

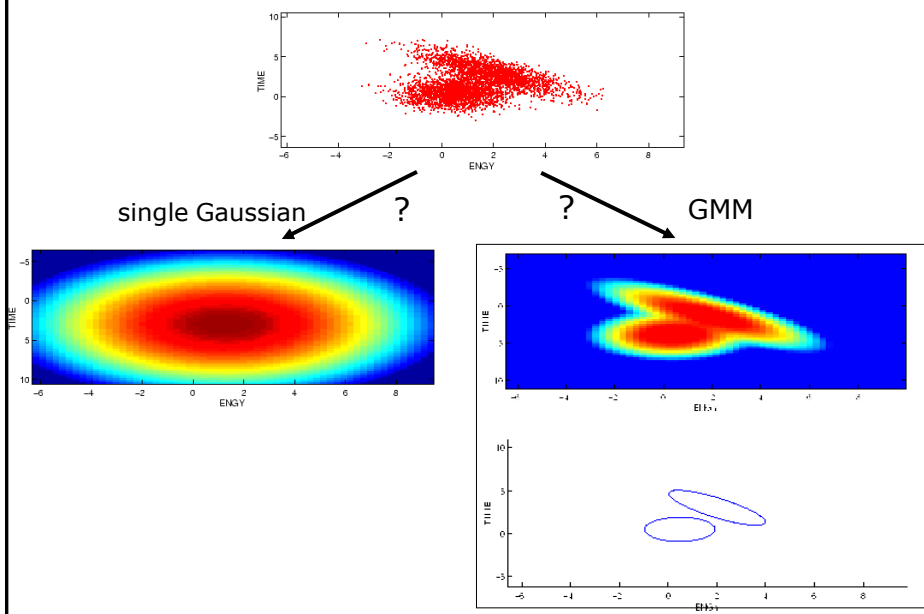
# Density Estimation

1

- Parametric techniques
  - Maximum Likelihood
  - Maximum A Posteriori
  - Bayesian Inference
  - **Gaussian Mixture Models (GMM)**
    - EM-Algorithm
- Non-parametric techniques
  - Histogram
  - Parzen Windows
  - k-nearest-neighbor rule

# GMM Applications

2



## GMM Applications

3

### Density estimation

Observed data from a complex but unknown probability distribution.

Can we describe this data with a few parameters ?

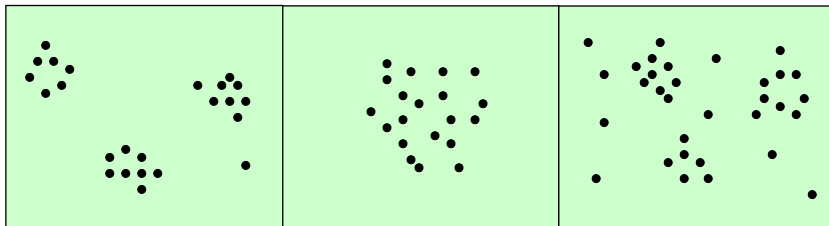
Which (new) samples are unlikely to come from this unknown distribution (Outlier detection )?

## GMM Applications

4

### Clustering

Observations from  $K$  classes. Each class produces samples from a multivariate normal distribution.  
Which observations belong to which class ?



Sometimes  
easy

Sometimes  
impossible

Often possible  
but not clear-cut

## GMM: Definition

5

- Mixture models are linear combinations of densities:

$$p(x | \Theta) = \sum_{i=1}^K c_i p(x | \Theta_i)$$

$$\text{with } \sum_{i=1}^K c_i = 1 \quad , \quad \int_x p(x | \Theta_i) dx = 1$$

- Capable of approximating almost any complex and irregularly shaped distributions ( *K might get big* )!
- For Gaussian mixtures:  
 $\Theta_i = \{\mu_i, \Sigma_i\}, \quad \Rightarrow \quad p(x | \Theta_i) \equiv N(\mu_i, \Sigma_i)$

## Sampling a GMM

6

- How to generate a random variable according to a known GMM

$$p(x) = \sum_i^K c_i N(\mu_i, \Sigma_i)$$

Assume that each data point is generated according to the following recipe:

1. Pick a component ( $i \in [1..K]$ ) at random.  
Choose component  $i$  with probability  $c_i$ .
2. Sample data point  $\sim N(\mu_i, \Sigma_i)$ .

