



Signal Processing

Lecture 15: Expectation Maximization & Mixture Models

Konstantinos Chatzilygeroudis - costashatz@upatras.gr

Department of Electrical and Computer Engineering
University of Patras

Template made by Panagiotis Papagiannopoulos



Motivation: Why Latent-Variable Models?

Many signal models involve **hidden / latent variables** that we do not observe directly.

Example (mixture / multi-regime signals):

y_i is generated by one of K sources or regimes $\Rightarrow z_i \in \{1, \dots, K\}$.

Complete-data model:

$$p(y_i, z_i \mid \theta) = p(z_i \mid \theta) p(y_i \mid z_i, \theta).$$

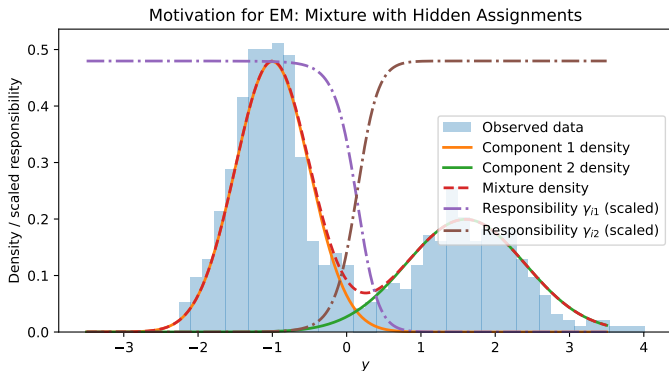
But we only observe y_i (N samples):

$$p(\mathbf{y} \mid \theta) = \prod_{i=1}^N \sum_{k=1}^K p(z_i = k \mid \theta) p(y_i \mid z_i = k, \theta).$$

Key difficulty:

$$\log p(\mathbf{y} \mid \theta) = \sum_i \log \sum_k (\cdot) \quad (\text{log of sum is hard to optimize/compute}).$$

Motivation: Why Latent-Variable Models? (2)



Expectation Maximization idea: alternate between

- **E-step:** infer soft assignments $p(z_i = k \mid y_i, \theta)$,
- **M-step:** update θ using these assignments.

Latent-Variable Models

We assume that the data is generated with the help of **unobserved (latent) variables**.

Observed data:

$$\mathbf{Y} = \{y_i\}_{i=1}^N$$

Latent variables:

$$\mathbf{Z} = \{z_i\}_{i=1}^N \quad (\text{e.g., cluster label, regime index, source ID}).$$

Generative structure:

$$p(\mathbf{Y}, \mathbf{Z} \mid \theta) = p(\mathbf{Z} \mid \theta) p(\mathbf{Y} \mid \mathbf{Z}, \theta).$$

Examples in signal processing

- Mixture / multi-regime signals: z_i = which regime generated y_i .
- Clustering features: z_i = cluster assignment.
- Speech / audio: z_i = phoneme/state producing a frame.

Complete-Data vs Incomplete-Data Likelihood

Complete-data likelihood (if we knew \mathbf{Z}):

$$p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta}) = \prod_{i=1}^N p(z_i \mid \boldsymbol{\theta}) p(y_i \mid z_i, \boldsymbol{\theta}).$$

Incomplete-data likelihood (what we actually have):

$$p(\mathbf{Y} \mid \boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta}) \quad (\text{or } \int d\mathbf{Z} \text{ if continuous}).$$

Key point:

Latent variables \mathbf{Z} are missing \Rightarrow we must marginalize them out.

Goal: estimate parameters by maximizing the incomplete-data likelihood

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{Y} \mid \boldsymbol{\theta}) \quad \Leftrightarrow \quad \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{Y} \mid \boldsymbol{\theta}).$$

Why the Likelihood is Hard

The log-likelihood for incomplete data becomes

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}) = \log \sum_{\mathbf{Z}} p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta}).$$

Problem: log of a sum

$$\log \sum_{\mathbf{Z}} (\cdot) \neq \sum_{\mathbf{Z}} \log(\cdot).$$

Consequences

- The objective is often **non-convex**.
- Direct maximization is typically **intractable**.
- Gradients involve ratios of sums of exponentials (unstable / messy).

We need a trick to handle the hidden \mathbf{Z} .

