

# COMS 4771

## Clustering

Nakul Verma

# Supervised Learning

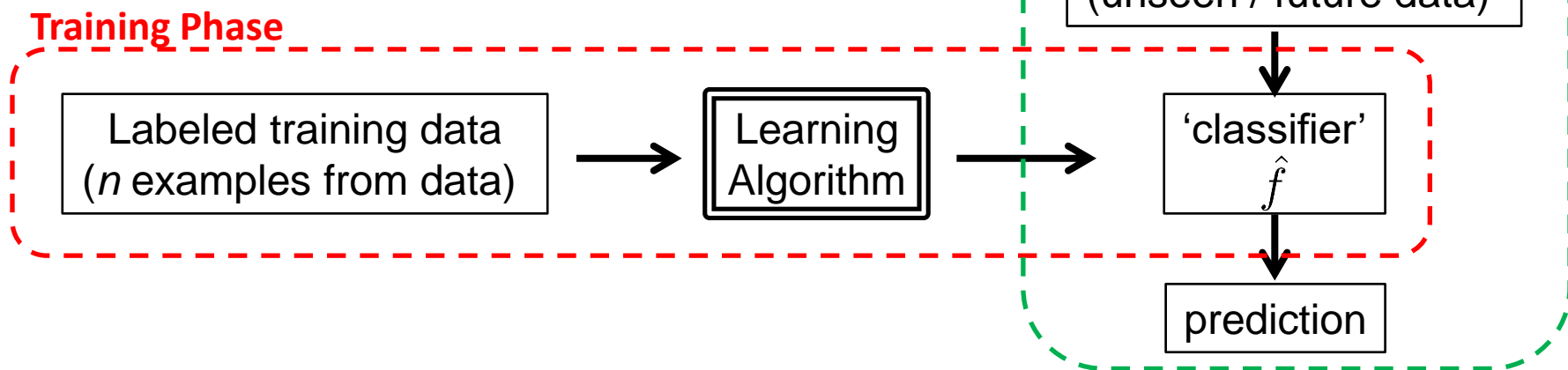
Data:  $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots \in \mathcal{X} \times \mathcal{Y}$

**Supervised learning**

Assumption: there is a (relatively simple) function  $f^* : \mathcal{X} \rightarrow \mathcal{Y}$   
such that  $f^*(\vec{x}_i) = y_i$  for most  $i$

Learning task: given  $n$  examples from the data, find an approximation  $\hat{f} \approx f^*$

Goal:  $\hat{f}$  gives mostly correct prediction  
on unseen examples



# Unsupervised Learning

Data:  $\vec{x}_1, \vec{x}_2, \dots \in \mathcal{X}$

**Unsupervised learning**

Assumption: there is an underlying structure in  $\mathcal{X}$

Learning task: discover the structure given  $n$  examples from the data

Goal: come up with the summary of the data using the discovered structure

*Partition the data into meaningful structures*

clustering

*Find a low-dimensional representation that retains important information, and suppresses irrelevant/noise information*

