

1

The SJSU logo is located in the top left corner of the slide. It consists of the letters 'SJSU' in a large, blue, stylized font, followed by 'SAN JOSÉ STATE UNIVERSITY' in a smaller, orange, sans-serif font.

## Agenda

- Recap of K-Means
- Mixture Models
- Gaussian Mixture Models (GMM)

2

## K-Means Clustering

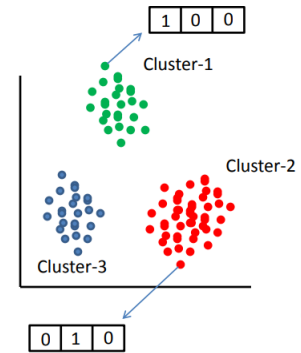
- K-Means is one of the most popular clustering algorithms.

- Input:

- Observations/data points ( $N$ ):  $x_i \quad \forall i \in \{1, \dots, N\}$
- # of clusters:  $k$

- Output:

- Cluster Assignments:  $w_{ij}$
- Cluster Centroids:  $c_j \quad \forall j \in \{1, \dots, k\}$



1-of- $k$  representation for cluster assignment.

- Objective Function: It aims to minimize the within-cluster sum of squares (WCSS):

3

## K-Means Clustering

- K-Means is a minimization problem with the objective function:

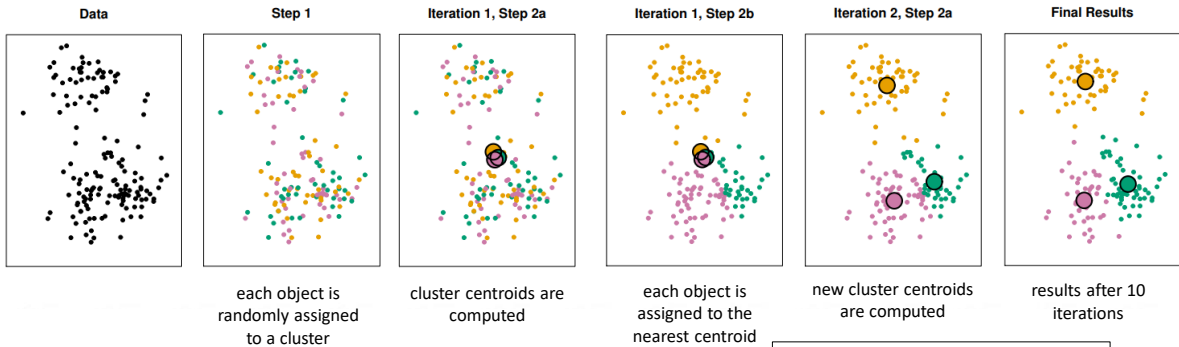
$$J = \sum_{j=1}^k \sum_{x_i \in C_j} (x_i - c_j)^2 = \sum_{j=1}^k \sum_{i=1}^N w_{ij} (x_i - c_j)^2 \quad w_{ij} = \begin{cases} 1 & \text{if } x_i \in C_j \\ 0 & \text{if } x_i \notin C_j \end{cases}$$

