



Machine learning

Parametric Models Part II Expectation-Maximization and Mixture Density Estimation

Mohammad-Reza A. Dehaqani

dehaqani@ut.ac.ir

Missing Features



- Suppose that we have a **Bayesian classifier** that uses the **feature vector** \mathbf{x} but a subset \mathbf{x}_g of \mathbf{x} are **observed** and the values for the remaining features \mathbf{x}_b are **missing**
- How can we make a decision?

$$g_i(\mathbf{x}) = p(\mathbf{x}|\omega_i)P(\omega_i)$$

 - Throw away the observations with missing values.
 - Or, substitute \mathbf{x}_b by their average $\overline{x_b}$ in the training data, and use $\mathbf{x} = (\mathbf{x}_g, \overline{x_b})$.
 - Or, **marginalize the posterior over the missing** features, and use the resulting posterior

$$P(w_i|\mathbf{x}_g) = \frac{\int P(w_i|\mathbf{x}_g, \mathbf{x}_b) p(\mathbf{x}_g, \mathbf{x}_b) d\mathbf{x}_b}{\int p(\mathbf{x}_g, \mathbf{x}_b) d\mathbf{x}_b}.$$

Marginal likelihood



- A **likelihood function** in which some **parameter** variables have been marginalized

$$p(\mathbf{X}|\alpha) = \int_{\theta} p(\mathbf{X}|\theta) p(\theta|\alpha) d\theta$$

θ has been **marginalized out** (integrated out)

$$\mathcal{L}(\psi; \mathbf{X}) = p(\mathbf{X}|\psi) = \int_{\lambda} p(\mathbf{X}|\lambda, \psi) p(\lambda|\psi) d\lambda$$

Expectation-Maximization



- Expectation–maximization (EM) algorithm is an **iterative** method to find **maximum likelihood** or **maximum a posteriori (MAP) estimates of parameters** in statistical models, from training data.
- The model depends on **unobserved latent** variables.
- It allows **learning of parameters** when some training patterns have **missing features (partial observation)**.

Applications of the EM algorithm



1. **Learning** when the **data is incomplete** or has missing values.
 2. **Optimizing** a likelihood (or posterior) function that is **analytically intractable** but can be **simplified** by assuming the existence of and values **for additional but missing** (or hidden) parameters.
- The **second** problem is more **common** in pattern recognition applications

General framework



- Assume that the **observed data \mathbf{X}** is generated by some distribution.
- Assume that a complete dataset $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ exists as a **combination** of the **observed but incomplete data \mathbf{X}** and the **missing data \mathbf{Y}** .
- The **observations in \mathbf{Z}** are assumed to be i.i.d. from the joint density

$$p(\mathbf{z}|\Theta) = p(\mathbf{x}, \mathbf{y}|\Theta) = p(\mathbf{y}|\mathbf{x}, \Theta)p(\mathbf{x}|\Theta).$$

New likelihood function



- We can define a new likelihood function

$$L(\Theta|\mathcal{Z}) = L(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta)$$

called the **complete-data likelihood** where $L(\theta|X)$ is referred to as the **incomplete-data likelihood**.

- The **EM** algorithm:
 - First, finds the **expected value** of the **complete-data log-likelihood** using the **current parameter** estimates (**expectation** step).
 - Then, **maximizes this expectation** (**maximization** step).
- Applying ML on the $E_{\mathbf{Y}|\mathbf{X}, \theta^t} \{\log[L(\theta|\mathbf{Z})]\}$
- Maximum likelihood with **partial observation**

Expected value of the $LL(\theta|Z)$ w.r.t. the unknown data Y ; **E-step**

- Define

$$Q(\Theta, \Theta^{(i-1)}) = E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta^{(i-1)}]$$

as the **expected value** of the **complete-data log-likelihood w.r.t.** the **unknown data Y given** the **observed data X** and the **current parameter estimates $\Theta^{(i-1)}$** ; $E_{\mathbf{Y}|\mathbf{X}, \Theta^{(i-1)}} \{\log[L(\theta|Z)]\}$