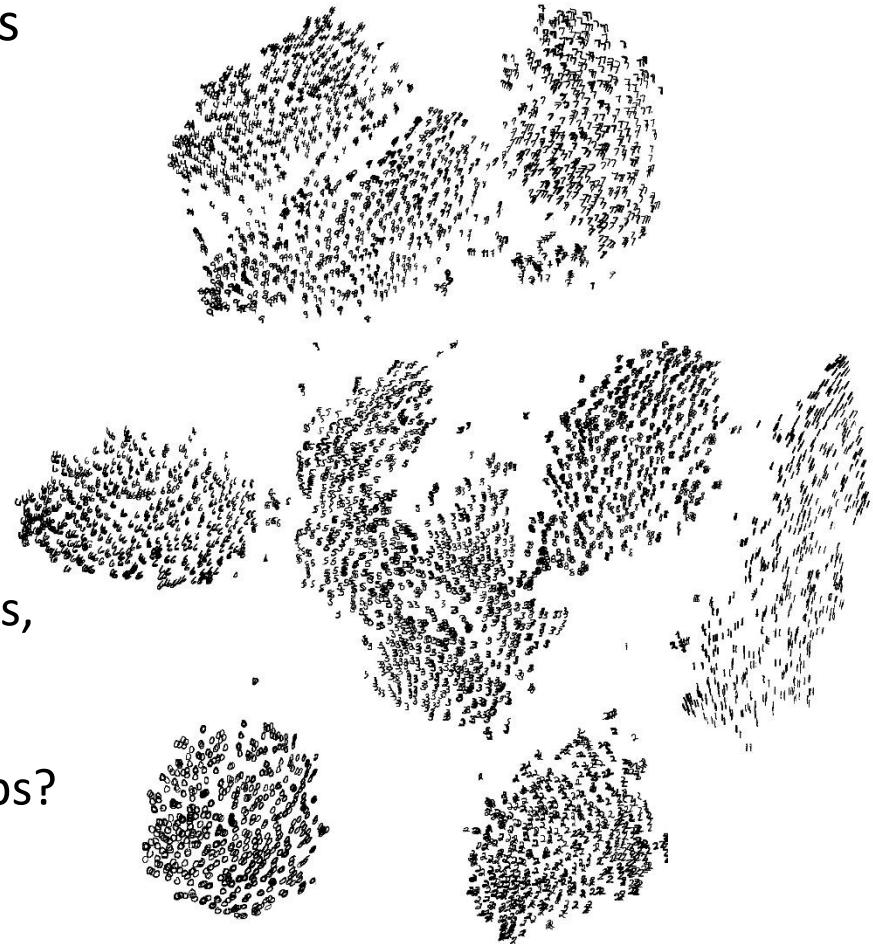
Dimensionality  
reduction

**Let's take a closer look using an example...**

# Example: Handwritten digits revisited

Handwritten digit data, but with no labels

0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9

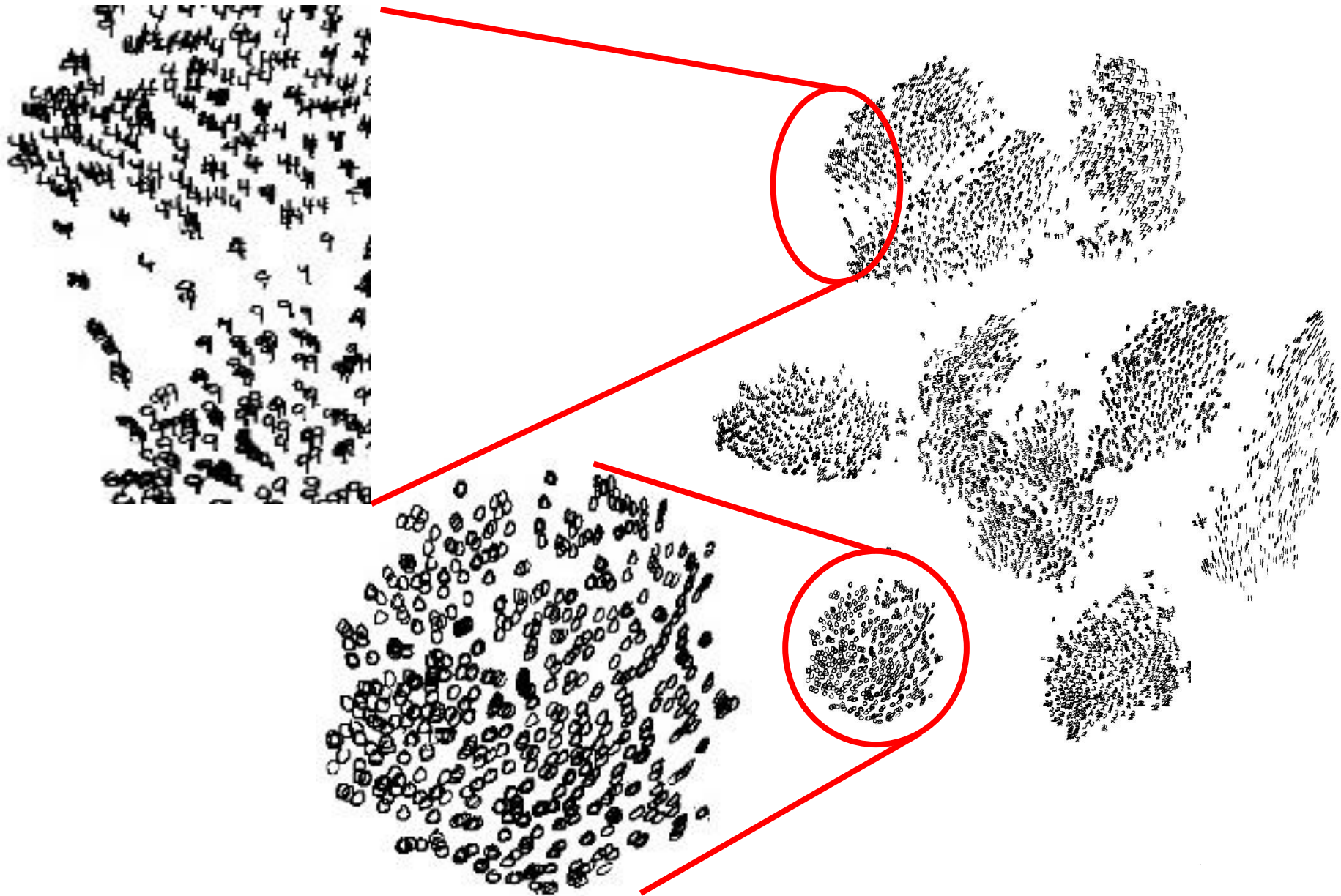


What can we do?

- Suppose know that there are 10 groupings, can we *find the groups*?
- What if we don't know there are 10 groups?
- How can we *discover/explore* other structure in such data?

*A 2D visualization of digits dataset*

# Handwritten digits visualization

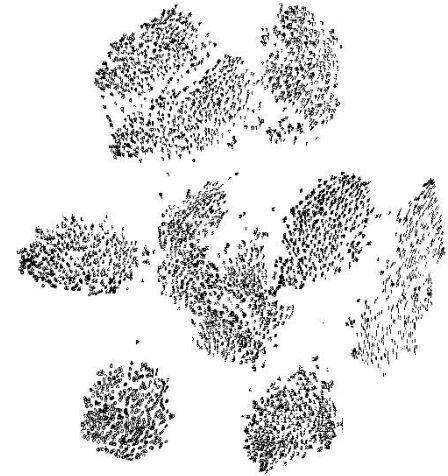


# Grouping The Data, aka Clustering

Data:  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathcal{X}$

Given: known target number of groups  $k$

Output: Partition  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$  into  $k$  groups.



