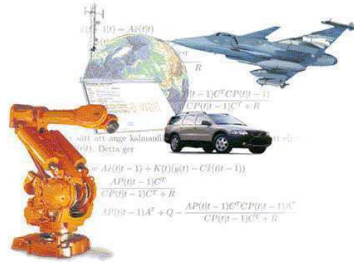


Machine Learning, Lecture 6 Expectation Maximization (EM) and clustering



Thomas Schön

Division of Automatic Control
Linköping University
Linköping, Sweden.

Email: schon@isy.liu.se,
Phone: 013 - 281373,
Office: House B, Entrance 27.

Outline lecture 6

2(35)

1. Summary of lecture 5
2. Expectation Maximization (EM)
 - General derivation
 - Example - identification of a linear state-space model
 - Example - identification of a Wiener system
3. Gaussian mixtures
 - Standard construction
 - Equivalent construction using latent variables
 - ML estimation using EM
4. Connections to the K -means algorithm for clustering

(Chapter 9)

The Gaussian mixture model and the K -means algorithm will be finished during lecture 7 on Wednesday.

Summary of lecture 5

3(35)

A **Gaussian process** is a collection of random variables, any finite number of which have a joint Gaussian distribution.

By assuming that the considered system is a Gaussian process, predictions can be made by computing the conditional distribution $p(y(x^*) | \text{all the observations})$, $y(x^*)$ being the output for which we seek a prediction. This regression approach is referred to as **Gaussian process regression**.

The **support vector machine** (SVM) is a discriminative classifier that gives the maximum margin decision boundary.

Latent variables – example

4(35)

A **latent variable** is a variable that is not directly observed. Other common names are hidden variables, unobserved variables or missing data.

An example of a latent variable is the state x_t in a state space model.

Consider the following linear Gaussian state space (LGSS) model

$$\begin{aligned} x_{t+1} &= \theta x_t + v_t, \\ y_t &= \frac{1}{2} x_t + e_t, \end{aligned} \quad \begin{pmatrix} v_t \\ e_t \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix} \right).$$

The **Expectation Maximization (EM)** algorithm computes ML estimates of unknown parameters in probabilistic models involving latent variables.

Strategy: Use *structure* inherent in the probabilistic model to separate the original ML problem into *two closely linked subproblems*, each of which is hopefully in some sense more tractable than the original problem.

EM focus on the joint log-likelihood function of the observed variables X and the latent variables $Z \triangleq \{z_1, \dots, z_N\}$,

$$\ell_{\theta}(X, Z) = \ln p_{\theta}(X, Z).$$

Algorithm 1 Expectation Maximization (EM)

1. **Initialise:** Set $i = 1$ and choose an initial θ_1 .
2. **While** not converged **do:**

- (a) **Expectation (E) step:** Compute

$$\begin{aligned} Q(\theta, \theta_i) &= E_{\theta_i} [\ln p_{\theta}(Z, X | X)] \\ &= \int \ln p_{\theta}(Z, X) p_{\theta_i}(Z | X) dZ \end{aligned}$$

- (b) **Maximization (M) step:** Compute

$$\theta_{i+1} = \arg \max_{\theta} Q(\theta, \theta_i)$$

- (c) $i \leftarrow i + 1$

Consider the following scalar LGSS model

$$\begin{aligned} x_{t+1} &= \theta x_t + v_t, \\ y_t &= \frac{1}{2} x_t + e_t, \end{aligned} \quad \begin{pmatrix} v_t \\ e_t \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix} \right).$$

The initial state is fully known ($x_1 = 0$) and the true θ -parameter is given by $\theta^* = 0.9$.

The identification problem is now to determine the parameter θ on the basis of the observations $Y = \{y_1, \dots, y_N\}$, using the EM algorithm.

The latent variables Z are given by the states
 $Z = X \triangleq \{x_1, \dots, x_{N+1}\}$.

