**Clustering**

Classification

$$\mathcal{X} = \{(x_i, y_i)\}_{i=1}^{N}$$

→ class labels

↳ data points

**Binary Classification**

$y_i \in \{0, 1\}$ or $y_i \in \{-1, +1\}$

**Multiclass Classification**

$y_i \in \{1, 2, \ldots, K\}$

Clustering

$$\mathcal{X} = \{x_i\}_{i=1}^{N}$$

**NO CLASS LABELS !**

**PARAMETRIC CLASSIFICATION**

- We assumed that each class follows a certain density

$$p(x \mid y = c)$$

- We estimated the parameters

$p(x \mid y=1) \quad P(y=1) \quad \ldots \ldots \quad p(x \mid y=k) \quad P(y=K)$

↓ ↓ ↓ ↓

$\hat{\mu}_1, \hat{\Sigma}_1 \qquad \hat{P}(y=1) \qquad \hat{\mu}_k, \hat{\Sigma}_k \qquad \hat{P}(y=K)$
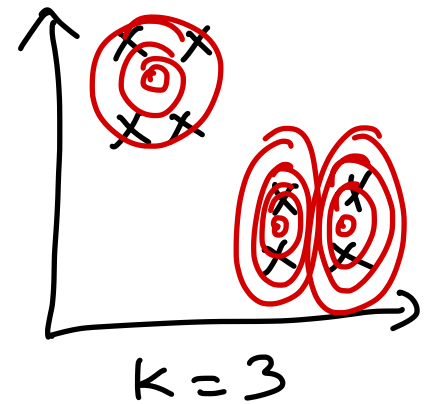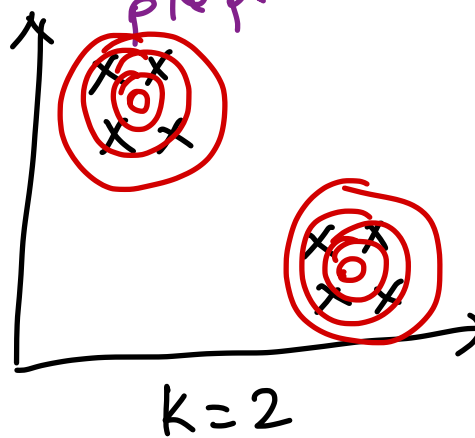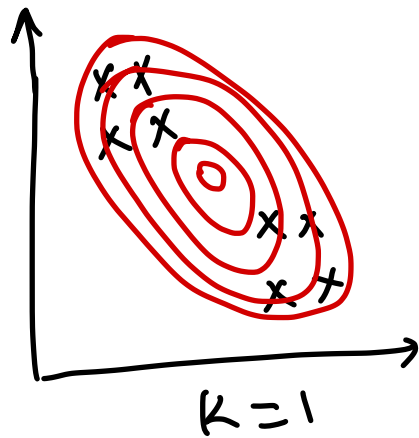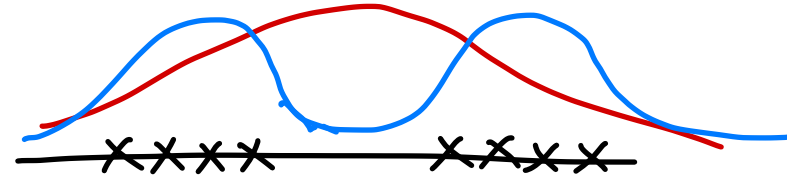
$$P(y=c \mid x) = ?$$

# Mixture Densities  K different clusters (unknown)

$C_k$ = cluster #$k$

$$p(x) = \sum_{k=1}^{K} \underbrace{p(x|C_k)}_{\substack{\text{component} \\ \text{density}}} \underbrace{P(C_k)}_{\substack{\text{mixture} \\ \text{proportions}}}$$

K = # of components
(clusters)
(groups)

$$\Phi = \{ \hat{P}(C_k), \hat{\mu}_k, \hat{\Sigma}_k \}_{i=1}^{k}$$

$$y_{ik} = \begin{cases} 1 & \text{if } x_i \text{ belongs to component/cluster/group } k. \\ 0 & \text{otherwise} \end{cases}$$

