# UVA CS 4501:
# Machine Learning

# Lecture 24: Unsupervised Clustering (III)

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

# Where are we ? ➜
# major sections of this course

❑ Regression (supervised)

❑ Classification (supervised)

    ❑ Feature selection

❑ Unsupervised models

    ➡ ❑ Dimension Reduction (PCA)

    ❑ Clustering (K-means, GMM/EM, Hierarchical )

❑ Learning theory

❑ Graphical models

$$X_1 \quad X_2 \quad X_3$$

$s_1$
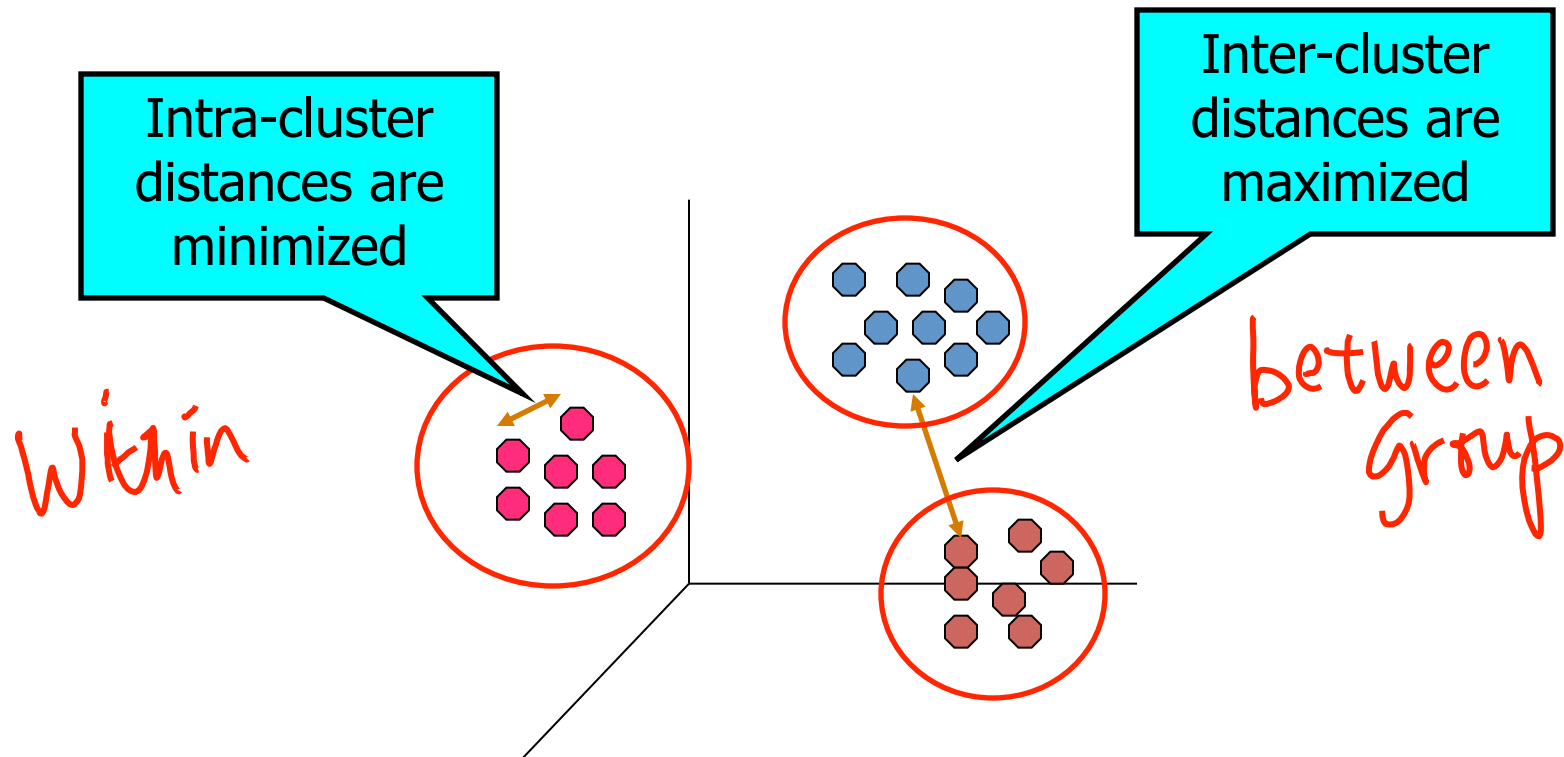
$s_2$

$s_3$

$s_4$

$s_5$

$s_6$

# An unlabeled Dataset X

a data matrix of $n$ observations on $p$ variables $x_1, x_2, \ldots x_p$

Unsupervised learning = learning from raw (unlabeled, unannotated, etc) data, as opposed to supervised data where a classification label of examples is given

- **Data**/*points/instances/examples/samples/records*: [ rows ]
- **Features**/*attributes/dimensions/independent variables/covariates/predictors/regressors*: [ columns]

# What is clustering?

- Find groups (clusters) of data points such that data points in a group will be similar (or related) to one another and different from (or unrelated to) the data points in other groups
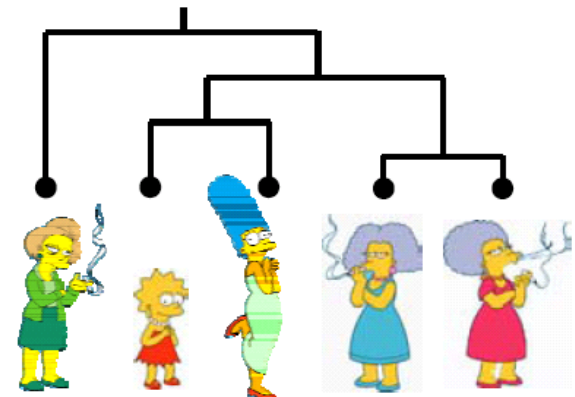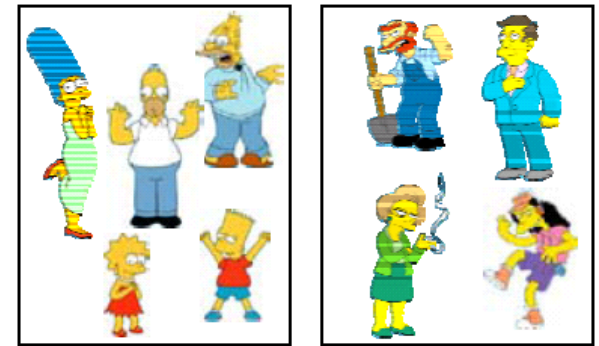


Intra-cluster distances are minimized

Inter-cluster distances are maximized

Within

between Group

# **Roadmap:** clustering

- Definition of "groupness"

- Definition of "similarity/distance"

- Representation for objects

- How many clusters?

- Clustering Algorithms

  ➡ ■ Partitional algorithms

  - Hierarchical algorithms

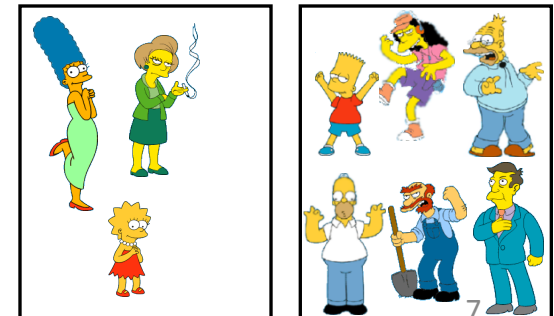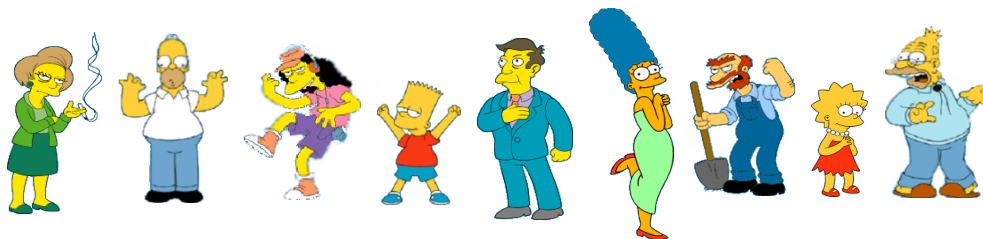- Formal foundation and convergence

# Clustering Algorithms

- Partitional algorithms
  - Usually start with a random (partial) partitioning
  - Refine it iteratively
    - K means clustering
    - Mixture-Model based clustering

- Hierarchical algorithms
  - Bottom-up, agglomerative
  - Top-down, divisive

# (2) Partitional Clustering

- Nonhierarchical

- Construct a partition of *n* objects into a set of *K* clusters

- User has to specify the desired number of clusters K.

# Other partitioning Methods

- Partitioning around medoids (PAM): instead of averages, use multidim medians as centroids (cluster "prototypes"). Dudoit and Freedland (2002).

# Other partitioning Methods

- Partitioning around medoids (PAM): instead of averages, use multidim medians as centroids (cluster "prototypes"). Dudoit and Freedland (2002).

- Self-organizing maps (SOM): add an underlying "topology" (neighboring structure on a lattice) that relates cluster centroids to one another. Kohonen (1997), Tamayo et al. (1999).

# Other partitioning Methods

- Partitioning around medoids (PAM): instead of averages, use multidim medians as centroids (cluster "prototypes"). Dudoit and Freedland (2002).

- Self-organizing maps (SOM): add an underlying "topology" (neighboring structure on a lattice) that relates cluster centroids to one another. Kohonen (1997), Tamayo et al. (1999).

- Fuzzy k-means: allow for a "gradation" of points between clusters; soft partitions. Gash and Eisen (2002).

# Other partitioning Methods

$C_j \in$ train Set

- Partitioning around medoids (PAM): instead of averages, use multidim medians as centroids (cluster "prototypes"). Dudoit and Freedland (2002).

- Self-organizing maps (SOM): add an underlying "topology" (neighboring structure on a lattice) that relates cluster centroids to one another. Kohonen (1997), Tamayo et al. (1999).

- Fuzzy k-means: allow for a "gradation" of points between clusters; soft partitions. Gash and Eisen (2002).