Note the difference in notation compared to Bishop! The observations are denoted Y and the latent variables are denoted by

The expectation (E) step:

$$Q(\theta, \theta_i) \triangleq \mathbf{E}_{\theta_i} \{ \ln p_{\theta}(X, Y) | Y \} = \int \ln p_{\theta}(X, Y) p_{\theta_i}(X | Y) dX.$$

Let us start investigating $\ln p_{\theta}(X, Y)$. Using conditional probabilities we have,

$$\begin{aligned} p_{\theta}(X, Y) &= p_{\theta}(x_{N+1}, X_N, y_N, Y_{N-1}) \\ &= p_{\theta}(x_{N+1}, y_N | X_N, Y_{N-1}) p_{\theta}(X_N, Y_{N-1}), \end{aligned}$$

According to the Markov property we have

$$p_{\theta}(x_{N+1}, y_N | X_N, Y_{N-1}) = p_{\theta}(x_{N+1}, y_N | x_N),$$

resulting in

$$p_{\theta}(X, Y) = p_{\theta}(x_{N+1}, y_N | x_N) p_{\theta}(X_N, Y_{N-1}).$$

Repeated use of the above ideas straightforwardly yields

$$p_{\theta}(X, Y) = p_{\theta}(x_1) \prod_{t=1}^N p_{\theta}(x_{t+1}, y_t | x_t).$$

According to the model, we have

$$p_{\theta} \left(\begin{pmatrix} x_{t+1} \\ y_t \end{pmatrix} | x_t \right) = \mathcal{N} \left(\begin{pmatrix} x_{t+1} \\ y_t \end{pmatrix}; \begin{pmatrix} \theta \\ 1/2 \end{pmatrix} x_t, \begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix} \right).$$

The resulting Q -function is

$$\begin{aligned} Q(\theta, \theta_i) &\propto -\mathbf{E}_{\theta_i} \left\{ \sum_{t=1}^N x_t^2 | Y \right\} \theta^2 + 2\mathbf{E}_{\theta_i} \left\{ \sum_{t=1}^N x_t x_{t+1} | Y \right\} \theta \\ &= -\varphi \theta^2 + 2\psi \theta, \end{aligned}$$

where we have defined

$$\varphi \triangleq \sum_{t=1}^N \mathbf{E}_{\theta_i} \{ x_t^2 | Y \}, \quad \psi \triangleq \sum_{t=1}^N \mathbf{E}_{\theta_i} \{ x_t x_{t+1} | Y \}.$$

There exist explicit expressions for these expected values.

The maximization (M) step,

$$\theta_{i+1} = \arg \max_{\theta} Q(\theta, \theta_i).$$

simply amounts to solving the following quadratic problem,

$$\theta_{i+1} = \arg \max_{\theta} -\varphi \theta^2 + 2\psi \theta.$$

The solution is given by

$$\theta_{i+1} = \frac{\psi}{\varphi}.$$

Algorithm 2 EM – example 1

1. **Initialise:** Set $i = 1$ and choose an initial θ_1 .
2. **While** not converged **do:**

(a) **Expectation (E) step:** Compute

$$\varphi = \sum_{t=1}^N \mathbf{E}_{\theta_i} \{ x_t^2 | Y \}, \quad \psi = \sum_{t=1}^N \mathbf{E}_{\theta_i} \{ x_t x_{t+1} | Y \}.$$

(b) **Maximization (M) step:** Find the next iterate according to

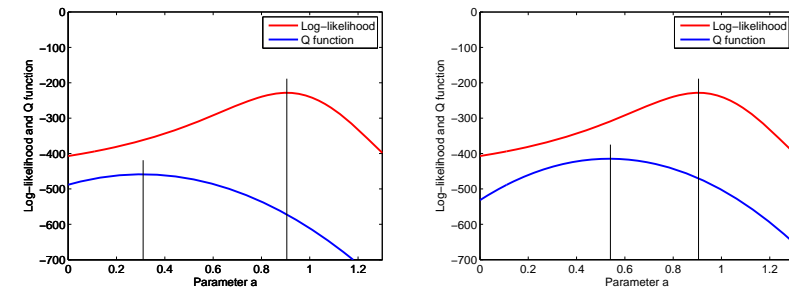
$$\theta_{i+1} = \frac{\psi}{\varphi}.$$

- (c) If $|L_{\theta_i}(Y) - L_{\theta_{i-1}}(Y)| \geq 10^{-6}$, update $i := i + 1$ and return to step 2, otherwise terminate.
-

- Different number of samples N used.
- Monte Carlo studies, each using 1000 realisations of data.
- Initial guess $\theta_0 = 0.1$.

