

# Introduction to Machine Learning

## Expectation Maximization

Andres Mendez-Vazquez

June 5, 2018

# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# Outline

## 1 Introduction

- **Maximum-Likelihood**
- Expectation Maximization
- Examples of Applications of EM

## 2 Incomplete Data

- Introduction
- Using the Expected Value
- Analogy

## 3 Derivation of the EM-Algorithm

- Hidden Features
  - Proving Concavity
- Using the Concave Functions for Approximation
- From The Concave Function to the EM
- The Final Algorithm
- Notes and Convergence of EM

## 4 Finding Maximum Likelihood Mixture Densities

- The Beginning of The Process
- Bayes' Rule for the components
  - Mixing Parameters
- Maximizing  $Q$  using Lagrange Multipliers
  - In Our Case
- Example on Mixture of Gaussian Distributions
- The EM Algorithm



# Maximum-Likelihood

We have a density function  $p(\mathbf{x}|\Theta)$

Assume that we have a data set of size  $N$ ,  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

- This data is known as **evidence**.

We assume in addition that:

The vectors are independent and identically distributed (i.i.d.) with distribution  $p$  under parameter  $\theta$ .



# Maximum-Likelihood

We have a density function  $p(\mathbf{x}|\Theta)$

Assume that we have a data set of size  $N$ ,  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

- This data is known as **evidence**.

We assume in addition that

The vectors are independent and identically distributed (i.i.d.) with distribution  $p$  under parameter  $\theta$ .



# What Can We Do With The Evidence?

We may use the Bayes' Rule to estimate the parameters  $\theta$

$$p(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (1)$$

Or, given a new observation  $\bar{x}$

$$p(\bar{x}|\mathcal{X}) \quad (2)$$

i.e. to compute the probability of the new observation being supported by the evidence  $\mathcal{X}$ .

Thus:

The former represents parameter estimation and the latter data prediction.



# What Can We Do With The Evidence?

We may use the Bayes' Rule to estimate the parameters  $\theta$

$$p(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (1)$$

Or, given a new observation  $\tilde{x}$

$$p(\tilde{x}|\mathcal{X}) \quad (2)$$

i.e. to compute the probability of the new observation being supported by the evidence  $\mathcal{X}$ .

Notes

The former represents parameter estimation and the latter data prediction.



# What Can We Do With The Evidence?

We may use the Bayes' Rule to estimate the parameters  $\theta$

$$p(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (1)$$

Or, given a new observation  $\tilde{x}$

$$p(\tilde{x}|\mathcal{X}) \quad (2)$$

i.e. to compute the probability of the new observation being supported by the evidence  $\mathcal{X}$ .

Thus

The former represents parameter estimation and the latter data prediction.





# Focusing First on the Estimation of the Parameters $\theta$

We can interpret the Bayes' Rule

$$p(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (3)$$

interpreted as

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \quad (4)$$

Thus, we want

$$\text{likelihood} = P(\mathcal{X}|\Theta)$$



# Focusing First on the Estimation of the Parameters $\theta$

We can interpret the Bayes' Rule

$$p(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (3)$$

Interpreted as

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior}}{\textit{evidence}} \quad (4)$$

Thus, we want

$$\textit{likelihood} = P(\mathcal{X}|\Theta)$$



# Focusing First on the Estimation of the Parameters $\theta$

We can interpret the Bayes' Rule

$$p(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (3)$$

Interpreted as

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior}}{\textit{evidence}} \quad (4)$$

Thus, we want

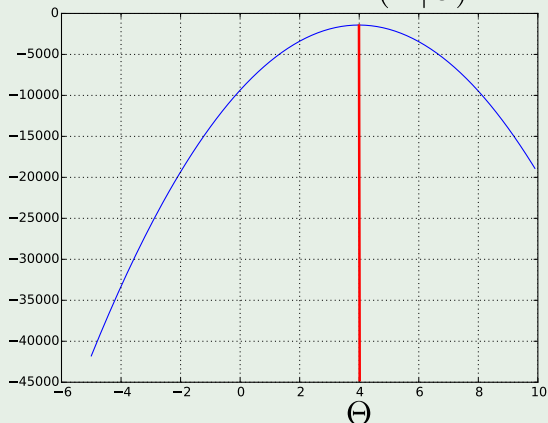
$$\textit{likelihood} = P(\mathcal{X}|\Theta)$$



## What we want...

We want to maximize the likelihood as a function of  $\theta$

$$\text{likelihood} = P(\mathcal{X}|\Theta)$$



# Maximum-Likelihood

We have

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \prod_{i=1}^N p(\mathbf{x}_i | \Theta) \quad (5)$$

Also known as the likelihood function.