- Mixture-based clustering: implemented through an EM (Expectation-Maximization)algorithm. This provides soft partitioning, and allows for modeling of cluster centroids and shapes. (Yeung et al. (2001), McLachlan et al. (2002))

$$M_{ij} \in \{1, 0\} \rightarrow [0, 1]$$

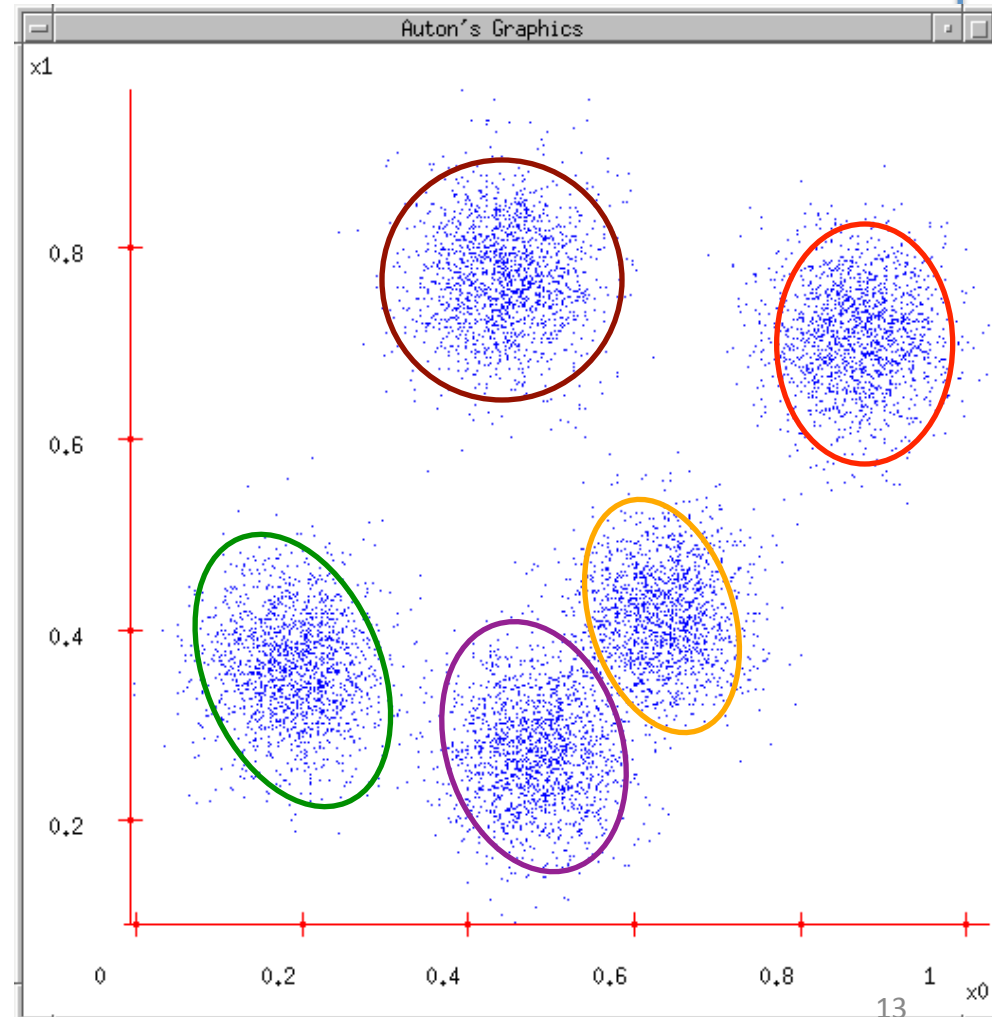# Partitional : Gaussian Mixture Model

- 1. Review of Gaussian Distribution
- 2. GMM for clustering : basic algorithm
- 3. GMM connecting to K-means
- 4. GMM examples
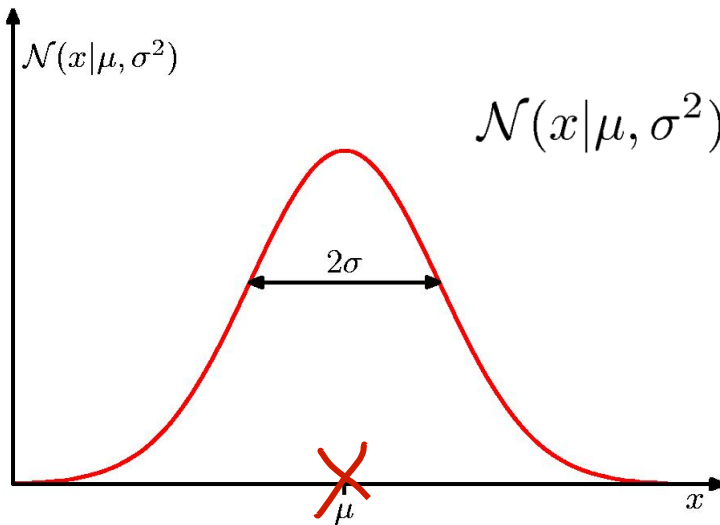- 5. Problems of GMM and K-means

# A Gaussian Mixture Model for Clustering

- Assume that data are generated from a mixture of Gaussian distributions

- For each Gaussian distribution
  - Center: $\mu_i$
  - covariance: $\sum_i$

- For each data point
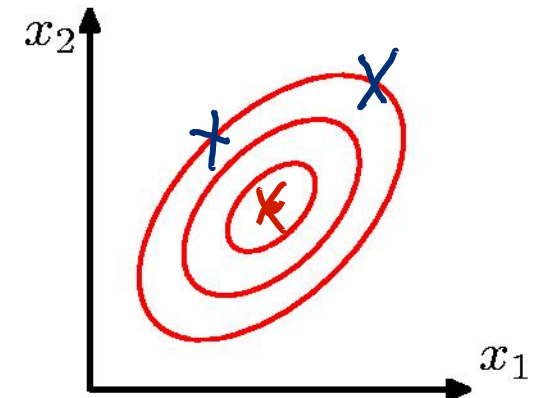  - Determine membership

$z_{ij}$ : if $x_i$ belongs to j-th cluster

# Gaussian Distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

$$X \sim \mathbf{N}\left(\mu, \sigma^2\right)$$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{P/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\}$$

Mean

Covariance Matrix

Courtesy: http://research.microsoft.com/~cmbishop/PRML/index.htm

# Example: the Bivariate Normal distribution

$$p(\vec{x}) = f\left(x_1, x_2\right) = \frac{1}{\left(2\pi\right)\left|\Sigma\right|^{1/2}} e^{-\frac{1}{2}\left(\vec{x}-\vec{\mu}\right)^T \Sigma^{-1}\left(\vec{x}-\vec{\mu}\right)}$$

with $\quad \vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ 2x1 $\quad$ and

$$\Sigma_{2\times 2} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$V(X_1) \qquad Cov(X_1, X_2)$
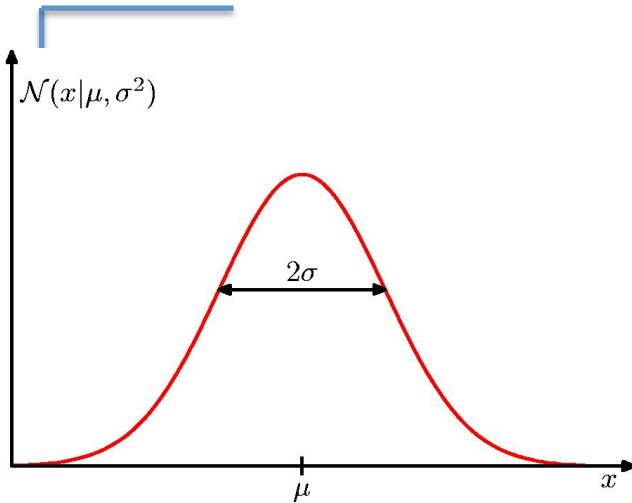
$V(X_2) \quad 2 \times 2$

$$\left|\Sigma\right| = \sigma_{11}\sigma_{22} - \sigma_{12}^2 = \sigma_1^2\sigma_2^2\left(1-\rho^2\right)$$

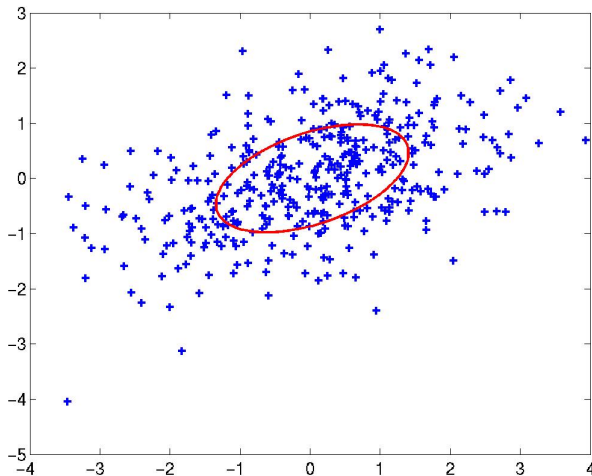# Scatter Plots of data from the bivariate Normal distribution

# How to Estimate Gaussian: MLE

$\mathcal{N}(x|\mu, \sigma^2)$



• In the 1D Gaussian case, we simply set the mean and the variance to the sample mean and the sample variance:

$$\overline{\mu} = \frac{1}{n}\sum_{i=1}^{n} x_i \qquad \overline{\sigma}^2 = \frac{1}{n}\sum_{i=1}^{n}\left(x_i - \overline{\mu}\right)^2$$