In the end, we might not know which data points came from which component (unless someone kept track during the sampling process)!

# Learning a GMM

7

## Recall ML-estimation

We have:

A density function  $p(\cdot; \Theta)$  governed by a set of unknown parameters  $\Theta$ .

A data set of size  $N$  drawn from this distribution

$X = \{x_1, \dots, x_N\}$

We wish:

to obtain the parameters best explaining data  $X$   
by maximizing the log-likelihood function:

$$L(\Theta) = \ln p(X; \Theta)$$

$$\Theta^* = \arg \max_{\Theta} L(\Theta)$$

# Learning a GMM

8

- For a single Gaussian distribution this is simple to solve. We have an analytical solution.
- Unfortunately for many problems (including GMM) it is not possible to find analytical expressions.
  - Resort to classical optimization techniques ?
  - Possible but there is a better way:  
EM - Algorithm (Expectation-Maximization)

## Expectation Maximization ( **EM** )

9

- General method for finding ML-estimates in the case of incomplete or missing data (GMM's are one application).
- Usually used when:
  - the observation is actually incomplete; some values are missing from the data set.
  - the likelihood function is analytically intractable but can be simplified by assuming the existence of additional but missing (so-called hidden/latent) parameters.

The latter technique is used for GMMs. Think of each data point as having a hidden label specifying the component it belongs to. These component labels are the latent parameters.

## General EM procedure

10

### The EM setting

Observed data set (*incomplete*):  $X$

Assume a complete data set exists:  $Z = (X, Y)$

$Z$  has a joint density function:

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

Define the *complete-data* log-likelihood function:

$$L(\Theta | Z) = L(\Theta | X, Y) = \ln p(X, Y | \Theta)$$

Our aim is to find a  $\Theta$  that maximizes this function.

## General EM procedure

11

- But: We cannot simply maximize  $L(\Theta | X, Y) = \ln p(X, Y | \Theta)$  because  $Y$  is not known.
  - $L(\Theta | X, Y)$  is in fact a random variable:
    - $Y$  can be assumed to come from some distribution  $f(y | X, \Theta)$
    - That is,  $L(\Theta | X, Y)$  can be interpreted as a function where  $X$  and  $\Theta$  are constant and  $Y$  is a random variable.
- The EM will compute a new, auxiliary function, based on  $L$ , that can be maximized instead.
- Let's assume we already have a reasonable estimate for the parameters:  $\Theta^{(i-1)}$ .

## General EM procedure

12

- EM uses **an auxiliary function**:

$$Q(\Theta, \Theta^{(i-1)}) = E \left[ \ln p(X, Y | \Theta) \mid X, \Theta^{(i-1)} \right]$$

How to read this:

- $X$  and  $\Theta^{(i-1)}$  are constants,
- $\Theta$  is a simple variable (the function argument),
- $Y$  is a random variable governed by distribution  $f$ .
- The task is to rewrite  $Q$  and perform some calculations to make it a **fully determined** function.
- $Q$  is the **expected value** of the complete-data log-likelihood w.r.t. to missing data ( $Y$ ), observed data ( $X$ ) and current parameter estimates ( $\Theta^{(i-1)}$ ).

————— This is called the **E-step** (*expectation-step*) —————

## General EM procedure

13

- $Q$  can be rewritten by means of the *marginal* distribution  $f$ :

If  $y$  is a continuous random variable:

$$Q(\Theta, \Theta^{(i-1)}) = E \left[ \ln p(X, Y | \Theta) \mid X, \Theta^{(i-1)} \right]$$
$$= \int \ln p(X, \mathbf{y} | \Theta) \cdot f(\mathbf{y} | X, \Theta^{(i-1)}) dy$$

If  $y$  is a discrete random variable:

$$Q(\Theta, \Theta^{(i-1)}) = E \left[ \ln p(X, Y | \Theta) \mid X, \Theta^{(i-1)} \right]$$
$$= \sum_y \ln p(X, \mathbf{y} | \Theta) f(\mathbf{y} | X, \Theta^{(i-1)})$$

Think of this as  
the expected  
value of a  
function of  $Y$   
 $E[g(Y)]$

Evaluate  $f(y | X, \Theta^{(i-1)})$ , using the current estimate  $\Theta^{(i-1)}$ .