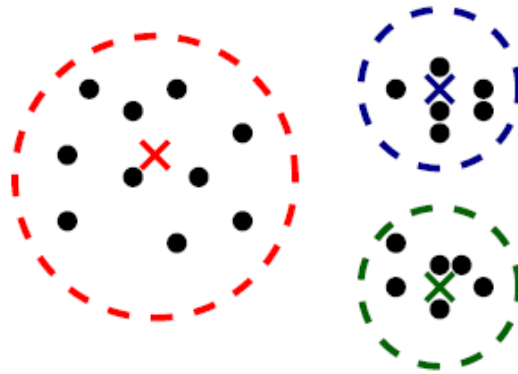
**This is called the clustering problem,  
also known as unsupervised classification, or quantization**

# k-means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Idea:

find a set of representatives  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  such that data is **close to** some representative



Optimization:

$$\text{minimize}_{c_1, \dots, c_k} \left[ \sum_{i=1}^n \min_{j=1, \dots, k} \|\vec{x}_i - \vec{c}_j\|^2 \right]$$

How do we optimize this?

*Unfortunately this is NP-hard  
Even for  $d=2$  and  $k=2$*

*How do we solve for  
 $d=1$  or  $k=1$  case?*

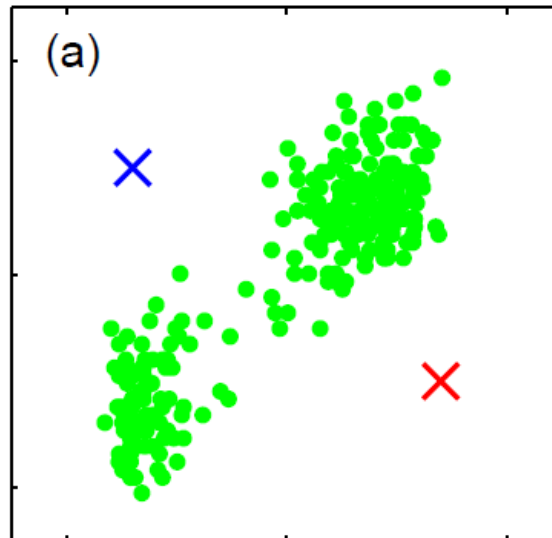
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:





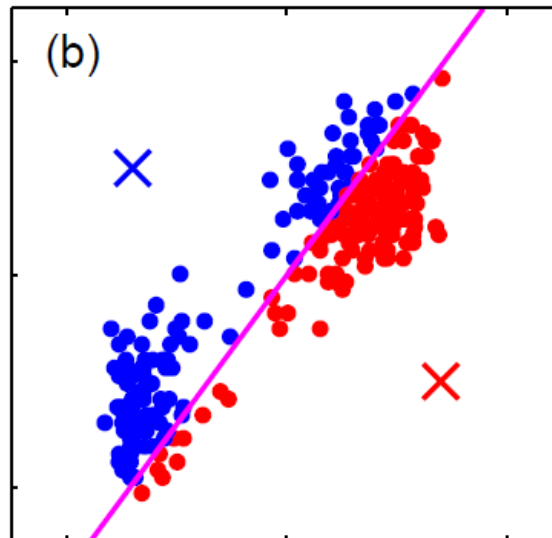
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:



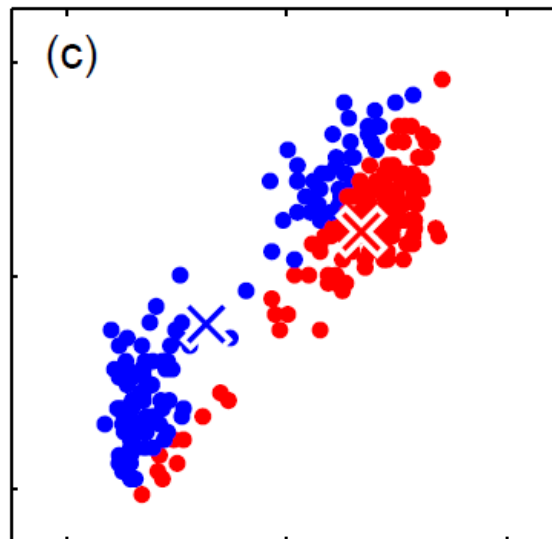
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:



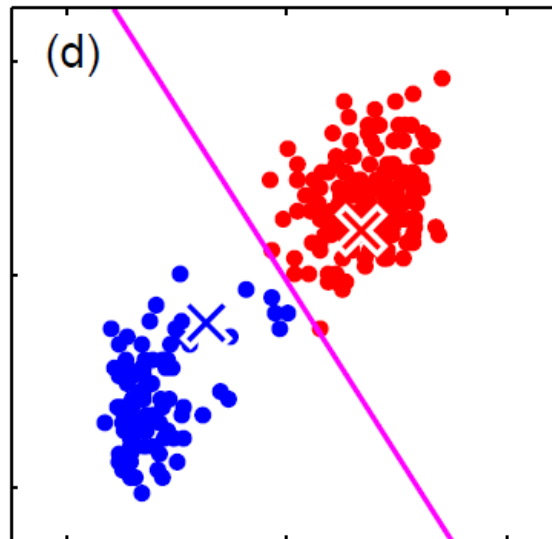
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:



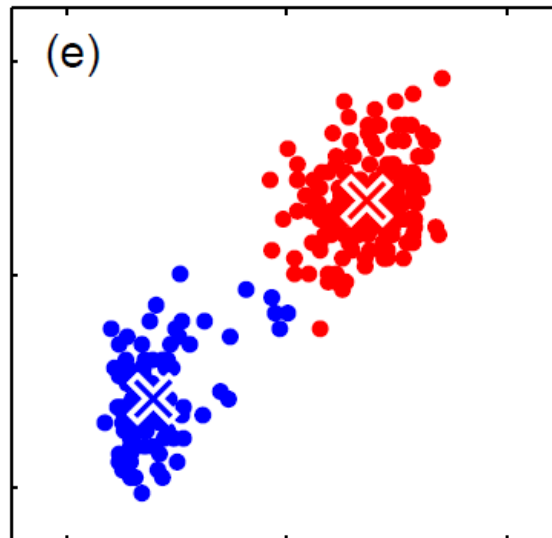
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:



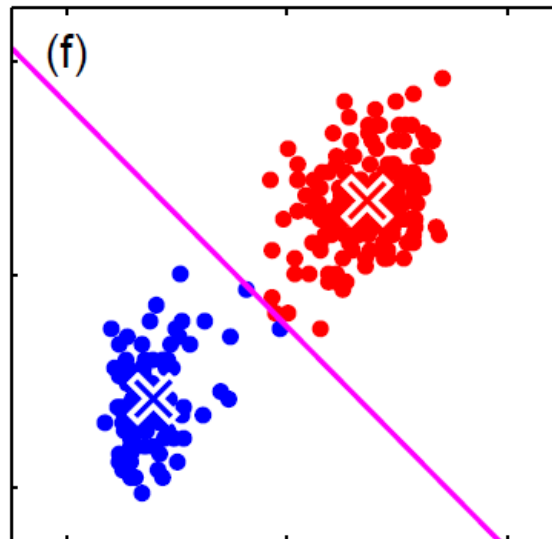
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:



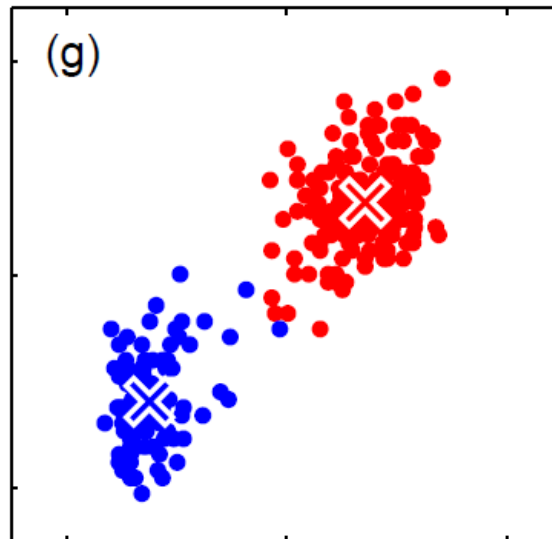
# Algorithm to approximate $k$ -means

Given: data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$ , and intended number of groupings  $k$

Alternating optimization algorithm:

- Initialize cluster centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (say randomly)
- Repeat till no more changes occur
  - Assign data to its closest center (this creates a partition) (assume centers are fixed)
  - Find the optimal centers  $\vec{c}_1, \vec{c}_2, \dots, \vec{c}_k$  (assuming the data partition is fixed)

Demo:



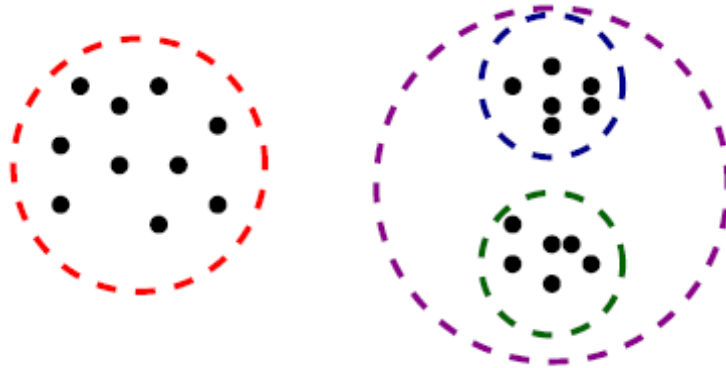
# *k*-means

Some properties of this alternating updates algorithm:

- The approximation can be arbitrarily bad, compared to the best cluster assignment!
- Performance quality heavily dependent on the initialization!

*k*-means:

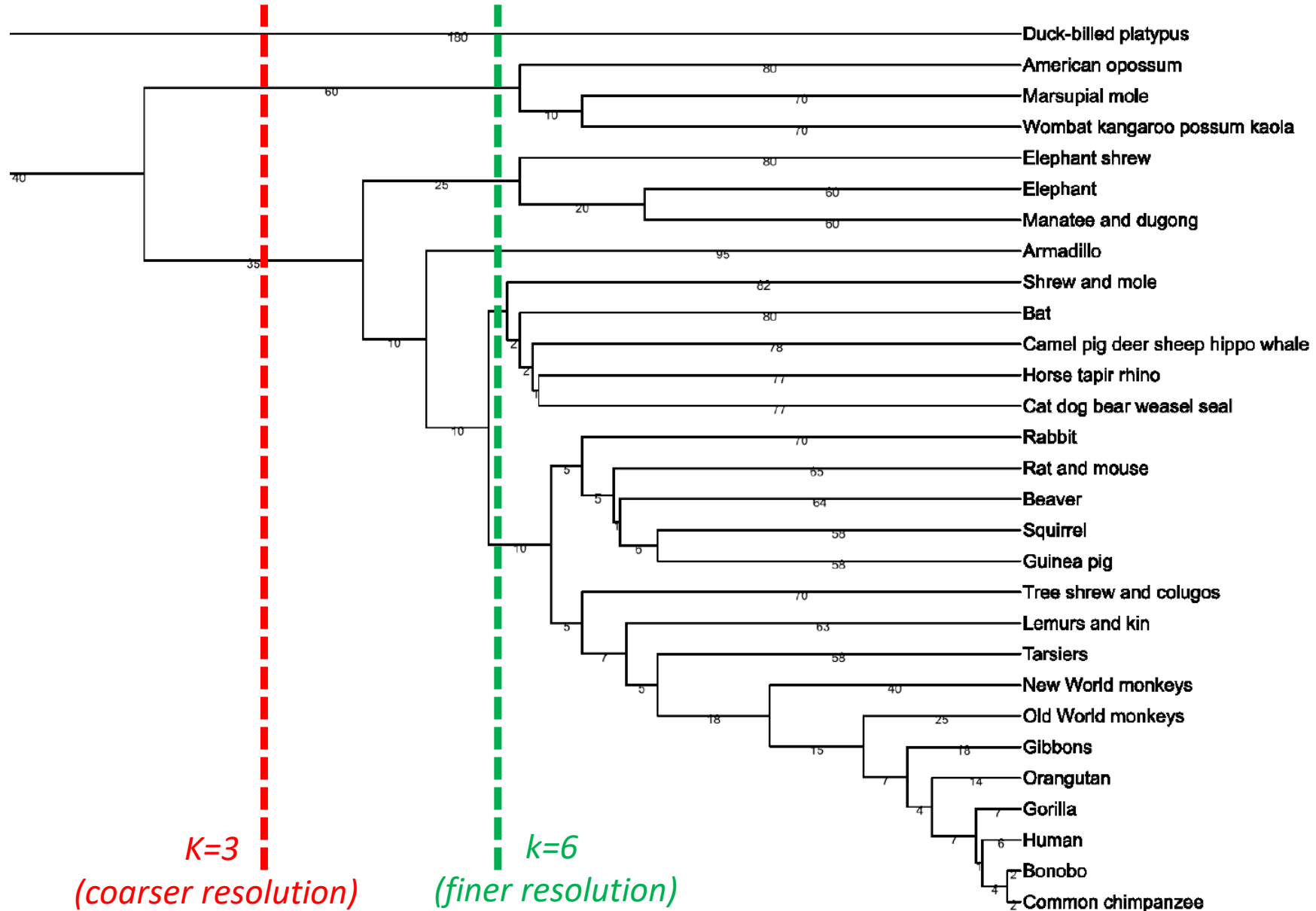
- How to select *k*?



*is the right  $k=2$  or  $k=3$ ?*

**Solution:** encode clustering  
for all values of  $k$ !  
(hierarchical clustering)

# Example: Clustering Without Committing to $k$





# Hierarchical Clustering

Two approaches:

Top Down (divisive):

- Partition data into two groups (say, by k-means, with  $k=2$ )
- Recurse on each part
- Stop when cannot partition data anymore (ie single points left)

Bottom Up (agglomerative):

- Start by each data sample as its own cluster (so initial number of clusters is  $n$ )
- Repeatedly merge “closest” pair of clusters
- Stop when only one cluster is left

# Clustering via Probabilistic Mixture Modeling

Alternative way to cluster data:

Given:  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$  and number of intended number of clusters  $k$ .

Assume a joint probability distribution  $(X, C)$  over the joint space  $\mathbf{R}^d \times [k]$

$$C \sim \begin{pmatrix} \pi_1 \\ \vdots \\ \pi_k \end{pmatrix} \quad \text{Discrete distribution over the clusters } P[C=i] = \pi_i$$

$$X|C = i \sim \text{Some multivariate distribution, e.g. } N(\vec{\mu}_i, \Sigma_i)$$

Parameters:  $\theta = (\pi_1, \vec{\mu}_1, \Sigma_1, \dots, \pi_k, \vec{\mu}_k, \Sigma_k)$  **looks familiar?**

Modeling assumption data  $(x_1, c_1), \dots, (x_n, c_n)$  i.i.d. from  $\mathbf{R}^d \times [k]$

BUT only get to see partial information:  $x_1, x_2, \dots, x_n$  ( $c_1, \dots, c_n$  hidden!)

# Gaussian Mixture Modeling (GMM)

Given:  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$  and  $k$ .

Assume a joint probability distribution  $(X, C)$  over the joint space  $\mathbf{R}^d \times [k]$

$$C \sim \begin{pmatrix} \pi_1 \\ \vdots \\ \pi_k \end{pmatrix} \quad X|C = i \sim N(\vec{\mu}_i, \Sigma_i) \quad \text{Gaussian Mixture Model}$$
$$\theta = (\pi_1, \vec{\mu}_1, \Sigma_1, \dots, \pi_k, \vec{\mu}_k, \Sigma_k)$$

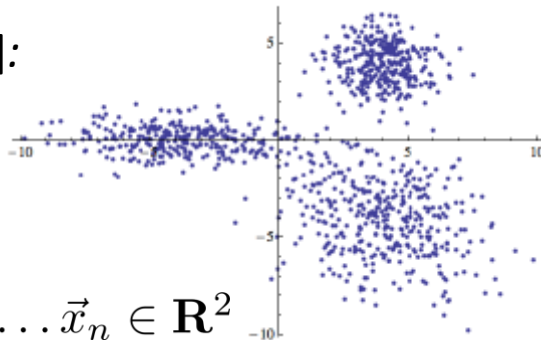
$$P[\vec{x} | \theta] = \sum_{j=1}^k \pi_j \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu}_j)^\top \Sigma_j^{-1} (\vec{x} - \vec{\mu}_j) \right\}$$