A Concrete Example: Mixture Likelihood

Suppose each y_i comes from one of K components:

$$p(z_i = k) = \pi_k, \quad p(y_i \mid z_i = k, \boldsymbol{\theta}) = p_k(y_i \mid \boldsymbol{\theta}_k).$$

Marginal likelihood per sample:

$$p(y_i \mid \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_k(y_i \mid \boldsymbol{\theta}_k).$$

Dataset log-likelihood:

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k p_k(y_i \mid \boldsymbol{\theta}_k) \right).$$

Expectation-Maximization (EM): Core Idea

We want to maximize the incomplete-data log-likelihood:

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}) = \log \sum_{\mathbf{Z}} p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta}).$$

Difficulty: the hidden \mathbf{Z} is inside a log-sum.

EM key idea:

- Introduce a tractable lower bound on $\log p(\mathbf{Y} \mid \boldsymbol{\theta})$.
- Alternate between:
 - estimating latent variables (E-step),
 - updating parameters (M-step).

EM Derivation I: Introduce an Auxiliary Distribution

Let $q(\mathbf{Z})$ be *any* distribution over latent variables.

Rewrite the likelihood:

$$\log p(\mathbf{Y} | \boldsymbol{\theta}) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{Y}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})}.$$

Interpretation:

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) (\cdot) = \mathbb{E}_q[(\cdot)].$$

So

$$\log p(\mathbf{Y} | \boldsymbol{\theta}) = \log \mathbb{E}_q \left[\frac{p(\mathbf{Y}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} \right].$$

EM Derivation I: Introduce an Auxiliary Distribution

Let $q(\mathbf{Z})$ be *any* distribution over latent variables.

Rewrite the likelihood:

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta})}{q(\mathbf{Z})}.$$

Interpretation:

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) (\cdot) = \mathbb{E}_q[(\cdot)].$$

So

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}) = \log \mathbb{E}_q \left[\frac{p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta})}{q(\mathbf{Z})} \right].$$

Jensen's Inequality:

- **Convex case (Jensen):** If φ is convex and X integrable, then

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)].$$

- **Concave case:** If φ is concave and X integrable, then the inequality reverses:

$$\varphi(\mathbb{E}[X]) \geq \mathbb{E}[\varphi(X)],$$

with equality iff X is a.s. constant or φ is affine on the support of X .

EM Derivation II: Jensen's Inequality (Lower Bound)

Apply Jensen's inequality (since $\log(\cdot)$ is concave):

$$\log \mathbb{E}_q[f(\mathbf{Z})] \geq \mathbb{E}_q[\log f(\mathbf{Z})].$$

Here, $f(\mathbf{Z}) = \frac{p(\mathbf{Y}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}$. Therefore:

$$\log p(\mathbf{Y} | \theta) \geq \mathbb{E}_q \left[\log \frac{p(\mathbf{Y}, \mathbf{Z} | \theta)}{q(\mathbf{Z})} \right] = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{Y}, \mathbf{Z} | \theta)}{q(\mathbf{Z})}.$$

Define the lower bound:

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{Y}, \mathbf{Z} | \theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z}).$$

Two terms:

- Expected complete-data log-likelihood.
- Entropy of $q(\mathbf{Z})$.

EM Derivation III: Tightness via KL Divergence

We can prove that:

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z} \mid \mathbf{Y}, \boldsymbol{\theta})).$$

Since $\text{KL}(\cdot \parallel \cdot) \geq 0$:

$$\mathcal{L}(q, \boldsymbol{\theta}) \leq \log p(\mathbf{Y} \mid \boldsymbol{\theta}).$$

Bound becomes tight when

$$q(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{Y}, \boldsymbol{\theta}).$$

KL divergence: measures how different q is from p (i.e., expected information loss using q instead of p).

$$\text{KL}(q \parallel p) = \mathbb{E}_q \left[\log \frac{q(\mathbf{Z})}{p(\mathbf{Z})} \right] = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{q(\mathbf{Z})}{p(\mathbf{Z})} \geq 0,$$

with equality iff $q = p$ (a.e.).

Expectation-Maximization Algorithm

EM performs coordinate ascent on $\mathcal{L}(q, \theta)$:

E-step (update latent distribution):

$$q^{(t+1)}(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{Y}, \theta^{(t)}).$$

Tightens the bound by minimizing KL: sets q to the posterior given our current parameters, $\theta^{(t)}$, and data \mathbf{Y} .

M-step (update parameters):

$$\theta^{(t+1)} = \arg \max_{\theta} \mathbb{E}_{q^{(t+1)}} [\log p(\mathbf{Y}, \mathbf{Z} \mid \theta)].$$

Maximizes the bound w.r.t. parameters.

