

# Statistical Computing

## The Expectation-Maximization Algorithm III Mixture Model for Multi-dimensional Gaussians

Uwe Menzel, 2018

[uwe.menzel@matstat.de](mailto:uwe.menzel@matstat.de)

[www.matstat.org](http://www.matstat.org)

## Multi-dimensional Gaussians

- Very similar to the 1-D case, for details see part II
- $d$  = dimension of the sample space
- $\boldsymbol{\mu}_k$ : vector of means (of length  $d$ )
- $\mathbf{x}_i$ : one observation, vector of length  $d$
- $\boldsymbol{\Sigma}_k$ : covariance matrix ( $d \times d$  matrix):

$$\boldsymbol{\Sigma}_k = \begin{pmatrix} Cov(\mathbf{x}_1, \mathbf{x}_1) & Cov(\mathbf{x}_1, \mathbf{x}_2) & \dots & Cov(\mathbf{x}_1, \mathbf{x}_d) \\ Cov(\mathbf{x}_2, \mathbf{x}_1) & Cov(\mathbf{x}_2, \mathbf{x}_2) & \dots & Cov(\mathbf{x}_2, \mathbf{x}_d) \\ \dots & \dots & \dots & \dots \\ Cov(\mathbf{x}_d, \mathbf{x}_1) & Cov(\mathbf{x}_d, \mathbf{x}_2) & \dots & Cov(\mathbf{x}_d, \mathbf{x}_d) \end{pmatrix}$$

In two dimensions, this translates to:

$$\boldsymbol{\Sigma}_k = \begin{pmatrix} Cov(\mathbf{x}_1, \mathbf{x}_1) & Cov(\mathbf{x}_1, \mathbf{x}_2) \\ Cov(\mathbf{x}_2, \mathbf{x}_1) & Cov(\mathbf{x}_2, \mathbf{x}_2) \end{pmatrix} = \begin{pmatrix} \sigma_x^2 & \rho \cdot \sigma_x \cdot \sigma_y \\ \rho \cdot \sigma_x \cdot \sigma_y & \sigma_y^2 \end{pmatrix}$$

$$Cov(\mathbf{x}, \mathbf{x}) = \sigma_x^2 \quad \rho = \rho(\mathbf{x}, \mathbf{y}) = \frac{Cov(\mathbf{x}, \mathbf{y})}{\sigma_x \sigma_y} \quad \text{correlation coefficient}$$

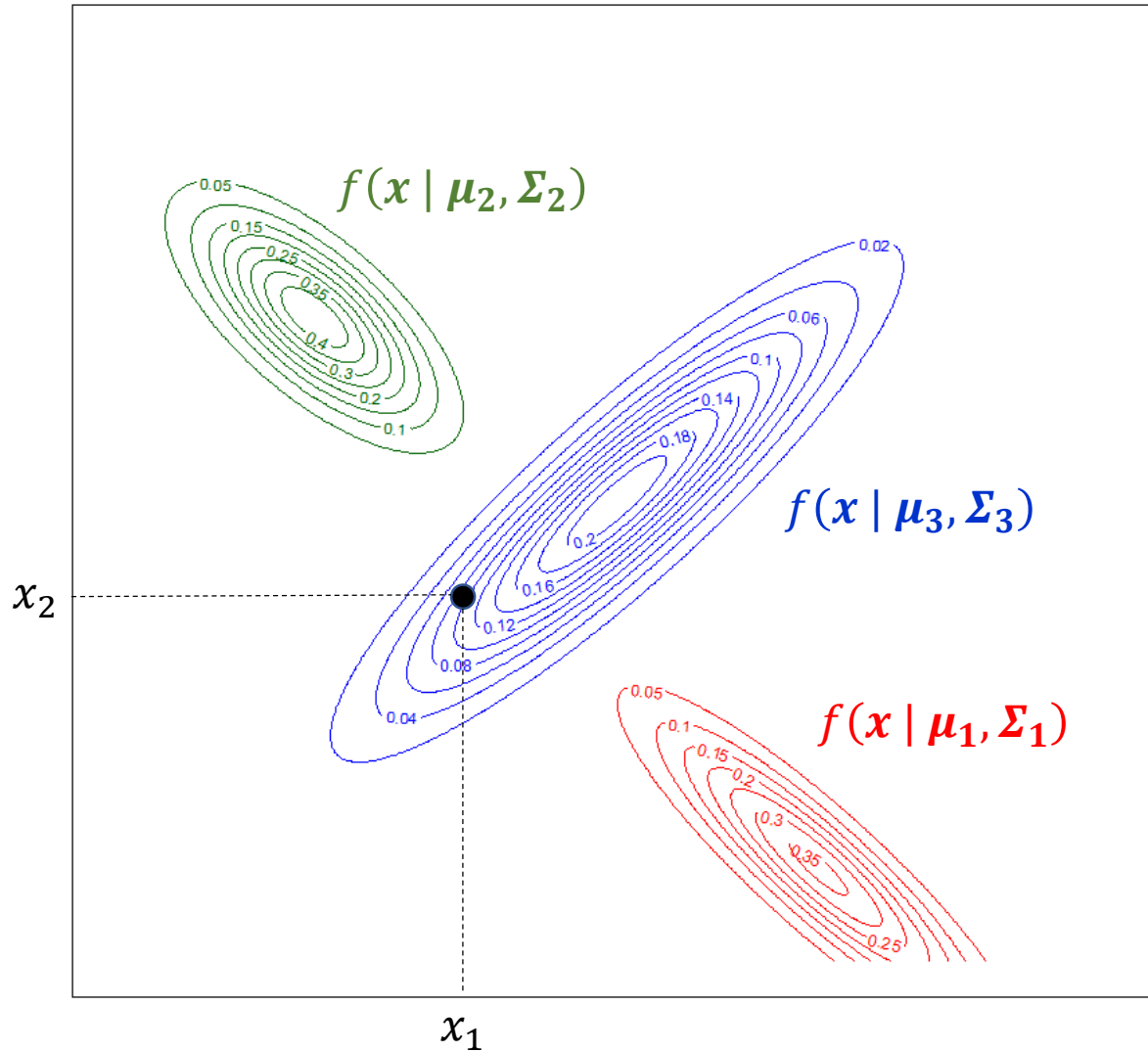
## Two-step experiment

Assume we have  $K$   $d$ -dimensional Gaussians with densities  $f_k(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , where  $k = (1, 2, \dots, K)$ . We carry out the same two-step experiment as we did in the one-dimensional case (see part II):