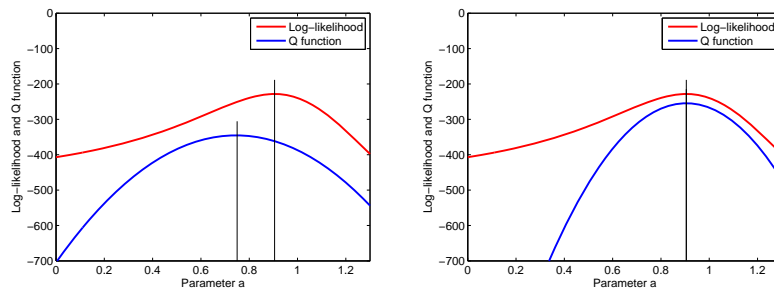
N	100	200	500	1000	2000	5000	10000
$\hat{\theta}$	0.8716	0.8852	0.8952	0.8978	0.8988	0.8996	0.8998

No surprise, since ML is asymptotically efficient.



(a) Iteration 1

(b) Iteration 2



(c) Iteration 3

(d) Iteration 11

All details (including MATLAB code) are provided in

Thomas B. Schön, **An Explanation of the Expectation Maximization Algorithm**. Division of Automatic Control, Linköping University, Sweden, *Technical Report nr: LITH-ISY-R-2915*, August 2009.

users.isy.liu.se/rt/schon/Publications/SchonEM2009.pdf



A general state space model (SSM) consists of a Markov process $\{x_t\}_{t \geq 1}$ and a measurement process $\{y_t\}_{t \geq 1}$, related according to

$$\begin{aligned} x_{t+1} | x_t &\sim f_{\theta,t}(x_{t+1} | x_t, u_t), \\ y_t | x_t &\sim h_{\theta,t}(y_t | x_t, u_t), \\ x_1 &\sim \mu_{\theta}(x_1). \end{aligned}$$

Identification problem: Find θ based on $\{u_{1:T}, y_{1:T}\}$.

According to the above, the first step is to compute the Q -function

$$Q(\theta, \hat{\theta}_k) = \mathbf{E}_{\theta_k} \{ \ln p_{\theta}(Z, Y) | Y \}$$



Applying $\mathbf{E}_{\theta_k}\{\cdot | Y\}$ to

$$\begin{aligned} \ln p_{\theta}(X, Y) &= \ln p_{\theta}(Y | X) + \ln p_{\theta}(X) \\ &= \ln p_{\theta}(x_1) + \sum_{t=1}^{N-1} \ln p_{\theta}(x_{t+1} | x_t) + \sum_{t=1}^N \ln p_{\theta}(y_t | x_t). \end{aligned}$$

This results in $Q(\theta, \theta_k) = I_1 + I_2 + I_3$, where

$$\begin{aligned} I_1 &= \int \ln p_{\theta}(x_1) p_{\theta_k}(x_1 | Y) dx_1, \\ I_2 &= \sum_{t=1}^{N-1} \int \int \ln p_{\theta}(x_{t+1} | x_t) p_{\theta_k}(x_{t+1}, x_t | Y) dx_t dx_{t+1}, \\ I_3 &= \sum_{t=1}^N \int \ln p_{\theta}(y_t | x_t) p_{\theta_k}(x_t | Y) dx_t. \end{aligned}$$

This leads us to a nonlinear state smoothing problem, which we can solve using a particle smoother (PS).