↳ cluster/component/group membership
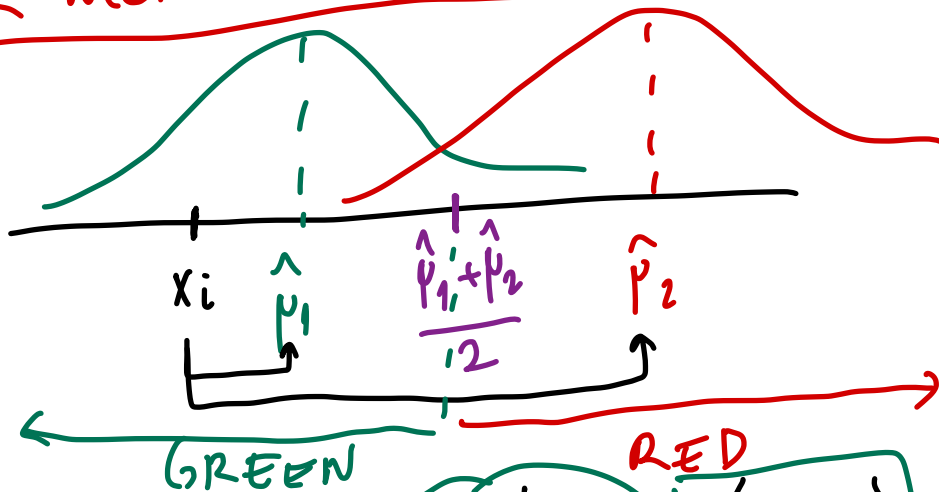
WE DO NOT KNOW "$y_{ik}$" VALUES APRIORI!

<u>Iterative Algorithm:</u>

STEP ① : Estimate the cluster memberships $(\hat{y}_{i:k})$

STEP ② : Estimate the parameters.

$$\hat{P}(C_k) = \frac{\sum_{i=1}^{N} \hat{y}_{ik}}{N}$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^{N} \hat{y}_{ik} \cdot x_i}{\sum_{i=1}^{N} \hat{y}_{ik}}$$

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^{N} \hat{y}_{ik}(x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T}{\sum_{i=1}^{N} \hat{y}_{ik}}$$

<u>K-MEANS CLUSTERING</u>



$x_i \quad \hat{\mu}_1 \qquad \frac{\hat{\mu}_1 + \hat{\mu}_2}{2} \qquad \hat{\mu}_2$

GREEN

RED

$$P(y=1|x) = \frac{\boxed{P(x|y=1)}\,\boxed{P(y=1)}}{P(x)} \approx$$

$$\hat{\sigma}_1^2 = \hat{\sigma}_2^2$$

$$P(y=2|x) = \frac{\boxed{P(x|y=2)}\,\boxed{P(y=2)}}{P(x)}$$

$$\exp\left[-\frac{(x_i - \hat{\mu}_1)^2}{2\hat{\sigma}_1^2}\right] \cdot \frac{1}{\sqrt{2\pi\hat{\sigma}_1^2}}$$

$$\exp\left[-\frac{(x_i - \hat{\mu}_2)^2}{2\hat{\sigma}_2^2}\right] \cdot \frac{1}{\sqrt{2\pi\hat{\sigma}_2^2}}$$

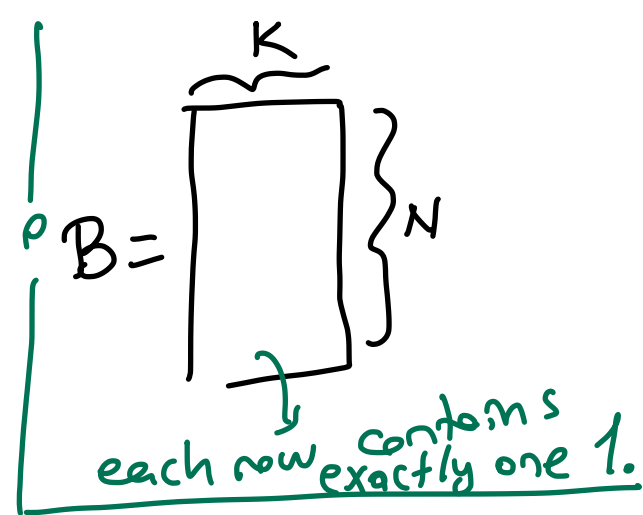$$\boxed{\hat{\sigma}_1^2 = \hat{\sigma}_2^2} \quad \Leftarrow$$

$$\|x_i - \hat{\mu}_1\|_2 \quad \|x_i - \hat{\mu}_2\|_2 \quad \cdots \quad \|x_i - \hat{\mu}_k\|_2$$

assume that 2nd distance is minimum.