# The p-multivariate Normal distribution

$$< X_1, X_2 \cdots, X_p > \sim N\left(\vec{\mu}, \Sigma\right)$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} \; p \times 1$$

$$\mu_i = \frac{1}{n} \sum_{\partial=1}^{N} X_{\partial}^{(i)}$$

$\in \{1, 2, \cdots, p\}$

i-th feature

∂-th Sample

$\in \{1, 2, \cdots, N\}$

$$\sum = \Sigma = \begin{bmatrix} Var(X_1) & & Cov(X_i, X_j) \\ & \ddots & \\ & & Var(X_p) \end{bmatrix}$$
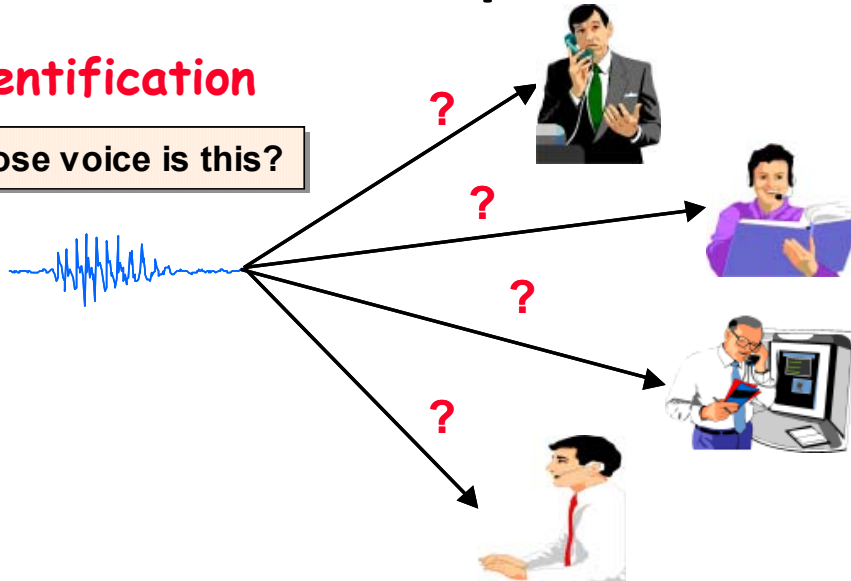
$P \times P$

# Partitional : Gaussian Mixture Model

- 1. Review of Gaussian Distribution
- 2. GMM for clustering : basic algorithm
- 3. GMM connecting to K-means
- 4. GMM examples
- 5. Problems of GMM and K-means

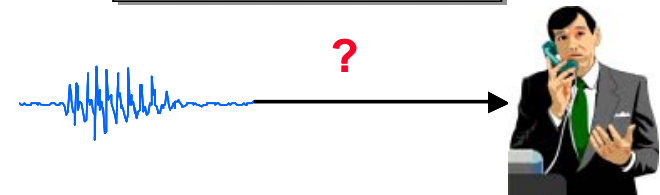# Application:
# Three Speaker Recognition Tasks

**Identification**

Whose voice is this?

? ? ? ?

**Verification/Authentication/ Detection**
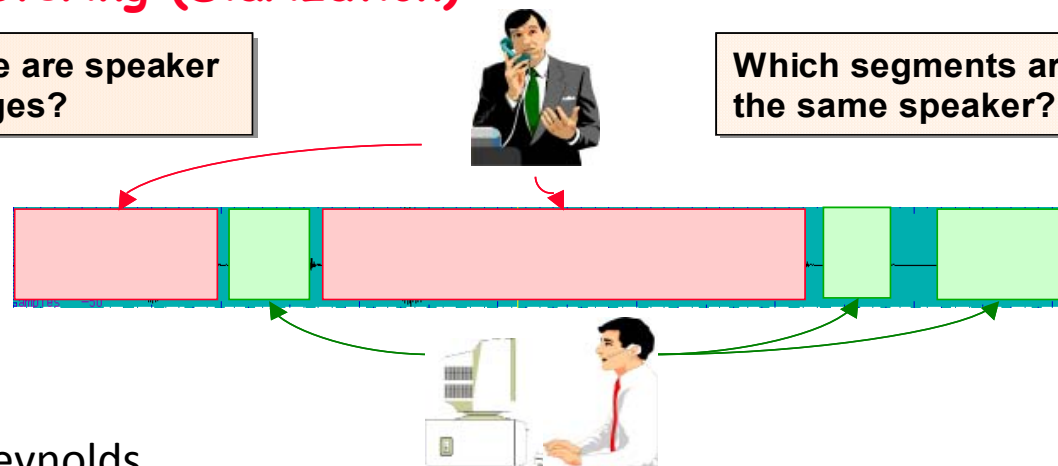
Is this Bob's voice?

?

**Segmentation and Clustering (Diarization)**

Where are speaker changes?

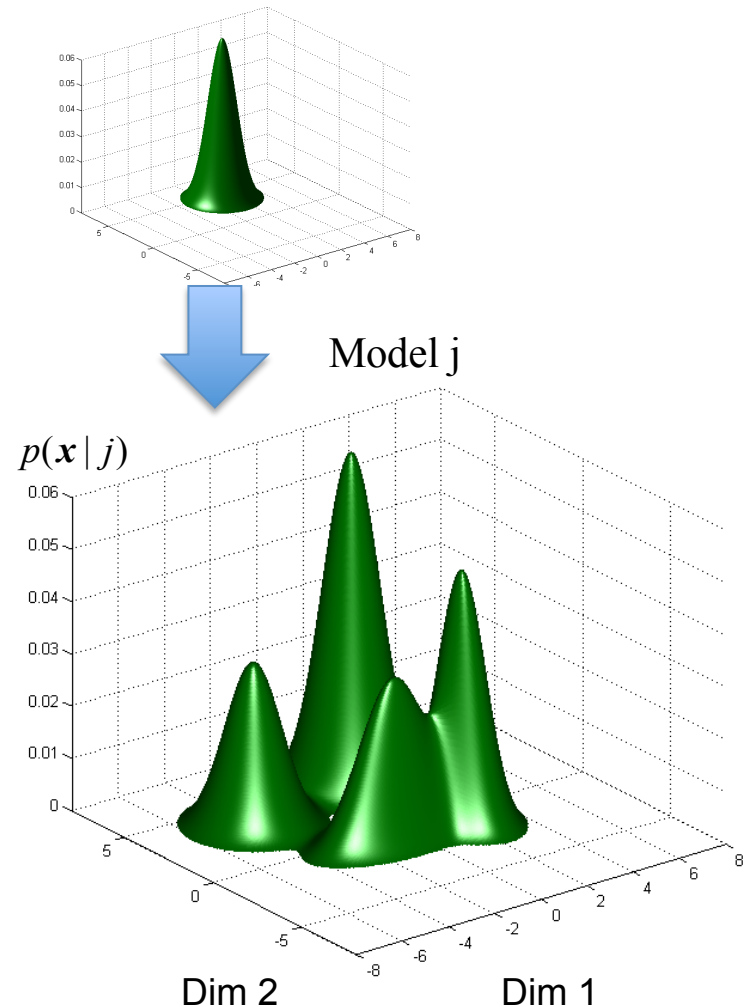Which segments are from the same speaker?

slide from Douglas Reynolds
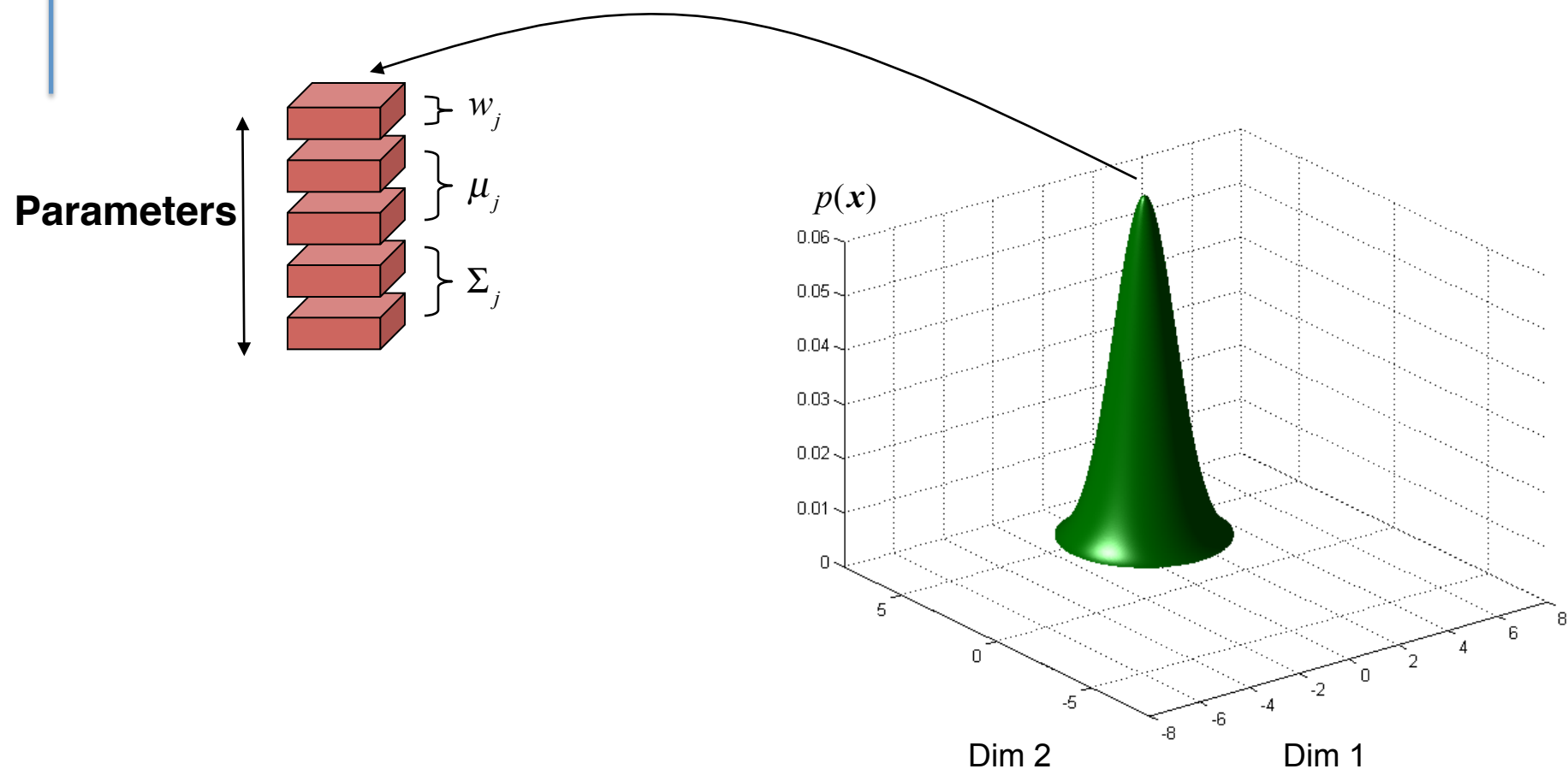
# Application : GMMs for speaker recognition

- A Gaussian mixture model (GMM) represents as the weighted sum of multiple Gaussian distributions