The PS provides us with the following approximation of the joint smoothing density

$$p(X | Y) \approx \frac{1}{M} \sum_{t=1}^M \delta(X - X^i),$$

which allows for the following approximations of the marginal smoothing densities that we need,

$$\begin{aligned} p_{\theta_k}(x_t | Y) &\approx \hat{p}_{\theta_k}(x_t | Y) = \frac{1}{M} \sum_{i=1}^M \delta(x_t - x_t^i), \\ p_{\theta_k}(x_{t:t+1} | Y) &\approx \hat{p}_{\theta_k}(x_{t:t+1} | Y) = \frac{1}{M} \sum_{i=1}^M \delta(x_{t:t+1} - x_{t:t+1}^i). \end{aligned}$$

Inserting the above approximations into the integrals straightforwardly yields the approximation we are looking for,

$$\begin{aligned} \hat{I}_1 &= \int \ln p_{\theta}(x_1) \sum_{i=1}^M \frac{1}{M} \delta(x_1 - x_1^i) dx_1 = \frac{1}{M} \sum_{i=1}^M \ln p_{\theta}(x_1^i), \\ \hat{I}_2 &= \sum_{t=1}^{N-1} \int \int \ln p_{\theta}(x_{t+1} | x_t) \sum_{i=1}^M \frac{1}{M} \delta(x_{t:t+1} - x_{t:t+1}^i) dx_{t:t+1} \\ &= \frac{1}{M} \sum_{t=1}^{N-1} \sum_{i=1}^M \ln p_{\theta}(x_{t+1}^i | x_t^i), \\ \hat{I}_3 &= \sum_{t=1}^N \int \ln p_{\theta}(y_t | x_t) \sum_{i=1}^M \frac{1}{M} \delta(x_t - x_t^i) dx_t = \frac{1}{M} \sum_{t=1}^N \sum_{i=1}^M \ln p_{\theta}(y_t | x_t^i) \end{aligned}$$

It is straightforward to make use of the approximation of the Q -function just derived in order to compute gradients of the Q -function,

$$\frac{\partial}{\partial \theta} \hat{Q}(\theta, \theta_k) = \frac{\partial \hat{I}_1}{\partial \theta} + \frac{\partial \hat{I}_2}{\partial \theta} + \frac{\partial \hat{I}_3}{\partial \theta}$$

For example (the other two terms are treated analogously),

$$\begin{aligned} \hat{I}_3 &= \frac{1}{M} \sum_{t=1}^N \sum_{i=1}^M \ln p_{\theta}(y_t | x_t^i), \\ \frac{\partial \hat{I}_3}{\partial \theta} &= \frac{1}{M} \sum_{t=1}^N \sum_{i=1}^M \frac{\partial \ln p_{\theta}(y_t | x_t^i)}{\partial \theta} \end{aligned}$$

With these gradients in place there are many algorithms that can be used in order to solve the maximization problem, we employ BFGS.

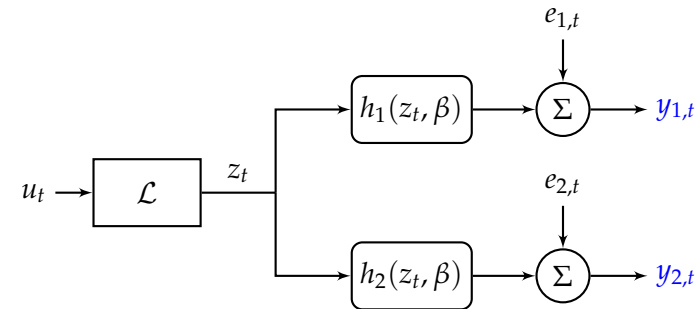
Algorithm 3 Nonlinear System Identification Using EM

1. **Initialise:** Set $i = 1$ and choose an initial θ_1 .
2. **While** not converged **do:**
 - (a) **Expectation (E) step:** Run a FFBS PS and compute

$$\hat{Q}(\theta, \theta_k) = \hat{I}_1(\theta, \theta_k) + \hat{I}_2(\theta, \theta_k) + \hat{I}_3(\theta, \theta_k)$$

- (b) **Maximization (M) step:** Compute $\theta_{k+1} = \arg \max_{\theta} \hat{Q}(\theta, \theta_k)$ using an off-the-shelf numerical optimization algorithm.
- (c) $k \leftarrow k + 1$

Thomas B. Schön, Adrian Wills and Brett Ninness. *System Identification of Nonlinear State-Space Models. Automatica*, 47(1):39-49, January 2011.

