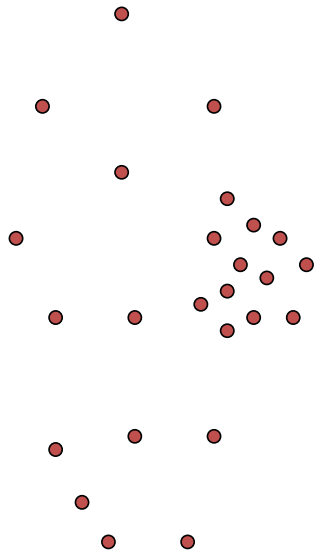
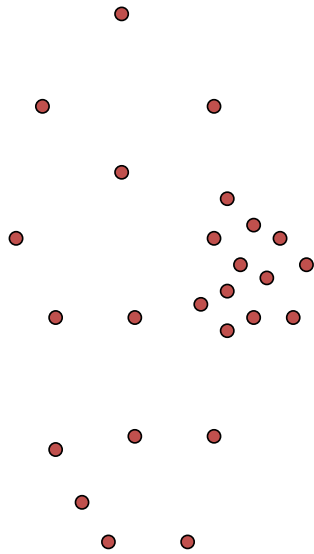


## (One) bad case for K-means



- Clusters may overlap
- Some clusters may be “wider” than others
- Clusters may not be linearly separable

## (One) bad case for K-means



- Clusters may overlap
- Some clusters may be “wider” than others
- Clusters may not be linearly separable

# Partitioning Algorithms

- K-means
  - **hard assignment**: each object belongs to only one cluster
- Mixture modeling
  - **soft assignment**: probability that an object belongs to a cluster

*Generative approach*

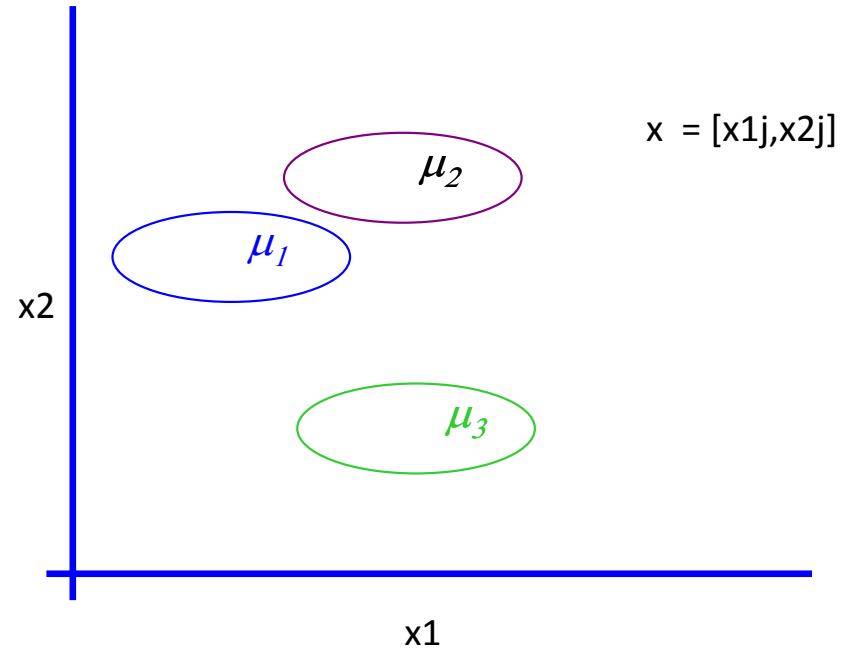
# Gaussian Mixture Model

Mixture of K Gaussian distributions: (Multi-modal distribution)

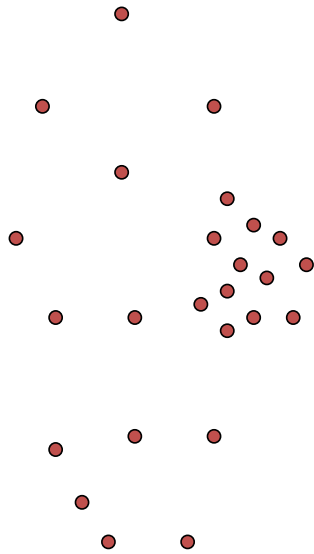
$$p(x|y=i) \sim N(\mu_i, \sigma^2 I)$$

$$p(x) = \sum_i p(x|y=i) P(y=i)$$

↓                      ↓  
**Mixture**           **Mixture**  
**component**       **proportion**