- The expected value can be computed as

$$E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta^{(i-1)}] = \int \log p(\mathcal{X}, \mathbf{y}|\Theta) p(\mathbf{y}|\mathcal{X}, \Theta^{(i-1)}) d\mathbf{y}.$$

y has been marginalized out

M-step



- Then, the expectation can be **maximized** by finding optimum values for the new parameters Θ as

$$\Theta^{(i)} = \arg \max_{\Theta} Q(\Theta, \Theta^{(i-1)}).$$

- These two steps are **repeated iteratively** where each iteration is **guaranteed to increase** the log-likelihood.
- The **EM algorithm** is also guaranteed to converge to **a local maximum** of the likelihood function.
- Q and $L(\theta)$ **behave similarly**; so we run optimization on Q
- Usually looking for **analytical solution** in M-step
- EM can be considered **as Quasi-static** solution for parameter estimation



Convergence properties of EM

- The solution depends on the **initial estimate** θ_0
- At each iteration, a value of θ is computed so that the likelihood function **does not decrease**.
- The algorithm is guaranteed to be **stable** (i.e., does not **oscillate**).
- There is **no guarantee** that it will converge to a **global** maximum.



- **Instead** of maximizing $Q(\Theta, \Theta^{(i-1)})$, the **Generalized Expectation-Maximization** algorithm finds some set of parameters $\Theta^{(i)}$ that satisfy

$$Q(\Theta^{(i)}, \Theta^{(i-1)}) > Q(\Theta, \Theta^{(i-1)})$$

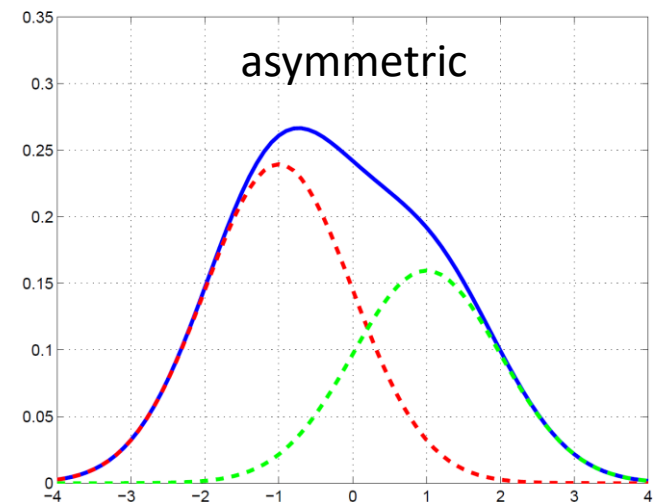
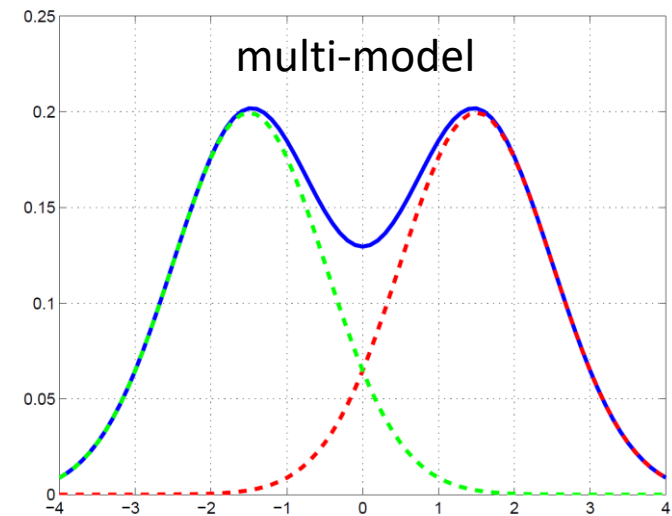
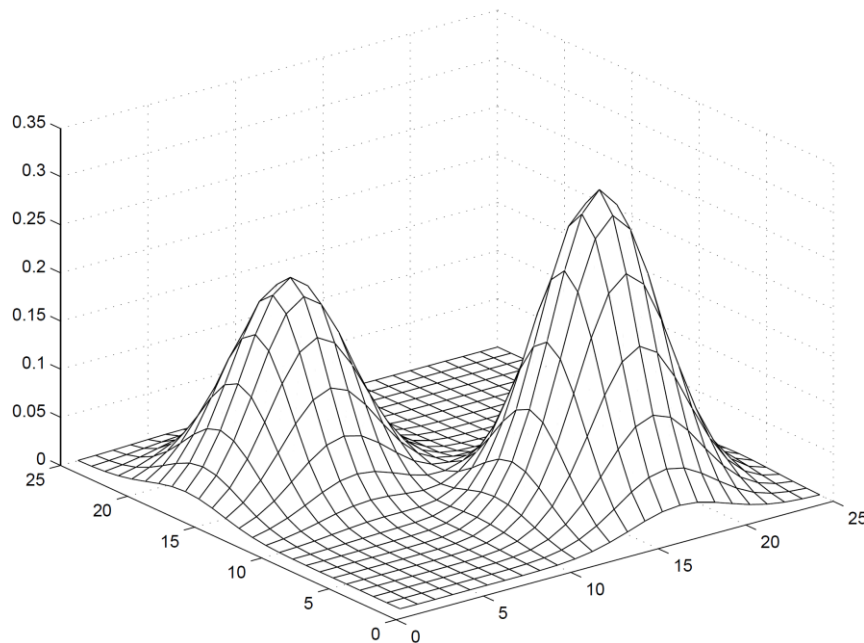
at **each iteration**.