1. Choose a Gaussian  $f_k$  randomly with some probability  $\alpha_k$ . This can be described by a multinomially distributed random variable  $Z$  with sample space  $\Omega_Z = \{1, 2, \dots, K\}$  and probability mass function  $P(Z = k) = \alpha_k$ . The  $\alpha_k$  are constrained by  $\sum \alpha_k = 1$  and  $\alpha_k > 0$  for all  $k$ .
2. Generate a sample  $\mathbf{x}$  from the above chosen distribution  $f_k$ . The vector  $\mathbf{x}$  is an observation of a  $d$ -dimensional normally distributed random variable  $\mathbf{X}$  with parameters  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$ , i.e.  $\mathbf{X} \sim MVN(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ .

The figure on the next page illustrates this experiment: In step 1, cluster 3 was chosen. Then the vector  $\mathbf{x} = (x_1, x_2)$  is generated from the distribution  $\mathbf{X} \sim MVN(\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$ .

# Contour plot for 2-D Gaussian mixture



# Probability Density Function for multidimensional Gaussian

The experiment includes a discrete ( $Z$ ) and a  $d$ -dimensional continuous ( $\mathbf{X}$ ) random variable. The (mixed) **joint density of  $\mathbf{X}$  and  $\mathbf{Z}$**  can be written:

$$\begin{aligned} f_{\mathbf{X},Z}(\mathbf{x}, Z = k) &= P(Z = k) \cdot f_{\mathbf{X}|Z}(\mathbf{x}|Z = k) && f_{\mathbf{X}|Z}: \text{conditional} \\ & && \text{probability} \\ &= \alpha_k \cdot f_k(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

where  $f_k$  is the probability density function for a  $d$ -dimensional Gaussian:

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \cdot \sqrt{|\boldsymbol{\Sigma}_k|}} \cdot \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

Here,  $|\boldsymbol{\Sigma}_k|$  is the determinant of the covariance matrix of the  $k^{\text{th}}$  component (cluster).

# Probability Density Function for two-dimensional Gaussian

In the 2-D case, the joint density of  $\mathbf{x} = (x, y)$  simplifies to

$$\begin{aligned} f_{\mathbf{k}}(\mathbf{x}) &= f_{\mathbf{k}}(x, y) = \\ &= \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \cdot \exp \left\{ -\frac{(x-\mu_x)^2}{2\sigma_x^2} + \rho \cdot \frac{(y-\mu_y)(x-\mu_x)}{\sigma_x\sigma_y} - \frac{(y-\mu_y)^2}{2\sigma_y^2} \right\} \end{aligned}$$

because  $|\boldsymbol{\Sigma}_{\mathbf{k}}| = \sigma_x^2\sigma_y^2 \cdot (1-\rho^2)$  and

$$\begin{aligned} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{k}})^T \cdot (\boldsymbol{\Sigma}_{\mathbf{k}})^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{k}}) &= \\ &= \frac{1}{(1-\rho^2)} \cdot \left\{ \frac{(x-\mu_x)^2}{\sigma_x^2} - 2 \cdot \rho \frac{(y-\mu_y)(x-\mu_x)}{\sigma_x\sigma_y} + \frac{(y-\mu_y)^2}{\sigma_y^2} \right\} \end{aligned}$$

See the appendix for a derivation of the two-dimensional expressions.

## Gaussian mixture

We have to find parameters  $\boldsymbol{\theta}$  that maximize  $f_X(\mathbf{x} | \boldsymbol{\theta})$ . As in the 1-D case, an expression for the density  $f_X(\mathbf{x} | \boldsymbol{\theta})$  that incorporates the latent variables  $Z$  can be found by using the **law of total probability**:

$$f_X(\mathbf{x} | \boldsymbol{\theta}) = \sum_{k=1}^K \underbrace{f_{\mathbf{X}|Z=k}(\mathbf{x} | Z = k)}_{f_k(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \cdot \underbrace{P(Z = k)}_{\alpha_k} = \sum_{k=1}^K \alpha_k \cdot f_k(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where  $K$  is the number of clusters.

The density  $f_X$  can be seen as a superposition of multiple probability density functions (Gaussians):

$$f_X(\mathbf{x} | \boldsymbol{\theta}) = \sum_{k=1}^K \alpha_k \cdot f_k(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

## Maximum Likelihood for a Gaussian mixture

If we have multiple independent observations  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T$ , the likelihood is the product of the density for the individual observations:

$$L(\boldsymbol{\theta}) = \prod_{i=1}^N f_X(\mathbf{x}_i | \boldsymbol{\theta}) = \prod_{i=1}^N \sum_{k=1}^K \alpha_k \cdot f_k(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Note that each  $\mathbf{x}_i$  is a  $d$ -dimensional vector here. We aim at calculating the  $\boldsymbol{\theta}$  that maximizes  $L(\boldsymbol{\theta})$ :

