \frac{1}{p} \sum_{i=1}^p y_i$ .



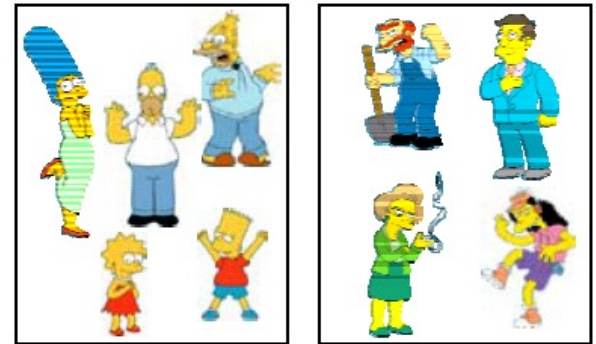


# Clustering Algorithms

- **Partition algorithms**

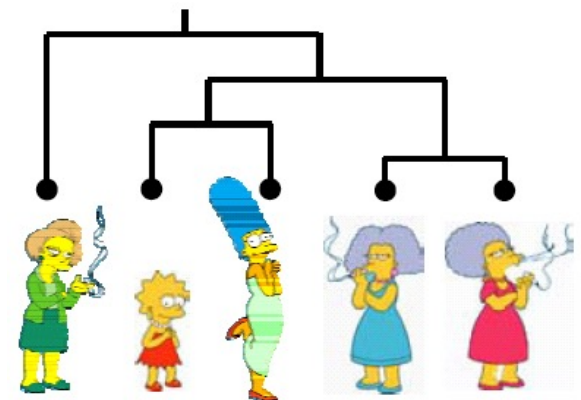
- K means clustering ✓
- Mixture-Model based clustering ✓

$p(x)$



- **Hierarchical algorithms**

- Single-linkage
- Average-linkage
- Complete-linkage
- Centroid-based



# Partitioning Algorithms

- Partitioning method: Construct a partition of  $n$  objects into a set of  $K$  clusters
- Given: a set of objects and the number  $K$
- Find: a partition of  $K$  clusters that optimizes the chosen partitioning criterion
  - Globally optimal: exhaustively enumerate all partitions
  - Effective heuristic method: K-means algorithm