Mixing weight

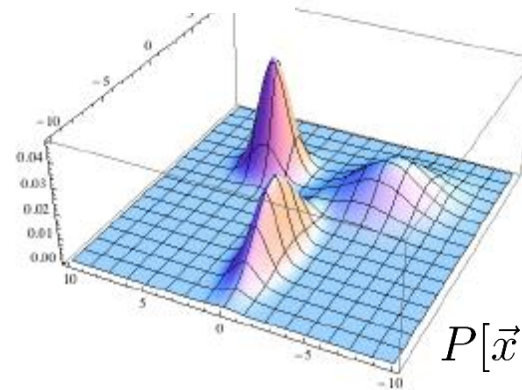
Mixture component

*(this is called a mixture model)*

Example in  $\mathbf{R}^2 \times [3]$ :



$\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^2$



$P[\vec{x} | \theta]$

# GMM: Parameter Learning

$$P[\vec{x} | \theta] = \sum_{j=1}^k \pi_j \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu}_j)^\top \Sigma_j^{-1} (\vec{x} - \vec{\mu}_j) \right\}$$

$$\theta = (\pi_1, \vec{\mu}_1, \Sigma_1, \dots, \pi_k, \vec{\mu}_k, \Sigma_k)$$

*So... how to learn the parameters  $\theta$ ?*

MLE approach:

Given data  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$  i.i.d.

$$\theta_{\text{MLE}} := \arg \max_{\theta} \sum_{i=1}^n \ln P[\vec{x}_i | \theta]$$

$$= \arg \max_{\theta} \sum_{i=1}^n \ln \left[ \sum_{j=1}^k \pi_j \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left\{ -\frac{1}{2} (\vec{x}_i - \vec{\mu}_j)^\top \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j) \right\} \right]$$

*ummm.... now what?*

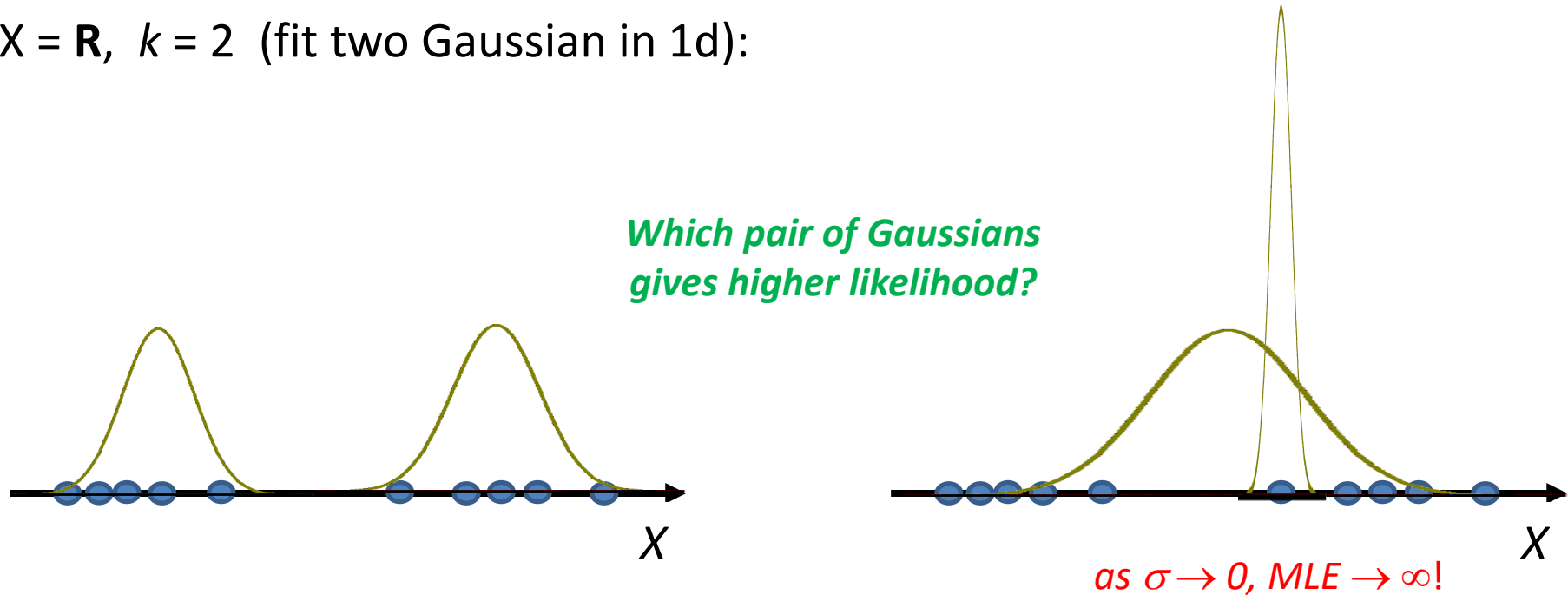
Cannot really simplify further!

# GMM: Maximum Likelihood

MLE for Mixture modeling (like GMMs) is NOT a convex optimization problem

In fact **Maximum** Likelihood Estimate for GMMs is degenerate!

$X = \mathbf{R}$ ,  $k = 2$  (fit two Gaussian in 1d):

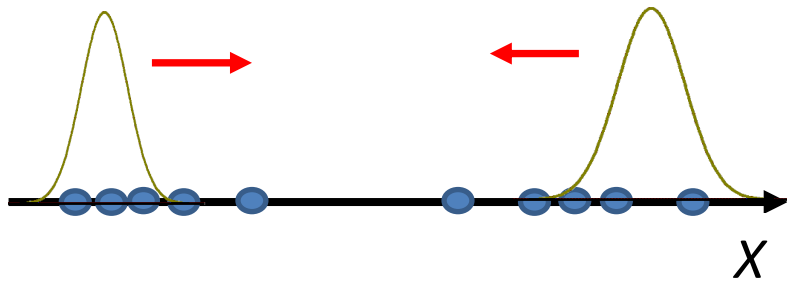


Aside: why doesn't this occur when fitting one Gaussian?