- Each Gaussian state i has a
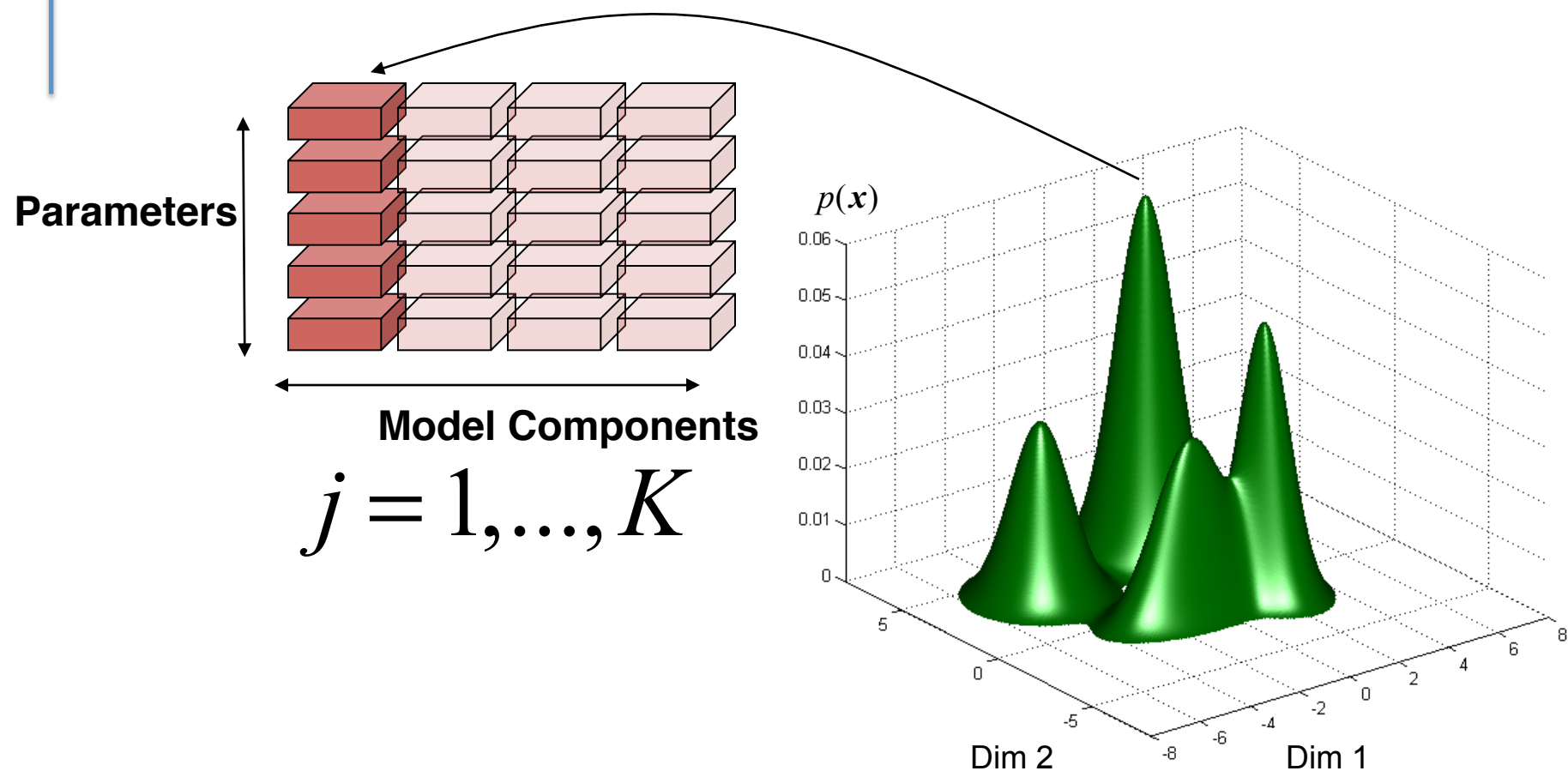  - Mean
  - Covariance $\mu_i$
  - Weight $\Sigma_j$

Model j

$p(\boldsymbol{x} \mid j)$

Dim 2          Dim 1

$$w_j \equiv p(\mu = \mu_j)$$

# Recognition Systems
# Gaussian Mixture Models



**Parameters**

$w_j$

$\mu_j$

$\Sigma_j$

$p(x)$

Dim 2

Dim 1

# Recognition Systems
# Gaussian Mixture Models

**Parameters**

**Model Components**

$$j = 1, ..., K$$

$p(x)$

Dim 2

Dim 1

Dr. Yanjun Qi / UVA CS

# Learning a Gaussian Mixture

- Probability Model

A Gaussian mixture model (GMM) represents as the weighted sum of multiple Gaussian distributions

$$p(\vec{x} = \vec{x}_i)$$

$$= \sum_j p(\vec{x} = \vec{x}_i, \vec{\mu} = \vec{\mu}_j)$$

Total low of probability

$$= \sum_j p(\vec{\mu} = \vec{\mu}_j) p(\vec{x} = \vec{x}_i \mid \vec{\mu} = \vec{\mu}_j)$$

Chain rule

$$= \sum_j p(\vec{\mu} = \vec{\mu}_j) \frac{1}{(2\pi)^{p/2} \left| \Sigma_j \right|^{1/2}} e^{-\frac{1}{2}\left(\vec{x}-\vec{\mu}_j\right)^T \Sigma_j^{-1}\left(\vec{x}-\vec{\mu}_j\right)}$$

# Learning a Gaussian Mixture
## (when assuming with known shared covariance)

$$p(\vec{x} = \vec{x}_i)$$

$$= \sum_{\mu_j} p(\vec{x} = \vec{x}_i, \vec{\mu} = \vec{\mu}_j)$$

$$= \sum_{j} p(\vec{\mu} = \vec{\mu}_j) p(\vec{x} = \vec{x}_i \mid \vec{\mu} = \vec{\mu}_j)$$

$$= \sum_{j} p(\vec{\mu} = \vec{\mu}_j) \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x}_i - \vec{\mu}_j)^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}_j)}$$

Assuming

# Max Log-likelihood of Observed Data Samples

☐ Log-likelihood of data   $log\, p(x_1, x_2, x_3, ..., x_n) =$

$$\log \prod_{i=1..n} \sum_{j=1..K} p(\vec{\mu} = \vec{\mu}_j) \frac{1}{\left(2\pi\right)^{p/2} \left|\Sigma_j\right|^{1/2}} e^{-\frac{1}{2}\left(\vec{x}_i - \vec{\mu}_j\right)^T \Sigma_j^{-1}\left(\vec{x}_i - \vec{\mu}_j\right)}$$

26

Apply MLE to find

$$\left\{ \{ p(\vec{\mu} = \mu_j) \}, j = 1...K \right\}$$

optimal Gaussian parameters   $\{ \vec{\mu}_j, \Sigma_j, j = 1...K \}$

# Expectation-Maximization for training GMM

- Start:
  - "Guess" the centroid and covariance for each of the K clusters
  - "Guess" the proportion of clusters, e.g., uniform prob 1/K

- Loop
  - For each point, revising its proportions belonging to each of the K clusters
  - For each cluster, revising both the mean (centroid position) and covariance (shape)

# Learning a Gaussian Mixture
## (when assuming with known shared covariance)

**E-Step**

[Bayes Rule]

assignment.  Soft

$$E[z_{ij}] = p(\vec{\mu} = \mu_j \mid x = x_i)$$

$$= \frac{p(x = x_i \mid \mu = \mu_j) p(\mu = \mu_j)}{\sum_{s=1}^{k} p(x = x_i \mid \mu = \mu_s) p(\mu = \mu_s)}$$

$$= \frac{\dfrac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x}_i - \vec{\mu}_j)^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}_j)} \; p(\mu = \mu_j)}{\sum_{s=1}^{k} \dfrac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x}_i - \vec{\mu}_s)^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}_s)} \; p(\mu = \mu_s)}$$

# Learning a Gaussian Mixture
## (when assuming with known shared covariance)

$$m_{ij} = \left\{ \begin{array}{l} 0 \\ 1 \end{array} \right.$$

$$E[z_{ij}] = p(\mu = \mu_j \mid x = x_i)$$

**E-Step**

Soft assignment
$p(\mu = \mu_j \mid x = x_i)$

How $x_i$ belongs
in proportion
to cluster $\{1, 2, \cdots, k\}$

VS. $m_{ij}$ Hard
assignment in
K-means

$$= \frac{p(x = x_i \mid \mu = \mu_j)\, p(\mu = \mu_j)}{\displaystyle\sum_{s=1}^{k} p(x = x_i \mid \mu = \mu_s)\, p(\mu = \mu_s)}$$

$$= \frac{\dfrac{1}{\left(2\pi\right)^{p/2}\left|\Sigma\right|^{1/2}} e^{-\frac{1}{2}\left(\vec{x}-\vec{\mu}_j\right)^T \Sigma^{-1}\left(\vec{x}-\vec{\mu}_j\right)}\, p(\mu = \mu_j)}{\displaystyle\sum_{s=1}^{k} \dfrac{1}{\left(2\pi\right)^{p/2}\left|\Sigma\right|^{1/2}} e^{-\frac{1}{2}\left(\vec{x}-\vec{\mu}_s\right)^T \Sigma^{-1}\left(\vec{x}-\vec{\mu}_s\right)}\, p(\mu = \mu_s)}$$

# Learning a Gaussian Mixture

**when assuming with known shared covariance**

**M-Step**

$$\mu_j^{(t+1)} \leftarrow \frac{1}{\sum_{i=1}^{n} E[z_{ij}]} \sum_{i=1}^{n} E[z_{ij}]^{(t)} x_i$$

$$p(\mu = \mu_j)^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^{n} E[z_{ij}]^{(t)}$$

Covariance: $\Sigma_j$ (j: 1 to K) can also be derived in the M-step under a full setting

# Learning a Gaussian Mixture

when assuming with known shared covariance

$$k \text{ mean} \Rightarrow \text{ centroid} = \frac{1}{N_j} \sum_{i=1}^{} m_{ij} x_i$$

**M-Step**

$$\mu_j^{(t+1)} \leftarrow \frac{1}{\sum_{i=1}^{n} E[z_{ij}]} \sum_{i=1}^{n} E[z_{ij}]^{(t)} x_i$$

$$\underbrace{E[z_{ij}]} \rightarrow [0, 1]$$

$$\sum_{j=1}^{k} E[z_{ij}] = 1$$

$$p(\mu = \mu_j)^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^{n} E[z_{ij}]^{(t)}$$