Now  $Q$  is fully determined and we can use it!

## General EM procedure

14

- In a **second** step  $Q$  is used to obtain a better set of parameters  $\Theta$ :

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

This is called the **M-step** (*maximization-step*)

- Both E- and M-steps are repeated until convergence.
  - In each E-Step, we find a new auxiliary function  $Q$
  - In each M-Step, we find a new parameter set  $\Theta$

# General EM algorithm

15

## Summary of the general EM algorithm (see also Bishop, p.440)

1. Choose an initial setting for the parameters  $\Theta^{(i-1)}$ .
2. E-step: evaluate  $f(y | X, \Theta^{(i-1)})$ ,  
 plug it into  $Q(\Theta, \Theta^{(i-1)}) = \int f(y | X, \Theta^{(i-1)}) \ln p(X, Y | \Theta) dy$   
 to obtain a fully determined auxiliary function
3. M-step: evaluate  $\Theta^{(i)}$  given by  $\Theta^{(i)} = \arg \max_{\Theta} Q(\Theta, \Theta^{(i-1)})$
4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let  $\Theta^{(i-1)} \leftarrow \Theta^{(i)}$  and return to step 2.

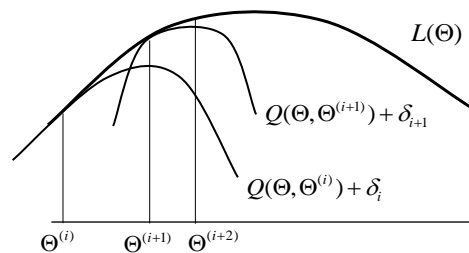
# General EM Illustration

16

## Iterative majorisation

Aim of EM: Find local maximum of function  $L(\Theta)$  by using auxiliary function  $Q(\Theta, \Theta^{(i)})$ .

How does this work?



- Q touches  $L$  at point  $[\Theta^{(i)}, L(\Theta^{(i)})]$  and lies everywhere below  $L$ .
- Maximize auxiliary function.
- The position of the maximum  $\Theta^{(i+1)}$  gives a value of  $L$  which is greater than in the previous iteration.
- Repeat this scheme with new auxiliary function until convergence.



## General EM Summary

17

- Iterative algorithm for ML-estimation of systems with hidden/missing values.
- Calculates expectation for hidden values based on observed data and joint distribution.
- Slow but guaranteed convergence.
- May get „stuck“ in local maximum.
- There is no general EM implementation. The details of both steps depend very much on the particular application.

## Application: EM for Mixture Models

18

- Our probabilistic model is now:

$$p(x | \Theta) = \sum_{i=1}^M c_i p_i(x | \theta_i)$$

with parameters:  $\Theta = (c_1, \dots, c_M, \theta_1, \dots, \theta_M)$

such that:  $\sum_{i=1}^M c_i = 1$

- That is, we have  $M$  component densities  $p_i$  (of the same family) combined through  $M$  mixing coefficients  $c_i$ .

## EM for Mixture Models

19

- The incomplete-data log-likelihood becomes (remember we assume  $X$  is iid):

$$L(\Theta | X) = \ln \prod_{i=1}^N p(x_i | \Theta) = \sum_{i=1}^N \ln \left( \sum_{j=1}^M c_j p_j(x_i | \theta_j) \right)$$

- Difficult to optimize with log of sum
- Now let's try the EM-trick:
  - Consider  $X$  as incomplete.
  - Introduce unobserved data  $Y = \{y_i\}_{i=1}^N$  whose values indicate which component of the MM generated each data item.
  - That is,  $y_i \in 1, \dots, M$  and  $y_i = k$  if the  $i$ -th sample stems from the  $k$ -th component.