Because multiplication of quantities  $p(\mathbf{x}_i | \Theta) \leq 1$  can be problematic

$$\mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(\mathbf{x}_i | \Theta) = \sum_{i=1}^N \log p(\mathbf{x}_i | \Theta) \quad (6)$$



# Maximum-Likelihood

We have

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \Theta) = \prod_{i=1}^N p(\mathbf{x}_i | \Theta) \quad (5)$$

Also known as the likelihood function.

Because multiplication of quantities  $p(\mathbf{x}_i | \Theta) \leq 1$  can be problematic

$$\mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(\mathbf{x}_i | \Theta) = \sum_{i=1}^N \log p(\mathbf{x}_i | \Theta) \quad (6)$$



# Maximum-Likelihood

We want to find a  $\Theta$

$$\Theta^* = \operatorname{argmax}_{\Theta} \mathcal{L}(\Theta | \mathcal{X}) \quad (7)$$

The classic method

$$\frac{\partial \mathcal{L}(\Theta | \mathcal{X})}{\partial \theta_i} = 0 \quad \forall \theta_i \in \Theta \quad (8)$$



# Maximum-Likelihood

We want to find a  $\Theta^*$

$$\Theta^* = \operatorname{argmax}_{\Theta} \mathcal{L}(\Theta | \mathcal{X}) \quad (7)$$

The classic method

$$\frac{\partial \mathcal{L}(\Theta | \mathcal{X})}{\partial \theta_i} = 0 \quad \forall \theta_i \in \Theta \quad (8)$$





# What happened if we have incomplete data

## Data could have been split

①  $\mathcal{X}$  = observed data or **incomplete** data

②  $\mathcal{Y}$  = unobserved data



# What happened if we have incomplete data

## Data could have been split

- 1  $\mathcal{X}$  = observed data or **incomplete** data
- 2  $\mathcal{Y}$  = unobserved data

For this type of problems

We have the famous Expectation Maximization (EM)



# What happened if we have incomplete data

Data could have been split

- 1  $\mathcal{X}$  = observed data or **incomplete** data
- 2  $\mathcal{Y}$  = unobserved data

For this type of problems

We have the famous Expectation Maximization (EM)



# Outline

## 1 Introduction

- Maximum-Likelihood
- **Expectation Maximization**
- Examples of Applications of EM

## 2 Incomplete Data

- Introduction
- Using the Expected Value
- Analogy

## 3 Derivation of the EM-Algorithm

- Hidden Features
  - Proving Concavity
- Using the Concave Functions for Approximation
- From The Concave Function to the EM
- The Final Algorithm
- Notes and Convergence of EM

## 4 Finding Maximum Likelihood Mixture Densities

- The Beginning of The Process
- Bayes' Rule for the components
  - Mixing Parameters
- Maximizing  $Q$  using Lagrange Multipliers
  - In Our Case
- Example on Mixture of Gaussian Distributions
- The EM Algorithm



# The Expectation Maximization

## The EM algorithm

It was first developed by Dempster et al. (1977).



# The Expectation Maximization

## The EM algorithm

It was first developed by Dempster et al. (1977).

## Its popularity comes from the fact

It can estimate an underlying distribution when data is incomplete or has missing values.



# The Expectation Maximization

## The EM algorithm

It was first developed by Dempster et al. (1977).

## Its popularity comes from the fact

It can estimate an underlying distribution when data is incomplete or has missing values.

## Two main applications

- 1 When missing values exists.
- 2 When a likelihood function can be simplified by assuming extra parameters that are missing or hidden.



# The Expectation Maximization

## The EM algorithm

It was first developed by Dempster et al. (1977).

## Its popularity comes from the fact

It can estimate an underlying distribution when data is incomplete or has missing values.

## Two main applications

- 1 When missing values exists.
- 2 When a likelihood function can be simplified by assuming extra parameters that are **missing** or **hidden**.





