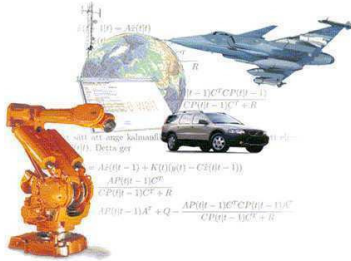


Machine Learning, Lecture 3 Expectation Maximization (EM) and Clustering



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 3

2(40)

1. Summary of Lecture 2
2. Expectation Maximization
 - 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)

Summary of Lecture 2 (I/IV)

3(40)

Using the posterior mean $m_N = \beta S_N \Phi^T T$, the linear regression model is given by

$$y(x, m_N) = m_N^T \phi(x) = \beta \phi(x)^T S_N \Phi^T T = \sum_{n=1}^N \underbrace{\beta \phi(x)^T S_N \phi(x_n)}_{k(x, x_n)} t_n$$

Hence, the mean of the predictive distribution at a point x is a linear combination of the training target variables t_n ,

$y(x, m_N) = \sum_{n=1}^N k(x, x_n) t_n$, where $k(x, x') = \beta \phi(x)^T S_N \phi(x')$ is called the *equivalent kernel*.

An alternative approach to regression is to directly make use of a localized kernel, rather than starting from basis functions. This leads to the so called *kernel methods*.

Summary of Lecture 2 (II/IV)

4(40)

Investigated one linear **discriminant** (a function that takes an input and assigns it to one of K classes) method in detail (least squares)

Modelled each class as $y_k(x) = w_k^T x + w_{k,0}$ and solved the LS problem, resulting in $\hat{w} = (X^T X)^{-1} X^T T$

Showed how probabilistic **generative models** could be built for classification using the strategy,

1. Model $p(x | C_k)$ (a.k.a. class-conditional density)
2. Model $p(C_k)$
3. Use ML to find the parameters in $p(x | C_k)$ and $p(C_k)$.
4. Use Bayes' rule to find $p(C_k | x)$

The “direct” method called **logistic regression** was introduced. Start by stating the model

$$p(C_1 | \phi) = \sigma(w^T \phi) = \frac{1}{1 + e^{-w^T \phi}},$$

which results in a log-likelihood function according to

$$L(w) = -\ln p(T | w) = -\sum_{n=1}^N (t_n \ln(y_n) + (1 - t_n) \ln(1 - y_n)),$$

where $y_n = p(C_1 | \phi) = \sigma(w^T \phi)$. Note that this is a nonlinear, **but** concave function of w .

Hence, we can easily find the global minimum using Newton’s method (resulting in an algorithm known as IRLS).

The likelihood function for logistic regression is

$$p(T | w) = \prod_{n=1}^N \sigma(w^T \phi_n)^{t_n} (1 - \sigma(w^T \phi_n))^{1-t_n}$$

Hence, computing the posterior density $p(w | T) = \frac{p(T|w)p(w)}{p(T)}$ is intractable and we considered the **Laplace approximation** for solving this.

The Laplace approximation is a simple (local) approximation that is obtained by fitting a Gaussian centered around the (MAP) mode of the distribution.

Let us define the following concepts,

- **True positives (tp):** Items correctly classified as belonging to the class.
- **True negatives (tn):** Items correctly classified as not belonging to the class.
- **False positives (fp):** Items incorrectly classified as belonging to the class. In other words, a false alarm.
- **False negatives (fn):** Items that are not classified in the correct class, even though they should have. In other words, a missed detection.

$$\text{Precision} = \frac{\text{tp}}{\text{tp} + \text{fp}} \quad \text{Recall} = \frac{\text{tp}}{\text{tp} + \text{fn}}$$

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 scalar state-space 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 strategy underlying the EM algorithm is to separate the original ML problem into **two linked problems**, each of which is hopefully easier to solve than the original problem.

This separation is accomplished by exploiting the **structure** inherent in the probabilistic model.