$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L(\boldsymbol{\theta})$$

i.e. we search for the parameter (vector)  $\boldsymbol{\theta}$  that makes the observed data most likely. Often, it is more convenient to maximize the logarithm of  $L(\boldsymbol{\theta})$ :

$$l(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta}) = \sum_{i=1}^N \log \left[ \sum_{k=1}^K \alpha_k \cdot f(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]$$

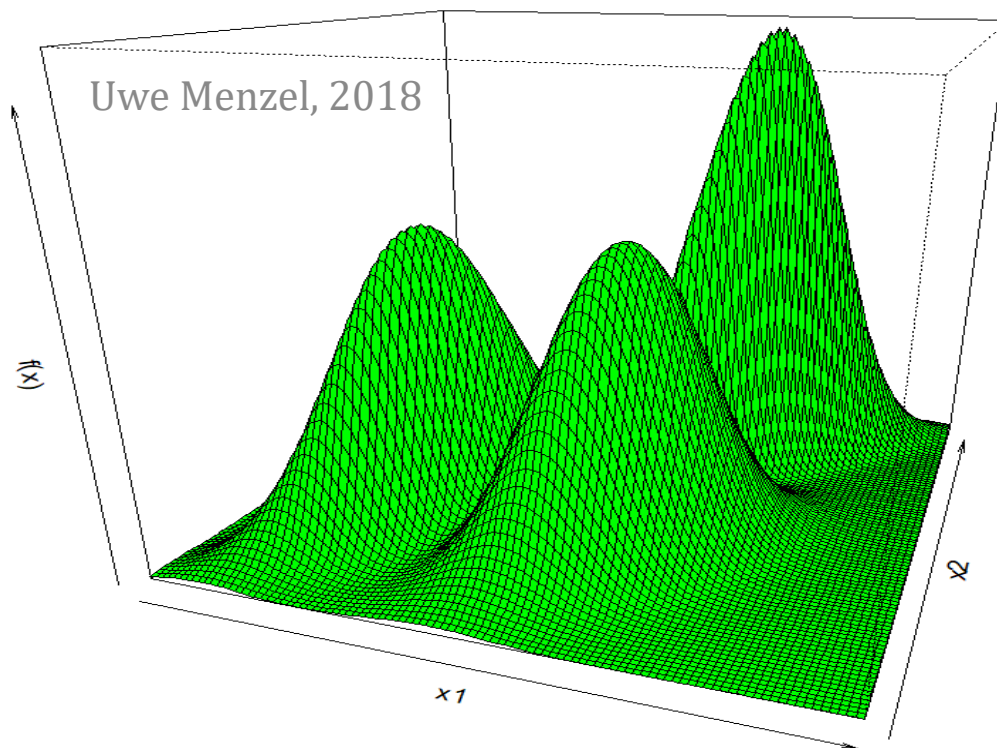


# Mixture of Gaussians: E-step

1. **Initialize:**  $\theta_t = 1^{\text{st}}$  guess for the set of parameters

2. **E-step:** calculate  $P(Z_i = k \mid \mathbf{X}_i = \mathbf{x}_i, \theta_t)$ .

This is the probability that  $Z_i$  is equal to  $k$ , i.e. the probability that an observation  $\mathbf{x}_i$  originates from the  $k^{\text{th}}$  Gaussian, given that observation ( $\mathbf{x}_i$ ) and all parameters  $\theta_t = \{\alpha_k^t, \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t\}$ . Since the parameters can be considered as given (we made a 1<sup>st</sup> guess) , we know the exact positions and shapes of the  $d$ -dimensional Gaussians and it's superposition, as shown in the figure for  $d = 2$ .



## Mixture of Gaussians: E-step

2. **E-step**: calculate  $P(Z_i = k | \mathbf{X} = \mathbf{x}_i, \boldsymbol{\theta}_t)$ .

Knowledge of  $\mathbf{x}_i$  and  $\boldsymbol{\theta}_t$  enables us to calculate the conditional probability using **Bayes theorem** :

$$\begin{aligned} P(Z_i = k | \mathbf{X} = \mathbf{x}_i, \boldsymbol{\theta}_t) &= \frac{f_{\mathbf{X}|Z}(\mathbf{x}_i | Z_i = k, \boldsymbol{\theta}_t) \cdot P(Z_i = k | \boldsymbol{\theta}_t)}{f_{\mathbf{X}}(\mathbf{x}_i | \boldsymbol{\theta}_t)} \\ &= \frac{f_k(\mathbf{x}_i | \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t) \cdot \alpha_k^t}{\sum_k \alpha_k^t \cdot f_k(\mathbf{x}_i | \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)} = \omega_{ik} \end{aligned}$$

The last ratio is labelled  $\omega_{ik}$  and often named "degree of membership" (of observation  $\mathbf{x}_i$  to component  $k$ ). The  $\omega_{ik}$  are known numbers since they are calculated based on the known  $\boldsymbol{\theta}_t = \{\alpha_k^t, \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t\}$ .

$$\omega_{ik} = \frac{\alpha_k^t \cdot f_k(\mathbf{x}_i | \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)}{\sum_k \alpha_k^t \cdot f_k(\mathbf{x}_i | \boldsymbol{\mu}_k^t, \boldsymbol{\Sigma}_k^t)} \quad \sum_{k=1}^K \omega_{ik} = 1$$

## Degree of membership: $\omega_{ik}$

It is not easy to illustrate the geometric meaning of the  $\omega_{ik}$  in the multi-dimensional case. In two dimensions, having two clusters, we might see the figure below as a cross section perpendicular to the  $x - y$  plane.