## (One) bad case for K-means



- Clusters may overlap
- Some clusters may be “wider” than others
- Clusters may not be linearly separable

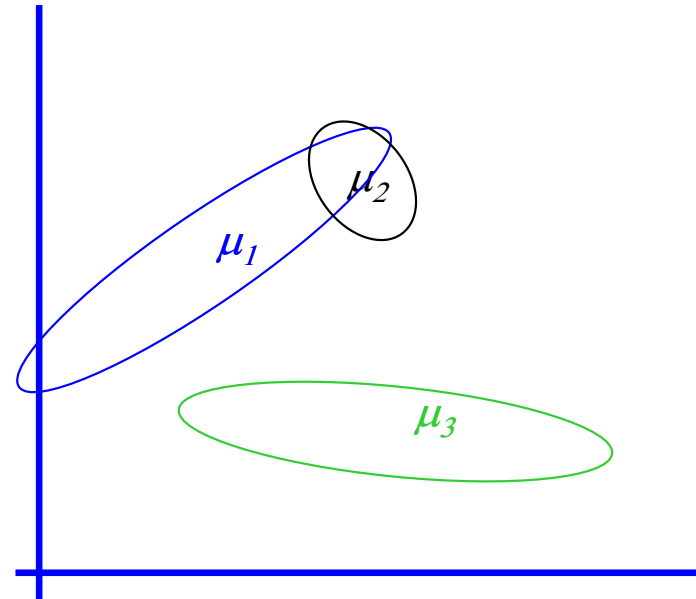
# General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

$$p(x|y=i) \sim \mathcal{N}(\mu_i, \Sigma_i)$$

$$p(x) = \sum_i p(x|y=i) P(y=i)$$

↓                      ↓  
**Mixture**           **Mixture**  
**component**       **proportion**



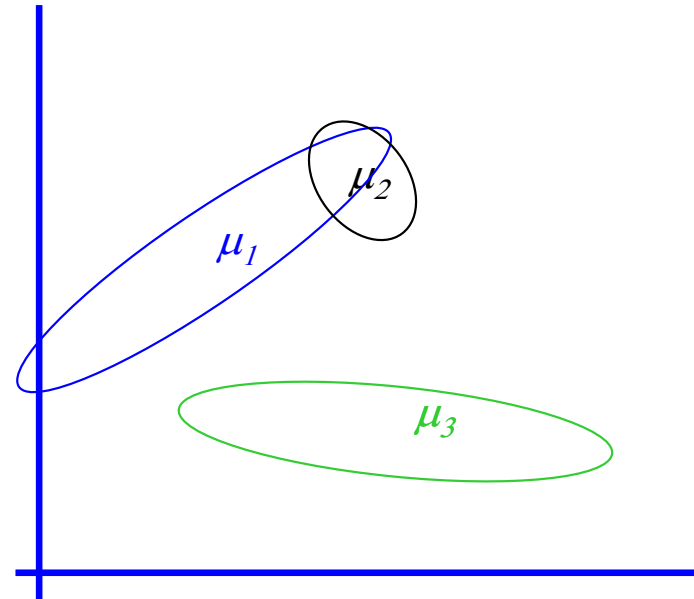
# General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

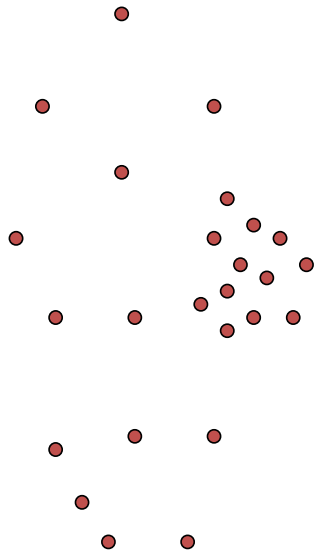
- There are  $k$  components
- Component  $i$  has an associated mean vector  $\mu_i$
- Each component generates data from a Gaussian with mean  $\mu_i$  and covariance matrix  $\Sigma_i$

Each data point is generated according to the following recipe:

- 1) Pick a component at random:  
Choose component  $i$  with probability  $P(y=i)$
- 2) Datapoint  $x \sim N(\mu_i, \Sigma_i)$



## (One) bad case for K-means



- Clusters may overlap
- Some clusters may be “wider” than others
- Clusters may not be linearly separable