# GMM: (local) Maximum Likelihood

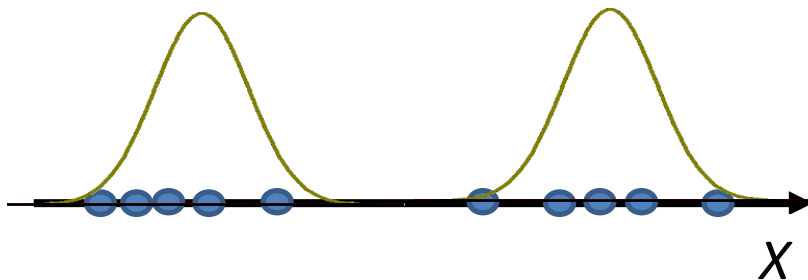
So, can we make any progress?

Observation: even though a global MLE maximizer is not appropriate, several local maximizers are desirable!



*An example  
non-maximized  
likelihood*

*(do a few steps of gradient ascent)*



*Reaches a desirable  
local maximum!*

*A better algorithm for finding good parameters: Expectation Maximization (EM)*

# Expectation Maximization (EM) Algorithm

Similar in spirit to the alternating update for  $k$ -means algorithm

Idea:

- Initialize the parameters arbitrarily
- Given the current setting of parameters find the best (soft) assignment of data samples to the clusters (**Expectation-step**)
- Update all the parameters with respect to the current (soft) assignment that maximizes the likelihood (**Maximization-step**)
- Repeat until no more progress is made.

# EM for GMM

Initialize  $\theta = (\pi_1, \vec{\mu}_1, \Sigma_1, \dots, \pi_k, \vec{\mu}_k, \Sigma_k)$  arbitrarily

**Expectation-step:** For each  $i \in \{1, \dots, n\}$  and  $j \in \{1, \dots, k\}$  compute the assignment  $w_j^{(i)}$  of data  $x_i$  to cluster  $j$

$$w_j^{(i)} := \frac{\pi_j \sqrt{\det(\Sigma_j^{-1})} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu}_j)^\top \Sigma_j^{-1}(\vec{x}_i - \vec{\mu}_j)\right)}{\sum_{j'=1}^k \pi_{j'} \sqrt{\det(\Sigma_{j'}^{-1})} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu}_{j'})^\top \Sigma_{j'}^{-1}(\vec{x}_i - \vec{\mu}_{j'})\right)}$$

**Maximization-step:** Maximize the log-likelihood of the parameters

$$n_j := \sum_{i=1}^n w_j^{(i)} \left( \begin{array}{l} \text{Effective number of points} \\ \text{assigned to cluster } j \end{array} \right) \quad \pi_j := \frac{n_j}{n}$$

$$\vec{\mu}_j := \frac{1}{n_j} \sum_{i=1}^n w_j^{(i)} \vec{x}_i \quad \Sigma_j := \frac{1}{n_j} \sum_{i=1}^n w_j^{(i)} (\vec{x}_i - \vec{\mu}_j)(\vec{x}_i - \vec{\mu}_j)^\top$$

Why?



# EM for GMM

Calculation for updating  $\mu_j$

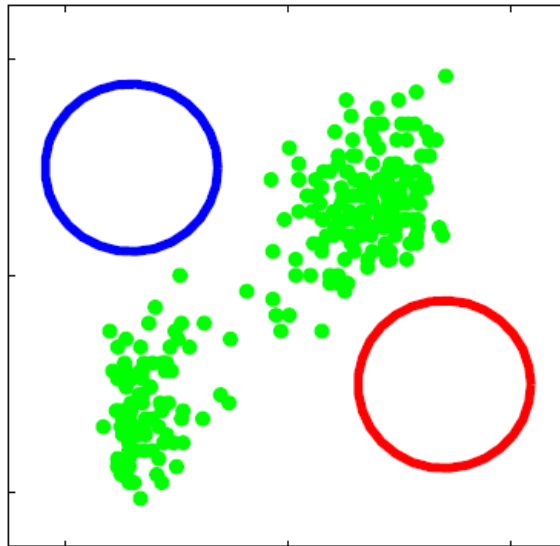
Recall, log likelihood: 
$$\sum_{i=1}^n \ln \left[ \sum_{j=1}^k \pi_j \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left\{ -\frac{1}{2} (\vec{x}_i - \vec{\mu}_j)^\top \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j) \right\} \right]$$

To maximize w.r.t.  $\mu_j$

$$\begin{aligned} & \left( \frac{\partial}{\partial \vec{\mu}_j} \right) \left( \sum_{i=1}^n \ln \left[ \sum_{j=1}^k \pi_j \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left\{ -\frac{1}{2} (\vec{x}_i - \vec{\mu}_j)^\top \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j) \right\} \right] \right) \\ &= \sum_{i=1}^n w_j^{(i)} \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j) \quad \left( \begin{array}{l} \text{since} \\ \frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{s})^\top \mathbf{W} (\mathbf{x} - \mathbf{s}) = -2\mathbf{W}(\mathbf{x} - \mathbf{s}) \end{array} \right) \end{aligned}$$