- We need to find  $w_{ij}$  and  $c_j \rightarrow$  minimize  $J$

$$c_k = \frac{\sum_{i=1}^N w_{ik} x_i}{\sum_{i=1}^N w_{ij}}$$

4

## K-Means Example & Algorithm



### Iterative Algorithm For K-Means

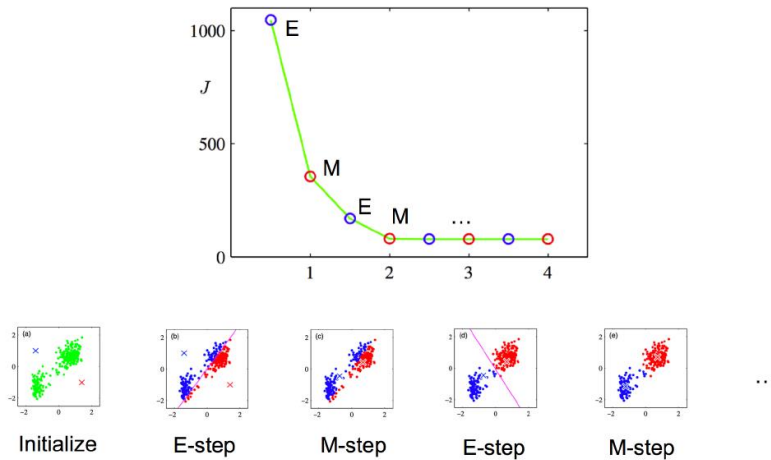
- Initialize  $k$  centroids
- Repeat till **convergence**
  - Calculate  $w_{ij}$
  - Update  $c_j = \frac{\sum_{i=1}^N w_{ij} x_i}{\sum_{i=1}^N w_{ij}}$

Expectation Step (E-Step)

Maximization Step (M-Step)

5

## K Means as an E-M Algorithm



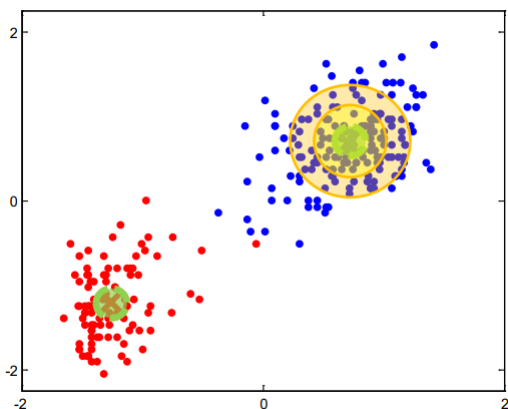
6

## Pros & Cons

- **Good**
  - Simple to implement
  - Fast
- **Bad**
  - Local minima
  - Model only “spherical” clusters
  - Sensitive to the features scale
  - Number of clusters  $K$  to be chosen in advance
  - Cluster assignments are “hard”, not probabilistic => Gaussian Mixture Model

7

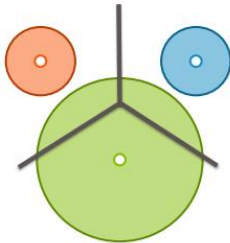
## Problem with K-Means



- K-Means makes hard guesses for cluster assignment.
- For some cases our model may not be sure about exact cluster assignment.
- Can we make this probabilistic (probability that the  $n^{\text{th}}$  observation belongs to the  $k^{\text{th}}$  cluster) ?

8

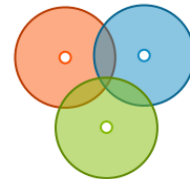
## Different Failure Modes of K-Means



disparate cluster sizes



different  
shaped/oriented  
clusters



overlapping clusters

9

## Mixture Models



10

## Hard Clustering vs. Soft Clustering

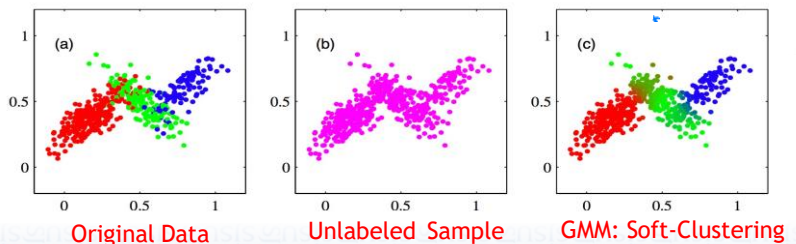
- Hard Clustering

- Every object  $i$  is assigned to one cluster  $j$  (e.g., k-means)
- $w_i = \{0, 1\}$  and  $\sum_j w_{ij} = 1$

$$z_i = \arg \min_j \|x_i - \mu_j\|^2$$

- Soft Clustering

- Every object  $i$  is assigned with a probability to different clusters
- $w_i = [0, 1]$  and  $\sum_j w_{ij} = 1$



