



TEXAS A&M UNIVERSITY
Engineering

ECEN 758 Data Mining and Analysis: Lecture 8, Gaussian Mixture Models

Joshua Peeples, Ph.D.

Assistant Professor

Department of Electrical and Computer Engineering

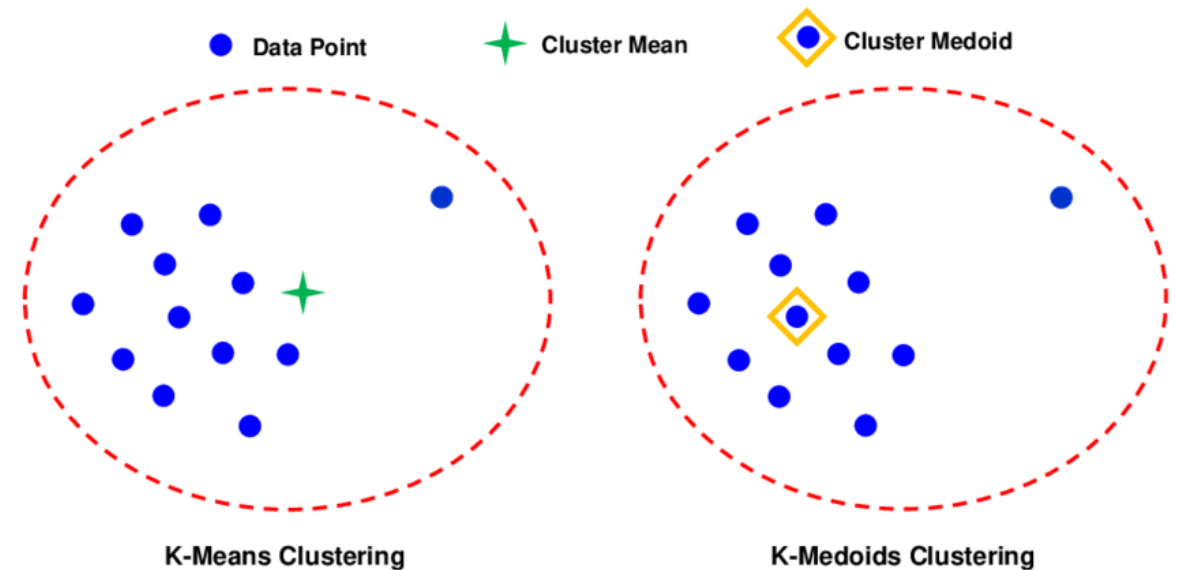
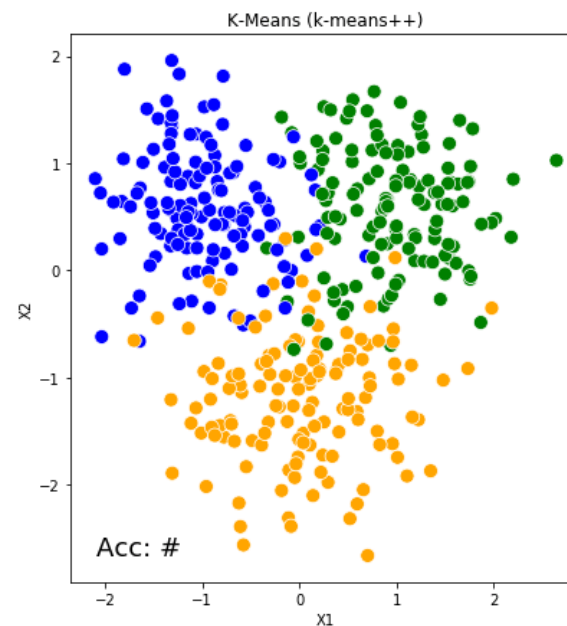
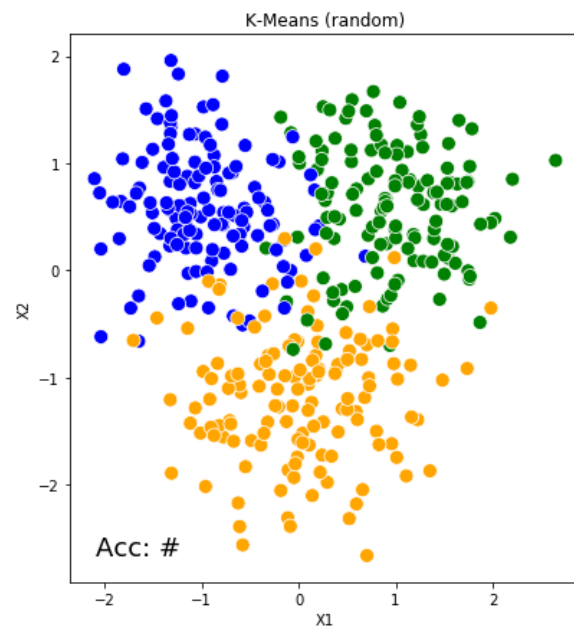
- Assignment #1 grades available
 - Please revise any grade discrepancies within a week (COB, 09/23)
 - Email Dr. Peeples (do not contact Grader) and/or stop by office hours
- Assignment #2 will be released this Wednesday (09/18)
 - Please upload submission as single PDF
 - Please share Python code (e.g., Jupyter Notebooks, Google Colab)

Assignment 1 Observations



- Make sure to clearly label your figures and tables
 - Communication is important
- Please have a clear discussion
 - Formal writing (i.e., no contractions)
- Show your work
 - State equations and show your steps
- Disclose if you use AI (e.g., ChatGPT, Copilot) for code development
 - Do not use for your discussions
- Ask questions if clarification is needed

- Representative Clustering II



- Gaussian Mixture Models
- Reading: ZM Chapter 13

- We will discuss several variants of clustering
 - **Representative-based Clustering**
 - Hierarchical Clustering
 - Density-Based Clustering



What disadvantages of k-means?

k-Means Disadvantages



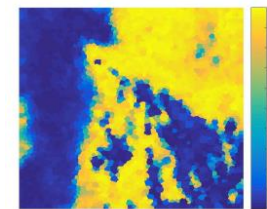
- Linear boundaries between clusters
- Only uses Euclidean distance
 - Assumes spherical clusters
 - Sensitive to outliers
- Non-symmetrical clusters
- Initialization
- Batch processing
- Selecting number of clusters (k)
- **“Crisp”/Hard clustering**

k-Means Disadvantage: “Crisp”/Hard Clustering

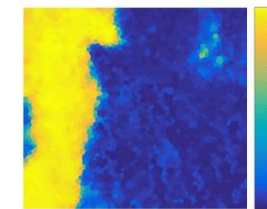


- Points can only “belong” to one cluster
- Different applications may require “soft” clustering
 - Points may belong to more than one group