$$\binom{n}{K} \approx n^K$$

# K-Means

## Algorithm

**Input** – Desired number of clusters,  $k$

**Initialize** – the  $k$  cluster centers (randomly if necessary)

**Iterate** –

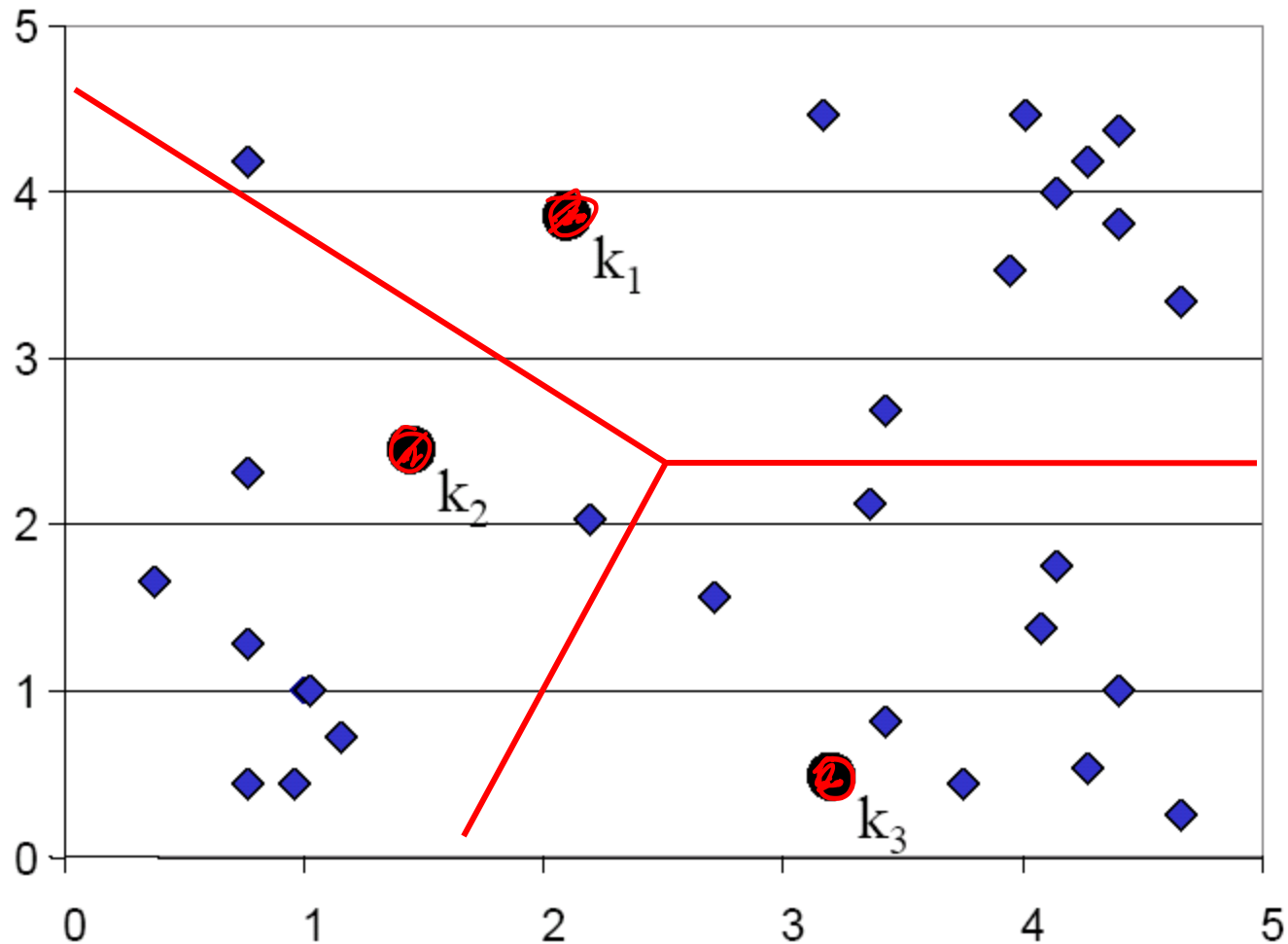
1. Assign points to the nearest cluster centers
2. Re-estimate the  $k$  cluster centers (aka the **centroid** or **mean**), by assuming the memberships found above are correct.