- **Convergence** will not be as rapid as the EM algorithm but it allows **greater flexibility** to choose **computationally simpler steps**.

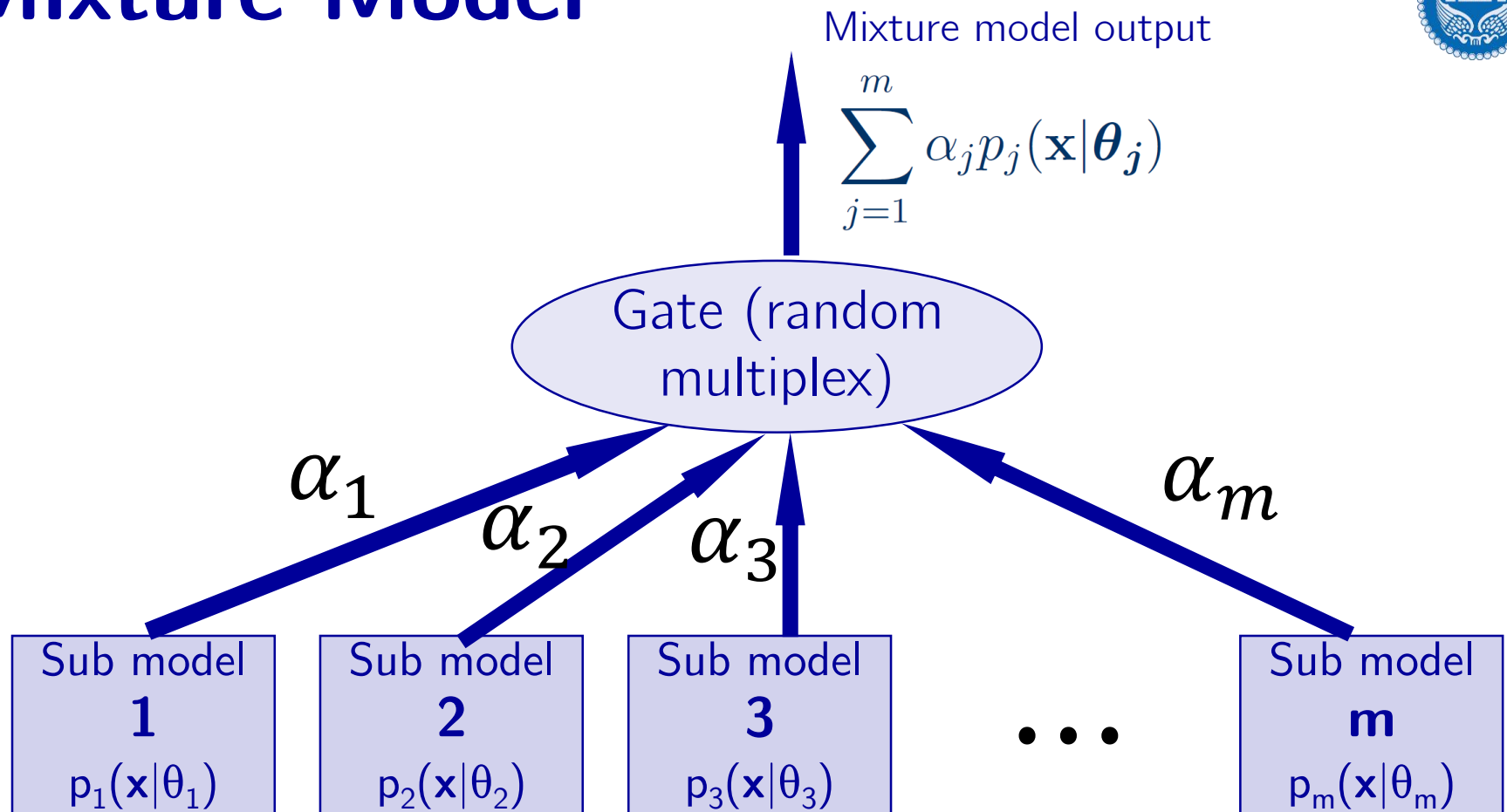
Gaussians do not well model ! multi-model and asymmetric data