11

## Latent Variables

- Latent variables are “hidden” or unobserved variables that influenced the observed data.
- In clustering, latent variables  $z$  often represent the cluster to which an observation or data point belongs:
  - cluster assignment for K-Means
  - probability of belonging to a cluster for GMMs and other probabilistic models

12

## Mixture Models

$$p(z, \mathbf{x}) = \underbrace{p(z)}_{\text{prior probability of assignment}} \underbrace{p(\mathbf{x} | z)}_{\text{likelihood}}$$

- Prior probability encodes our belief about the latent variable  $z$  (cluster assignment) before observing any data.

$$p(z = k) = \pi_k \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$$

- Likelihood captures the probability of the data  $\mathbf{x}$  given a cluster assignment  $z$ .
- Recall that the law of total probability  $\rightarrow$  the marginal probability of the data  $\mathbf{x}$ :

$$p(\mathbf{x}) = \sum_{k=1}^K p(z = k) p(\mathbf{x} | z = k) = \sum_{k=1}^K \pi_k p(\mathbf{x} | z = k)$$

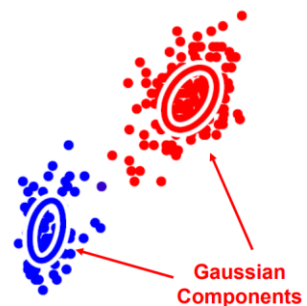
13

## Gaussian Mixture Model

- One of the most common mixtures are over Gaussians:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k p(\mathbf{x} | z = k) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mu_k, \Sigma_k)$$

- Unlike K-Means, assignment (probabilities) aren't just distance.



14

## Soft GMM Assignments (Responsibilities)

- Using Bayes' rule, we have the posterior probability/inference for a given data point  $\mathbf{x}_i$ :

$$p(z_i = k | \mathbf{x}_i) = \frac{p(z_i = k) \cdot p(\mathbf{x}_i | z_i = k)}{p(\mathbf{x})}$$

- For Gaussian mixtures, this is:

$$p(z_i = k | \mathbf{x}_i) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k)}{p(\mathbf{x})}$$

- In mixture modeling, this is called the responsibility (it is how "responsible" cluster  $k$  is for data point  $\mathbf{x}_i$ ):

$$\gamma(z_i = k) = p(z_i = k | \mathbf{x}_i) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i; \mu_j, \Sigma_j)}$$

15

## Gaussian Mixture Model Summary

- Mixture model is a weighted combination of component distributions:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k p(\mathbf{x} | z = k)$$

- Bayes' rule gives the posterior probability/inference of assignment (responsibility):

$$p(z_i = k | \mathbf{x}_i) = \frac{p(z_i = k) \cdot p(\mathbf{x}_i | z_i = k)}{p(\mathbf{x})}$$

- A GMM uses Gaussian component distributions with responsibilities:

$$p(z_i = k | \mathbf{x}_i) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i; \mu_j, \Sigma_j)}$$

16



## Parameter Estimation for GMMs

- For data  $\mathbf{x}$  that's  $N$ -dimensional:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

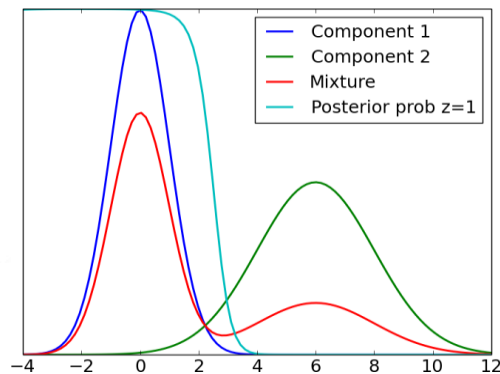
mixture probabilities or proportions  
 $\pi_k$   
 $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$   
 $\boldsymbol{\mu}_k$ :  $N$ -dimensional vector of mean parameters  
 $\boldsymbol{\Sigma}_k$ :  $N \times N$ -dimensional matrix of covariance parameters



- For  $K$  components, we need to estimate (or train):  $O(K + KN + KN^2)$  parameters!!!