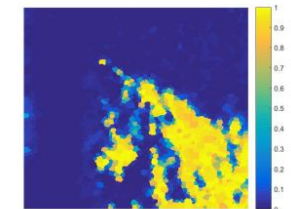
Input Image



(h) FLICM Cluster 1



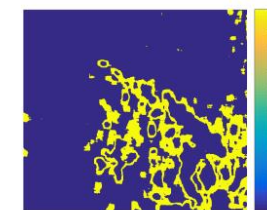
(i) FLICM Cluster 2



(j) FLICM Cluster 3



(k) K-Means Cluster 1



(l) K-Means Cluster 2



(m) K-Means Cluster 3



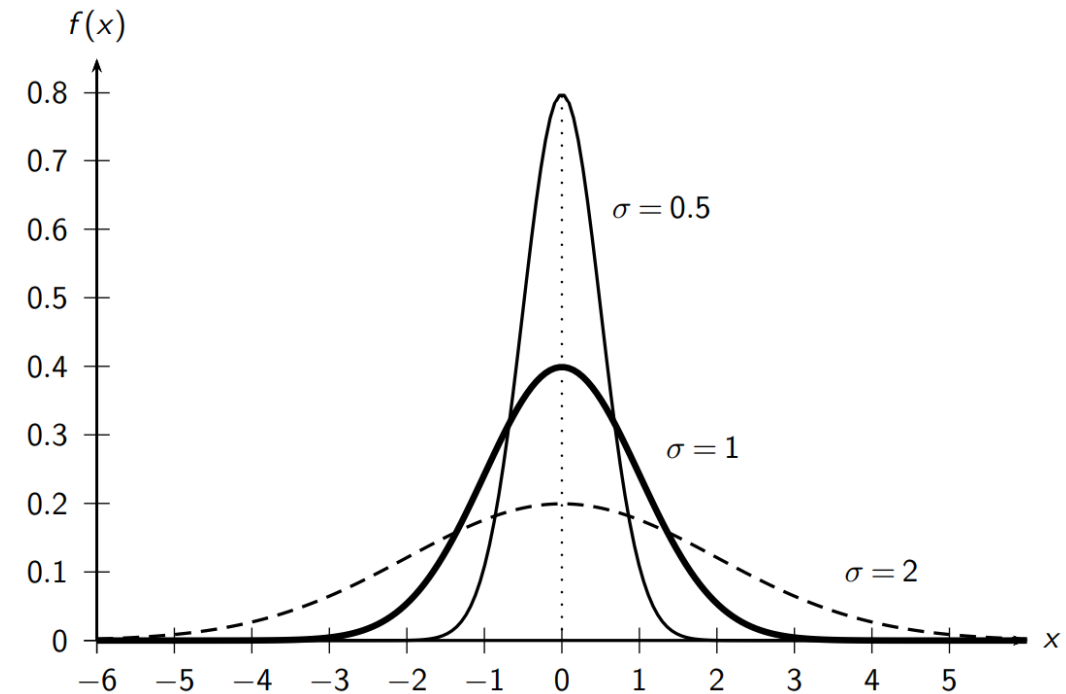
Gaussian Mixture Models

Gaussian/Normal Distribution (1D)



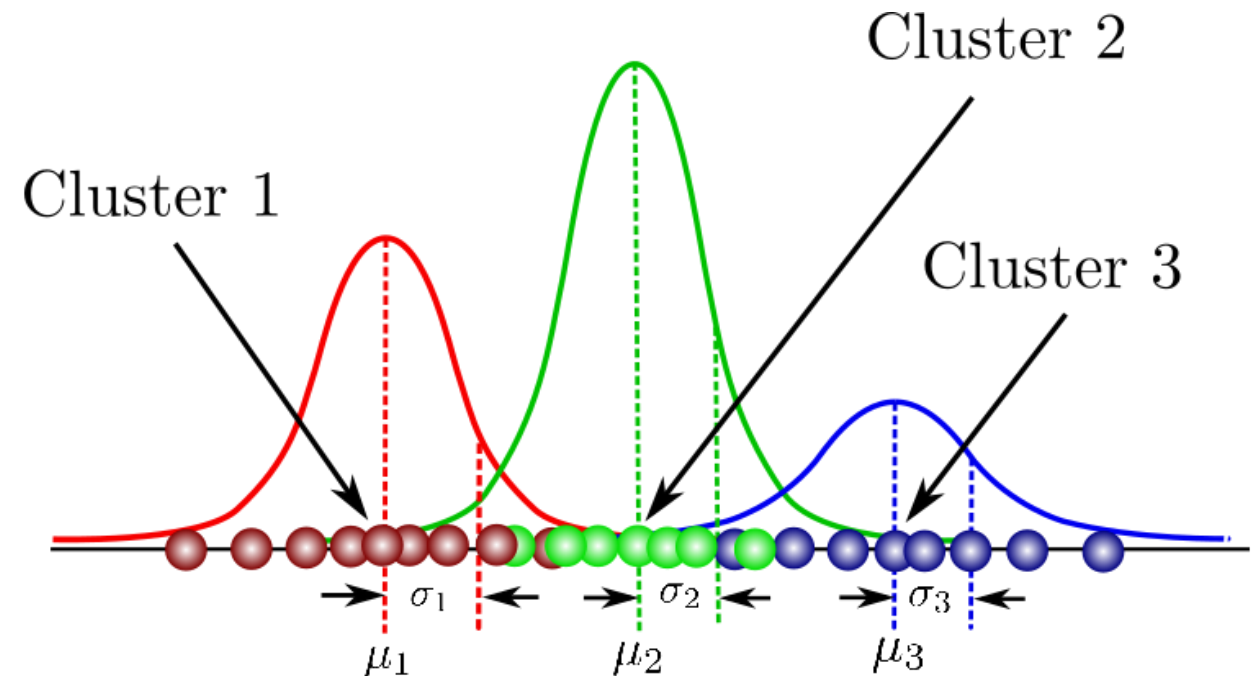
TEXAS A&M UNIVERSITY
Engineering

- Two parameters, mean (μ) and variance (σ^2)
- Probability density decreases exponentially as a function of the distance from mean
- Maximum value when $x = \mu$



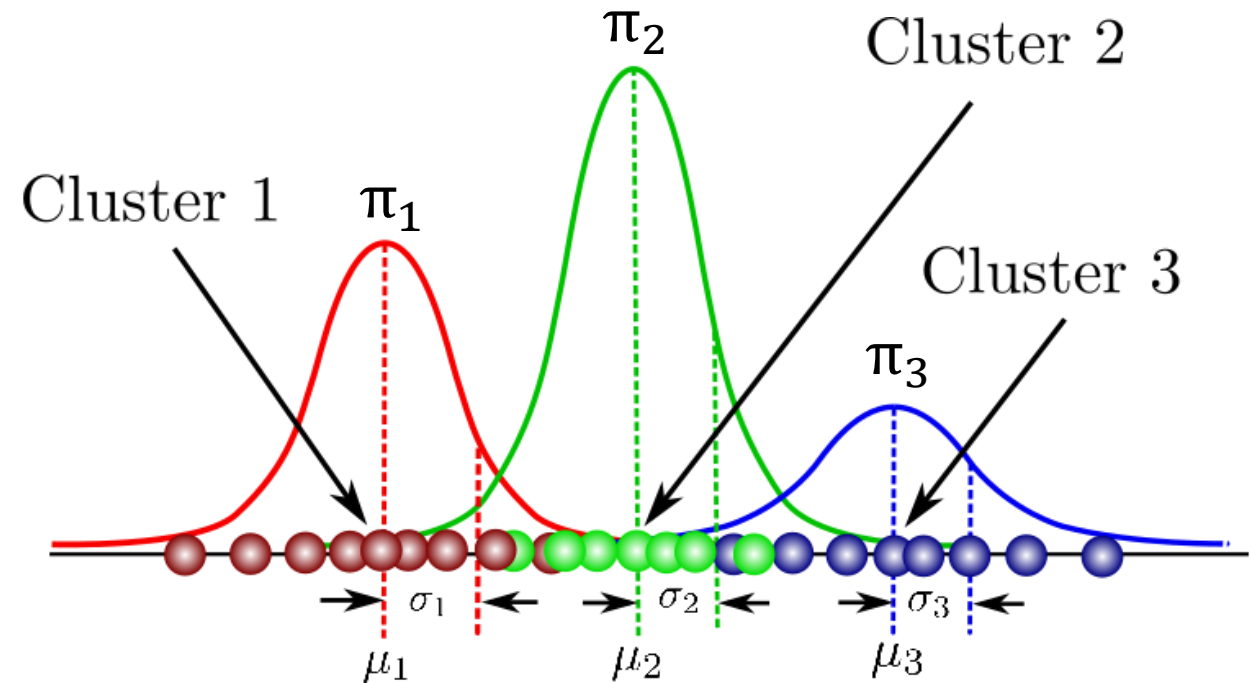
$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}$$