The **key idea** is to consider the joint log-likelihood function of both the observed variables X and the latent variables Z ,

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

Algorithm (*Expectation Maximization (EM)*)

1. Set $i = 0$ and choose an initial guess θ_0 .
2. **Expectation (E) step:** Compute

$$\begin{aligned} Q(\theta, \theta_i) &= \mathbf{E}_{\theta_i} \{ \ln p_{\theta}(Z, X) \mid X \} \\ &= \int \ln p_{\theta}(Z, X) p_{\theta_i}(Z \mid X) dZ. \end{aligned}$$

3. **Maximization (M) step:** Compute

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

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

Consider the following linear scalar state-space 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, $x_1 = 0$. Finally, 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 given by X .

The resulting Q -function is

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

where we have defined

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

There exists explicit expressions for these expected values.

The maximization (M) step:

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

Hence, the M step simply amounts to solving the following quadratic problem,

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

which results in

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



Algorithm

1. Set $i = 0$ and choose an initial guess θ_0 .
2. **Expectation (E) step:** Compute

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

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

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

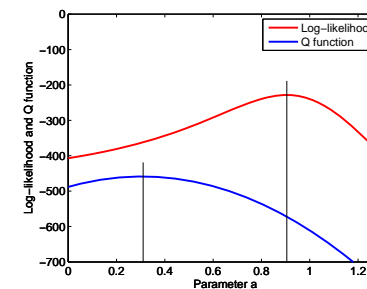
4. If $|L_{\theta_i}(Y) - L_{\theta_{i-1}}(Y)| \geq 10^{-6}$, update $i := i + 1$ and return to step 2, otherwise terminate.



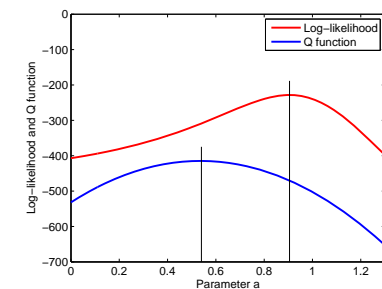
- Different number of samples N used.
- Monte Carlo studies, each using 1000 realisations of data.
- Initial guess $\theta_0 = 0.1$.

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.

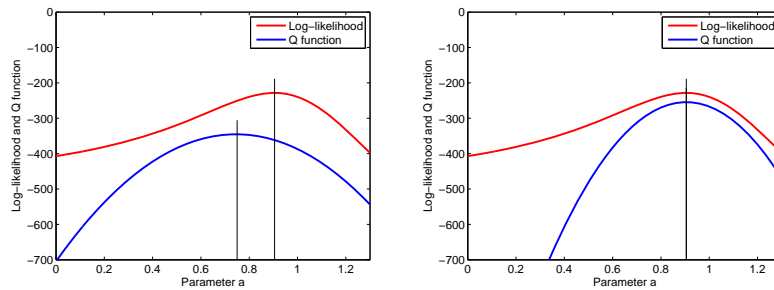


(a) Iteration 1



(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.

<http://www.rti.isy.liu.se/~schon/Publications/SchonEM2009.pdf>



Now, consider the problem of estimating the parameters θ in

$$\begin{aligned}x_{t+1} &= f_t(x_t, \theta) + v_t(\theta), \\ y_t &= h_t(x_t, \theta) + e_t(\theta).\end{aligned}$$

Commonly also written as

$$\begin{aligned}x_{t+1} &\sim p_\theta(x_{t+1} | x_t), \\ y_t &\sim p_\theta(y_t | x_t).\end{aligned}$$

According to the above, the first step is to compute the \mathcal{Q} -function

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



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

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

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

$$I_1 = \int \log p_\theta(x_1) p_{\theta_k}(x_1 | Y) dx_1,$$

$$I_2 = \sum_{t=1}^{N-1} \int \int \log p_\theta(x_{t+1} | x_t) p_{\theta_k}(x_{t+1}, x_t | Y) dx_t dx_{t+1},$$

$$I_3 = \sum_{t=1}^N \int \log p_\theta(y_t | x_t) p_{\theta_k}(x_t | Y) dx_t.$$



This leads us to a nonlinear state smoothing problem, which we can solve using particle smoothers.