$$\hat{y}_{i1} = 0 \quad \hat{y}_{i2} = 1 \quad \hat{y}_{i3} = 0 \quad \cdots \quad \hat{y}_{ik} = 0$$

$$\text{Error} = \sum_{i=1}^{N} \sum_{k=1}^{K} \widehat{\left(b_{ik}\right)} \| x_i - \widehat{\mu_k} \|_2^2$$

$$b_{ik} = \begin{cases} 1 & \text{if } \| x_i - \hat{\mu}_k \|_2 = \min_{c=1}^{K} \| x_i - \hat{\mu}_c \|_2 \\ 0 & \text{otherwise} \end{cases}$$

$$B = \begin{bmatrix} \phantom{x} \\ \phantom{x} \\ \phantom{x} \end{bmatrix} \Big\} N$$

each row contains exactly one 1.

MIP $\begin{bmatrix} \text{minimize} \sum_{i=1}^{N} \sum_{k=1}^{K} b_{ik} \| x_i - \hat{\mu}_k \|_2^2 \\ \text{with respect to: } \hat{\mu}_1, \hat{\mu}_2, \ldots, \hat{\mu}_K, \{b_{ik}\}_{i=1, k=1}^{N, K} \end{bmatrix}$

— Initialize $\hat{\mu}_1, \hat{\mu}_2, \ldots, \hat{\mu}_K$ randomly

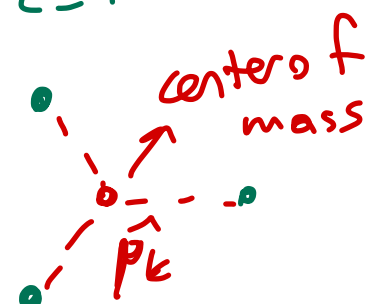— Repeat

E-STEP $a \to$ $\begin{bmatrix} \text{for all } x_i : \\ b_{ik} = \begin{cases} 1 & \text{if } \| x_i - \hat{\mu}_k \|_2 = \min_{c=1}^{K} \| x_i - \hat{\mu}_c \|_2 \\ 0 & \text{otherwise} \end{cases} \end{bmatrix}$
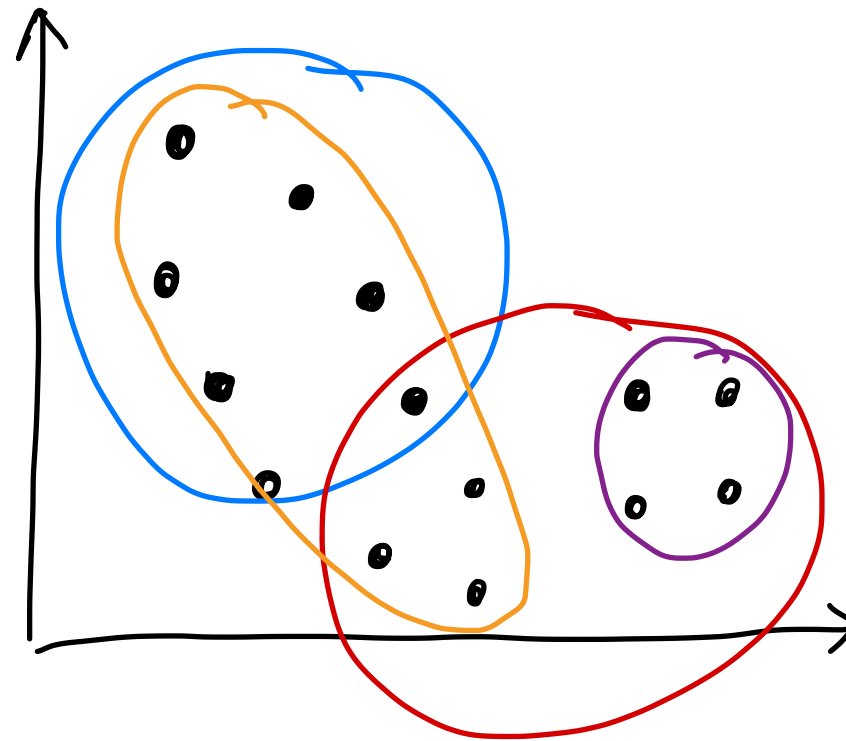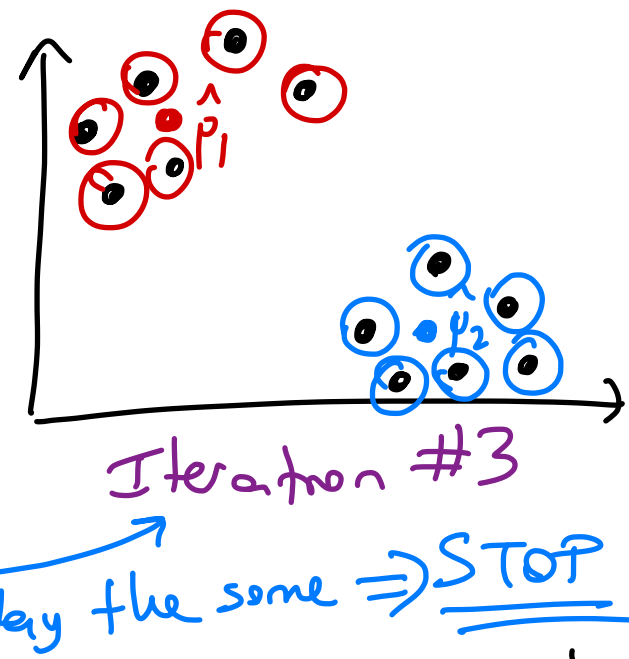
M-STEP $b \to$ $\begin{bmatrix} \text{for all } \hat{\mu}_k : \\ \hat{\mu}_k = \dfrac{\sum_{i=1}^{N} b_{ik} x_i}{\sum_{i=1}^{N} b_{ik}} \end{bmatrix}$