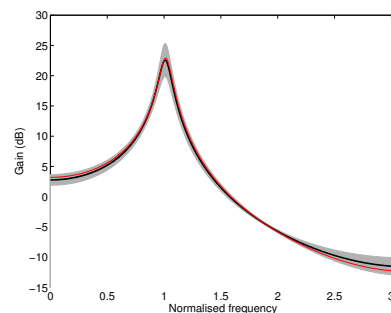


$$x_{t+1} = \begin{pmatrix} A & B \end{pmatrix} \begin{pmatrix} x_t \\ u_t \end{pmatrix}, \quad u_t \sim \mathcal{N}(0, Q),$$

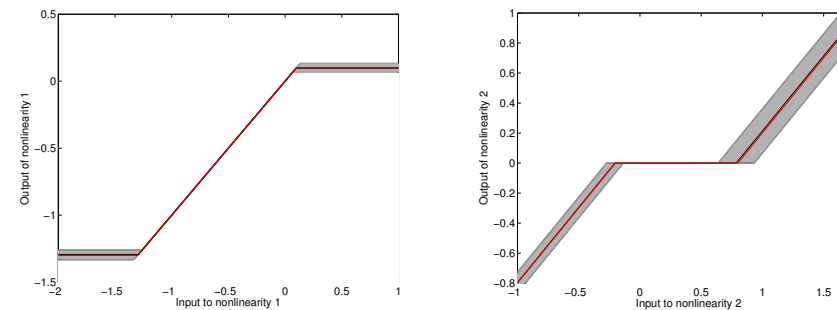
$$z_t = Cx_t, \quad y_t = h(z_t, \beta) + e_t, \quad e_t \sim \mathcal{N}(0, R).$$

Identification problem: Find $A, B, C, \beta, Q,$ and R based on $\{y_{1,1:T}, y_{2,1:T}\}$ using EM.

- Second order LGSS model with complex poles.
- Employ the EM-PS with $M = 100$ particles.
- EM-PS was terminated after 100 iterations.
- Results obtained using $T = 1000$ samples.
- The plots are based on 100 realizations of data.
- Nonlinearities (dead-zone and saturation) shown on next slide.



Bode plot of estimated mean (black), true system (red) and the result for all 100 realisations (gray).



Estimated mean (black), true static nonlinearity (red) and the result for all 100 realisations (gray).

Adrian Wills, Thomas B. Schön, Lennart Ljung and Brett Ninness. *Identification of Hammerstein-Wiener Models. Automatica*, 49(1): 70-81, January 2013.

A linear superposition of Gaussians

$$p(x) = \sum_{k=1}^K \underbrace{\pi_k}_{p(k)} \underbrace{\mathcal{N}(x_n | \mu_k, \Sigma_k)}_{p(x_n|k)}$$

is called a **Gaussian mixture (GM)**. The mixture coefficients π_k satisfies

$$\sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1.$$

Interpretation: The density $p(x | k) = \mathcal{N}(x | \mu_k, \Sigma_k)$ is the probability of x , given that component k was chosen. The probability of choosing component k is given by the prior probability $p(k)$.



Consider the following GM,

$$p(x) = \underbrace{0.3}_{\pi_1} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 4 \\ 4.5 \end{pmatrix}}_{\mu_1}, \underbrace{\begin{pmatrix} 1.2 & 0.6 \\ 0.6 & 0.5 \end{pmatrix}}_{\Sigma_1}\right) + \underbrace{0.5}_{\pi_2} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 8 \\ 1 \end{pmatrix}}_{\mu_2}, \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{\Sigma_2}\right) + \underbrace{0.2}_{\pi_3} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 9 \\ 8 \end{pmatrix}}_{\mu_3}, \underbrace{\begin{pmatrix} 0.6 & 0.5 \\ 0.5 & 1.5 \end{pmatrix}}_{\Sigma_3}\right)$$

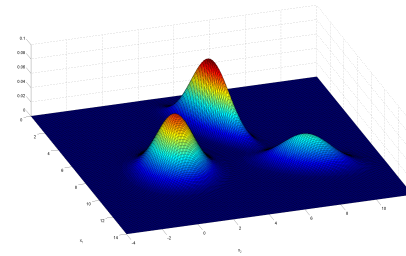


Figure: Probability density function.