Covariance: $\Sigma_j$ (j: 1 to K) will also be derived in the M-step under a full setting

# M-step for Estimating unknown Covariance Matrix (more general, details in EM-Extra lecture)

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^{n} E[z_{ij}]^{(t)}(x_i - \mu_j^{(t+1)})(x_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^{n} E[z_{ij}]^{(t)}}$$

for small Trainset too many parameters to estimate

$j = 1, \ldots, K$

$\Sigma_j \Rightarrow O(P^2/2)$

$\Sigma_j \Rightarrow O(KP^2/2)$

$M_j \Rightarrow O(KP + K)$    $p(u=u_j)$

$E[z_{ij}] \; O(Kn)$

# Recap: Expectation-Maximization for training GMM

- Start:
  - "Guess" the centroid and covariance for each of the K clusters
  - "Guess" the proportion of clusters, e.g., uniform prob 1/K

- Loop
  - For each point, revising its proportions belonging to each of the K clusters
  - For each cluster, revising both the mean (centroid position) and covariance (shape)

4/26/18

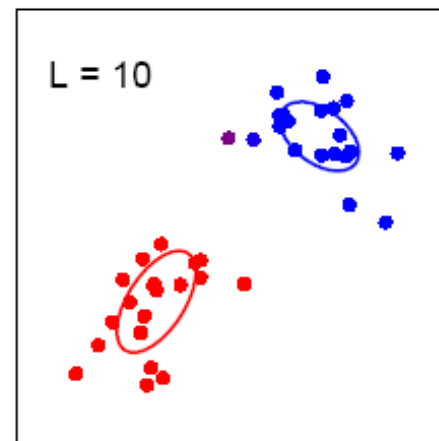each cluster, revising both the mean (centroid position) and covariance (shape)
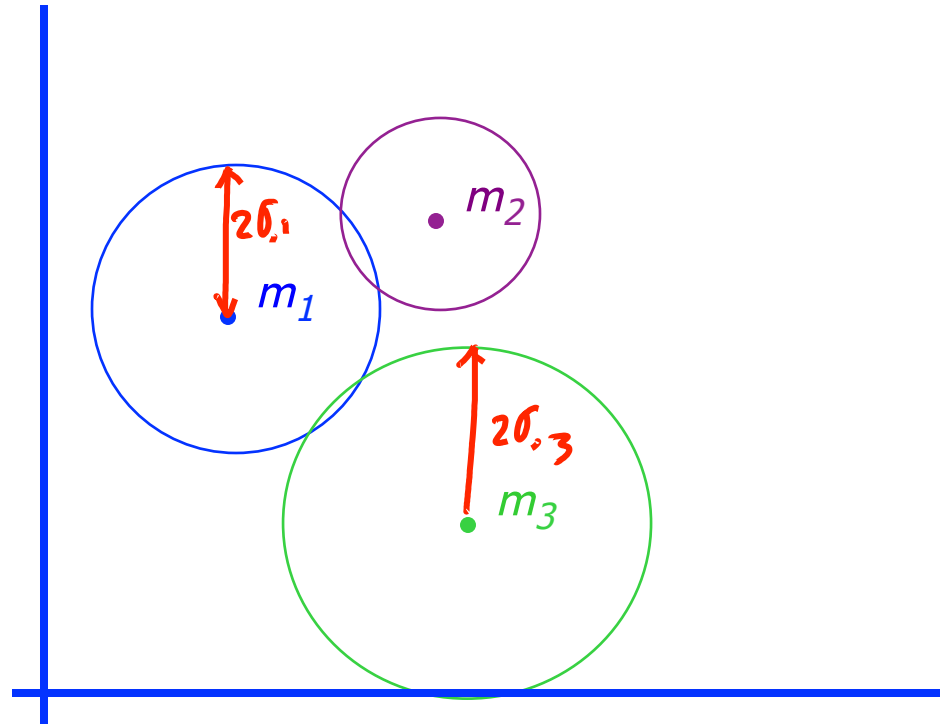
# The Simplest GMM assumption

- Each component generates data from a Gaussian with

  - mean $\mu_i$

  - Shared covariance matrix $\sigma^2 \boldsymbol{I}$



$$\Sigma_j = \Sigma = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$
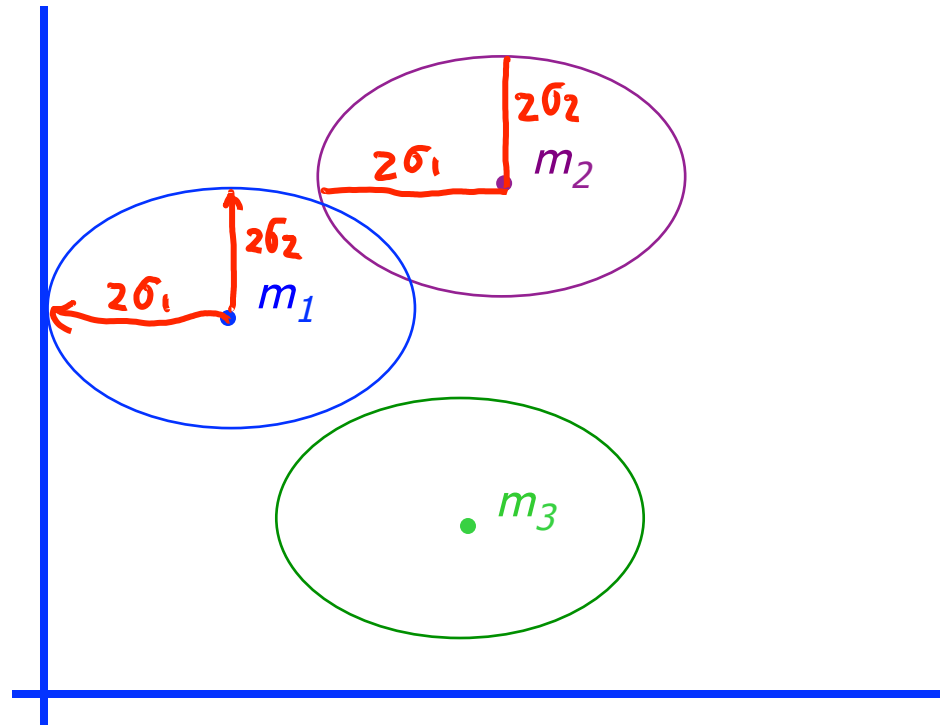
# A Simple GMM assumption

- Each component generates data from a Gaussian with

  - mean $\mu_i$

  - Cluster-specific covariance matrix as $\sigma_j^2 \boldsymbol{I}$



$$\sum_j = \sigma_j^2 I = \begin{bmatrix} \delta_j^2 & 0 \\ 0 & \delta_j^2 \end{bmatrix}$$
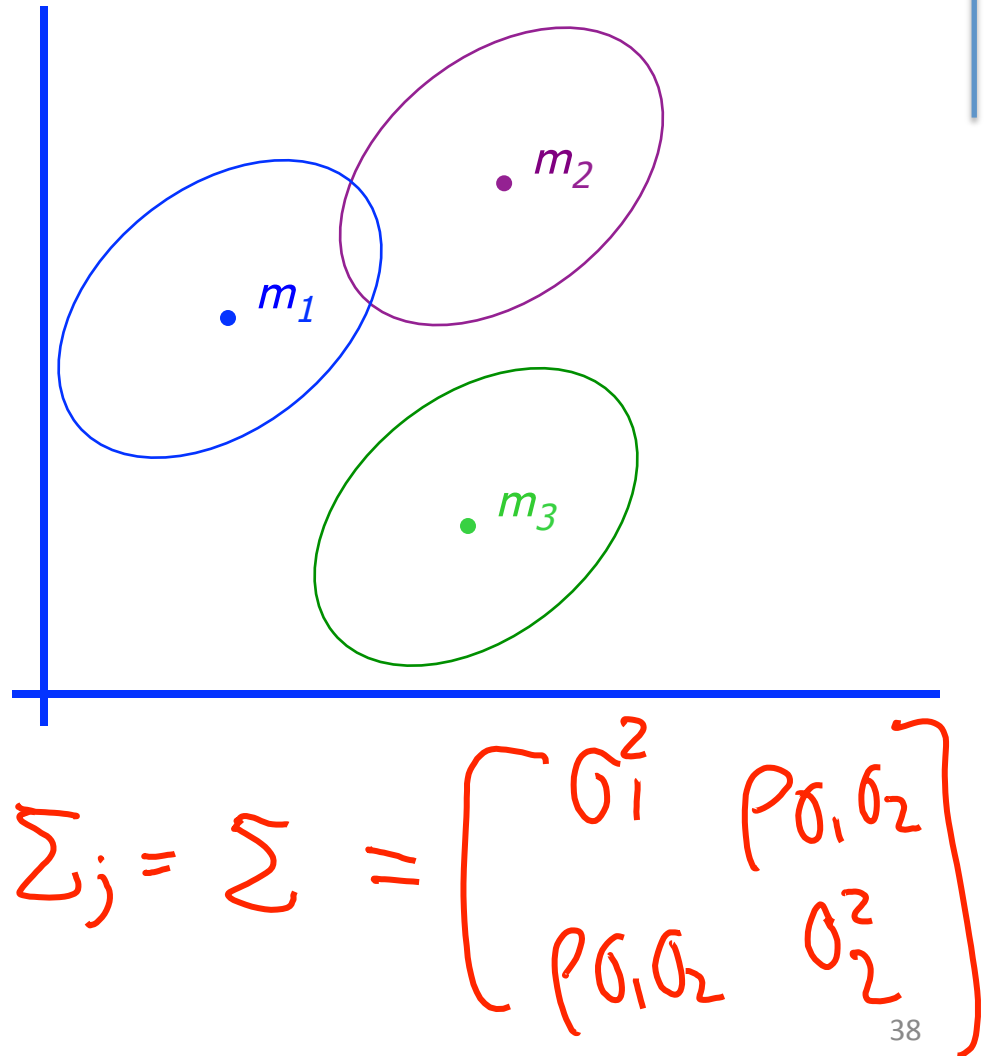
# Another Simple GMM assumption