$$\vec{\mu}_k = \frac{1}{C_k} \sum_{i \in C_k} \vec{x}_i$$

**Termination** –

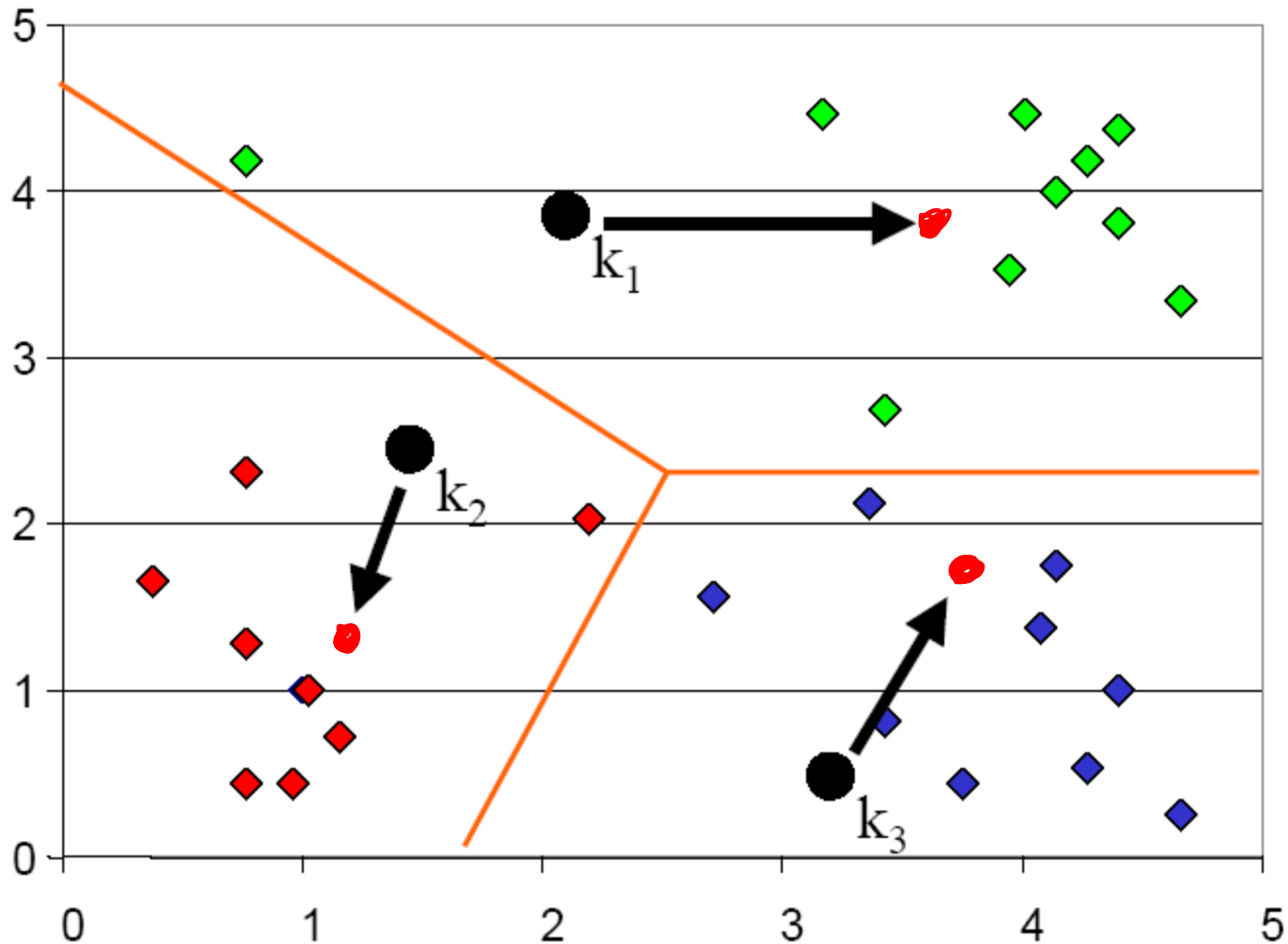
If none of the objects changed membership in the last iteration, exit.  
Otherwise go to 1.