# Outline

## 1 Introduction

- Maximum-Likelihood
- Expectation Maximization
- **Examples of Applications of EM**

## 2 Incomplete Data

- Introduction
- Using the Expected Value
- Analogy

## 3 Derivation of the EM-Algorithm

- Hidden Features
  - Proving Concavity
- Using the Concave Functions for Approximation
- From The Concave Function to the EM
- The Final Algorithm
- Notes and Convergence of EM

## 4 Finding Maximum Likelihood Mixture Densities

- The Beginning of The Process
- Bayes' Rule for the components
  - Mixing Parameters
- Maximizing  $Q$  using Lagrange Multipliers
  - In Our Case
- Example on Mixture of Gaussian Distributions
- The EM Algorithm



# Clustering

Given a series of data sets

Given the fact that Radial Gaussian Functions are Universal Approximators

- Samples  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  are the visible parameters
- The Gaussian distributions generating each of the samples are the hidden parameters

Then, we model the cluster as a mixture of Gaussians



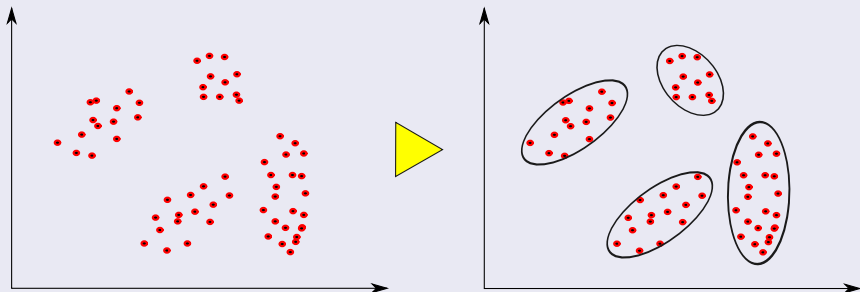
# Clustering

Given a series of data sets

Given the fact that Radial Gaussian Functions are Universal Approximators

- Samples  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  are the visible parameters
- The Gaussian distributions generating each of the samples are the hidden parameters

Then, we model the cluster as a mixture of Gaussian's



# Natural Language Processing

## Unsupervised induction of probabilistic context-free grammars

Here given a series of words  $o_1, o_2, o_3, \dots$  and normalized Context-Free Grammar

- We want to know the probabilities of each rule  $P(i \rightarrow jk)$

### This

- Here the you have two variables:
  - ▶ The Visible Ones: The sequence of words
  - ▶ The Hidden Ones: The rule that produces the possible sequence  
 $o_i \rightarrow o_j$



# Natural Language Processing

## Unsupervised induction of probabilistic context-free grammars

Here given a series of words  $o_1, o_2, o_3, \dots$  and normalized Context-Free Grammar

- We want to know the probabilities of each rule  $P(i \rightarrow jk)$

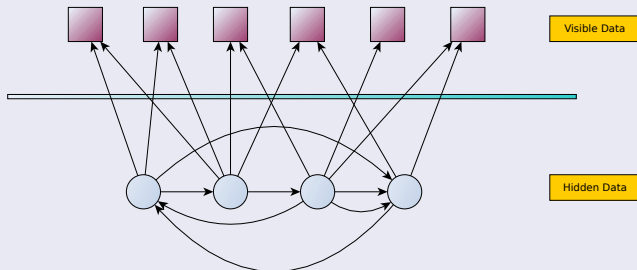
## Thus

- Here the you have two variables:
  - ▶ The Visible Ones: The sequence of words
  - ▶ The Hidden Ones: The rule that produces the possible sequence
$$o_i \rightarrow o_j$$



# Natural Language Processing

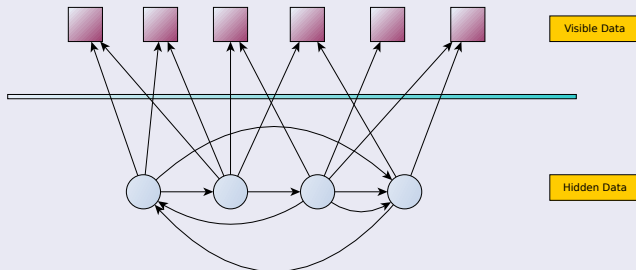
## Baum-Welch Algorithm for Hidden Markov Models



- Hidden Variables: The circular nodes producing the data
- Visible Variables: The square nodes representing the samples.

# Natural Language Processing

## Baum-Welch Algorithm for Hidden Markov Models



### Here

- Hidden Variables: The circular nodes producing the data
- Visible Variables: The square nodes representing the samples.

# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - **Introduction**
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# Incomplete Data

We assume the following

Two parts of data

①  $\mathcal{X}$  = observed data or incomplete data

②  $\mathcal{Y}$  = unobserved data



# Incomplete Data

We assume the following

Two parts of data

①  $\mathcal{X}$  = observed data or **incomplete** data

②  $\mathcal{Y}$  = unobserved data

Thus

$$\mathcal{Z} = (\mathcal{X}, \mathcal{Y}) = \text{Complete Data} \quad (9)$$



# Incomplete Data

We assume the following

Two parts of data

- 1  $\mathcal{X}$  = observed data or **incomplete** data
- 2  $\mathcal{Y}$  = unobserved data

Thus

$$\mathcal{Z} = (\mathcal{X}, \mathcal{Y}) = \text{Complete Data} \quad (9)$$

Thus, we have the following probability

$$p(\mathcal{z}|\Theta) = p(\mathcal{x}, \mathcal{y}|\Theta) = p(\mathcal{y}|\mathcal{x}, \Theta) p(\mathcal{x}|\Theta) \quad (10)$$



# Incomplete Data

We assume the following

Two parts of data

- 1  $\mathcal{X}$  = observed data or **incomplete** data
- 2  $\mathcal{Y}$  = unobserved data

Thus

$$\mathcal{Z} = (\mathcal{X}, \mathcal{Y}) = \text{Complete Data} \quad (9)$$

Thus, we have the following probability

$$p(\mathcal{Z}|\Theta) = p(\mathcal{X}, \mathcal{Y}|\Theta) = p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \quad (10)$$



# Incomplete Data

We assume the following

Two parts of data

- 1  $\mathcal{X}$  = observed data or **incomplete** data
- 2  $\mathcal{Y}$  = unobserved data

Thus

$$\mathcal{Z} = (\mathcal{X}, \mathcal{Y}) = \text{Complete Data} \quad (9)$$

Thus, we have the following probability

$$p(\mathbf{z}|\Theta) = p(\mathbf{x}, \mathbf{y}|\Theta) = p(\mathbf{y}|\mathbf{x}, \Theta) p(\mathbf{x}|\Theta) \quad (10)$$



# New Likelihood Function

## The New Likelihood Function

$$\mathcal{L}(\Theta|\mathcal{Z}) = \mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta) \quad (11)$$

**Note:** The complete data likelihood.

This, we have

$$\mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta) = p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \quad (12)$$

# New Likelihood Function

## The New Likelihood Function

$$\mathcal{L}(\Theta|\mathcal{Z}) = \mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta) \quad (11)$$

**Note:** The complete data likelihood.

Thus, we have

$$\mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta) = p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \quad (12)$$

Did you notice?

- $p(\mathcal{X}|\Theta)$  is the likelihood of the observed data.

# New Likelihood Function

## The New Likelihood Function

$$\mathcal{L}(\Theta|\mathcal{Z}) = \mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta) \quad (11)$$

**Note:** The complete data likelihood.

## Thus, we have

$$\mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = p(\mathcal{X}, \mathcal{Y}|\Theta) = p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \quad (12)$$

## Did you notice?

- $p(\mathcal{X}|\Theta)$  is the likelihood of the observed data.
- $p(\mathcal{Y}|\mathcal{X}, \Theta)$  is the likelihood of the no-observed data under the observed data!!!



# Rewriting

This can be rewritten as

$$\mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = h_{\mathcal{X}, \Theta}(\mathcal{Y}) \quad (13)$$

This basically signify that  $\mathcal{X}, \Theta$  are constant and the only random part is  $\mathcal{Y}$ .

In addition

$$\mathcal{L}(\Theta|\mathcal{X}) \quad (14)$$

It is known as the incomplete-data likelihood function.



# Rewriting

This can be rewritten as

$$\mathcal{L}(\Theta|\mathcal{X}, \mathcal{Y}) = h_{\mathcal{X}, \Theta}(\mathcal{Y}) \quad (13)$$

This basically signify that  $\mathcal{X}, \Theta$  are constant and the only random part is  $\mathcal{Y}$ .

In addition

$$\mathcal{L}(\Theta|\mathcal{X}) \quad (14)$$

It is known as the incomplete-data likelihood function.