17

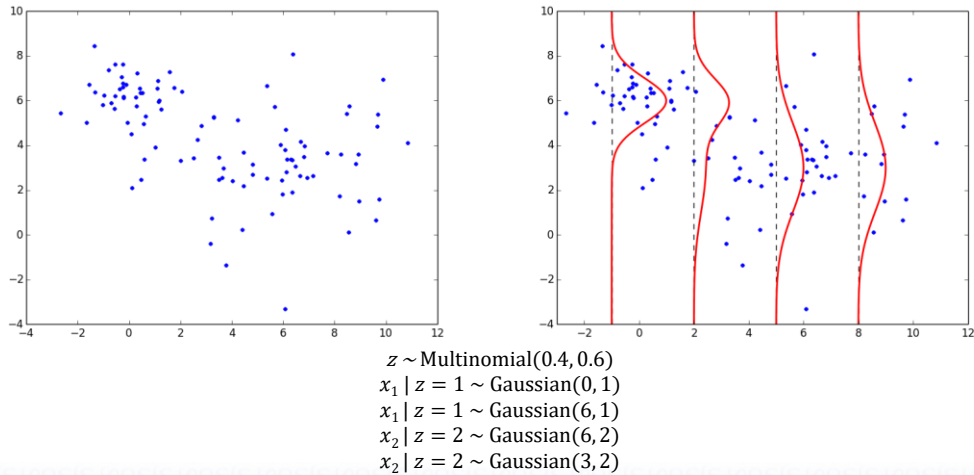
## Example: Mixture of 2 Gaussians



$z \sim \text{Multinomial}(0.7, 0.3)$   
 $\mathbf{x} | z = 1 \sim \text{Gaussian}(0, 1)$   
 $\mathbf{x} | z = 2 \sim \text{Gaussian}(6, 2)$

18

## Example: Mixture of Gaussians in 2D



19

## Parameter Estimation Using MLE

- We need to fit two sets of parameters:
  - The mixture probabilities  $\pi_k$
  - The mean  $\mu_k$  and standard deviation  $\Sigma_k$  for each component
- Recall that the likelihood of a single data point  $\mathbf{x}_i$  is given by:

$$p(\mathbf{x}_i) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k)$$

- For  $N$  points, the log likelihood is:

$$\mathcal{L}(\mu, \Sigma) = \log \prod_{i=1}^N p(\mathbf{x}_i) = \sum_{i=1}^N \log \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k) \right\}$$

20

## Parameter Estimation Using MLE

- We need to maximize  $\mathcal{L}(\mu, \Sigma)$  wrt  $\pi_k$ ,  $\mu_k$  and  $\Sigma_k$

$$\mathcal{L}(\mu, \Sigma) = \log \prod_{i=1}^N p(\mathbf{x}_i) = \sum_{i=1}^N \log \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k) \right\}$$

- This is highly non-convex → very difficult to optimize!!!
- Instead, we try to find the local minima by using an iterative algorithm → EM method (Expectation-Maximization)

21

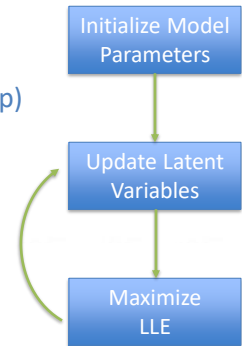
## Expectation Maximization (EM) Method

- A very powerful method for dealing with probabilistic models that involve latent/missing variables.
- Each iteration of the EM is guaranteed to maximize the data log likelihood.
- Guaranteed to converge to a local maxima.
- Sensitive to starting points.
- We have applied it to Gaussian Mixture Models, which can model any arbitrary shaped densities. Can be used for data density estimation aside from clustering