- Each component generates data from a Gaussian with

  - mean $\mu_i$

  - Shared covariance matrix as diagonal matrix



$$\Sigma_j = \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$
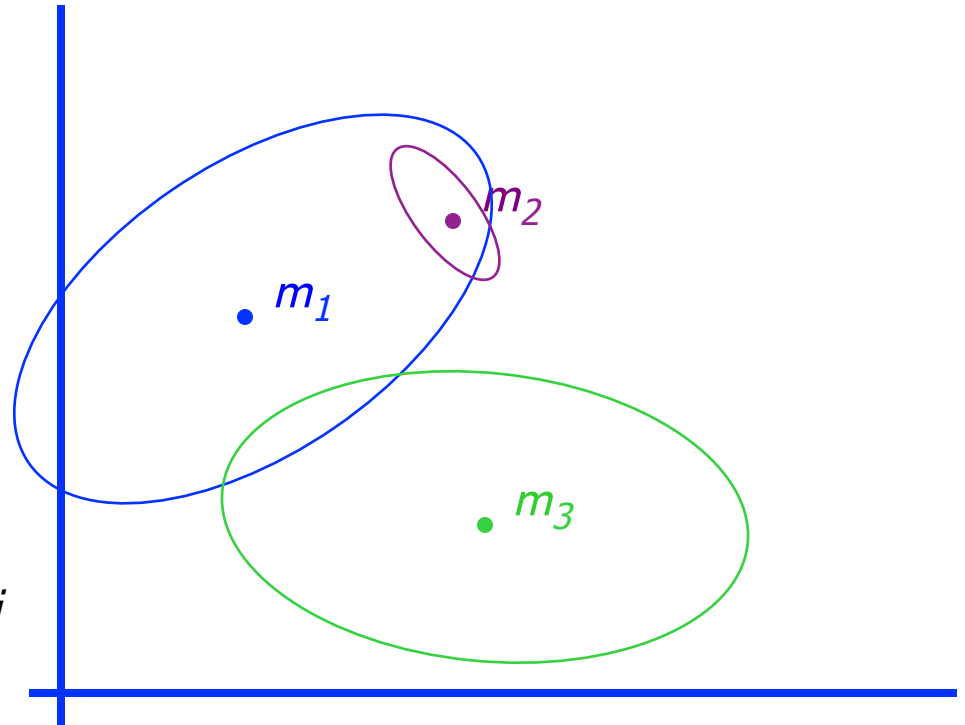
# A bit More General GMM assumption

- Each component generates data from a Gaussian with

  - mean $\mu_i$

  - Shared covariance matrix as full matrix

$$\Sigma_j = \Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$m_1$

$m_2$

$m_3$

# The General GMM assumption

- Each component generates data from a Gaussian with

  - mean $\mu_i$

  - covariance matrix $\Sigma_i$



$$\Sigma_j = \begin{bmatrix} \sigma_{1j} & Cov_j(x_1, x_2) \\ Cov_j(x_1, x_2) & \sigma_{2j} \end{bmatrix}$$

# Partitional : Gaussian Mixture Model

- 1. Review of Gaussian Distribution
- 2. GMM for clustering : basic algorithm
- 3. GMM connecting to K-means
- 4. GMM examples
- 5. Problems of GMM and K-means

# Recap: K-means iterative learning

$$\underset{\{\vec{C}_j, m_{i,j}\}}{\arg\min} \sum_{j=1}^{K} \sum_{i=1}^{n} m_{i,j} \left( \vec{x}_i - \vec{C}_j \right)^2$$

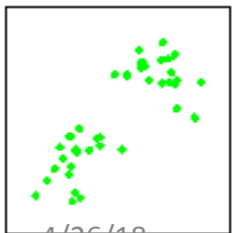Memberships $\{m_{i,j}\}$ and centers $\{C_j\}$ are correlated.

**E-Step**  Given centers $\{\vec{C}_j\}$, $m_{i,j} = \begin{cases} 1 & j = \underset{k}{\arg\min}(\vec{x}_i - \vec{C}_j)^2 \\ 0 & \text{otherwise} \end{cases}$

**M-Step**  Given memberships $\{m_{i,j}\}$, $\vec{C}_j = \dfrac{\sum_{i=1}^{n} m_{i,j} \vec{x}_i}{\sum_{i=1}^{n} m_{i,j}}$
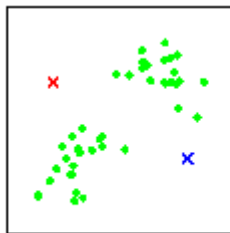
# Compare: K-means

- The EM algorithm for mixtures of Gaussians is like a "soft version" of the K-means algorithm.

- In the K-means "E-step" we do hard assignment:

- In the K-means "M-step" we update the means as the weighted sum of the data, but now the weights are 0 or 1:
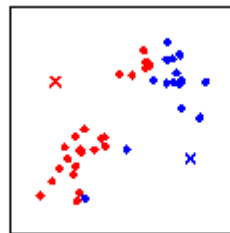


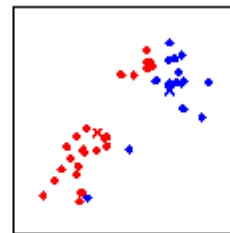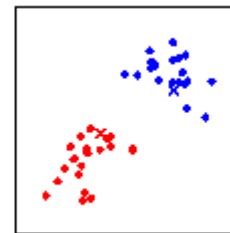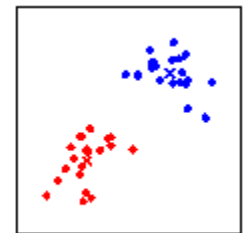(a)    (b)    (c)    (d)    (e)    (f)

K-means: $\arg\min\limits_{\{\vec{C}_j, m_{i,j}\}} \sum\limits_{j=1}^{K} \sum\limits_{i=1}^{n} m_{i,j} \left( \vec{x}_i - \vec{C}_j \right)^2$

$m_{ij} = \begin{cases} 0 \\ 1 \end{cases}$

GMM : $\sum\limits_{i} \log \prod\limits_{i=1}^{n} p(x = x_i) = \sum\limits_{i} \log \left[ \sum\limits_{\mu_j}^{j=1,..k} p(\mu = \mu_j) \frac{1}{(2\pi)|\Sigma|^{1/2}} e^{-\frac{1}{2}\left( \vec{x} - \vec{\mu}_j \right)^T \Sigma^{-1} \left( \vec{x} - \vec{\mu}_j \right)} \right]$

$i=1..n$

- K-Mean only detect spherical clusters.
- GMM can adjust its self to elliptic shape clusters.

# Partitional : Gaussian Mixture Model

- 1. Review of Gaussian Distribution
- 2. GMM for clustering : basic algorithm
- 3. GMM connecting to K-means
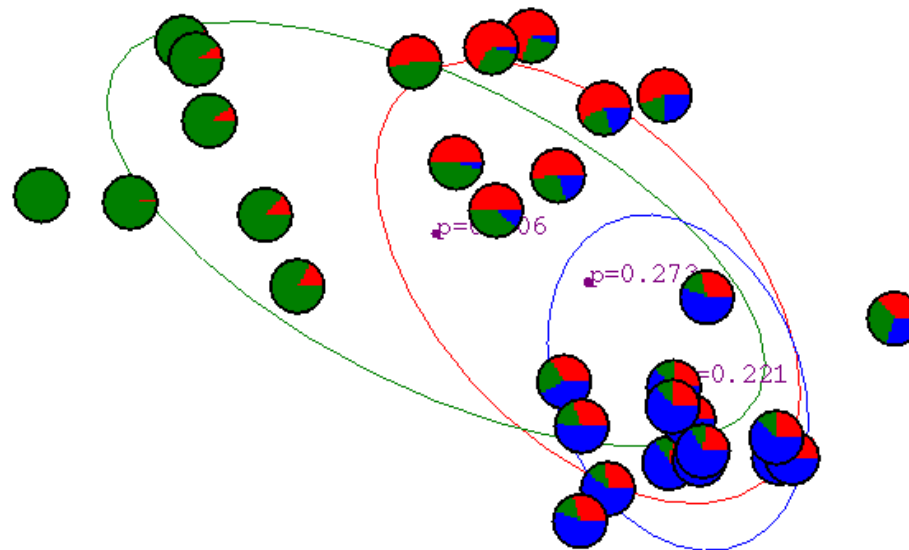- 4. GMM examples
- 5. Problems of GMM and K-means

# Gaussian Mixture Example: Start



$p(\mu_j | x_i)$

p=0.333

p=0.333    p=0.333
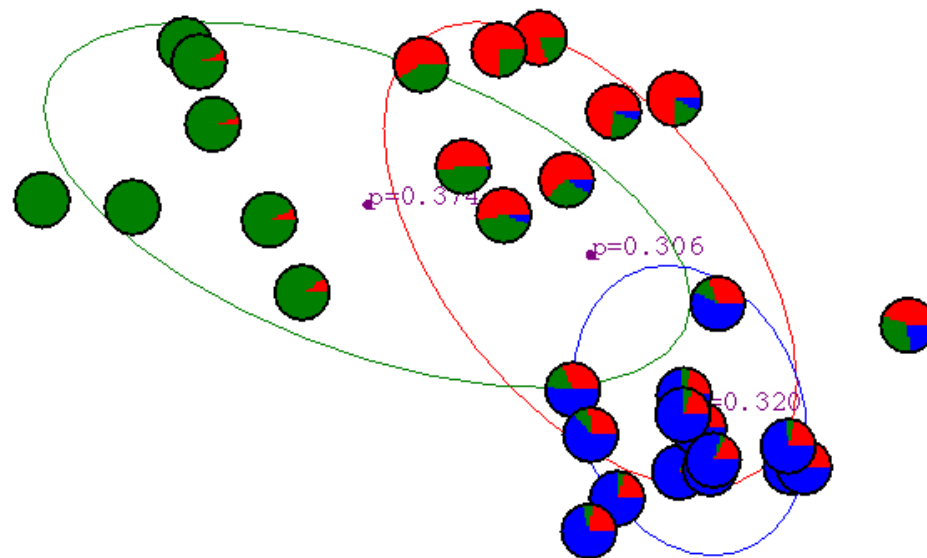
# After First Iteration