The performance of a **generative model** is highly dependent on the accuracy of the class-conditional PDFs, $p(x|\omega)$



Mixture Model



$$p(\mathbf{x}|\boldsymbol{\Theta}) = \sum_{j=1}^m \alpha_j p_j(\mathbf{x}|\boldsymbol{\theta}_j)$$

Mixture Densities



- A mixture model is a **linear combination** of m densities

$$p(\mathbf{x}|\Theta) = \sum_{j=1}^m \alpha_j p_j(\mathbf{x}|\theta_j)$$

where $\Theta = (\alpha_1, \dots, \alpha_m, \theta_1, \dots, \theta_m)$ such that $\alpha_j \geq 0$ and $\sum_{j=1}^m \alpha_j = 1$.

- $\alpha_1, \dots, \alpha_m$ are called the **mixing parameters**.
- $p_j(\mathbf{x}|\theta_j)$, $j = 1, \dots, m$ are called the **component densities**

The log-likelihood of mixture density



- Suppose that $X = \{x_1, \dots, x_n\}$ is a set of observations i.i.d. with distribution $p(x|\Theta)$.
- The log-likelihood function of Θ becomes

$$\log L(\Theta|\mathcal{X}) = \log \prod_{i=1}^n p(\mathbf{x}_i|\Theta) =$$

- We **cannot** obtain an **analytical** solution for Θ by simply setting the **derivatives** of $\log L(\Theta|X)$ to zero because of the **logarithm of the sum**.

Mixture Density Estimation via EM



- Consider \mathbf{X} as incomplete and define **hidden variables** $\mathcal{Y} = \{y_i\}_{i=1}^n$ where y_i corresponds to **which mixture component** generated the data vector \mathbf{x}_i . ($y_i \in \{1, 2, \dots, m\}$)
- In other words, $y_i = j$ if the i 'th **data vector** was generated by the j 'th **mixture component**.
- Then, the **log-likelihood** becomes:

$$\log L(\Theta | \mathcal{X}, \mathcal{Y}) = \log p(\mathcal{X}, \mathcal{Y} | \Theta)$$

$$\begin{aligned} &= \sum_{i=1}^n \log(p(\mathbf{x}_i | y_i, \theta_i) p(y_i | \theta_i)) \\ &= \sum_{i=1}^n \log(\alpha_{y_i} p_{y_i}(\mathbf{x}_i | \theta_{y_i})). \end{aligned}$$