- Model clusters as Gaussians
- “Soft” clustering approach
 - Assign probability of belonging to clustering
- Generative model



Mixtures of Gaussians (1D)

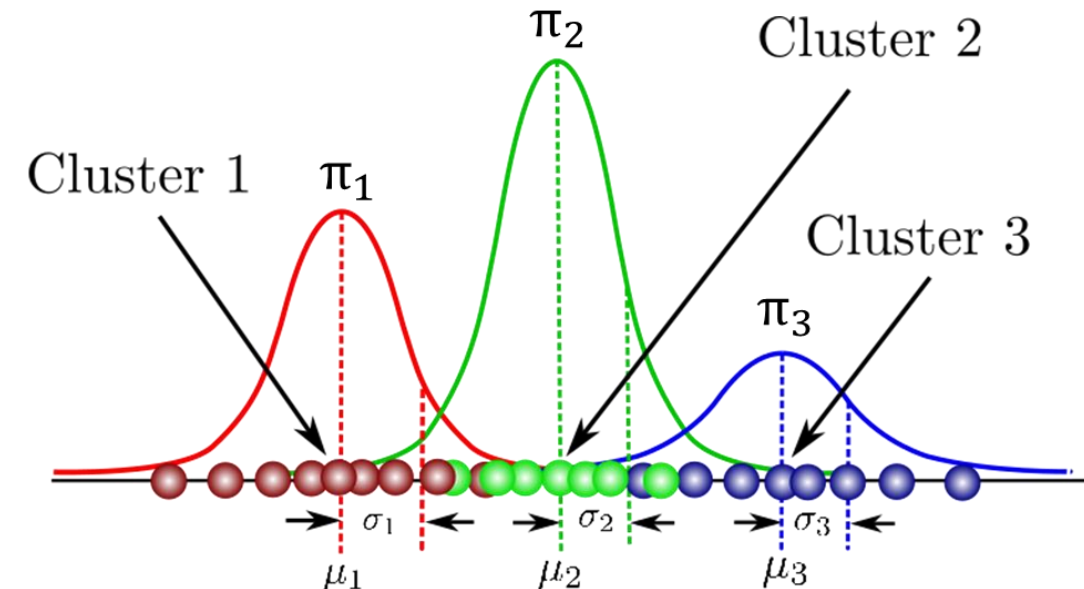
- Three parameters to describe clusters:
 - Mean (μ_k)
 - Variance (σ_k^2)
 - Mixture parameters (π_k)
 - Weights, “size”, prior probability
 - Sum to one constraint



Mixtures of Gaussians (1D)

- Three parameters to describe clusters:
 - Mean (μ_k)
 - Variance (σ_k^2)
 - Mixture parameters (π_k)
 - Weights, “size”, prior probability
- Probability distribution:

$$p(x) = \sum_{i=1}^k \pi_i \mathcal{N}(x | \mu_i, \sigma_i)$$



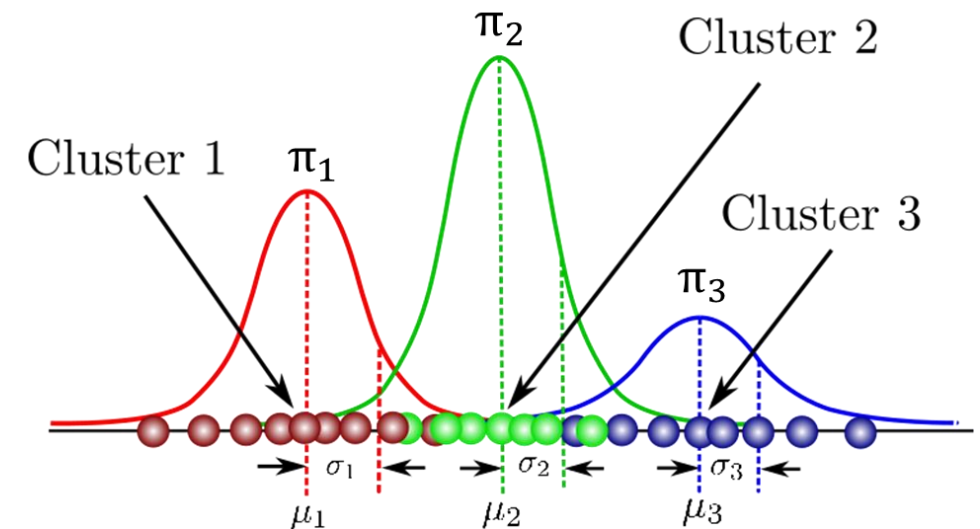
Mixtures of Gaussians (1D)

- Probability distribution:

$$p(x) = \sum_{i=1}^k \pi_i \mathcal{N}(x | \mu_i, \sigma_i)$$

- Select mixture component with probability π_k

$$p(z = k) = \pi_k$$



Mixtures of Gaussians (1D)

- Probability distribution:

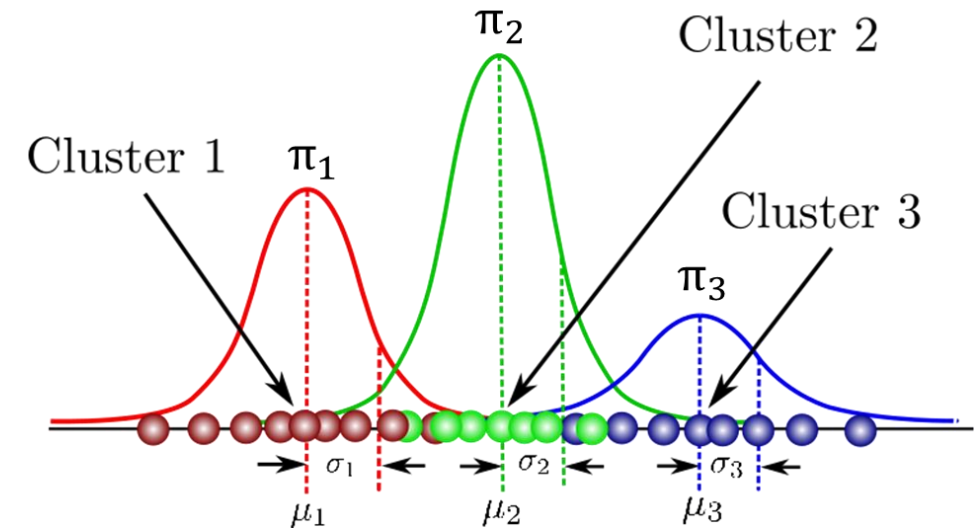
$$p(x) = \sum_{i=1}^k \pi_i \mathcal{N}(x|\mu_i, \sigma_i)$$