For each point, revising its proportions belonging to each of the K clusters



For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture

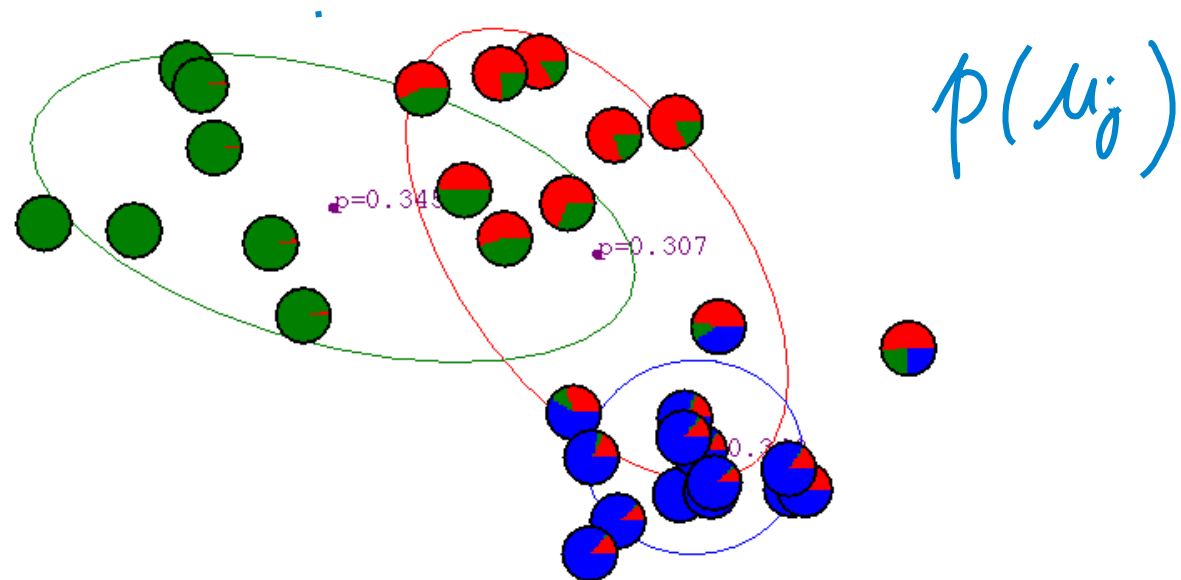4/26/18

# After 2nd Iteration

For each point, revising its proportions belonging to each of the K clusters



For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture

# After 3rd Iteration

For each point, revising its proportions belonging to each of the K clusters
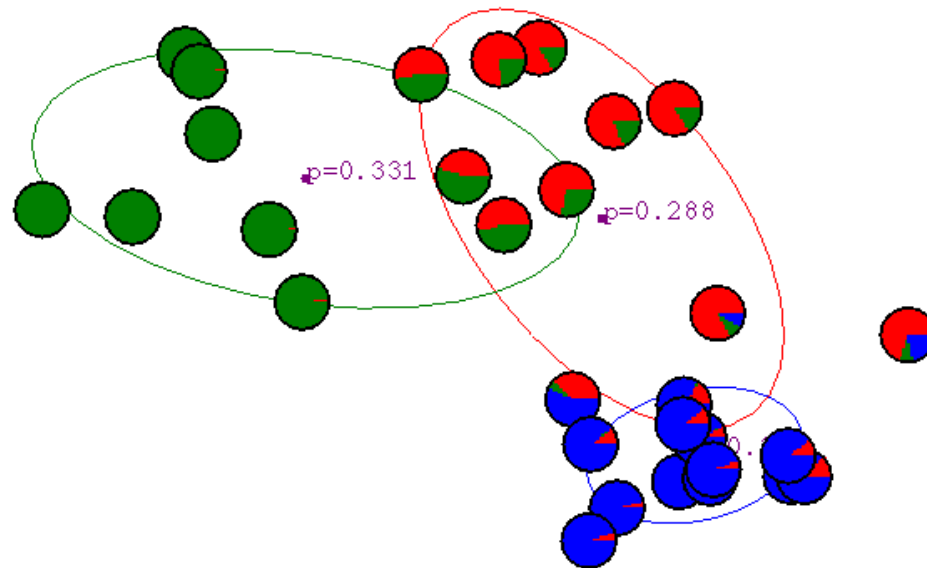


$p(\mu_j)$

p=0.345
p=0.307

For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture

# After 4th Iteration

For each point, revising its proportions belonging to each of the K clusters



p=0.331

p=0.288

For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture

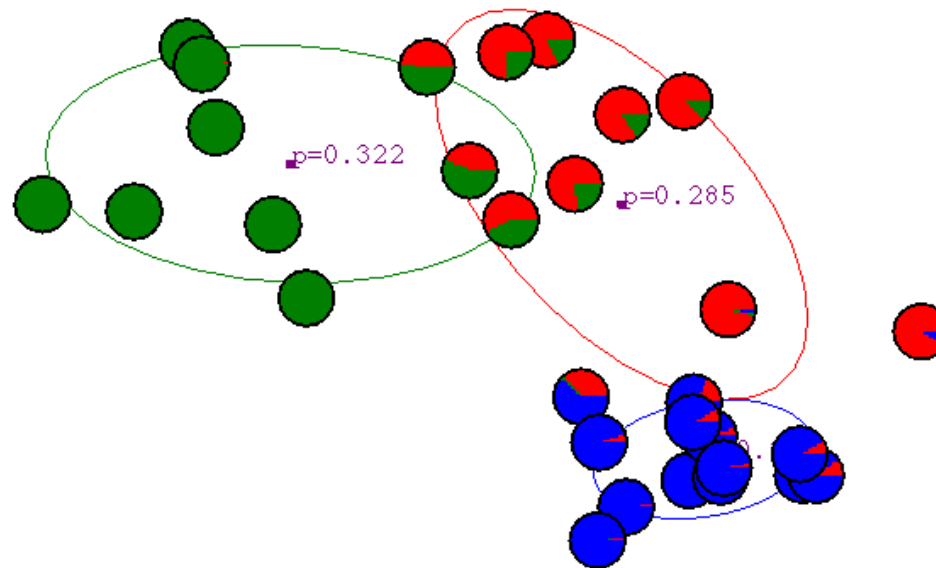4/26/18

# After 5th Iteration

For each point, revising its proportions belonging to each of the K clusters



For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture
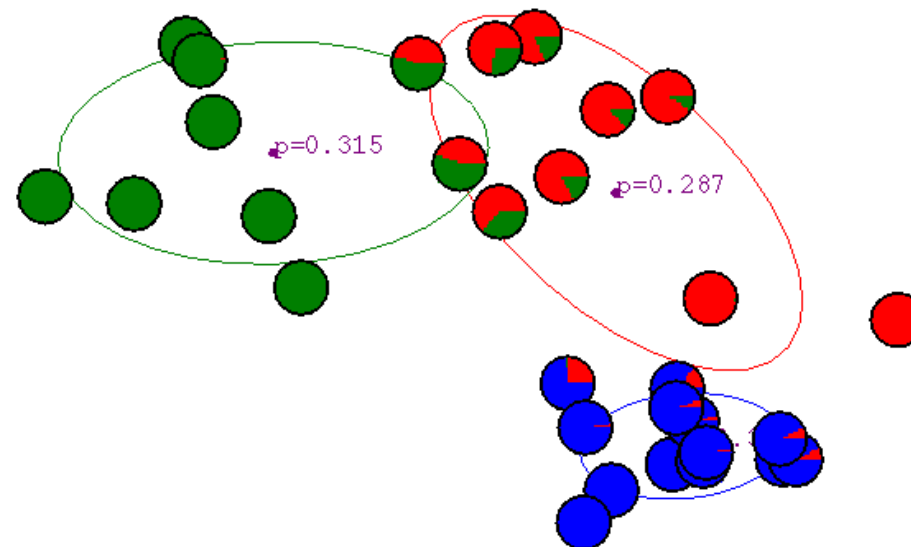
# After 6th Iteration