The current particle smoothers provides us with the following approximation of the joint smoothing density

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

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

$$p_{\theta_k}(x_t | Y) \approx \hat{p}_{\theta_k}(x_t | Y) = \frac{1}{M} \sum_{i=1}^M \delta(x_t - x_t^i),$$

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



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

$$\hat{I}_1 = \int \log p_\theta(x_1) \sum_{i=1}^M \frac{1}{M} \delta(x_1 - x_1^i) dx_1 = \frac{1}{M} \sum_{i=1}^M \log p_\theta(x_1^i),$$

$$\hat{I}_2 = \sum_{t=1}^{N-1} \int \int \log p_\theta(x_{t+1} | x_t) \sum_{i=1}^M \frac{1}{M} \delta(x_{t+1} - x_{t+1}^i) dx_{t+1}$$

$$= \frac{1}{M} \sum_{t=1}^{N-1} \sum_{i=1}^M \log p_\theta(x_{t+1}^i | x_t^i),$$

$$\hat{I}_3 = \sum_{t=1}^N \int \log p_\theta(y_t | x_t) \sum_{i=1}^M \frac{1}{M} \delta(y_t - y_t^i) dy_t = \frac{1}{M} \sum_{t=1}^N \sum_{i=1}^M \log p_\theta(y_t^i | x_t^i)$$

It is straightforward to make use of the approximation of the Q-functions 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 analogously),

$$\hat{I}_3 = \frac{1}{M} \sum_{t=1}^N \sum_{i=1}^M \log p_\theta(y_t | x_t^i),$$

$$\frac{\partial \hat{I}_3}{\partial \theta} = \frac{1}{M} \sum_{t=1}^N \sum_{i=1}^M \frac{\partial \log p_\theta(y_t | x_t^i)}{\partial \theta}$$

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

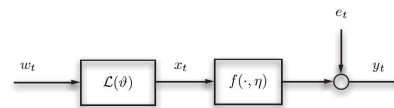
Algorithm (Nonlinear System Identification Using EM)

1. Set $i = 0$ and choose an initial guess θ_0 .
2. **Expectation (E) step:**
 1. Run the particle filter and the particle smoother
 2. Calculate an approximation of the Q-function $\hat{Q}_M(\theta, \hat{\theta}_i) = \hat{I}_1 + \hat{I}_2 + \hat{I}_3$.
3. **Maximization (M) step:** Compute

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

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

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



$$\zeta_{t+1} = A\zeta_t + w_t, \quad w_t \sim \mathcal{N}(0, Q),$$

$$x_t = C\zeta_t,$$

$$y_t = f(x_t, \eta) + e_t, \quad e_t \sim \mathcal{N}(0, R).$$

Figure: The blind (only the outputs y_t are available) Wiener (linear dynamic system followed by a static nonlinearity) problem.

$$\theta = [\vartheta^T, \eta^T, \text{vec}\{R\}^T],$$

$$\vartheta = [\text{vec}\{A\}^T, \text{vec}\{C\}^T, \text{vec}\{Q\}^T]^T.$$

Most existing work deals with special cases, with assumptions like;

1. Invertible nonlinearities,
2. no measurement noise,
3. no process noise.

We are not forced to do any of these assumptions.

Recall that $(Z = \Xi) \mathcal{Q}(\theta, \theta_i) \triangleq \int \ln p_\theta(\Xi, Y) p_{\theta_i}(\Xi | Y) d\Xi$.
 According to the model we have

$$\begin{aligned} \ln p_\theta(\Xi, Y) &= \ln p_\theta(Y | \Xi) + \ln p_\theta(\Xi) \\ &= \sum_{t=1}^N \ln p_\theta(y_t | \xi_t) + \sum_{t=1}^{N-1} \ln p_\theta(\xi_{t+1} | \xi_t), \end{aligned}$$

resulting in $\mathcal{Q}(\theta, \theta_i) = I_1 + I_2$,

$$\begin{aligned} I_1 &= \sum_{t=1}^N \int \int \ln p_\theta(\xi_{t+1} | \xi_t) p_{\theta_i}(\xi_{t+1}, \xi_t | Y) d\xi_t d\xi_{t+1}, \\ I_2 &= \sum_{t=1}^N \ln p_\theta(y_t | \xi_t) p_{\theta_i}(\xi_t | Y) d\xi_t, \end{aligned}$$

where $p_{\theta_i}(\xi_{t+1}, \xi_t | Y)$ and $p_{\theta_i}(\xi_t | Y)$ are provided by the particle smoother.