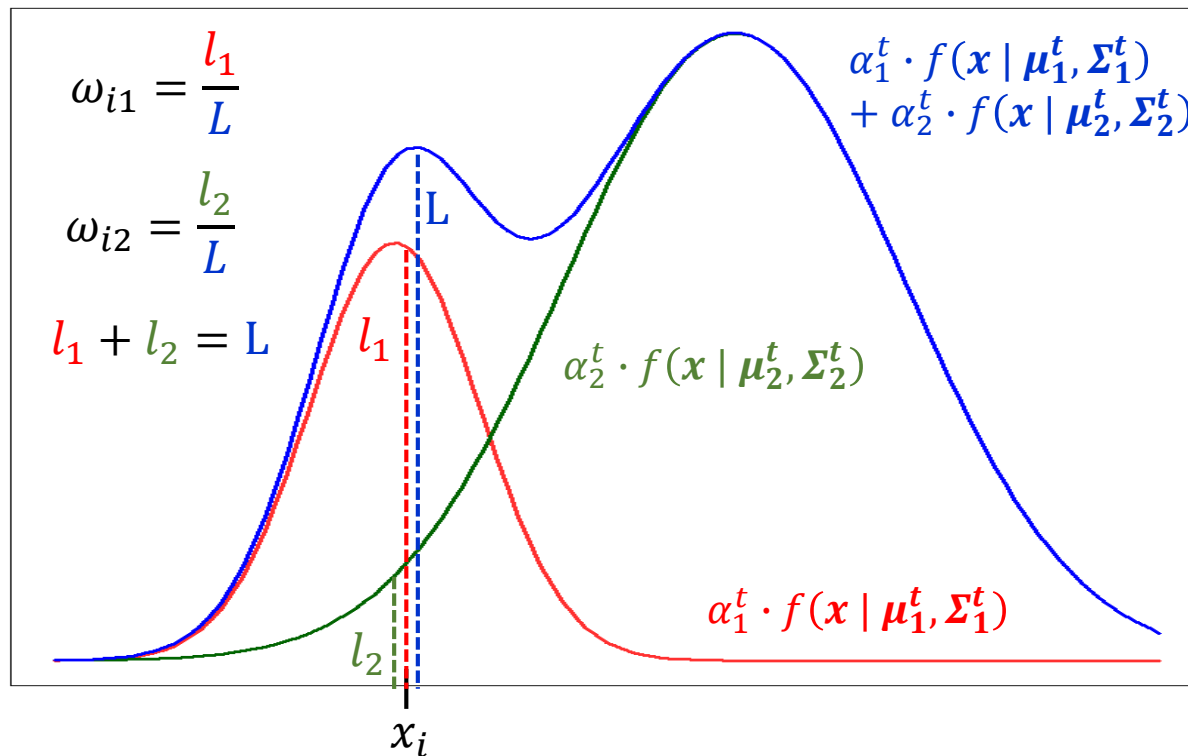


Illustration of the  $\omega_{ik}$  (dashed lines jittered around  $x_i$  for better visibility)

## Completion of the E-step

It remains to calculate:

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \underbrace{P(Z_i = k \mid \mathbf{X} = \mathbf{x}_i, \boldsymbol{\theta}_t)}_{\omega_{ik}} \cdot \log f_{\mathbf{X}, Z}(\mathbf{x}_i, Z_i = k \mid \boldsymbol{\theta})$$

$f_{\mathbf{X}, Z}(\mathbf{x}, Z) = \alpha_k \cdot f_k(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  (mixed) joint probability distribution

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \log \{ \underbrace{\alpha_k \cdot f_k(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{unknown parameters (depending on } \boldsymbol{\theta})} \}$$

$\omega_{ik}$  known (depending on  $\boldsymbol{\theta}_t$ )

The expression  $Q_1$  has to be maximized for the unknown  $\alpha_k, \mu_k, \sigma_k \rightarrow$  **M-step**.

## Mixture of multivariate Gaussians: M-step

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \log \{ \alpha_k \cdot f_k(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}$$

$Q_1$  has to be maximized for the unknown  $\alpha_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \{ \log \alpha_k + \log f_k(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}$$

$$\log f_k(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \left\{ \log \alpha_k - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}$$

The parameters  $\boldsymbol{\Sigma}_k, \boldsymbol{\mu}_k$  and  $\alpha_k$  in the curly brackets are the actual variables  $Q_1$  must be maximized for, while the  $\omega_{ik}$  contain only variables that have been specified by first guess, so that  $\omega_{ik}$  can be treated as a constant in this expression.

## Mixture of Gaussians: M-step

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \left\{ \log \alpha_k - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}$$

Maximization of  $Q_1$  with regard to  $\boldsymbol{\mu}_m$ :

(here, we have ignored the term  $-\frac{d}{2} \cdot \log 2\pi$  since it doesn't depend on any parameter)

$$\frac{\partial Q_1}{\partial \boldsymbol{\mu}_m} = -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_m} \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k)$$

$$= -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial \boldsymbol{\mu}_m} (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\boldsymbol{\Sigma}_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)$$

$$= \sum_{i=1}^N \omega_{im} \cdot (\boldsymbol{\Sigma}_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) = 0 \quad (\text{see Appendix for the last step})$$

$$\Rightarrow (\boldsymbol{\Sigma}_m)^{-1} \sum_{i=1}^N \omega_{im} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) = 0$$

## Mixture of Gaussians: M-step

$$(\Sigma_m)^{-1} \sum_{i=1}^N \omega_{im} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) = 0$$

$$\Rightarrow \sum_{i=1}^N \omega_{im} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) = 0 \quad (\Sigma \text{ is positive definite, covariance})$$

$$\Rightarrow \boldsymbol{\mu}_m = \frac{\sum_{i=1}^N \omega_{im} \cdot \mathbf{x}_i}{\sum_{i=1}^N \omega_{im}}$$