22

## Expectation Maximization Method

- Here's the overview:
  - Start with initial guesses for the model parameters:  $\pi_k, \mu_k, \Sigma_k$
  - Update latent variables based on our expectations (E-Step)
  - Update model parameters to maximize log likelihood estimates (M-Step)
  - Keep repeating steps 2 and 3 until changes in  $\mathcal{L}(\mu, \Sigma)$  is small



23

## Expectation Step (E-Step)

- This is similar to the cluster assignment step in K-Means, except that we update the (fractional) responsibilities here:

$$\gamma(z_i = k) = p(z_i = k | \mathbf{x}_i) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i; \mu_j, \Sigma_j)}$$

- This is the posterior inference that was discussed earlier.

24

## Maximization Step

- Define the effective number of points in cluster  $k$  by:

$$N_k = \sum_{j=1}^N \gamma(z_j = k)$$

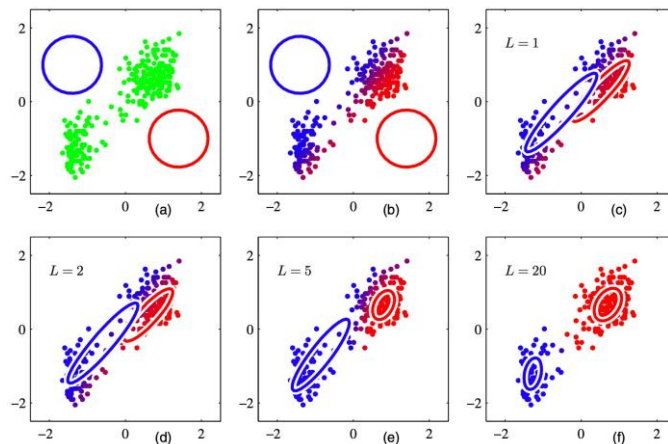
- Update model parameters  $\pi_k, \mu_k, \Sigma_k$  with updated responsibilities  $\gamma(z_i = k)$ :

$$\mu_k = \frac{1}{N_k} \sum_{j=1}^N \gamma(z_j = k) \mathbf{x}_j \quad \pi_k = \frac{N_k}{N}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{j=1}^N \gamma(z_j = k) (\mathbf{x}_j - \mu_k)(\mathbf{x}_j - \mu_k)^T$$

25

## Example: GMM with EM



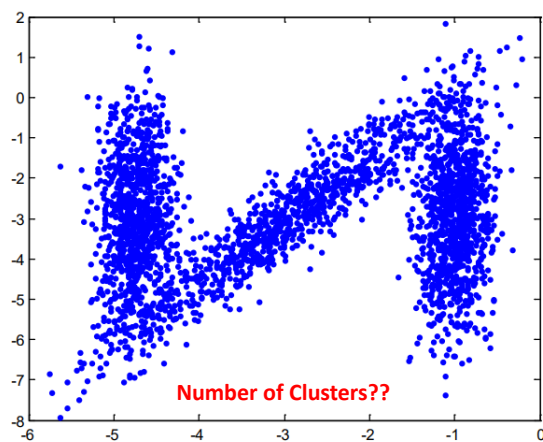
26

## Summary of GMM Components

- Observations or Data:  $\mathbf{x}_i \in \mathbb{R}^N$ ,  $i = \{1, 2, \dots, N\}$
- Hidden Cluster Labels  $z_i \in \{1, 2, \dots, N\}$ ,  $i = \{1, 2, \dots, N\}$
- Hidden Mixture Means  $\boldsymbol{\mu}_k \in \mathbb{R}^N$ ,  $k = \{1, 2, \dots, K\}$
- Hidden Mixture Covariances  $\boldsymbol{\Sigma}_k \in \mathbb{R}^{N \times N}$ ,  $k = \{1, 2, \dots, K\}$
- Hidden Mixture Probabilities  $\pi_k \in \mathbb{R}^N$ ,  $\sum_{j=1}^N \pi_k = 1$