# K-means Clustering: Step 1

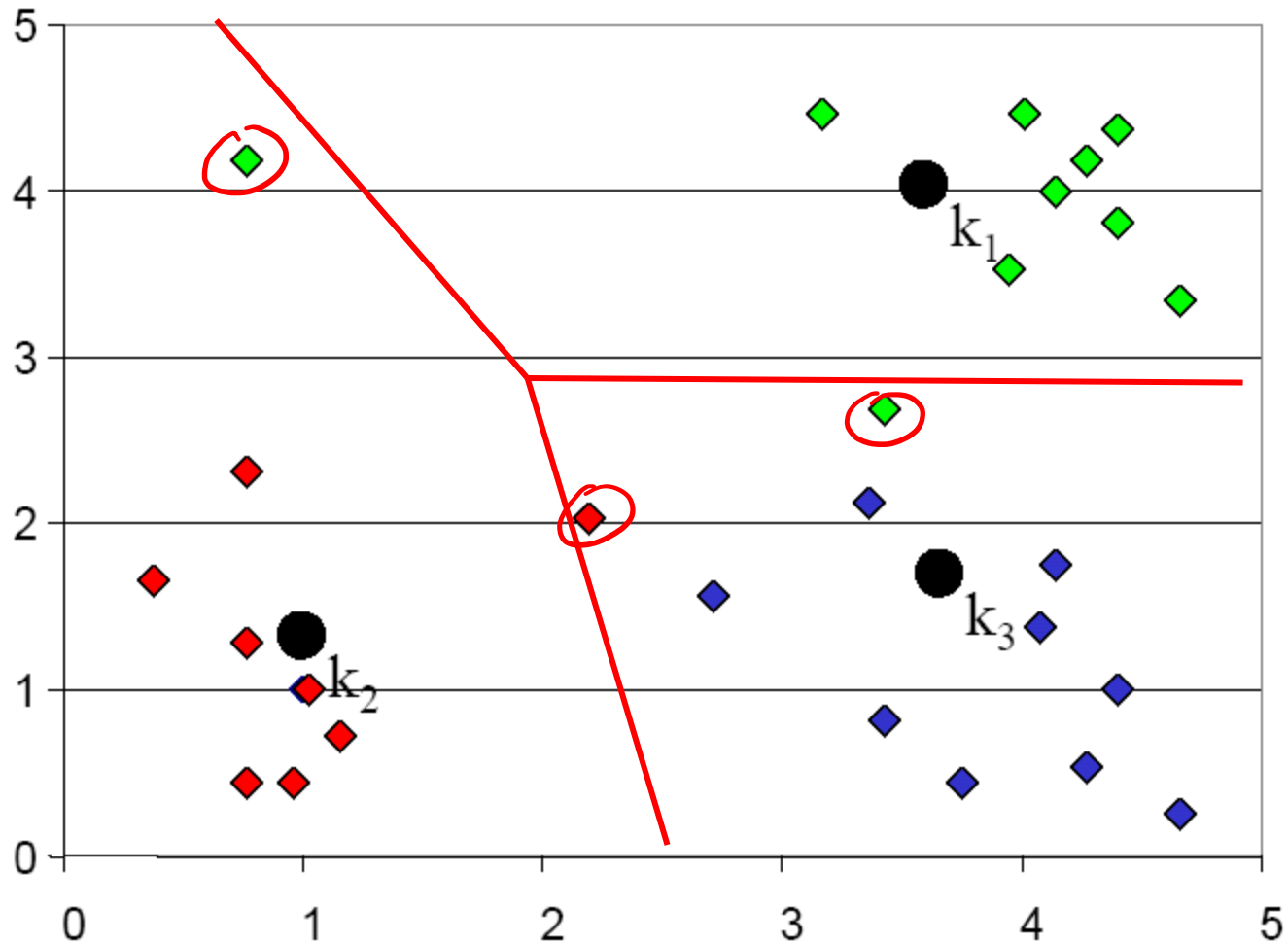


Voronoi diagram

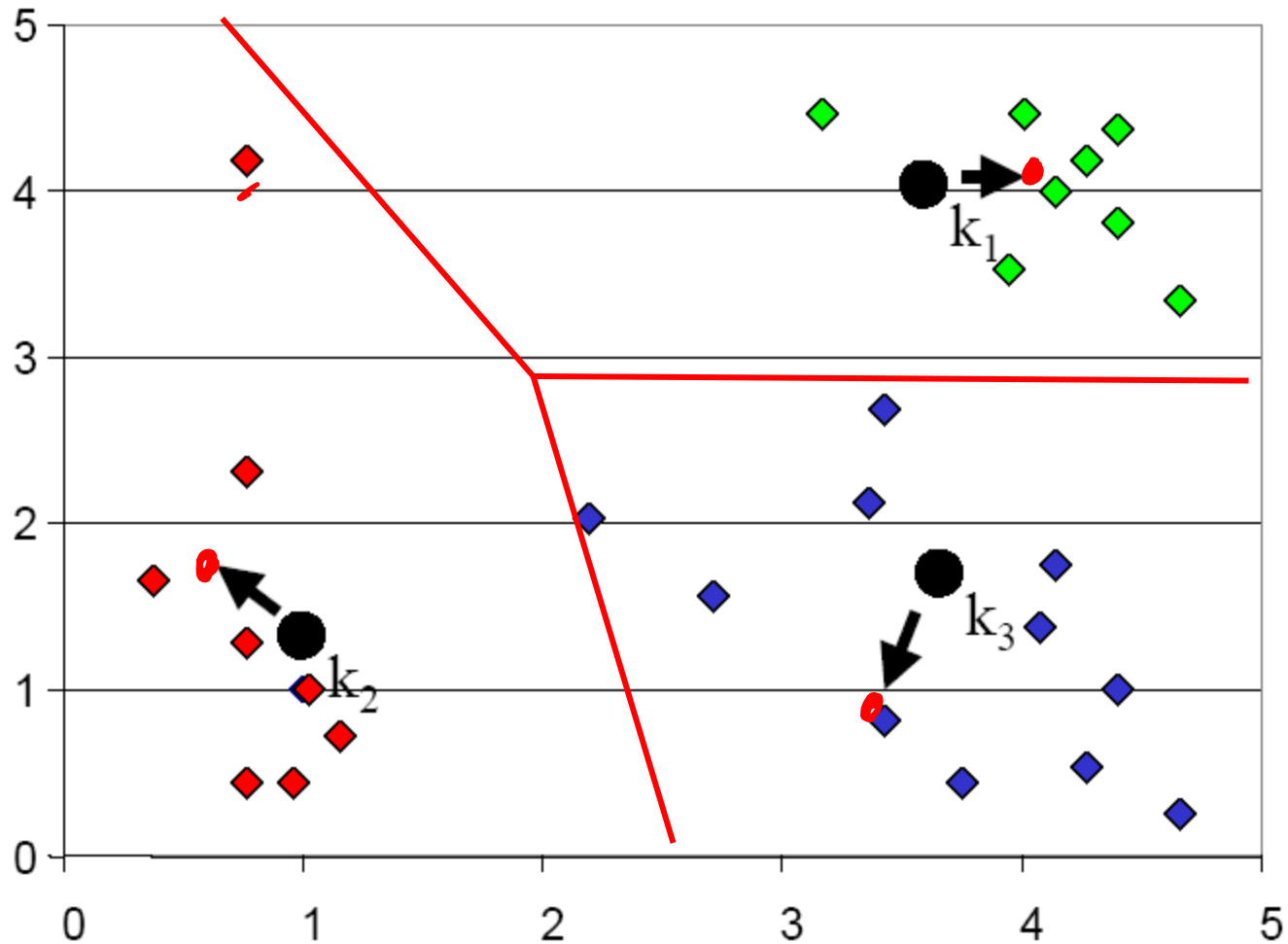
# K-means Clustering: Step 2



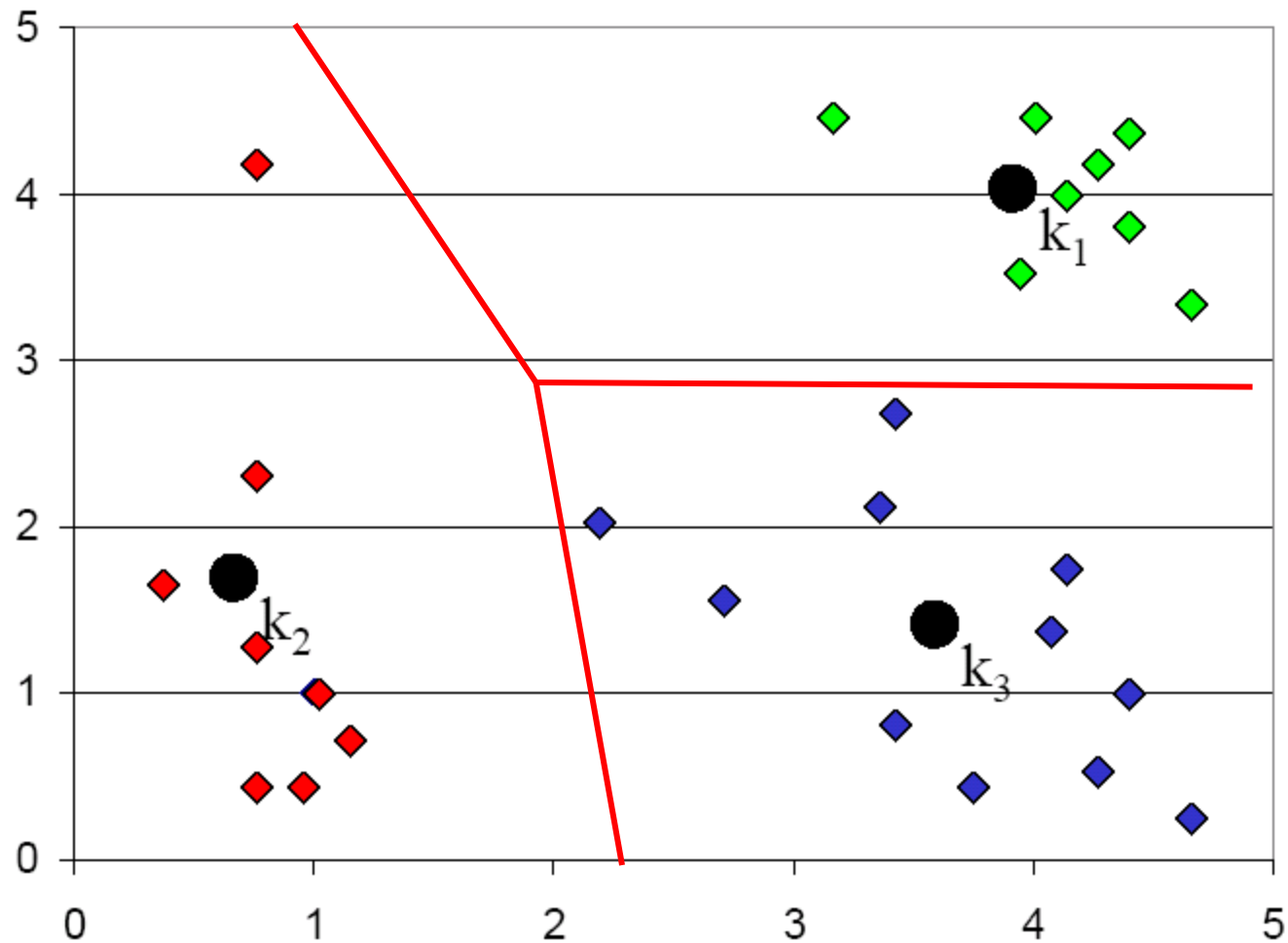
# K-means Clustering: Step 3



# K-means Clustering: Step 4



# K-means Clustering: Step 5





# K-means Recap ...

- Randomly initialize  $k$  centers
  - $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

# K-means Recap ...

- Randomly initialize  $k$  centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

Iterate  $t = 0, 1, 2, \dots$

- **Classify:** Assign each point  $j \in \{1, \dots, m\}$  to nearest center:

- $\underline{C^{(t)}}(j) \leftarrow \arg \min_{\underline{i=1, \dots, k}} \|\underline{\mu_i^{(t)}} - \underline{x_j}\|^2 \leftarrow$