## EM for Mixture Models

20

- **If** we knew the values of  $Y$ , the log likelihood would simplify to:

$$\begin{aligned} L(\Theta | X, Y) &= \ln p(X, Y | \Theta) \\ &= \sum_{i=1}^N \ln (p(x_i | y_i, \Theta) p(y_i | \Theta)) = \sum_{i=1}^N \ln (c_{y_i} p_{y_i}(x_i | \theta_{y_i})) \end{aligned}$$

Could apply standard optimization techniques

- **But** we don't know  $Y$ , so we follow the EM-procedure:
  1. Start with an initial guess of the mixture parameters:  
 $\Theta^g = (c_1^g, \dots, c_M^g, \theta_1^g, \dots, \theta_M^g)$
  2. Find an expression for the marginal density function of the unobserved data  $p(y | X, \Theta)$ :

## EM for Mixture Models

21

$$p(\mathbf{y} | \mathbf{X}, \Theta^g) = \prod_{i=1}^N p(y_i | x_i, \Theta^g)$$

$y_i$  is the (unknown) component label of data point  $x_i$ .

Using Bayes's rule, we get:

$$\begin{aligned} p(y_i | x_i, \Theta^g) &= \frac{p(x_i | y_i, \Theta^g) p(y_i)}{p(x_i | \Theta^g)} \\ &= \frac{c_{y_i} p_{y_i}(x_i | \theta_{y_i}^g)}{p(x_i | \Theta^g)} = \frac{c_{y_i} p_{y_i}(x_i | \theta_{y_i}^g)}{\sum_{k=1}^M c_k p_k(x_i | \theta_k^g)} \end{aligned}$$

- Using guessed parameters, we obtained the desired marginal density function.
- This can now be substituted in Q (i.e. in the E-step).

## EM for **Gaussian** Mixtures

22

- For our mixture model, the **E-step** is:

$$Q(\Theta, \Theta^g) = \sum_y \ln(L(\Theta | \mathbf{X}, \mathbf{y})) p(\mathbf{y} | \mathbf{X}, \Theta^g)$$

Substituted marginal hidden data density

- The **M-step** is to find a parameter set  $\Theta^{new}$  that maximizes  $Q$ .
- But for Gaussian mixtures, there is no need to deal with  $Q$  in the above form!

$$\Theta^{new} = \operatorname{argmax}_{\Theta} \sum_y L(\Theta | \mathbf{X}, \mathbf{y}) p(\mathbf{y} | \mathbf{X}, \Theta^g)$$

Here it is not necessary to deal with this directly

- Instead, a set of simple formulas for updating  $\Theta$  can be used.

## EM for Gaussian Mixtures

23

3. Compute parameters  $\Theta^{new}$ , using update formulas (perform E- and M-step simultaneously):

Update formulas

$$c_k^{new} = \frac{1}{N} \sum_{i=1}^N p(k | x_i, \Theta^g) \quad \mu_k^{new} = \frac{\sum_{i=1}^N x_i p(k | x_i, \Theta^g)}{\sum_{i=1}^N p(k | x_i, \Theta^g)}$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^N p(k | x_i, \Theta^g) (x_i - \mu_k^{new})(x_i - \mu_k^{new})^T}{\sum_{i=1}^N p(k | x_i, \Theta^g)}$$

Plug in the expression found in previous step (k = label of the k-th component)

These formulas are derived from  $Q(\Theta, \Theta^g)$ .

## EM for Mixture Models

24

### Derivation of the update formulas

Q in its initial form:

$$Q(\Theta, \Theta^g) = \sum_{\mathbf{y}} \ln(L(\Theta | \mathcal{X}, \mathbf{y})) p(\mathbf{y} | \mathcal{X}, \Theta^g)$$

Substituted marginal hidden data density

- After a lot of simplification we arrive at an equation where the  $c_k$  and  $\theta_k$  are expressed independently:

$$Q(\Theta, \Theta^g) = \sum_{k=1}^M \sum_{i=1}^N \ln(c_k) p(k | x_i, \Theta^g) + \sum_{k=1}^M \sum_{i=1}^N \ln(p_k(x_i | \theta_k)) p(k | x_i, \Theta^g)$$

Get formula for  $c_k$  from this part

Get formulas for  $\theta_k$  from this part

Formula for  $c_k$ , after further simplification

with  $c_k = \frac{1}{N} \sum_{i=1}^N p(k | x_i, \Theta^g)$

## EM for Gaussian Mixtures

25

- Formula for  $c_k$  (previous slide) is valid for any mixture model, not just Gaussian.
- Formulas for  $\theta_k$  will be specific to the Gaussian mixture.
- For a d-dimensional Gaussian component, use

$$\theta = (\mu, \Sigma) \quad p_k(x | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{d/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}$$

Plug this into the expression on the previous slide

- Take the derivatives of the resulting expression with respect to  $\mu_k$  and  $\Sigma_k$  (*very technical*).
- Set the derivatives to zero, then solve for  $\mu_k$  and  $\Sigma_k$ .  
→ The results are the update formulas for  $\mu_k^{\text{new}}$  and  $\Sigma_k^{\text{new}}$ .