Wiener system:

$$H(q) = \frac{q^{-1} + 0.1q^{-2} - 0.49q^{-3} + 0.01q^{-4}}{1 + 0.3676q^{-1} + 0.88746q^{-2} + 0.52406q^{-3} + 0.55497q^{-4}},$$

$$f(x) = \begin{cases} 0.3 & : x > 0.3, \\ x & : -0.2 \leq x \leq 0.3, \\ -0.2 & : x < -0.2. \end{cases}$$

$N = 1000$ samples used, $w_t \sim \mathcal{N}(0, 0.1), e_t \sim \mathcal{N}(0, 0.001)$.

$M = 200$ particles are used and 100 Monte Carlo runs.

Initialize the linear system using a subspace method.

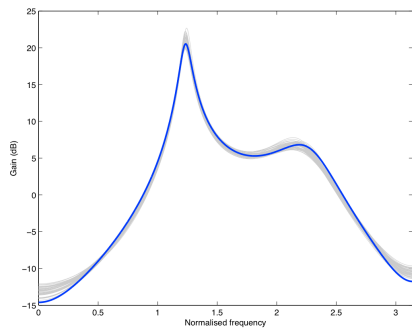


Figure: Bode plot of estimated (grey) and true (blue) systems.

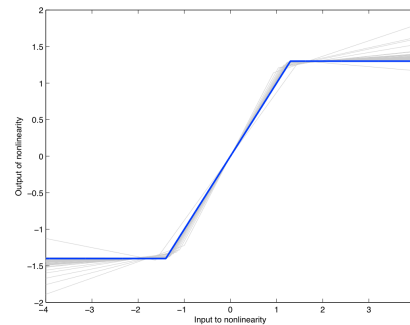


Figure: Estimated (grey) and true (blue) memoryless nonlinearity (saturation).

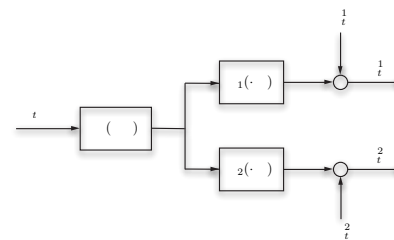


Figure: Block diagram of blind Wiener model with two outputs.

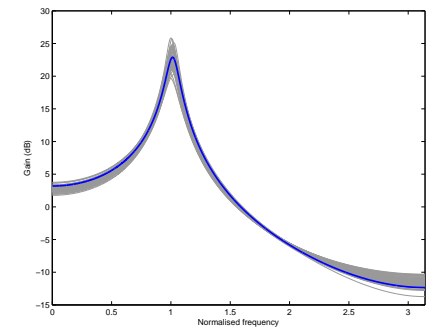


Figure: Bode plot of estimated (grey) and true (blue) systems.

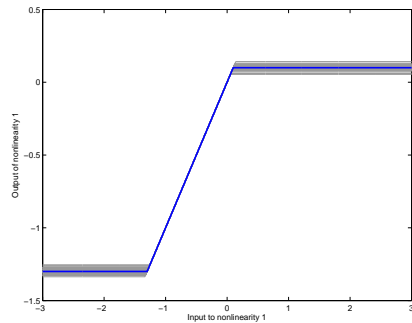


Figure: Estimated (grey) and true (blue) nonlinearity (saturation).