Initial parameters and latent variable distribution in EM Mixture Density

- Assume we have the **initial parameter** estimates

$$\Theta^{(g)} = (\alpha_1^{(g)}, \dots, \alpha_m^{(g)}, \boldsymbol{\theta}_1^{(g)}, \dots, \boldsymbol{\theta}_m^{(g)}).$$

- Compute

$$p(y_i | \mathbf{x}_i, \Theta^{(g)}) = \frac{\alpha_{y_i}^{(g)} p_{y_i}(\mathbf{x}_i | \boldsymbol{\theta}_{y_i}^{(g)})}{p(\mathbf{x}_i | \Theta^{(g)})} = \frac{\alpha_{y_i}^{(g)} p_{y_i}(\mathbf{x}_i | \boldsymbol{\theta}_{y_i}^{(g)})}{\sum_{j=1}^m \alpha_j^{(g)} p_j(\mathbf{x}_i | \boldsymbol{\theta}_j^{(g)})}$$

- and

$$p(\mathcal{Y} | \mathcal{X}, \Theta^{(g)}) = \prod_{i=1}^n p(y_i | \mathbf{x}_i, \Theta^{(g)}).$$

$Q(\Theta, \Theta^{(g)})$ in EM Mixture Density



$$\begin{aligned} Q(\Theta, \Theta^{(g)}) &= \sum_{\mathbf{y}} \log p(\mathcal{X}, \mathbf{y} | \Theta) p(\mathbf{y} | \mathcal{X}, \Theta^{(g)}) \\ &= \sum_{j=1}^m \sum_{i=1}^n \log(\alpha_j p_j(\mathbf{x}_i | \boldsymbol{\theta}_j)) p(j | \mathbf{x}_i, \Theta^{(g)}) \\ &= \sum_{j=1}^m \sum_{i=1}^n \log(\alpha_j) p(j | \mathbf{x}_i, \Theta^{(g)}) \\ &\quad + \sum_{j=1}^m \sum_{i=1}^n \log(p_j(\mathbf{x}_i | \boldsymbol{\theta}_j)) p(j | \mathbf{x}_i, \Theta^{(g)}). \end{aligned}$$

Mixture Density Estimation via EM



- We can maximize the two sets of summations for α_j and θ_j **independently** because they are **not related**.
- The estimate for α_j can be computed as

lagrange multiplier method

$$\hat{\alpha}_j = \frac{1}{n} \sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)})$$

- where

$$p(j|\mathbf{x}_i, \Theta^{(g)}) = \frac{\alpha_j^{(g)} p_j(\mathbf{x}_i|\boldsymbol{\theta}_j^{(g)})}{\sum_{t=1}^m \alpha_t^{(g)} p_t(\mathbf{x}_i|\boldsymbol{\theta}_t^{(g)})}.$$

It is a number and completely describe by $\theta^{(g)}$

Mixture of Gaussians



- We can obtain analytical expressions for θ_j for the special case of a **Gaussian mixture** where $\theta_j = (\mu_j, \Sigma_j)$ and

$$\begin{aligned} p_j(\mathbf{x}|\boldsymbol{\theta}_j) &= p_j(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \\ &= \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_j|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j) \right]. \end{aligned}$$

- Equating the **partial derivative** of $Q(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{(g)})$ with respect to $\boldsymbol{\mu}_j$ to zero gives

$$\begin{aligned} &+ \sum_{j=1}^m \sum_{i=1}^n \log(p_j(\mathbf{x}_i|\boldsymbol{\theta}_j)) p(j|\mathbf{x}_i, \boldsymbol{\Theta}^{(g)}). \\ \hat{\boldsymbol{\mu}}_j &= \frac{\sum_{i=1}^n p(j|\mathbf{x}_i, \boldsymbol{\Theta}^{(g)}) \mathbf{x}_i}{\sum_{i=1}^n p(j|\mathbf{x}_i, \boldsymbol{\Theta}^{(g)})}. \end{aligned}$$