For each point, revising its proportions belonging to each of the K clusters



For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture

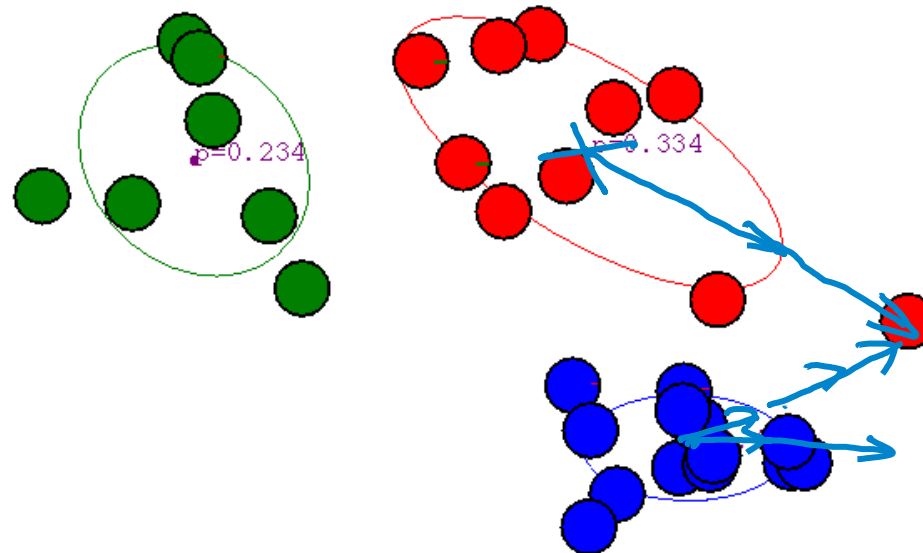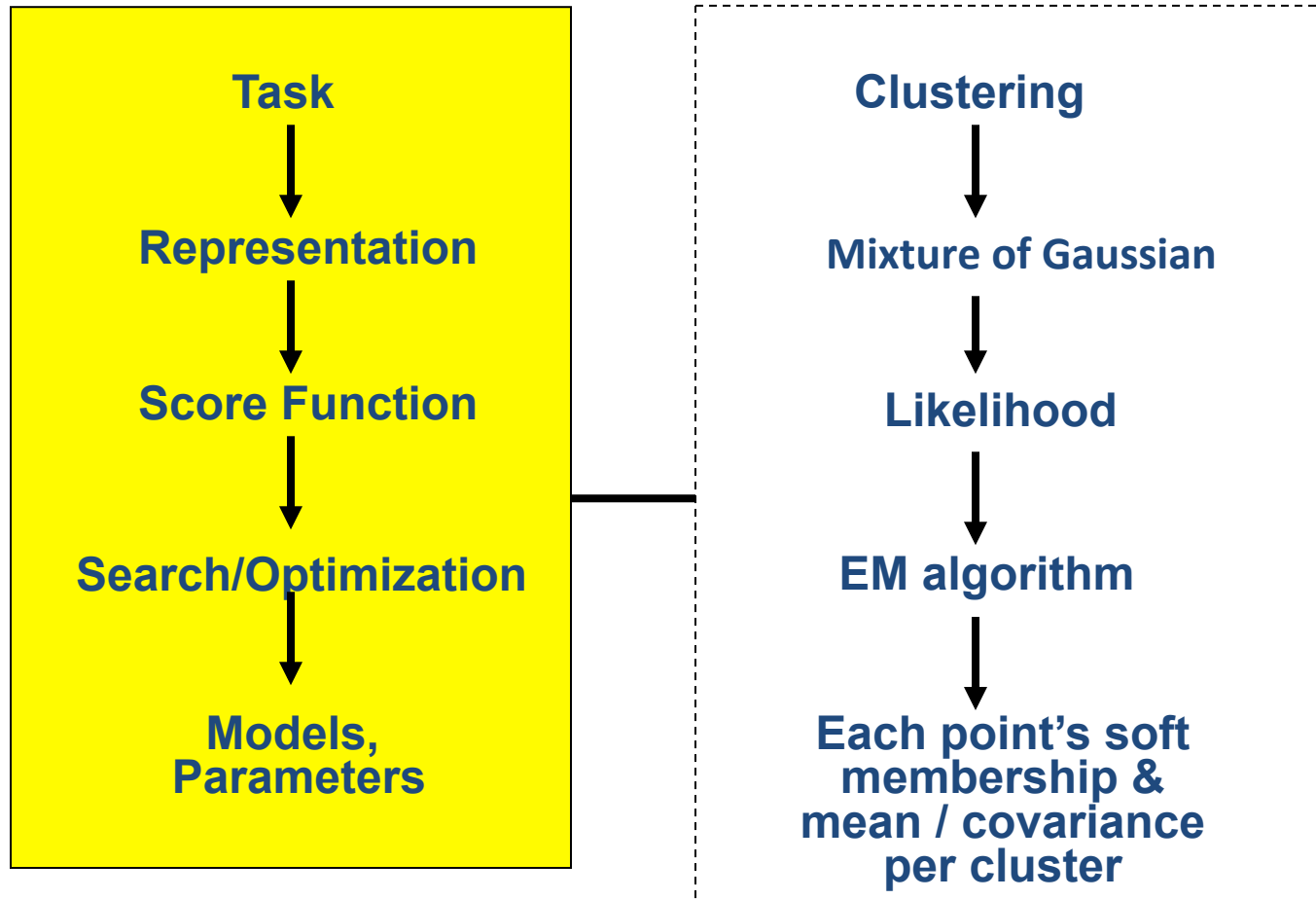4/26/18

# After 20th Iteration

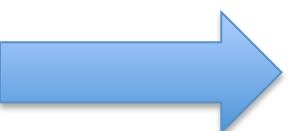For each point, revising its proportions belonging to each of the K clusters



For each cluster, revising its mean (centroid position), covariance (shape) and proportion in the mixture
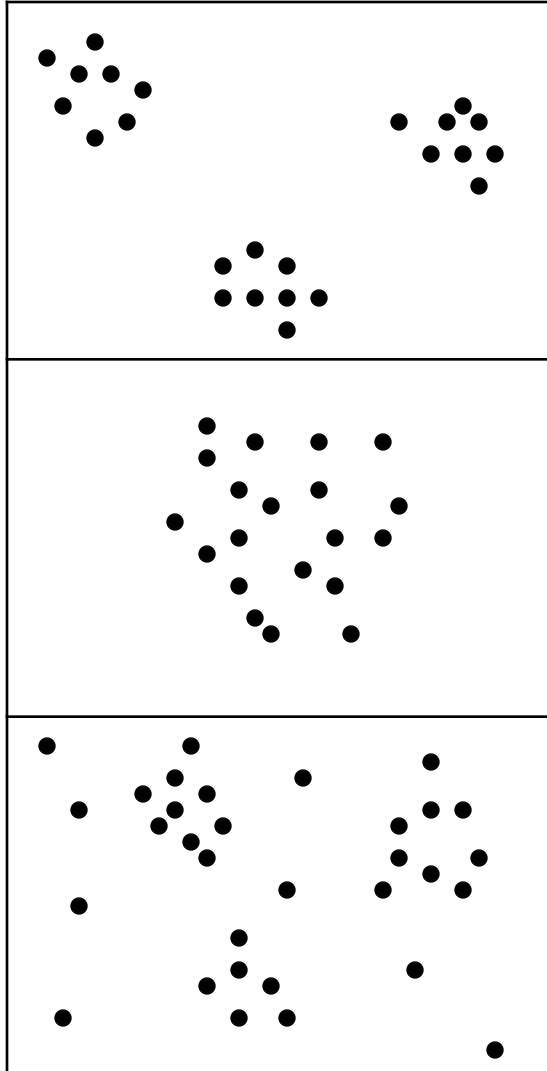
4/26/18

# (3) GMM Clustering

**Task**

↓

**Representation**

↓

**Score Function**

↓

**Search/Optimization**

↓

**Models, Parameters**

**Clustering**

↓

**Mixture of Gaussian**

↓

**Likelihood**

↓

**EM algorithm**

↓

**Each point's soft membership & mean / covariance per cluster**

$$\sum_i \log \prod_{i=1}^{n} p(x = x_i) = \sum_i \log \left[ \sum_{\mu_j} p(\mu = \mu_j) \frac{1}{(2\pi)\left|\Sigma_j\right|^{1/2}} e^{-\frac{1}{2}\left(\vec{x}-\vec{\mu}_j\right)^T \Sigma_j^{-1}\left(\vec{x}-\vec{\mu}_j\right)} \right]$$

4/26/18

53

# Partitional : Gaussian Mixture Model

- 1. Review of Gaussian Distribution
- 2. GMM for clustering : basic algorithm
- 3. GMM connecting to K-means
- 4. GMM examples
- 5. Problems of GMM and K-means

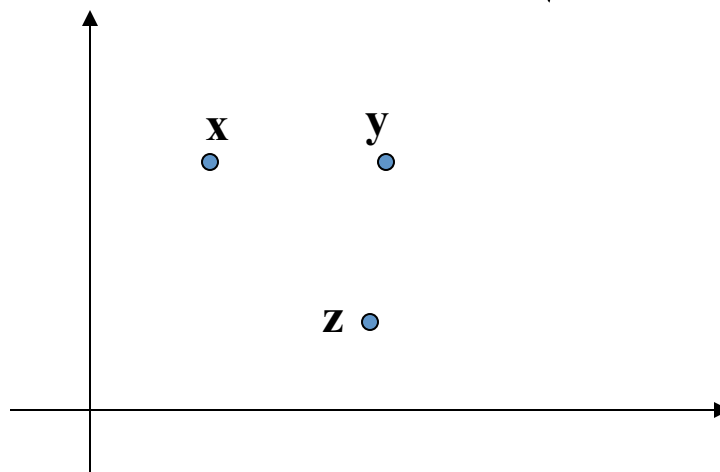# Unsupervised Learning: not as hard as it looks

Sometimes easy

Sometimes impossible

and sometimes
in between

# Problems (I)

- Both k-means and mixture models need to compute centers of clusters and explicit distance measurement
  - Given strange distance measurement, the center of clusters can be hard to compute

E.g.,

$$\left\| \vec{x} - \vec{x}' \right\|_\infty = \max\left( \left| x_1 - x_1' \right|, \left| x_2 - x_2' \right|, ..., \left| x_p - x_p' \right| \right)$$

**x**     **y**

$$\left\| \mathbf{x} - \mathbf{y} \right\|_\infty = \left\| \mathbf{x} - \mathbf{z} \right\|_\infty$$
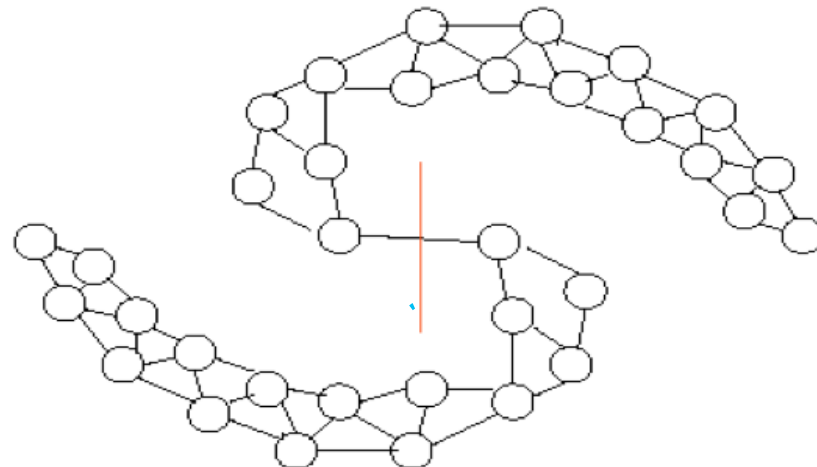
**z**

# Problem (II)

*tight*

- Both k-means and mixture models look for compact clustering structures

  – In some cases, connected clustering structures are more desirable
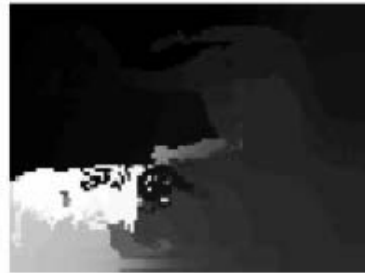
**Graph based clustering**

**e.g. MinCut, Spectral clustering**

# e.g. Image Segmentation through minCut



(a)    (b)    (c)

(d)    (e)    (f)

# **References**

❑ Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.

❑ Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides

❑ Big thanks to Prof. Ziv Bar-Joseph @ CMU for allowing me to reuse some of his slides

❑ clustering slides from Prof. Rong Jin @ MSU