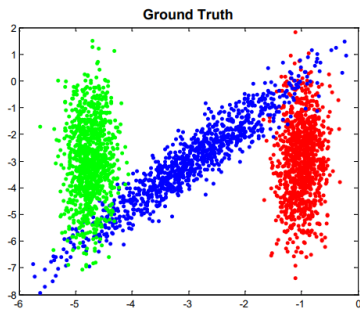
27

## Example: GMM vs K-Means

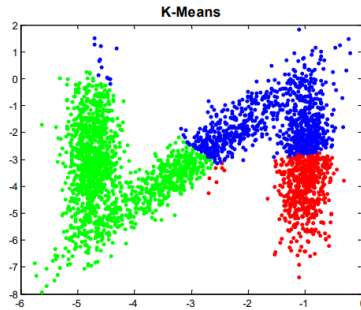


28

## Example

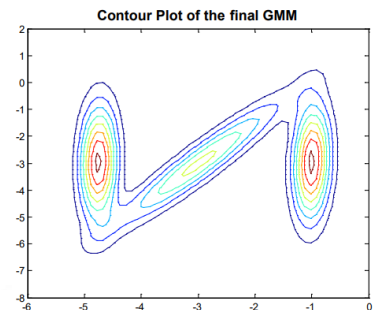


[-1,-3]  
[-3,-3]  
[-4.75,-3]



K-Means with 50 random starting points

[-1.0335,-4.057]  
[-1.5821, -1.6458]  
[-4.3681, -3.4009]



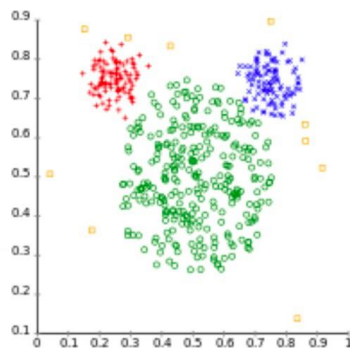
Soft Clustering using a 3 component Gaussian Mixture Model with random starting point

[-1.0006, -2.9663]  
[-2.9747, -2.9921]  
[-4.7488, -2.9717]

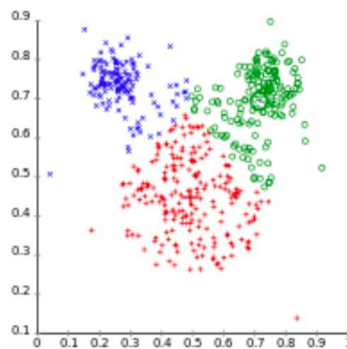
29

## Example

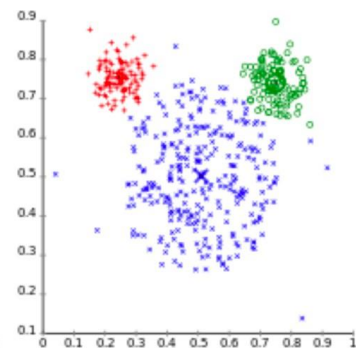
- Mouse Clusters...



Original Clusters



K-means



GMM

30

## Summary

- Mixture models are probabilistic models that represent data as a combination of multiple underlying distributions. Each underlying distribution models a specific cluster (or component) within the data.
- GMM is one type of mixture models by using Gaussian as the probability distribution.
- Model parameters for GMM include means, covariances, and mixing probability for each Gaussian component.
- GMMs are more flexible as it can model different mean/covariance for each cluster
- However, GMMs have (local) convergence issues, just like K-Means
- GMMs are also slower than K-Means due to the higher computation effort needed by the EM algorithm. K-Means is sometimes used to initialize the cluster means.