The new means are weighted means of the  $x_i$  (weighted with the degree of membership of each datapoint). This can be compared with the ML estimation of the mean for a single Gaussian:  $\mu = \frac{1}{N} \cdot \sum x_i$

This is similar to the 1-D case where we had  $\mu_m = \frac{\sum_i^N \omega_{im} \cdot x_i}{\sum_i^N \omega_{im}}$

## Mixture of Gaussians: M-step

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \left\{ \log \alpha_k - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}$$

**Maximization of  $Q_1$  with respect to  $\boldsymbol{\Sigma}_m$ :**

Instead of deriving for  $\boldsymbol{\Sigma}_m$ , we derive for  $(\boldsymbol{\Sigma}_m)^{-1}$ . This will automatically lead to an expression for  $\boldsymbol{\Sigma}_m$ . This strategy was already presented by Xavier Bourret Sicotte<sup>1</sup>.

$$\frac{\partial Q_1}{\partial (\boldsymbol{\Sigma}_m)^{-1}} = -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial (\boldsymbol{\Sigma}_m)^{-1}} \left\{ \log |\boldsymbol{\Sigma}_m| + (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\boldsymbol{\Sigma}_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$$

We obtain:

$$\frac{\partial Q_1}{\partial (\boldsymbol{\Sigma}_m)^{-1}} = \frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \left\{ \boldsymbol{\Sigma}_m - (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \right\} \quad \text{see Appendix}$$

<sup>1</sup> <https://stats.stackexchange.com/questions/351549/maximum-likelihood-estimators-multivariate-gaussian>



## Mixture of Gaussians: M-step

$$\frac{\partial Q_1}{\partial (\boldsymbol{\Sigma}_m)^{-1}} = \frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \left\{ \boldsymbol{\Sigma}_m - (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \right\}$$

$$\Rightarrow \frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \left\{ \boldsymbol{\Sigma}_m - (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \right\} = 0$$

$$\Rightarrow \boldsymbol{\Sigma}_m = \frac{\sum_{i=1}^N \omega_{im} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T}{\sum_{i=1}^N \omega_{im}}$$

The new estimate for the covariance matrix  $\boldsymbol{\Sigma}_m$  is a matrix where each element is a weighted mean of the squared distances between datapoints and mean  $\boldsymbol{\mu}_m$  (weighted with the degree of membership of each datapoint). This can be compared with the one-dimensional case, where we obtained for the variance:

$$\sigma_m^2 = \frac{\sum_i \omega_{im} \cdot (x_i - \mu_m)^2}{\sum_i \omega_{im}}$$

## Mixture of Gaussians: M-step

$$Q_1(\boldsymbol{\theta}, \boldsymbol{\theta}_t) = \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \left\{ \log \alpha_k - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \cdot (\boldsymbol{\Sigma}_k)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}$$

Maximization of  $Q_1$  with respect to  $\alpha_m$ :

$$\frac{\partial Q_1}{\partial \alpha_m} = - \frac{\partial}{\partial \alpha_m} \sum_{i=1}^N \sum_{k=1}^K \omega_{ik} \cdot \log \alpha_k$$

This has to be solved with the constraint  $\sum_k \alpha_k = 1$

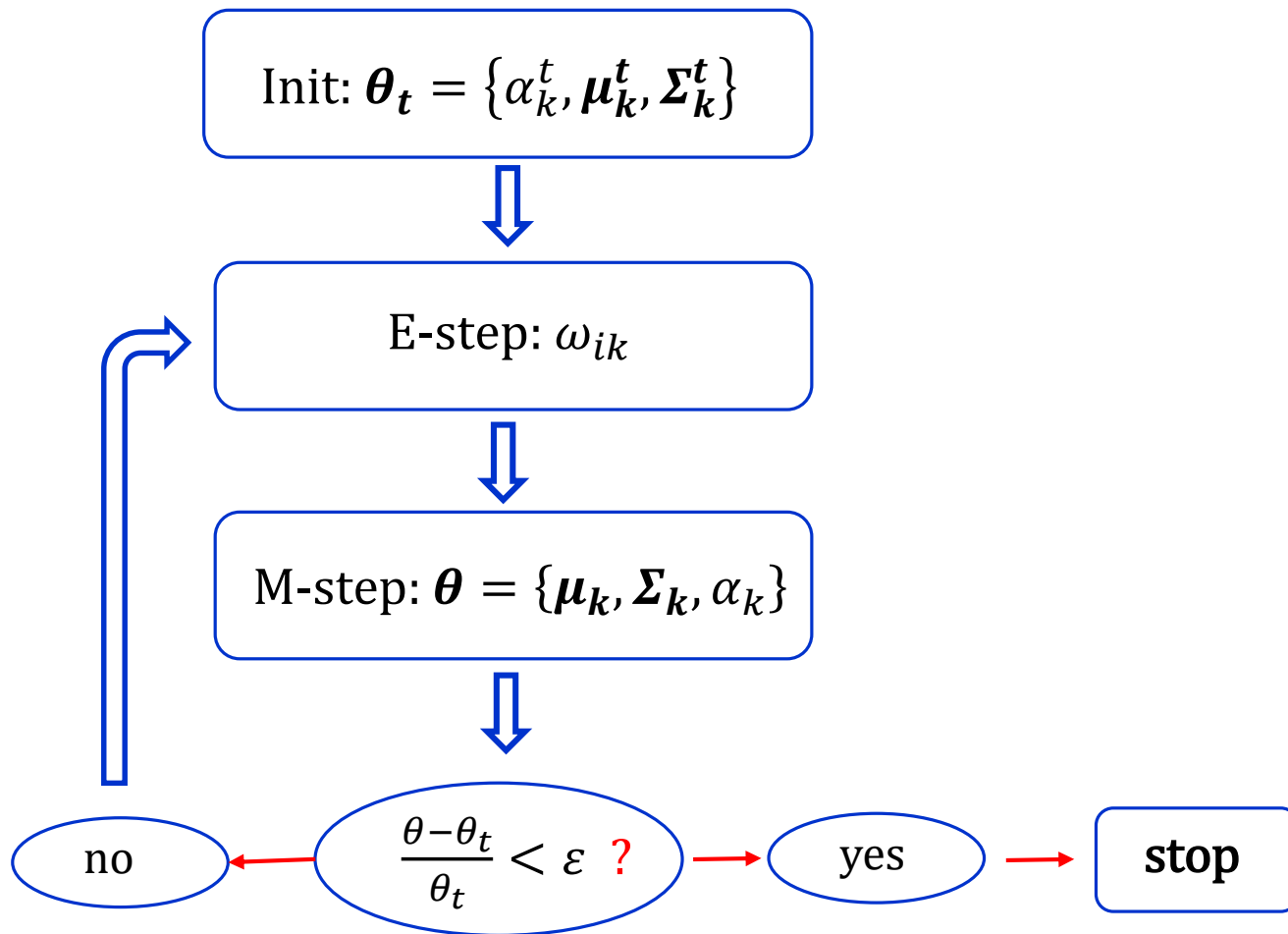
This expression is the same as for the one-dimensional case, so that we can directly adopt the solution (see part II):