# General GMM

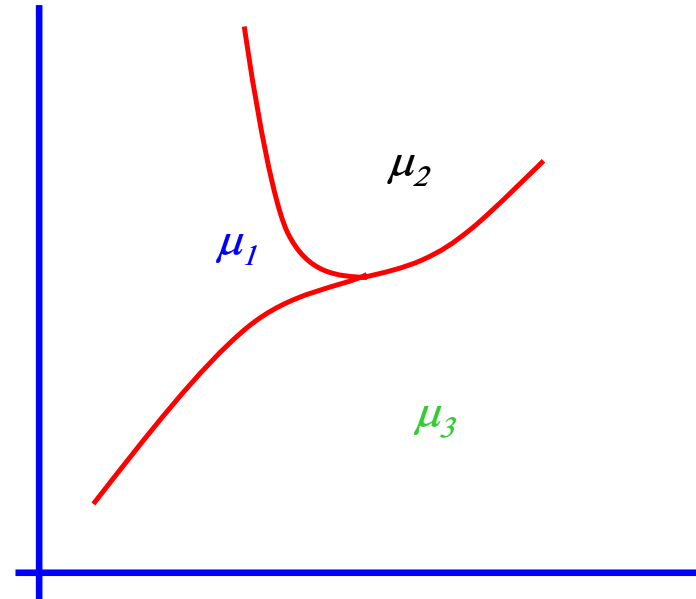
GMM – Gaussian Mixture Model (Multi-modal distribution)

$$p(x|y=i) \sim N(\mu_i, \Sigma_i)$$

Decision boundary when probabilities are equal:

$$\begin{aligned} \log \frac{P(y=i|x)}{P(y=j|x)} \\ &= \log \frac{p(x|y=i)P(y=i)}{p(x|y=j)P(y=j)} \\ &= \mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{w}^T \mathbf{x} \end{aligned}$$

Depend on  $\mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K, P(y=1), \dots, P(y=K)$



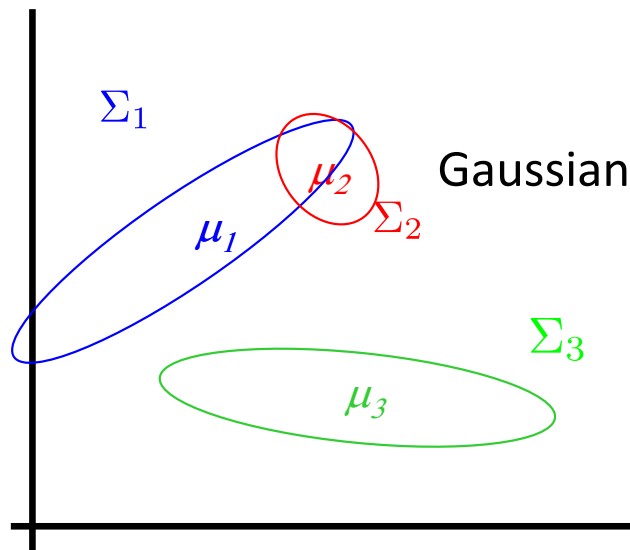
**“Quadratic Decision boundary”** – second-order terms don’t cancel out

# Learning General GMM

$$x_1, \dots, x_m \sim p(x) = \sum_{i=1}^k p(x|Y = i) P(Y = i)$$

↓  
**Mixture  
component**

↓  
**Mixture  
proportion,  $p_i$**



Gaussian mixture model

$$p(x|Y = i) \sim \mathcal{N}(\mu_i, \Sigma_i)$$

**Parameters:**  $\{p_i, \mu_i, \Sigma_i\}_{i=1}^K$

- How to estimate parameters? Max Likelihood  
But don't know labels  $Y$

# Learning General GMM

Maximize marginal likelihood:

$$\begin{aligned}\operatorname{argmax} \prod_j P(x_j) &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i, x_j) \\ &= \operatorname{argmax} \prod_j \sum_{i=1}^K P(y_j=i) p(x_j | y_j=i)\end{aligned}$$

$P(y_j=i) = P(y=i)$  Mixture component  $i$  is chosen with prob  $P(y = i)$

$$= \operatorname{arg max} \prod_{j=1}^m \sum_{i=1}^k P(y = i) \frac{1}{\sqrt{\det(\Sigma_i)}} \exp \left[ -\frac{1}{2} (x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i) \right]$$

How do we find the  $\mu_i, \Sigma_i$  s and  $P(y=i)$ s which give max. marginal likelihood?

\* Set  $\frac{\partial}{\partial \mu_i} \log \text{Prob} (\dots) = 0$  and solve for  $\mu_i$ 's. Non-linear not-analytically solvable

\* Use gradient descent: Doable, but often slow