## EM for Gaussian Mixture Models

26

**Summary of the algorithm for GMM** (see Bishop, p.438):

1. Initialize the parameters  $\Theta^{\text{old}} = (c_1 \dots c_M, \mu_1 \dots \mu_M, \Sigma_1 \dots \Sigma_M)$
2. E-step: evaluate the responsibilities of each component for all data points:

$$p(k | x_i, \Theta^{\text{old}}) = \frac{c_k p_k(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^M c_j p_j(x_i | \mu_j, \Sigma_j)}$$

Responsibility of the k-th component for the i-th data point

No need to compute  $Q(\Theta, \Theta^{(i-1)})$  explicitly!

## EM for Gaussian Mixture Models

27

3. E-step/M-step: Update the parameters

$$c_k^{new} = \frac{1}{N} \sum_{i=1}^N p(k | x_i, \Theta^{old}) \quad \mu_k^{new} = \frac{\sum_{i=1}^N p(k | x_i, \Theta^{old}) x_i}{\sum_{i=1}^N p(k | x_i, \Theta^{old})}$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^N p(k | x_i, \Theta^{old}) (x_i - \mu_k^{new})(x_i - \mu_k^{new})^T}{\sum_{i=1}^N p(k | x_i, \Theta^{old})}$$

4. Evaluate the log likelihood

$$\ln p(X | \Theta) = \sum_{i=1}^N \ln \left( \sum_{k=1}^M c_k p_k(x_i | \mu_k, \Sigma_k) \right)$$

and check it for convergence. If the convergence criterion is not satisfied, return to step 2.

## Relation to k-means

28

- Let  $c_k = 1/M$  and  $\Sigma_k = \sigma^2 I$
- k-means procedure:
  1. Random initialize M cluster centers.
  2. Assign each data point to a cluster according to the minimum distance criterion:

$$p(k | x_i) = \begin{cases} 1 & \text{if } \forall j \| x_i - \mu_k \| \leq \| x_i - \mu_j \| \\ 0 & \text{otherwise} \end{cases}$$

3. Re-calculate cluster centers:

$$\mu_k^{new} = \frac{\sum_{i=1}^N p(k | x_i) x_i}{\sum_{i=1}^N p(k | x_i)}$$

4. Go to step 2 until no change in cluster centers.

## Relation to k-means

29

- GMM is referred to as soft clustering
  - Probability  $p(k/x_i, \Theta)$  indicates the responsibility of the k-th component for the i-th observation (i.e. the posterior prob. that the i-th observation comes from the k-th component).
  - For each point  $x_i$ , GMM produces smooth posterior. From these, one can find cluster label for each  $x_i$ :  $C(x_i) = \underset{k}{\operatorname{argmax}} p(k/x_i, \Theta)$
- k-means is a hard clustering method
  - The responsibility of can only be 1 or 0 .

## GMM: Open questions

30

- How many components are required ?
  - Answer is highly problem dependent.
  - One possibility: Try different numbers, then choose model (number) which gives best performance on a validation data set .
- Which initial parameters to use ?
  - Same here: in general we don't know where to look for global maximum.
  - Obvious approaches:
    1. Perform k-means to obtain initial  $\mu$ 's.
    2. Try different random values and choose the ones which lead to maximal likelihood.

## GMM/EM Resources

31

- **J. A. Bilmes et al : A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models (1998)**
- **GMMBAYES - Gaussian Mixture Model Methods Matlab-Toolbox**  
<http://www.it.lut.fi/project/gmbayes/>
- **Gaussian Mixtures Demo Applet**  
<http://lcn.epfl.ch/tutorial/english/gaussian/html/index.html>