$$\alpha_m = \frac{1}{N} \cdot \sum_{i=1}^N \omega_{im}$$

Again, it is  $\sum_{m=1}^K \alpha_m = 1$  and  $\sum_{i=1}^N \sum_{m=1}^K \omega_{im} = N$

# Iteration

Now, we set  $\theta_t = \theta$  and repeat the E step. This continues until convergence is reached:



# Summary of EM for multidimensional Gaussian mixture

**Initialization:** 1<sup>st</sup> guess  $\theta_t = \{\alpha_k^t, \mu_k^t, \Sigma_k^t\}$

**E-step:**

$$\omega_{ik} = \frac{\alpha_k^t \cdot f_k(\mathbf{x}_i | \mu_k^t, \Sigma_k^t)}{\sum_k \alpha_k^t \cdot f_k(\mathbf{x}_i | \mu_k^t, \Sigma_k^t)}$$

**M-step:**

$$\mu_m = \frac{\sum_{i=1}^N \omega_{im} \cdot \mathbf{x}_i}{\sum_{i=1}^N \omega_{im}}$$

$$\Sigma_m = \frac{\sum_{i=1}^N \omega_{im} \cdot (\mathbf{x}_i - \mu_m) \cdot (\mathbf{x}_i - \mu_m)^T}{\sum_{i=1}^N \omega_{im}}$$

$$\alpha_m = \frac{1}{N} \cdot \sum_{i=1}^N \omega_{im}$$

$$\sum_{k=1}^K \omega_{ik} = 1$$

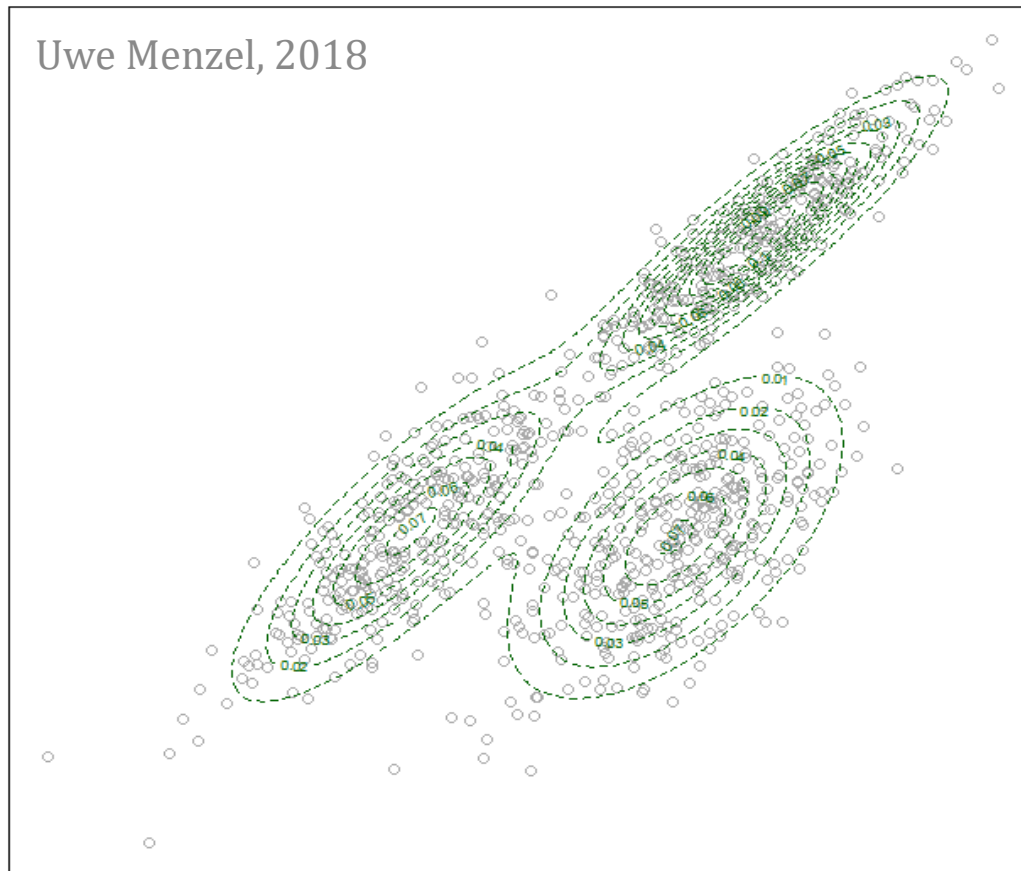
$$\sum_{i=1}^N \sum_{m=1}^K \omega_{im} = N$$

**Iterate** between E- and M-step until convergence, i.e. until  $\frac{\theta - \theta_t}{\theta_t} < \varepsilon$ .

Alternatively, check convergence of the likelihood.