Monotonic improvement

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}^{(t+1)}) \geq \log p(\mathbf{Y} \mid \boldsymbol{\theta}^{(t)}).$$

What EM guarantees

- Each iteration *does not decrease* the incomplete likelihood.
- Converges to a **local optimum** (not necessarily global).

Practical notes

- Initialization matters (often use k-means / random restarts).
- Stop when likelihood improvement is small.
- Non-convexity \Rightarrow multiple runs recommended.

Gaussian Mixture Models (GMM): Model

We model data $\{\mathbf{y}_i\}_{i=1}^N$, $\mathbf{y}_i \in \mathbb{R}^d$ as coming from a mixture of K Gaussians.

Latent assignment:

$$z_i \in \{1, \dots, K\}.$$

Mixing weights:

$$p(z_i = k \mid \boldsymbol{\theta}) = \pi_k, \quad \pi_k \geq 0, \quad \sum_{k=1}^K \pi_k = 1.$$

Component densities:

$$p(\mathbf{y}_i \mid z_i = k, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Parameters:

$$\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K.$$

GMM: Incomplete (Observed) Likelihood

Marginal likelihood per sample:

$$p(\mathbf{y}_i | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Incomplete-data log-likelihood:

$$\log p(\mathbf{Y} | \boldsymbol{\theta}) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right).$$

Hard part: log-sum prevents direct optimization \Rightarrow use EM.

GMM: Complete-Data Likelihood

If the assignments z_i were known:

$$p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta}) = \prod_{i=1}^N \prod_{k=1}^K \left(\pi_k \mathcal{N}(\mathbf{y}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)^{\mathbb{I}[z_i=k]}.$$

Complete-data log-likelihood:

$$\log p(\mathbf{Y}, \mathbf{Z} \mid \boldsymbol{\theta}) = \sum_{i=1}^N \sum_{k=1}^K \mathbb{I}[z_i = k] \left(\log \pi_k + \log \mathcal{N}(\mathbf{y}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right).$$

EM replaces $\mathbb{I}[z_i = k]$ by its posterior expectation, where $\mathbb{I}[z_i = k]$ is the **indicator function**:

$$\mathbb{I}[z_i = k] = \begin{cases} 1, & \text{if } z_i = k, \\ 0, & \text{otherwise.} \end{cases}$$

E-step: Responsibilities

E-step: compute posterior probability that sample i belongs to component k :

$$\gamma_{ik} \equiv p(z_i = k \mid \mathbf{y}_i, \boldsymbol{\theta}^{(t)}) = \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{y}_i \mid \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}{\sum_{j=1}^K \pi_j^{(t)} \mathcal{N}(\mathbf{y}_i \mid \boldsymbol{\mu}_j^{(t)}, \boldsymbol{\Sigma}_j^{(t)})}.$$

Interpretation:

- $\gamma_{ik} \in [0, 1]$, and $\sum_k \gamma_{ik} = 1$.
- “Soft” cluster assignment (not hard labels).

M-step: Update Mixing Weights and Means

Let the effective number of points in cluster k be:

$$N_k = \sum_{i=1}^N \gamma_{ik}.$$

Mixing weights:

$$\pi_k^{(t+1)} = \frac{N_k}{N}.$$

Means:

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_{ik} \mathbf{y}_i.$$

Interpretation: re-weighted averages using responsibilities.

M-step: Update Covariances

Covariances:

$$\mathbf{\Sigma}_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_{ik} (\mathbf{y}_i - \boldsymbol{\mu}_k^{(t+1)})(\mathbf{y}_i - \boldsymbol{\mu}_k^{(t+1)})^\top.$$

Interpretation: responsibility-weighted covariance within each component.

Numerical note: often regularize to avoid singularities:

$$\mathbf{\Sigma}_k \leftarrow \mathbf{\Sigma}_k + \epsilon \mathbf{I}.$$

GMM-EM Algorithm (Summary)

Initialize $\{\pi_k^{(0)}, \boldsymbol{\mu}_k^{(0)}, \boldsymbol{\Sigma}_k^{(0)}\}$.

Repeat until convergence:

- **E-step:** compute responsibilities γ_{ik} .
- **M-step:** update $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ using γ_{ik} .

Convergence check:

$$\log p(\mathbf{Y} \mid \boldsymbol{\theta}^{(t+1)}) - \log p(\mathbf{Y} \mid \boldsymbol{\theta}^{(t)}) < \delta.$$

Each iteration increases (or leaves unchanged) the incomplete likelihood.