- Select mixture component with probability π_k

$$p(z = k) = \pi_k$$

- Sample from that component's Gaussian

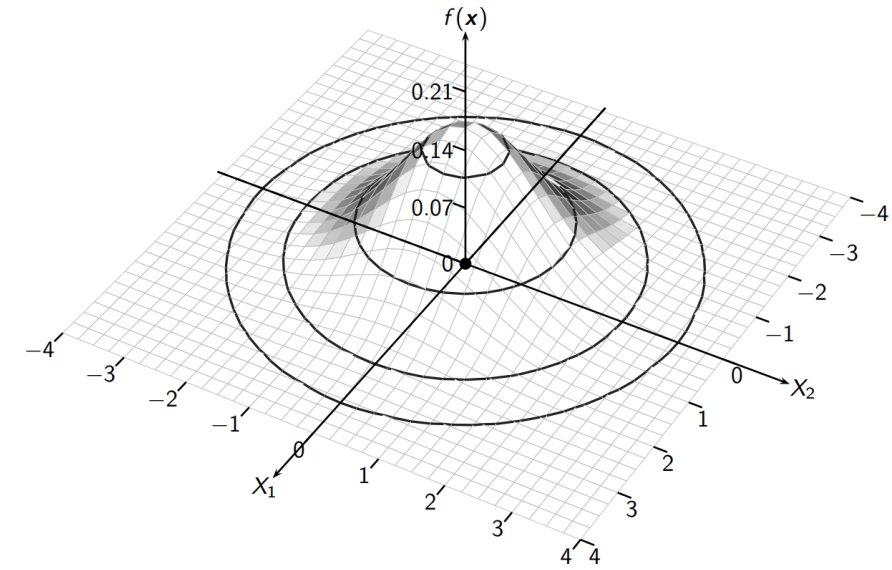
$$p(x|z = k) = \mathcal{N}(x|\mu_k, \sigma_k)$$





Gaussian Mixture Models: Multivariate

- Parameters: mean vector (μ) and covariance matrix (Σ)
- $|\Sigma|$ determinant of covariance matrix
- Numerator in exponential referred to as **Mahalanobis distance**



$$f(x|\mu, \Sigma) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|\Sigma|}} \exp \left\{ -\frac{(x - \mu)^T \Sigma^{-1} (x - \mu)}{2} \right\}$$

- Three parameters to describe clusters:
 - Mean vector (μ_i)
 - Covariance matrix (Σ_i^2)
 - Mixture parameters (π_i or $P(C_i)$)
 - Weights, “size”, prior probability
 - Sum to one constraint

$$\sum_{i=1}^k P(C_i) = 1.$$

i^{th} Cluster:

$$f_i(\mathbf{x}) = f(\mathbf{x}|\mu_i, \Sigma_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp \left\{ -\frac{(\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i)}{2} \right\}$$

Probability Density function of \mathbf{x} as GMM:

$$f(\mathbf{x}) = \sum_{i=1}^k f_i(\mathbf{x}) P(C_i) = \sum_{i=1}^k f(\mathbf{x}|\mu_i, \Sigma_i) P(C_i)$$



Gaussian Mixture Models Algorithm

GMM Algorithm: Objective



- Parameters of model represented as Θ

$$\theta = \{\mu_1, \Sigma_1, P(C_1), \dots, \mu_k, \Sigma_k, P(C_k)\}$$

- Maximum likelihood estimation (MLE)
- Usually maximize log-likelihood function

Likelihood:

$$P(\mathbf{D}|\theta) = \prod_{j=1}^n f(\mathbf{x}_j)$$

MLE:

$$\theta^* = \arg \max_{\theta} \{\ln P(\mathbf{D}|\theta)\}$$

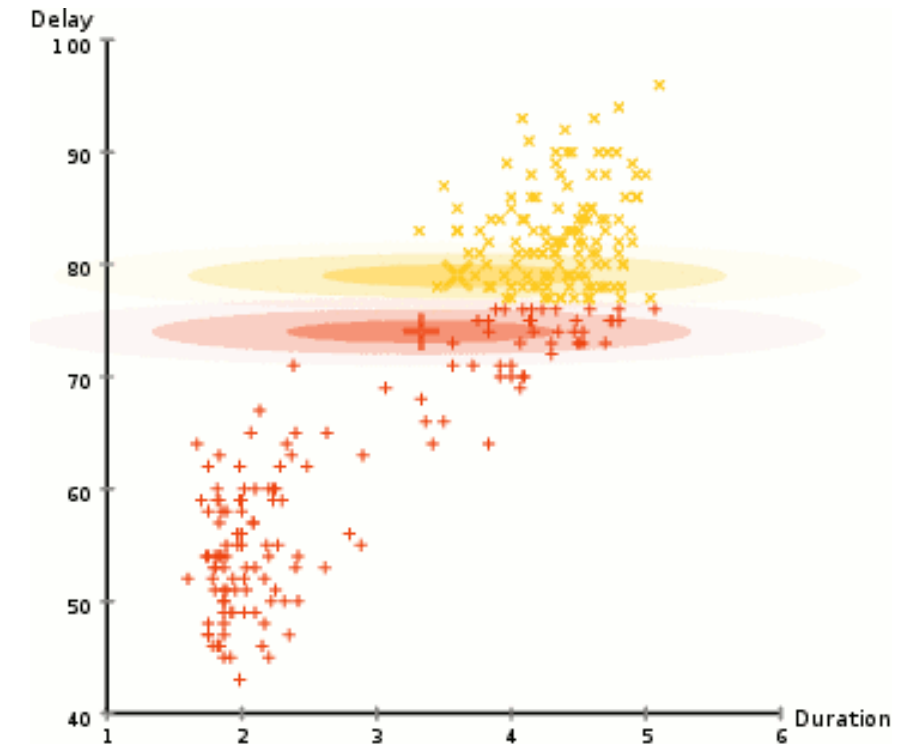
Log-likelihood:

$$\ln P(\mathbf{D}|\theta) = \sum_{j=1}^n \ln f(\mathbf{x}_j) = \sum_{j=1}^n \ln \left(\sum_{i=1}^k f(\mathbf{x}_j|\mu_i, \Sigma_i) P(C_i) \right)$$

GMM Algorithm: Objective



- Directly maximizing log-likelihood over Θ is hard
- Alternative approach: Expectation-Maximization (EM)
- Two steps:
 - Expectation: Assignment of points
 - Maximization: Estimation of parameters
- We will do a deep dive into EM next lecture!



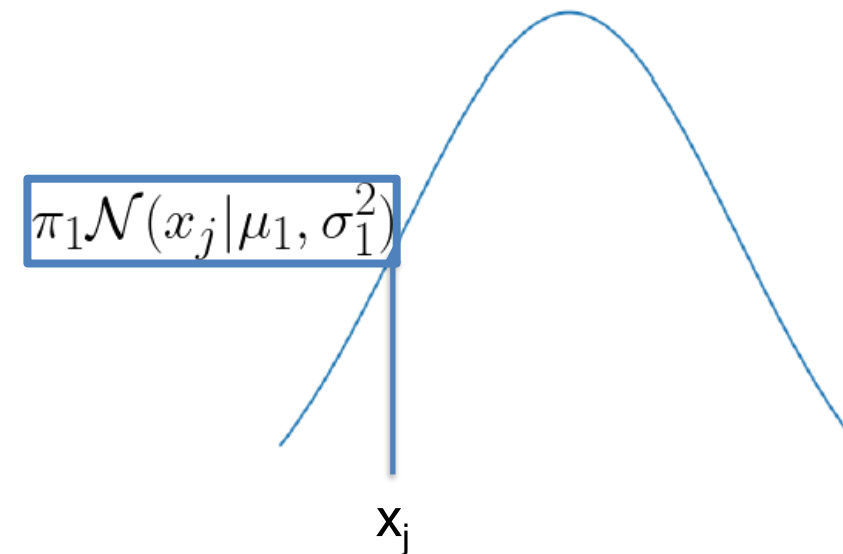


GMM Expectation-Maximization (1D)