## R-script, EM for 2-D Gaussian mixture

- In the R-code provided, we generate sample points for a two-dimensional Gaussian mixture by modelling the two-step experiment described above (grey points in the figure below)
- Starting with these data, we calculate the means, covariance matrices, and mixture weights using the EM algorithm presented here. The solution is indicated by the green contour lines.



# Appendix

## The Expectation-Maximization algorithm III

Uwe Menzel, 2018

[uwe.menzel@slu.se](mailto:uwe.menzel@slu.se) ; [uwe.menzel@matstat.de](mailto:uwe.menzel@matstat.de)

[www.matstat.org](http://www.matstat.org)

# Probability density for the 2-D Gaussian

**Probability density function for the d-dimensional case:**

$$f_{\mathbf{k}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \cdot \sqrt{|\boldsymbol{\Sigma}_{\mathbf{k}}|}} \cdot \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{k}})^T \cdot (\boldsymbol{\Sigma}_{\mathbf{k}})^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{k}}) \right\}$$

**Probability density function for the 2-dimensional case:**

$$\begin{aligned} f_{\mathbf{k}}(\mathbf{x}) &= f_{\mathbf{k}}(x, y) = \\ &= \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \cdot \exp \left\{ -\frac{(x-\mu_x)^2}{2\sigma_x^2} + \rho \cdot \frac{(y-\mu_y)(x-\mu_x)}{\sigma_x\sigma_y} - \frac{(y-\mu_y)^2}{2\sigma_y^2} \right\} \end{aligned}$$

**Determinant of the covariance matrix for the 2-D case:**

$$|\boldsymbol{\Sigma}_{\mathbf{k}}| = \det \begin{pmatrix} \sigma_x^2 & \rho \cdot \sigma_x\sigma_y \\ \rho \cdot \sigma_x\sigma_y & \sigma_y^2 \end{pmatrix} = \sigma_x^2\sigma_y^2 - \rho^2 \cdot \sigma_x^2\sigma_y^2 = \sigma_x^2\sigma_y^2 \cdot (1 - \rho^2)$$

## Inverse of the 2D covariance matrix

A **general expression** for the inverse of a 2-D matrix is:

$$\mathbf{A}^{-1} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{\det \mathbf{A}} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

Hence, the **inverse of the 2-D covariance matrix** is:

$$\Sigma^{-1} = \begin{pmatrix} \sigma_x^2 & \rho \cdot \sigma_x \sigma_y \\ \rho \cdot \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}^{-1} = \frac{1}{\sigma_x^2 \sigma_y^2 \cdot (1 - \rho^2)} \cdot \begin{pmatrix} \sigma_y^2 & -\rho \cdot \sigma_x \sigma_y \\ -\rho \cdot \sigma_x \sigma_y & \sigma_x^2 \end{pmatrix}$$

We still need the exponent  $(\mathbf{x} - \boldsymbol{\mu})^T \cdot (\boldsymbol{\Sigma})^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}) \Rightarrow$

**Note:** (index  $k$  suppressed to simplify notation)

$(\mathbf{x} - \boldsymbol{\mu})$  is a column vector: 2x1

$$(\mathbf{x} - \boldsymbol{\mu}) = \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix}$$

$(\mathbf{x} - \boldsymbol{\mu})^T$  is a row vector: 1x2

$$(\mathbf{x} - \boldsymbol{\mu})^T = (x - \mu_x, y - \mu_y)$$



## Exponent of the 2D probability density function

$$\begin{aligned}(\boldsymbol{\Sigma})^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}) &= \frac{1}{\sigma_x^2 \sigma_y^2 \cdot (1 - \rho^2)} \cdot \begin{pmatrix} \sigma_y^2 & -\rho \cdot \sigma_x \sigma_y \\ -\rho \cdot \sigma_x \sigma_y & \sigma_x^2 \end{pmatrix} \cdot \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix} \\ &= \frac{1}{\sigma_x^2 \sigma_y^2 \cdot (1 - \rho^2)} \cdot \begin{pmatrix} \sigma_y^2(x - \mu_x) - \rho \cdot \sigma_x \sigma_y(y - \mu_y) \\ -\rho \cdot \sigma_x \sigma_y(x - \mu_x) + \sigma_x^2(y - \mu_y) \end{pmatrix}\end{aligned}$$

It remains to multiply this with  $(\mathbf{x} - \boldsymbol{\mu})^T$  from the left side:

$$\begin{aligned}(\mathbf{x} - \boldsymbol{\mu})^T \cdot (\boldsymbol{\Sigma})^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}) &= \\ &= \frac{1}{\sigma_x^2 \sigma_y^2 \cdot (1 - \rho^2)} \begin{pmatrix} x - \mu_x \\ y - \mu_y \end{pmatrix}^T \cdot \begin{pmatrix} \sigma_y^2(x - \mu_x) - \rho \cdot \sigma_x \sigma_y(y - \mu_y) \\ -\rho \cdot \sigma_x \sigma_y(x - \mu_x) + \sigma_x^2(y - \mu_y) \end{pmatrix} \\ &= \frac{1}{\sigma_x^2 \sigma_y^2 \cdot (1 - \rho^2)} \cdot \{ \sigma_y^2(x - \mu_x)^2 - 2 \cdot \rho \cdot \sigma_x \sigma_y(y - \mu_y)(x - \mu_x) + \sigma_x^2(y - \mu_y)^2 \} \\ &= \frac{1}{(1 - \rho^2)} \cdot \left\{ \frac{(x - \mu_x)^2}{\sigma_x^2} - 2 \cdot \rho \frac{(y - \mu_y)(x - \mu_x)}{\sigma_x \sigma_y} + \frac{(y - \mu_y)^2}{\sigma_y^2} \right\}\end{aligned}$$