Initialization matters

- Initial centers, $\mu_k^{(0)}$, are very important!.
- Multiple random restarts help avoid bad local maxima.

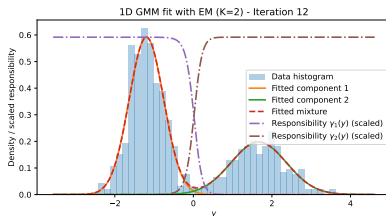
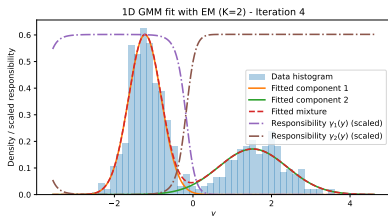
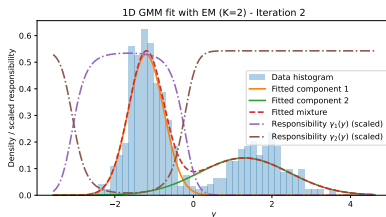
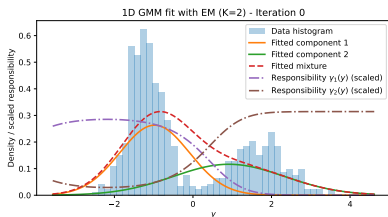
Choosing K

- BIC / AIC / cross-validation.
- Too small K : underfit; too large K : overfit.

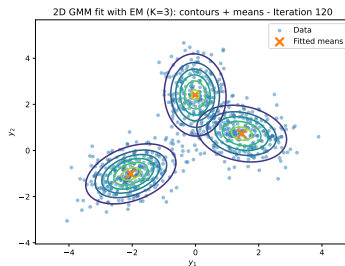
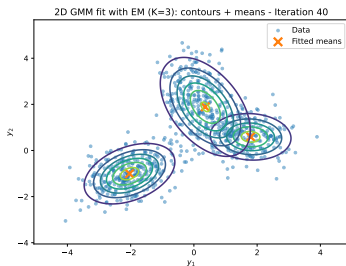
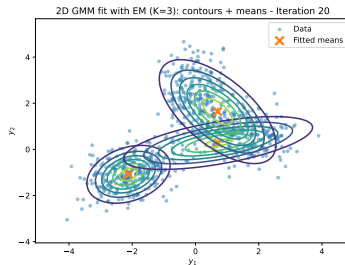
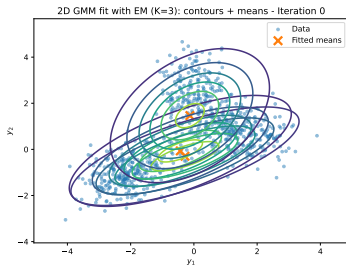
Covariance structure

- Full Σ_k : flexible but more parameters.
- Diagonal / spherical: faster, more stable.

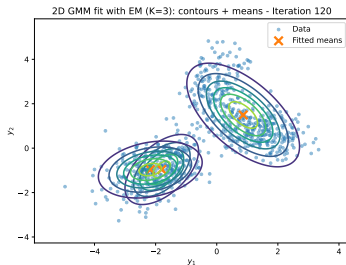
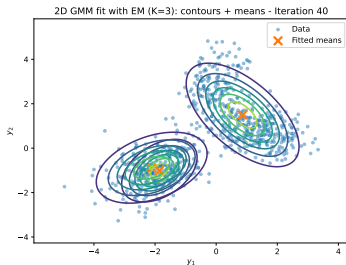
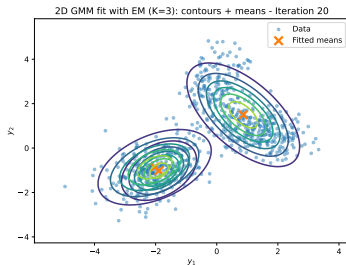
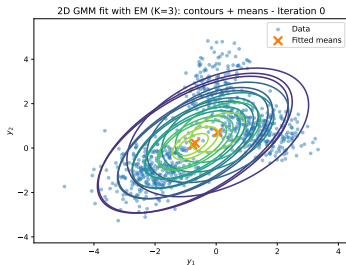
GMMs: 1D Example



GMMs: 2D Example



GMMs: 2D Example - Bad Initialization



Thank you

- **Any Questions?**
- **Office Hours:**
 - **Tue & Thu (09:00-11:00)**
 - 24/7 by email (costashatz@upatras.gr, subject: *ECE_SP_AM*)
- **Material and Announcements**



Laboratory of Automation & Robotics