# K-means Recap ...

- Randomly initialize  $k$  centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

Iterate  $t = 0, 1, 2, \dots$

- **Classify:** Assign each point  $j \in \{1, \dots, m\}$  to nearest center:

- $C^{(t)}(j) \leftarrow \arg \min_{i=1, \dots, k} \|\mu_i^{(t)} - x_j\|^2$  ←

- **Recenter:**  $\mu_i$  becomes centroid of its points: ←

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2 \quad i \in \{1, \dots, k\}$

- Equivalent to  $\mu_i \leftarrow$  average of its points!

# What is K-means optimizing?

- Potential function  $F(\mu, C)$  of centers  $\mu$  and point allocations  $C$ :

$$\begin{aligned} F(\mu, C) &= \sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2 \\ &= \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2 \end{aligned}$$

*Handwritten annotations:* Red arrows point to  $\mu$  and  $C$  in the first equation. Red underlines are under  $\mu_{C(j)}$  and  $x_j$ . A red box encloses the second equation, with a red arrow pointing from the first equation to it. A red '1' is written next to the inner sum, and a red note  $C(j)=i$  is written to the right. A red 'j' is written below the inner sum.

- Optimal K-means:

- $\min_{\mu} \min_C F(\mu, C)$

➤ Is the K-means objective convex?

# K-means algorithm

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$$

- **K-means algorithm:** (coordinate descent on F)

**(1)** Fix  $\mu$ , optimize C

**Expected** cluster assignment

**(2)** Fix C, optimize  $\mu$

**Maximum** likelihood for center

Similar to EM/Baum Welch algorithm for learning HMM parameters