- **Initialize cluster parameters**
- **Expectation (E-Step)**
 - For each data point, x_j
 - Compute cluster posterior probability
 - Compute probability with respect to C_i
 - Normalize to sum to one over clusters

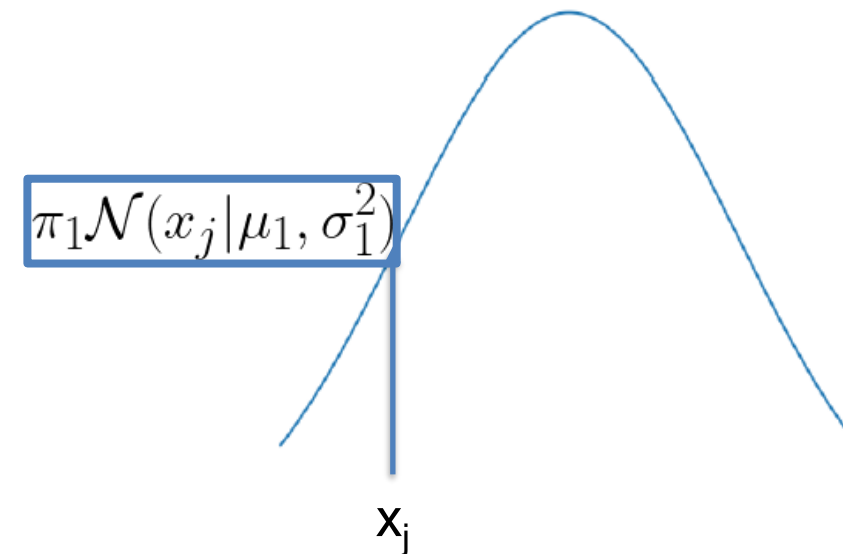
For each cluster:

$$f_i(x) = f(x|\mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left\{ -\frac{(x - \mu_i)^2}{2\sigma_i^2} \right\}$$



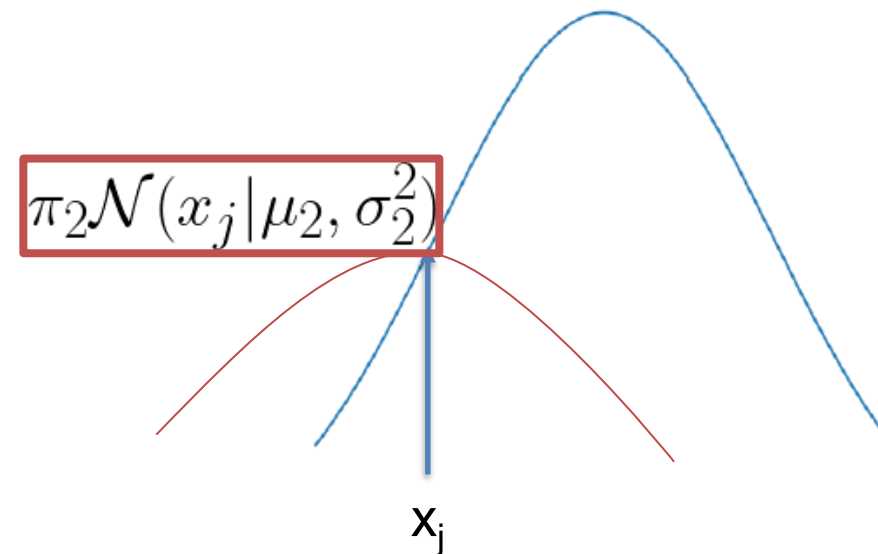
- Initialize cluster parameters
- Expectation (E-Step)
 - For each data point, x_j
 - **Compute cluster posterior probability**
 - Compute probability with respect to C_i
 - Normalize to sum to one over clusters

$$w_{ij} = P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$$



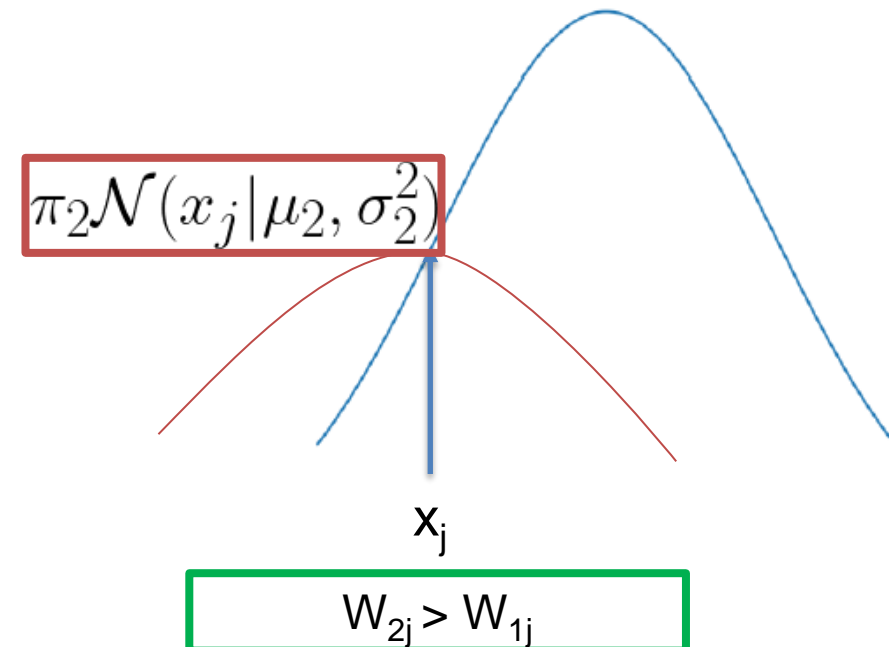
- Expectation (E-Step)
 - For each data point, x_j
 - **Compute cluster posterior probability**
 - Compute probability with respect to C_i
 - Normalize to sum to one over clusters

$$w_{ij} = P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$$



- Expectation (E-Step)
 - For each data point, x_j
 - Compute cluster posterior probability
 - Compute probability with respect to C_i
 - Normalize to sum to one over clusters
- **Higher probability will be assigned to Gaussian that is more likely**

$$w_{ij} = P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$$



- Maximization (M-Step)
 - Update parameters using (weighted) data points

$$w_{ij} = P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$$

Mean:

$$\mu_i = \frac{\sum_{j=1}^n w_{ij} \cdot x_j}{\sum_{j=1}^n w_{ij}}$$

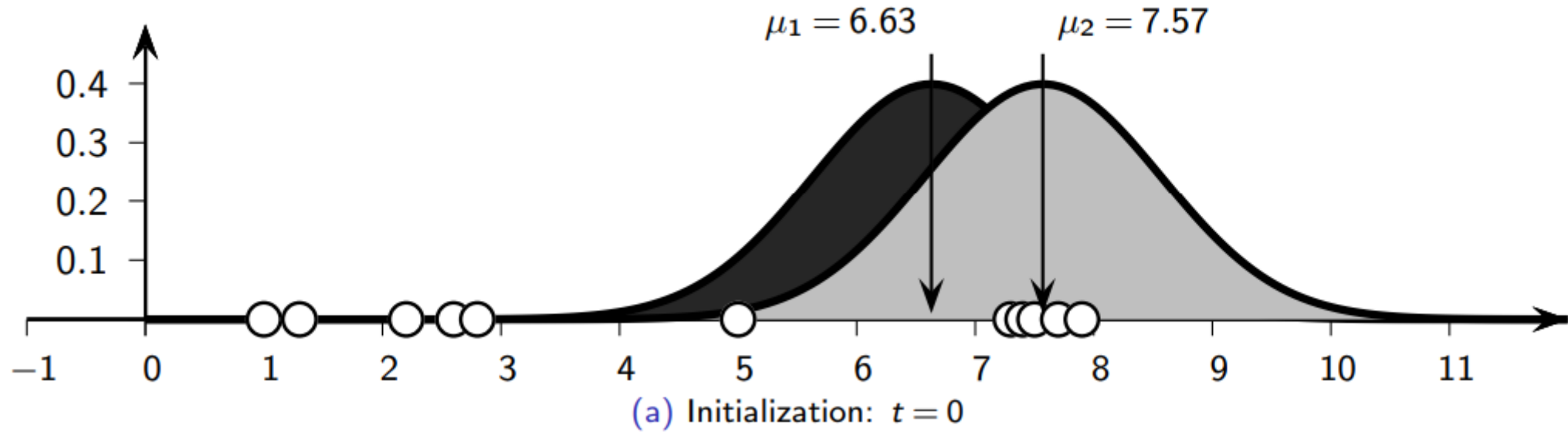
Variance:

$$\sigma_i^2 = \frac{\sum_{j=1}^n w_{ij} (x_j - \mu_i)^2}{\sum_{j=1}^n w_{ij}}$$

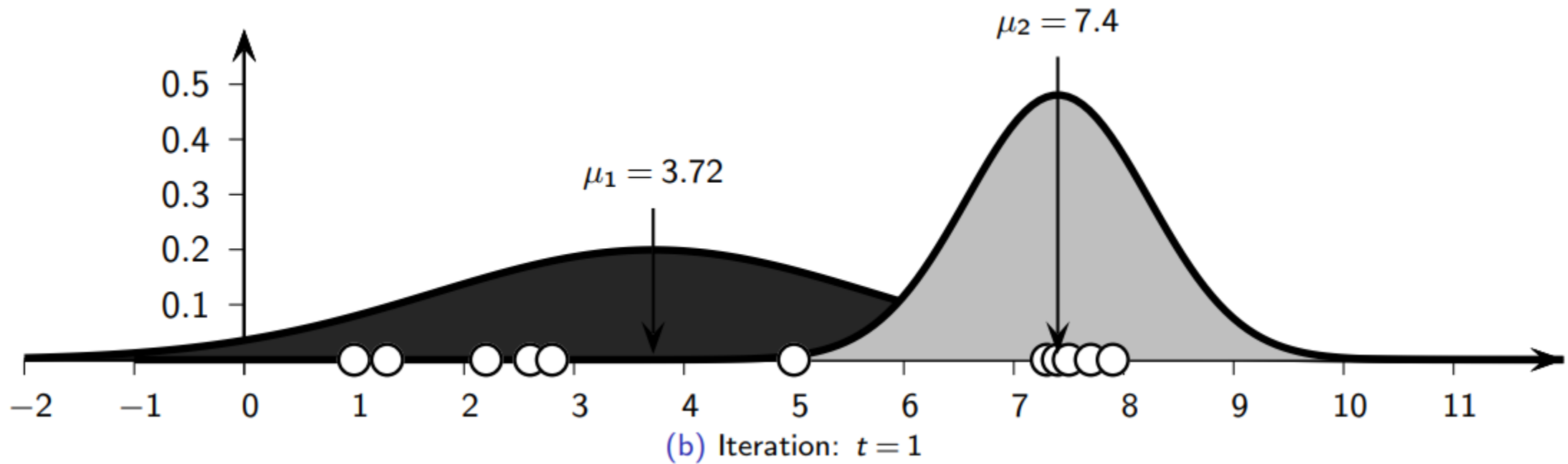
Mixture Weight/Prior Probability:

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{n}$$

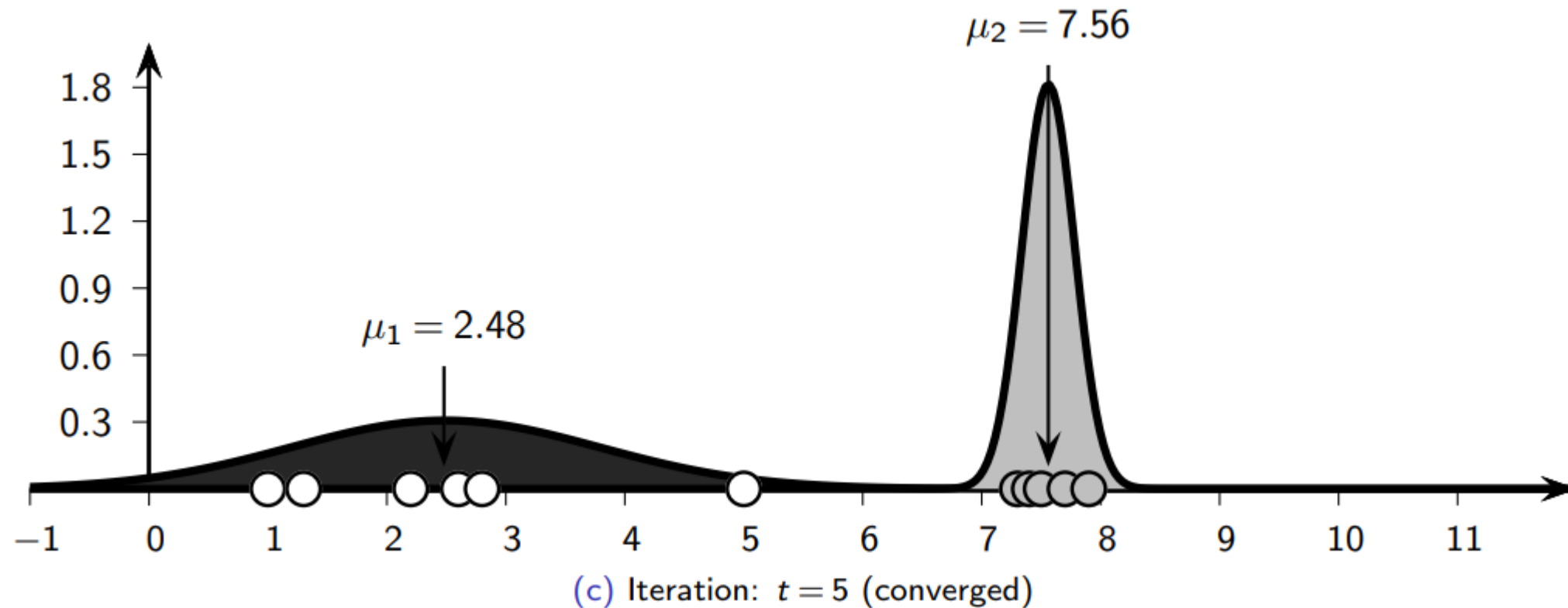
GMM EM 1D Example



GMM EM 1D Example



GMM EM 1D Example





GMM Expectation-Maximization (d-dimensions)

- Each cluster will have $d \times d$ covariance matrix
- Expensive to calculate and may be unreliable estimation
- Can use diagonal covariance
 - Assumes dimensions are independent