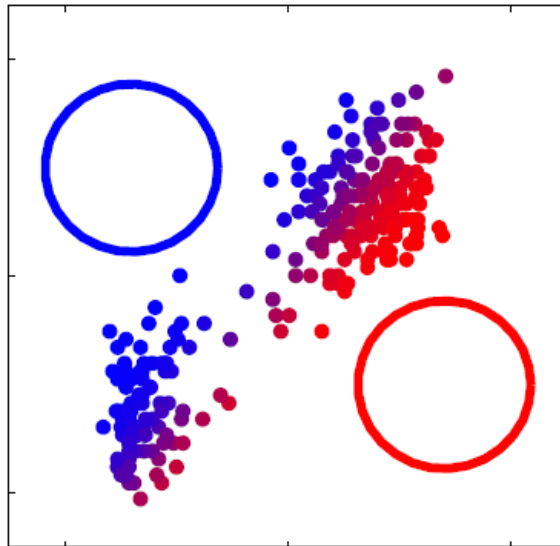
Therefore at stationarity 
$$\vec{\mu}_j := \frac{1}{n_j} \sum_{i=1}^n w_j^{(i)} \vec{x}_i$$

# EM for GMM in Action



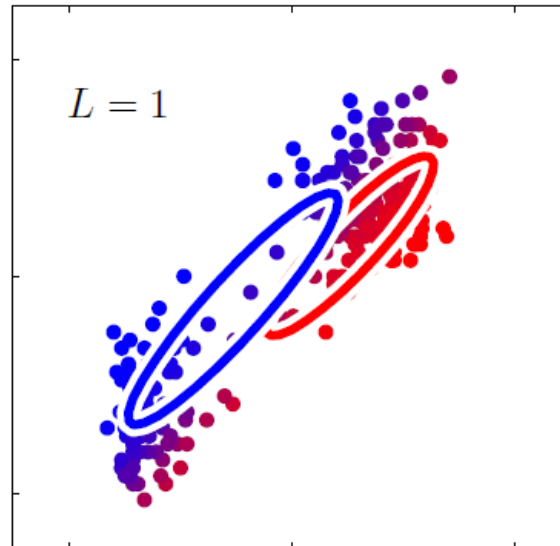
*Arbitrary  $\theta$   
assignment*

# EM for GMM in Action



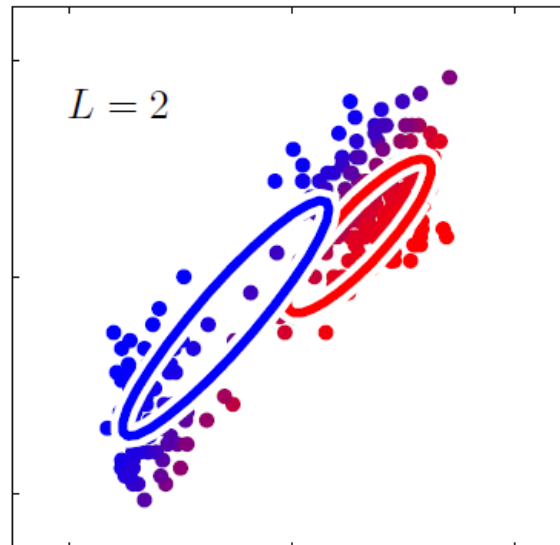
*E step: soft  
assignment of data*

# EM for GMM in Action



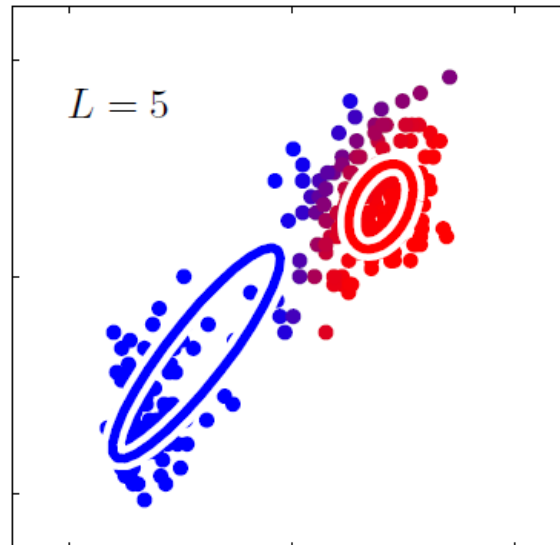
*M step: Maximize  
parameter estimate*

# EM for GMM in Action



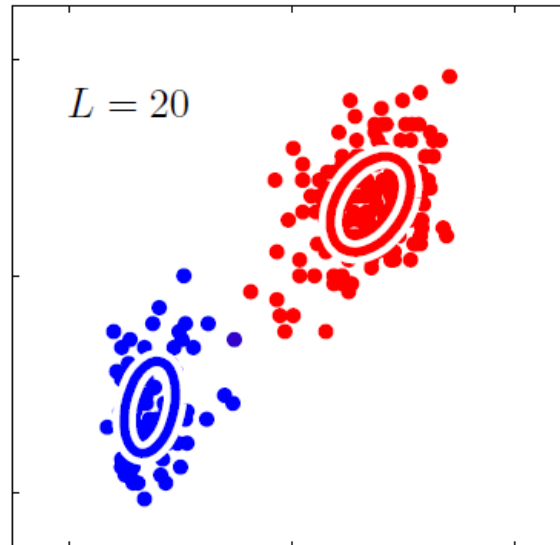
*After two rounds*

# EM for GMM in Action



*After five rounds*

# EM for GMM in Action



*After twenty rounds*

# What We Learned...

- Unsupervised Learning problems:  
    Clustering and Dimensionality Reduction
- K-means
- Hierarchical Clustering
- Gaussian Mixture Models
- EM algorithm



# Questions?

# Next time...

Dimension reduction!