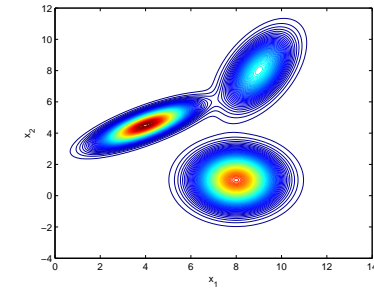


Figure: Contour plot.



Given N independent observations $\{x_n\}_{n=1}^N$, the log-likelihood function is given by

$$\ln p(X; \pi_{1:K}, \mu_{1:K}, \Sigma_{1:K}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right)$$

There is no closed form solution available (due to the sum inside the logarithm).

Let us now see how this problem can be separated into two simple problems using the EM algorithm.

First we introduce an **equivalent** construction of the Gaussian mixture by introducing a latent variable.



Based on

$$p(z_n) = \prod_{k=1}^K \pi_k^{z_{nk}} \quad \text{and} \quad p(x_n | z_n) = \prod_{k=1}^K \mathcal{N}(x_n | \mu_k, \Sigma_k)^{z_{nk}}$$

we have (for independent observations $\{x_n\}_{n=1}^N$)

$$p(X, Z) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(x_n | \mu_k, \Sigma_k)^{z_{nk}},$$

resulting in the following log-likelihood

$$\ln p(X, Z) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} (\ln \pi_k + \ln \mathcal{N}(x_n | \mu_k, \Sigma_k)). \quad (1)$$

Let us now use wishful thinking and assume that Z is known. Then, maximization of (1) is straightforward.



Algorithm 4 EM for Gaussian mixtures

- Initialise:** Initialize $\mu_k^1, \Sigma_k^1, \pi_k^1$ and set $i = 1$.
- While** not converged **do:**

- Expectation (E) step:** Compute

$$\gamma(z_{nk}) = \frac{\pi_k^i \mathcal{N}(x_n | \mu_k^i, \Sigma_k^i)}{\sum_{j=1}^K \pi_j^i \mathcal{N}(x_n | \mu_j^i, \Sigma_j^i)}, \quad n = 1, \dots, N, k = 1, \dots, K.$$

- Maximization (M) step:** Compute

$$\mu_k^{i+1} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n, \quad \pi_k^{i+1} = \frac{N_k}{N}, \quad N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$$\Sigma_k^{i+1} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k^{i+1})(x_n - \mu_k^{i+1})^T$$

- $i \leftarrow i + 1$



Consider the same Gaussian mixture as before,

$$p(x) = \underbrace{0.3}_{\pi_1} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 4 \\ 4.5 \end{pmatrix}}_{\mu_1}, \underbrace{\begin{pmatrix} 1.2 & 0.6 \\ 0.6 & 0.5 \end{pmatrix}}_{\Sigma_1}\right) + \underbrace{0.5}_{\pi_2} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 8 \\ 1 \end{pmatrix}}_{\mu_2}, \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{\Sigma_2}\right) + \underbrace{0.2}_{\pi_3} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 9 \\ 8 \end{pmatrix}}_{\mu_3}, \underbrace{\begin{pmatrix} 0.6 & 0.5 \\ 0.5 & 1.5 \end{pmatrix}}_{\Sigma_3}\right)$$

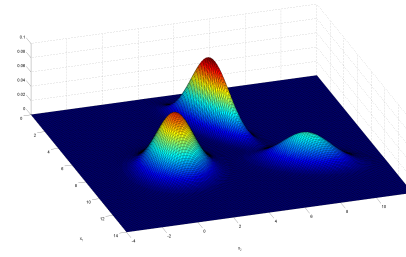
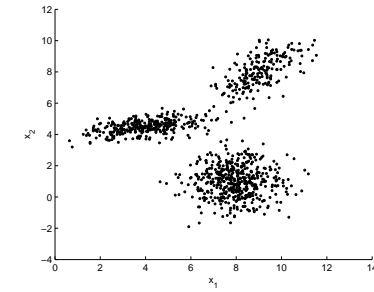


Figure: Probability density function.