## Derivative of the quadratic form $\mathbf{x}^T \mathbf{A} \mathbf{x}$

If the matrix  $\mathbf{A}$  is symmetric ( $\mathbf{A}^T = \mathbf{A}$ ) and independent of the vector  $\mathbf{a}$ , the following is generally valid:

$$\frac{\partial}{\partial \mathbf{a}} (\mathbf{a}^T \mathbf{A} \mathbf{a}) = 2\mathbf{A} \mathbf{a}$$

Therefore:

$$\frac{\partial}{\partial \underbrace{(\mathbf{x}_i - \boldsymbol{\mu}_m)}_{\mathbf{a}}} \left[ \underbrace{(\mathbf{x}_i - \boldsymbol{\mu}_m)^T}_{\mathbf{a}^T} \cdot \underbrace{(\boldsymbol{\Sigma}_m)^{-1}}_{\mathbf{A}} \cdot \underbrace{(\mathbf{x}_i - \boldsymbol{\mu}_m)}_{\mathbf{a}} \right] = 2 \cdot \underbrace{(\boldsymbol{\Sigma}_m)^{-1}}_{\mathbf{A}} \underbrace{(\mathbf{x}_i - \boldsymbol{\mu}_m)}_{\mathbf{a}}$$

Using the chain rule, we get:

$$\frac{\partial}{\partial \boldsymbol{\mu}_m} = \frac{\partial (\mathbf{x}_i - \boldsymbol{\mu}_m)}{\partial \boldsymbol{\mu}_m} \cdot \frac{\partial}{\partial (\mathbf{x}_i - \boldsymbol{\mu}_m)} = (-\mathbf{I}) \cdot \frac{\partial}{\partial (\mathbf{x}_i - \boldsymbol{\mu}_m)} \quad \text{\textit{I is the identity matrix}}$$

Finally, we get:

$$\frac{\partial}{\partial \boldsymbol{\mu}_m} \left[ (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\boldsymbol{\Sigma}_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right] = -2 \cdot (\boldsymbol{\Sigma}_m)^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_m)$$

## Derivation of $Q_1$ for $(\Sigma_m)^{-1}$

Instead of deriving w.r.t.  $\Sigma_m$

$$\frac{\partial Q_1}{\partial \Sigma_m} = -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial \Sigma_m} \left\{ \log |\Sigma_m| + (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$$

we derive w.r.t.  $(\Sigma_m)^{-1}$ , so that we have to calculate :

$$\frac{\partial Q_1}{\partial (\Sigma_m)^{-1}} = -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ \log |\Sigma_m| + (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$$

## Derivation of $Q_1$ for $(\Sigma_m)^{-1}$

$$\frac{\partial Q_1}{\partial (\Sigma_m)^{-1}} = -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ \log |\Sigma_m| + (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$$

Calculation of  $\frac{\partial}{\partial (\Sigma_m)^{-1}} \log |\Sigma_m|$        $|\Sigma_m|$  is the  
determinant of  $\Sigma_m$

Since  $|\mathbf{A}^{-1}| = 1/|\mathbf{A}|$  we can write

$$\frac{\partial}{\partial (\Sigma_m)^{-1}} \log |\Sigma_m| = -\frac{\partial}{\partial (\Sigma_m)^{-1}} \log |(\Sigma_m)^{-1}|$$

Furthermore, if  $\mathbf{A}$  is symmetric, we have:  $\frac{\partial}{\partial \mathbf{A}} \log |\mathbf{A}| = \mathbf{A}^{-1}$  (\*)

so that  $\frac{\partial}{\partial (\Sigma_m)^{-1}} \log |\Sigma_m| = \underline{\underline{-\Sigma_m}}$

(\*): see "The Matrix Cookbook" by Petersen & Pedersen  
<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>)

## Derivation of $Q_1$ for $(\Sigma_m)^{-1}$

$$\frac{\partial Q_1}{\partial (\Sigma_m)^{-1}} = -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ \log |\Sigma_m| + (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$$

Calculation of  $\frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$

In general we have :  $\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T$  (\*)

Using this, we get:

$$\frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\} = (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T$$

(\*) see “[The Matrix Cookbook](https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf)” by Petersen and Pedersen, eqn. (72)  
<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

## Derivation of $Q_1$ for $(\Sigma_m)^{-1}$

In **summary**, we obtained

$$\frac{\partial}{\partial (\Sigma_m)^{-1}} \log |\Sigma_m| = -\Sigma_m$$

$$\frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\} = (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T$$

so that the final result becomes

$$\begin{aligned} \frac{\partial Q_1}{\partial (\Sigma_m)^{-1}} &= -\frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \frac{\partial}{\partial (\Sigma_m)^{-1}} \left\{ \log |\Sigma_m| + (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \cdot (\Sigma_m)^{-1} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\} \\ &= \frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \left\{ \Sigma_m - (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \right\} \end{aligned}$$

$$\frac{\partial Q_1}{\partial (\Sigma_m)^{-1}} = \frac{1}{2} \sum_{i=1}^N \omega_{im} \cdot \left\{ \Sigma_m - (\mathbf{x}_i - \boldsymbol{\mu}_m) \cdot (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \right\}$$