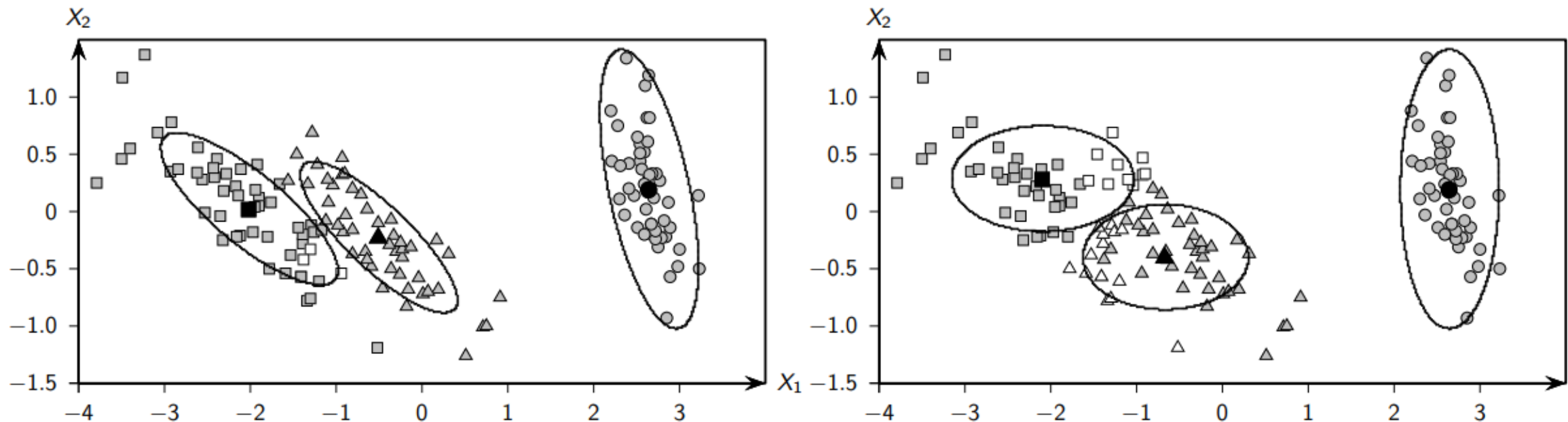
Full Covariance:

$$\Sigma_i = \begin{pmatrix} (\sigma_1^i)^2 & \sigma_{12}^i & \dots & \sigma_{1d}^i \\ \sigma_{21}^i & (\sigma_2^i)^2 & \dots & \sigma_{2d}^i \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1}^i & \sigma_{d2}^i & \dots & (\sigma_d^i)^2 \end{pmatrix}$$

Diagonal Covariance:

$$\Sigma_i = \begin{pmatrix} (\sigma_1^i)^2 & 0 & \dots & 0 \\ 0 & (\sigma_2^i)^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & (\sigma_d^i)^2 \end{pmatrix}$$

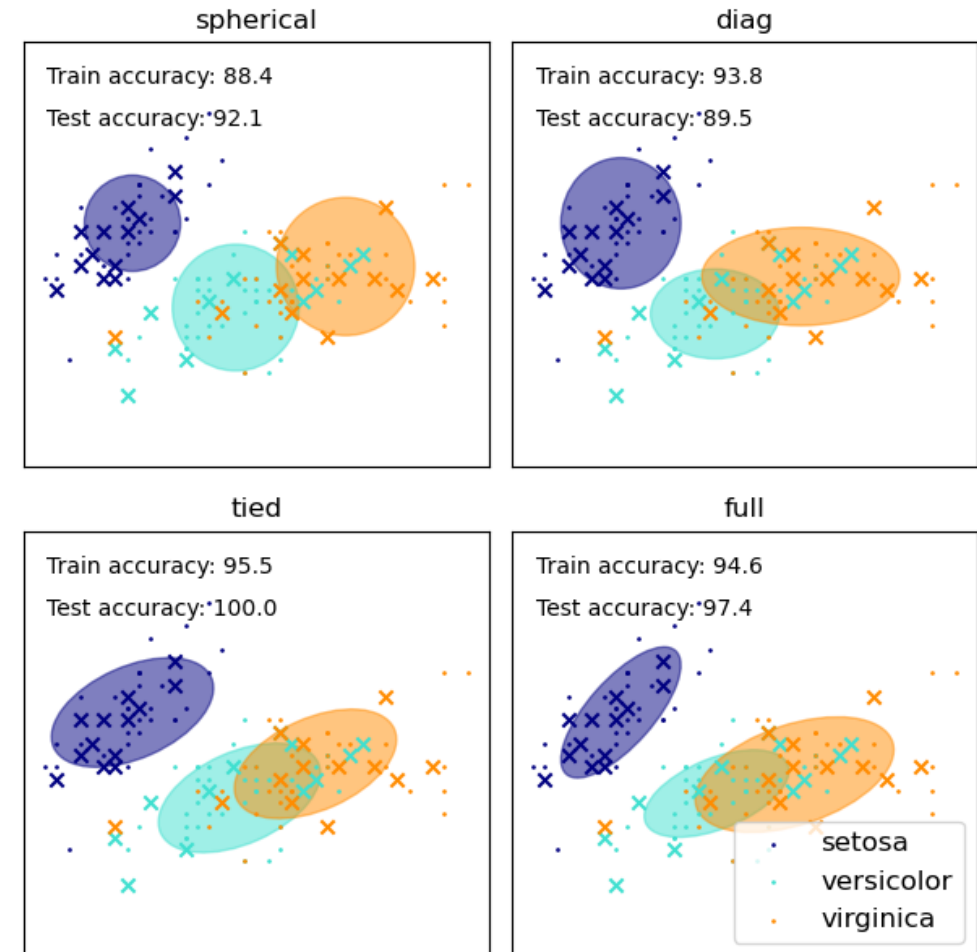
Full vs Diagonal



(a) Full covariance matrix ($t = 36$)

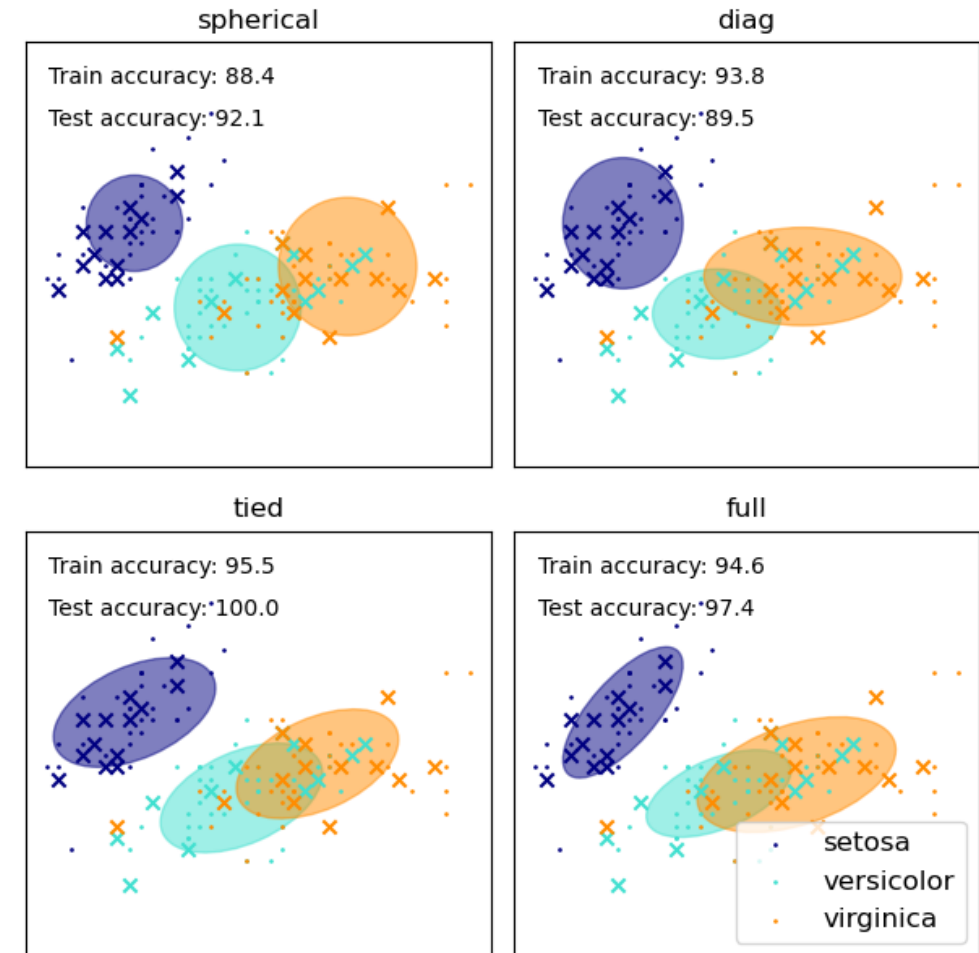
(b) Diagonal covariance matrix ($t = 29$)