Figure: $N = 1000$ samples from the Gaussian mixture $p(x)$.

- Apply the EM algorithm to estimate a Gaussian mixture with $K = 3$ Gaussians, i.e. use the 1000 samples to compute estimates of $\pi_1, \pi_2, \pi_3, \mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3$.
- 200 iterations.

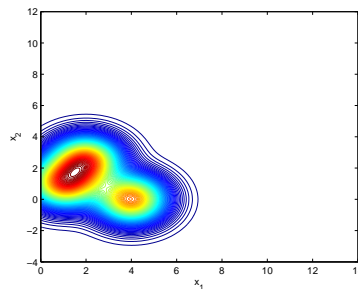


Figure: Initial guess.

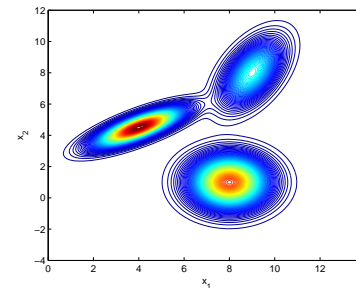


Figure: True PDF.

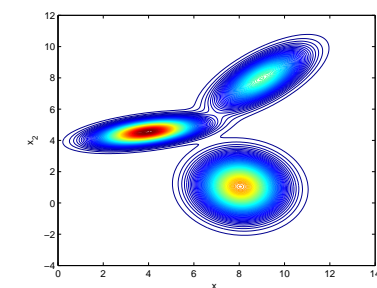


Figure: Estimate after 200 iterations of the EM algorithm.



Algorithm 5 K -means algorithm, a.k.a. Lloyd's algorithm

1. Initialize μ_k^1 and set $i = 1$.
2. Minimize J w.r.t. r_{nk} keeping $\mu_k = \mu_k^i$ fixed.

$$r_{nk}^{i+1} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j^i\|^2 \\ 0 & \text{otherwise} \end{cases}$$

3. Minimize J w.r.t. μ_k keeping $r_{nk} = r_{nk}^{i+1}$ fixed.

$$\mu_k^{i+1} = \frac{\sum_{n=1}^N r_{nk}^{i+1} x_n}{\sum_{n=1}^N r_{nk}^{i+1}}.$$

4. If not converged, update $i := i + 1$ and return to step 2.

The name K -means stems from the fact that in step 3 of the algorithm, μ_k is given by the mean of all the data points assigned to cluster k .

Note the **similarities** between the K -means algorithm and the EM algorithm for Gaussian mixtures!

K -means is deterministic with “hard” assignment of data points to clusters (no uncertainty), whereas EM is a probabilistic method that provides a “soft” assignment.

If the Gaussian mixtures are modeled using covariance matrices

$$\Sigma_k = \epsilon I, \quad k = 1, \dots, K,$$

it can be shown that the EM algorithm for a mixture of K Gaussian's is **equivalent** to the K -means algorithm, when $\epsilon \rightarrow \infty$.

Latent variable: A variable that is not directly observed. Sometimes also referred to as hidden variable or missing data.

Expectation Maximization (EM): The EM algorithm computes maximum likelihood estimates of unknown parameters in probabilistic models involving latent variables.

Jensen's inequality: States that if f is a convex function, then $\mathbf{E}(f(x)) \geq f(\mathbf{E}(x))$.

Clustering: Unsupervised learning, where a set of observations is divided into clusters. The observations belonging to a certain cluster are similar in some sense.

K -means algorithm (a.k.a. Lloyd's algorithm): A clustering algorithm that assigns N observations into K clusters, such that each observation belongs to the cluster with nearest (in the Euclidean sense) mean.