Mixture of Gaussian; Σ estimation



$$\Sigma_j = \sigma^2 \mathbf{I}$$

$$\hat{\sigma}^2 = \frac{1}{nd} \sum_{j=1}^m \sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)}) \|\mathbf{x}_i - \hat{\mu}_j\|^2$$

$$\Sigma_j = \sigma_j^2 \mathbf{I}$$

$$\hat{\sigma}_j^2 = \frac{\sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)}) \|\mathbf{x}_i - \hat{\mu}_j\|^2}{d \sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)})}$$

Mixture of Gaussian; Σ estimation



$$\Sigma_j = \text{diag}(\{\sigma_{jk}^2\}_{k=1}^d)$$

$$\hat{\sigma}_{jk}^2 = \frac{\sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)}) (\mathbf{x}_{ik} - \hat{\mu}_{jk})^2}{\sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)})}$$

$$\Sigma_j = \Sigma$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^m \sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)}) (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T$$

Mixture of Gaussian; general case



$\Sigma_j = \text{arbitrary}$

$$\hat{\Sigma}_j = \frac{\sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)}) (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T}{\sum_{i=1}^n p(j|\mathbf{x}_i, \Theta^{(g)})}$$

Mixture of Gaussians



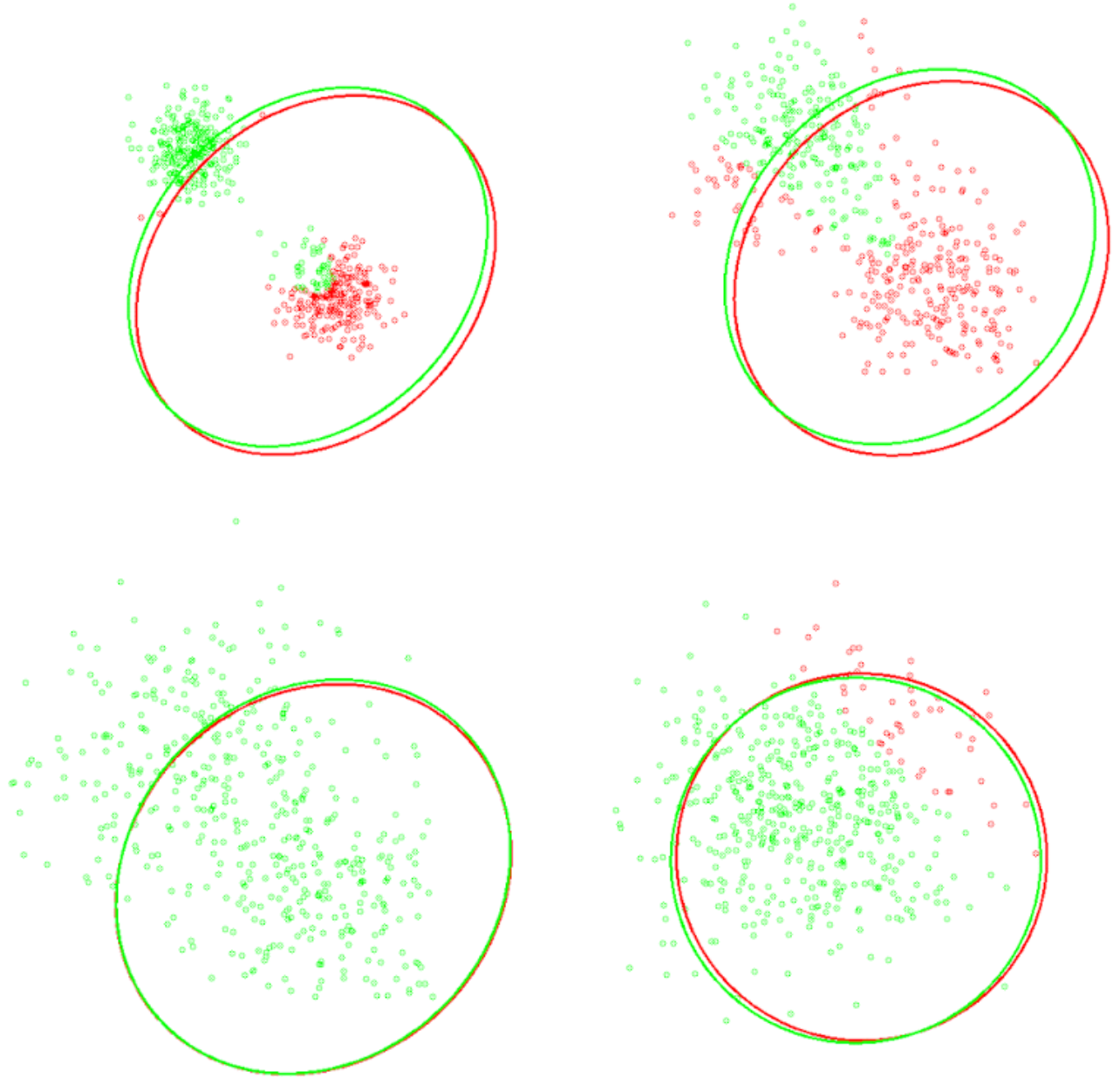
- Estimates for α_j , μ_j and Σ_j perform **both expectation and maximization steps simultaneously**.
- EM iterations proceed by using the **current estimates** as the initial estimates for the next iteration.
- The **priors are** computed from the proportion of examples belonging to each mixture component.
- The **means are** the component centroids.
- The **covariance matrices** are calculated as the sample covariance of the points associated with each component