- [Additional options](#) for covariance matrices include:
 - Spherical: Each cluster has a single variance (isotropic covariance)
 - Tied: All clusters share same covariance matrix



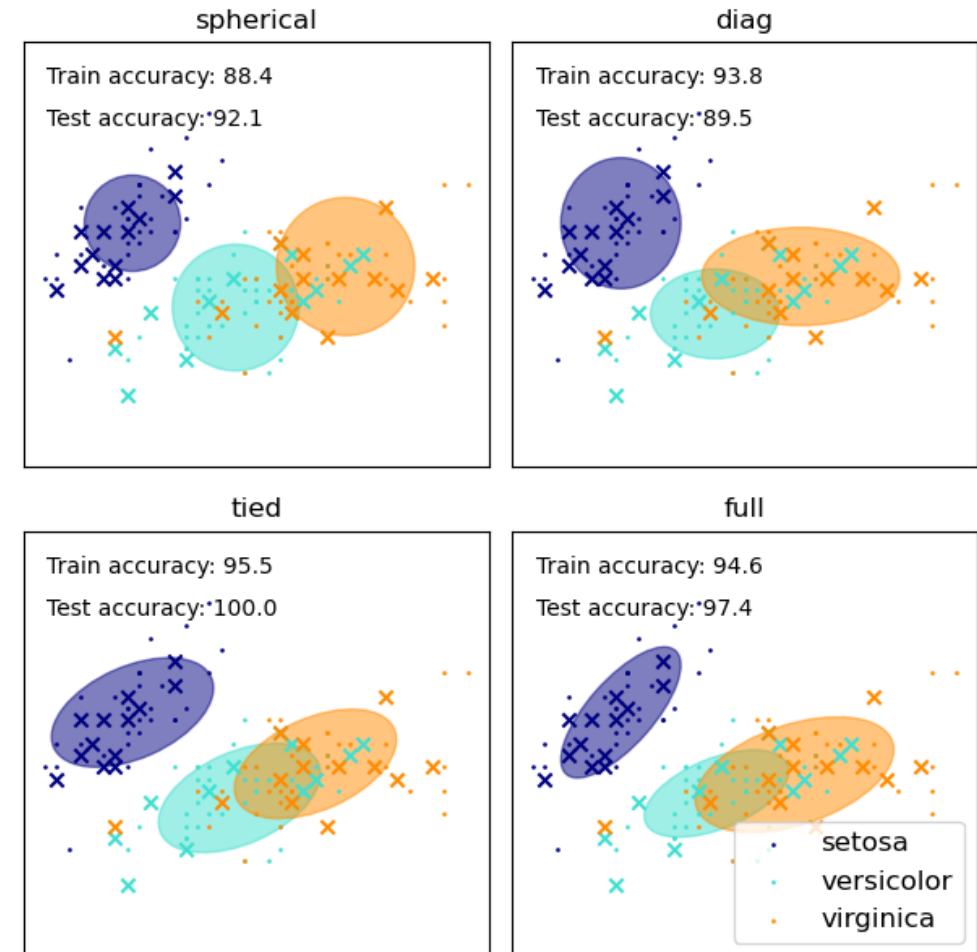
GMM Number of Parameters

- Given k clusters, n samples, and d features, what are the total number of parameters for a **tied covariance matrix** (i.e., all clusters share the same covariance matrix) GMM?
- Break into pairs
- 5 minutes for activity



GMM Number of Parameters

- Given k clusters, n samples, and d features, what are the total number of parameters for a **tied covariance matrix** (i.e., all clusters share the same covariance matrix) GMM?
- Solution:
 - $k*d + d^2 + k$
 - k mean vectors (d by 1), single covariance (d by d), and k mixture parameters (1×1 , scalars)



- Expectation step:

$$w_{ij} = P(C_i | \mathbf{x}_j) = \frac{f_i(\mathbf{x}_j) \cdot P(C_i)}{\sum_{a=1}^k f_a(\mathbf{x}_j) \cdot P(C_a)}$$

- Maximization step:

$$\mu_i = \frac{\sum_{j=1}^n w_{ij} \cdot \mathbf{x}_j}{\sum_{j=1}^n w_{ij}} \quad \Sigma_i = \frac{\sum_{j=1}^n w_{ij} (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T}{\sum_{j=1}^n w_{ij}} \quad P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{n}$$

- Each step maximizes log-likelihood
- Iterate until convergence
 - Set maximum iterations or set threshold for changes in parameters
 - May converge to local optima

MLE:

$$\theta^* = \arg \max_{\theta} \{\ln P(\mathbf{D}|\theta)\}$$

Log-likelihood:

$$\ln P(\mathbf{D}|\theta) = \sum_{j=1}^n \ln f(\mathbf{x}_j) = \sum_{j=1}^n \ln \left(\sum_{i=1}^k f(\mathbf{x}_j|\mu_i, \Sigma_i) P(C_i) \right)$$

GMM EM Algorithm Pseudocode



Expectation-Maximization (D, k, ϵ):

```
1  $t \leftarrow 0$ 
2 Randomly initialize  $\mu_1^t, \dots, \mu_k^t$ 
3  $\Sigma_i^t \leftarrow I, \forall i = 1, \dots, k$ 
4 repeat
5    $t \leftarrow t + 1$ 
6   for  $i = 1, \dots, k$  and  $j = 1, \dots, n$  do
7      $w_{ij} \leftarrow \frac{f(\mathbf{x}_j | \mu_i, \Sigma_i) \cdot P(C_i)}{\sum_{a=1}^k f(\mathbf{x}_j | \mu_a, \Sigma_a) \cdot P(C_a)}$  // posterior probability
8      $P^t(C_i | \mathbf{x}_j)$ 
9   for  $i = 1, \dots, k$  do
10     $\mu_i^t \leftarrow \frac{\sum_{j=1}^n w_{ij} \cdot \mathbf{x}_j}{\sum_{j=1}^n w_{ij}}$  // re-estimate mean
11     $\Sigma_i^t \leftarrow \frac{\sum_{j=1}^n w_{ij} (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T}{\sum_{j=1}^n w_{ij}}$  // re-estimate covariance
12    matrix
13     $P^t(C_i) \leftarrow \frac{\sum_{j=1}^n w_{ij}}{n}$  // re-estimate priors
14 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$ 
```


Next class



TEXAS A&M UNIVERSITY
Engineering

- Expectation-Maximization

INTEGRITY
EXCELLENCE LEADERSHIP



TEXAS A&M UNIVERSITY
Engineering

Thank You! Questions?
Joshua Peeples, Ph.D.
<https://www.joshpeeples.com/>
jpeeples@tamu.edu





TEXAS A&M UNIVERSITY
Engineering

Supplemental Slides

- [Gaussian Mixture Models and EM](#)
- [Gaussian Mixture Models Google Colab](#)