center of mass

$\hat{\mu}_k$

— Until convergence $\big[$ all $b_{ik}$'s stay the same $\big]$ or $\big[$ all $\mu_k$'s stay the same $\big]$

$K=2$

Iteration #1

Iteration #2

Iteration #3

$\hat{\mu}_1$    $\hat{\mu}_2$

$\mu$'s stay the same $\Rightarrow$ STOP

} assumed shared covariance

} assumed different covariances

# Expectation - Maximization Algorithm

$$X = \{x_i\}_{i=1}^{N} \qquad \text{log likelihood} \Rightarrow L(\Phi | x) = \log\left[\prod_{i=1}^{N} p(x_i | \Phi)\right]$$

$$\log L(\Phi | x) = \sum_{i=1}^{N} \log\left[\sum_{k=1}^{K} p(x_i | C_k) P(C_k)\right]$$

$$\underbrace{\qquad\qquad}_{\text{mixture densities}}$$

two sets of random variables

$$Z = \text{cluster memberships (hidden variables)}$$

$$\Phi = \text{parameters } [\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_K, \hat{\Sigma}_1, \hat{\Sigma}_2, \dots, \hat{\Sigma}_K]$$

E-STEP :
$$E\left[L_c(\Phi | x, z) \,\middle|\, x, \hat{\Phi}^{(t)}\right]$$

M-STEP :
$$\hat{\Phi}^{(t+1)} = \arg\max_{\Phi} E\left[L_c(\Phi | x, z) \,\middle|\, X, \hat{\Phi}^{(t)}\right]$$

$$\underline{\text{E-STEP:}} \quad h_{ik} = E\left[z_{ik} \mid \mathcal{X}, \Phi^{(t)}\right] = \frac{p(x_i \mid C_k, \Phi^{(t)}) \cdot P(C_k)}{\sum\limits_{c=1}^{K} p(x_i \mid C_c, \Phi^{(t)}) \, P(C_c)}$$

<span style="color:red">multivariate Gaussians</span>

$$h_{ik} \geqslant 0, \quad \sum_{k=1}^{K} h_{ik} = 1 \quad \forall i$$

$$H = \begin{array}{c} \overbrace{\phantom{\rule{2cm}{0pt}}}^{K} \\ \left. \rule{0pt}{2cm} \right\} N \end{array}$$

each row sums up to 1.

$$\underline{\text{M-STEP:}} \quad \hat{P}^{(t+1)}(C_k) = \frac{\sum\limits_{i=1}^{N} h_{ik}}{N}$$

$$\hat{\mu}_k^{(t+1)} = \frac{\sum\limits_{i=1}^{N} h_{ik} \cdot x_i}{\sum\limits_{i=1}^{N} h_{ik}}$$

$$\hat{\Sigma}_k^{(t+1)} = \frac{\sum\limits_{i=1}^{N} h_{ik} \left(x_i - \hat{\mu}_k^{(t+1)}\right) \left(x_i - \hat{\mu}_k^{(t+1)}\right)^T}{\sum\limits_{i=1}^{N} h_{ik}}$$