Mixture of Gaussians; an iterative algorithm

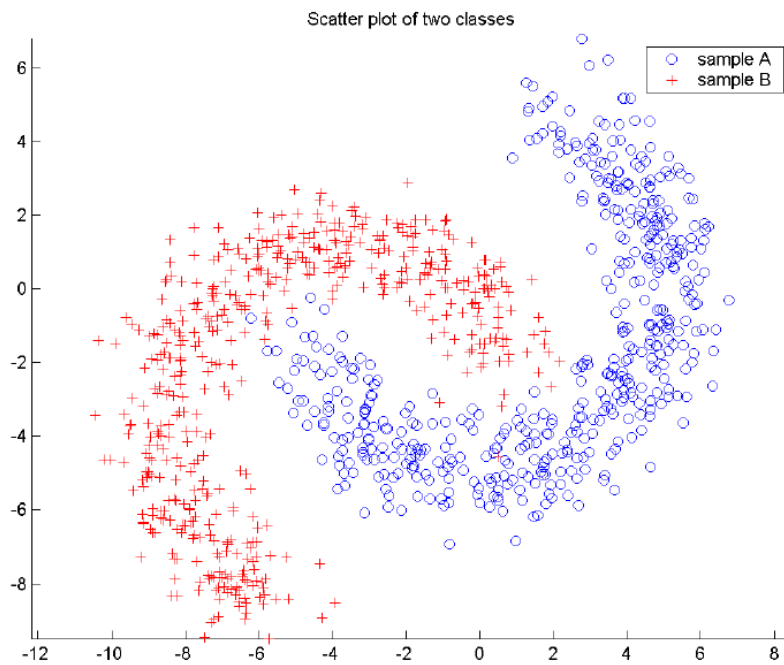


- The **number** of components in the mixture?
- The **initial** estimates for Θ ?
- When to **stop** the iterations?
 - Stop if the **change** in log-likelihood between two iterations is less **than a threshold**.
 - Or, use a threshold for the **number of iterations**

Example



2-D Bayesian classification examples where the data for each class come from a banana shaped distribution



(a) Scatter plot.