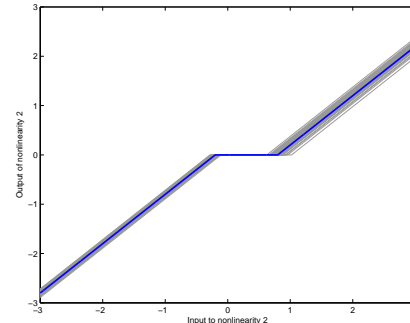


Figure: Estimated (grey) and true (blue) nonlinearity (dead band).

Adrian Wills, Thomas B. Schön, Lennart Ljung and Brett Ninness. Blind Identification of Wiener Models. *Proceedings of the 18th World Congress of the International Federation of Automatic Control (IFAC)*, Milan, Italy, September 2011. (accepted for publication)

A linear superposition of Gaussians

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

is called a Gaussian mixture. 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 Gaussian mixture

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

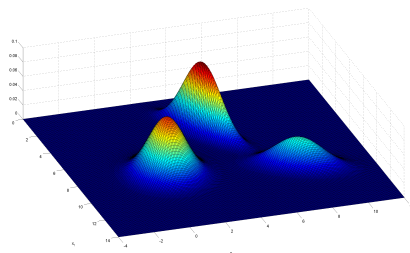


Figure: Probability density function.

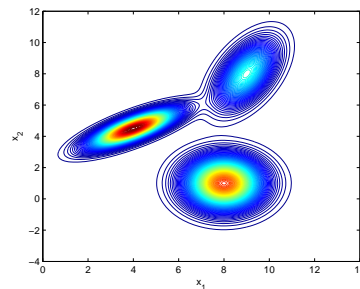


Figure: Contour plot.

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).

We will now see that 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 (EM for Gaussian Mixtures)

1. Initialize $\mu_k^1, \Sigma_k^1, \pi_k^1$ and set $i = 1$.
2. **Expectation (E) step:** Compute

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

3. **Maximization (M) step:** Compute

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

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

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

Example - EM for Gaussian Mixtures (I/III)

Consider the same Gaussian mixture as before,

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

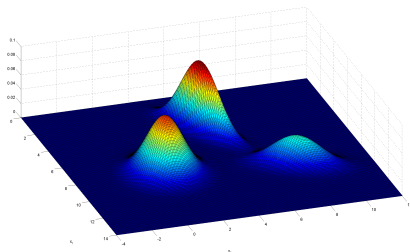


Figure: Probability density function.

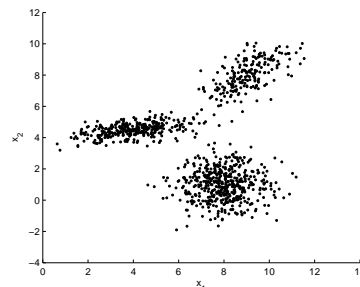


Figure: $N = 1000$ samples from the Gaussian mixture $p(x)$.

Example - EM for Gaussian Mixtures (II/III)

- 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.

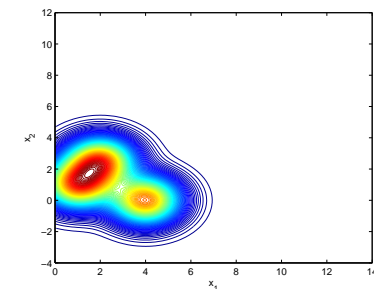


Figure: Initial guess.

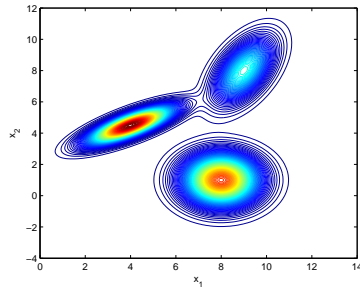


Figure: True PDF.

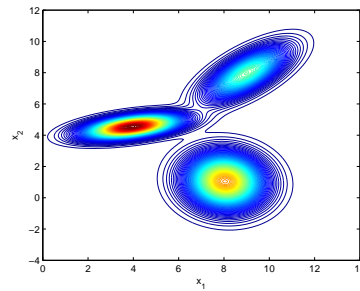


Figure: Estimate after 200 iterations of the EM algorithm.

Algorithm (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 is straightforward to show 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.