Thus

We can connect both incomplete-complete data equations by doing the following

$$\begin{aligned}\mathcal{L}(\Theta|\mathcal{X}) &= p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{X}, \mathcal{Y}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} \left( \prod_{i=1}^N p(x_i|\Theta) \right) p(\mathcal{Y}|\mathcal{X}, \Theta)\end{aligned}$$



Thus

We can connect both incomplete-complete data equations by doing the following

$$\begin{aligned}\mathcal{L}(\Theta|\mathcal{X}) &= p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{X}, \mathcal{Y}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} \left( \prod_{i=1}^N p(x_i|\Theta) \right) p(\mathcal{Y}|\mathcal{X}, \Theta)\end{aligned}$$



Thus

We can connect both incomplete-complete data equations by doing the following

$$\begin{aligned}\mathcal{L}(\Theta|\mathcal{X}) &= p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{X}, \mathcal{Y}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} \left( \prod_{i=1}^N p(x_i|\theta) \right) p(\mathcal{Y}|\mathcal{X}, \Theta)\end{aligned}$$



Thus

We can connect both incomplete-complete data equations by doing the following

$$\begin{aligned}\mathcal{L}(\Theta|\mathcal{X}) &= p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{X}, \mathcal{Y}|\Theta) \\ &= \sum_{\mathcal{Y}} p(\mathcal{Y}|\mathcal{X}, \Theta) p(\mathcal{X}|\Theta) \\ &= \sum_{\mathcal{Y}} \left( \prod_{i=1}^N p(x_i|\Theta) \right) p(\mathcal{Y}|\mathcal{X}, \Theta)\end{aligned}$$



## Problems

Normally, it is almost impossible to obtain a closed analytical solution for the previous equation.

## However

We can use the expected value of  $\log p(\mathcal{X}, \mathcal{Y} | \Theta)$ , which allows us to find an iterative procedure to approximate the solution.



# Remarks

## Problems

Normally, it is almost impossible to obtain a closed analytical solution for the previous equation.

## However

We can use the expected value of  $\log p(\mathcal{X}, \mathcal{Y} | \Theta)$ , which allows us to find an iterative procedure to approximate the solution.





# The function we would like to have

## The Q function

We want an estimation of the complete-data log-likelihood

$$\log p(\mathcal{X}, \mathcal{Y} | \Theta) \quad (15)$$

Based in the info provided by  $\mathcal{X}$ ,  $\Theta_{n-1}$  where  $\Theta_{n-1}$  is a previously estimated set of parameters at step  $n$ .

Think about the following, if we want to remove  $\mathcal{Y}$

$$\int [\log p(\mathcal{X}, \mathcal{Y} | \Theta)] p(\mathcal{Y} | \mathcal{X}, \Theta_{n-1}) d\mathcal{Y} \quad (16)$$

Remark: We integrate out  $\mathcal{Y}$  - Actually, this is the expected value of  $\log p(\mathcal{X}, \mathcal{Y} | \Theta)$ .

# The function we would like to have

## The Q function

We want an estimation of the complete-data log-likelihood

$$\log p(\mathcal{X}, \mathcal{Y} | \Theta) \quad (15)$$

Based in the info provided by  $\mathcal{X}$ ,  $\Theta_{n-1}$  where  $\Theta_{n-1}$  is a previously estimated set of parameters at step  $n$ .

Think about the following, if we want to remove  $\mathcal{Y}$

$$\int [\log p(\mathcal{X}, \mathcal{Y} | \Theta)] p(\mathcal{Y} | \mathcal{X}, \Theta_{n-1}) d\mathcal{Y} \quad (16)$$

**Remark:** We integrate out  $\mathcal{Y}$  - Actually, this is the expected value of  $\log p(\mathcal{X}, \mathcal{Y} | \Theta)$ .



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - **Using the Expected Value**
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



## Use the Expected Value

Then, we want an iterative method to guess  $\Theta$  from  $\Theta_{n-1}$

$$Q(\Theta, \Theta_{n-1}) = E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] \quad (17)$$



## Use the Expected Value

Then, we want an iterative method to guess  $\Theta$  from  $\Theta_{n-1}$

$$Q(\Theta, \Theta_{n-1}) = E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] \quad (17)$$

Take in account that

- 1  $\mathcal{X}, \Theta_{n-1}$  are taken as constants.
- 2  $\Theta$  is a normal variable that we wish to adjust.
- 3  $\mathcal{Y}$  is a random variable governed by distribution  $p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1})$  = marginal distribution of missing data.



## Use the Expected Value

Then, we want an iterative method to guess  $\Theta$  from  $\Theta_{n-1}$

$$Q(\Theta, \Theta_{n-1}) = E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] \quad (17)$$

Take in account that

- 1  $\mathcal{X}, \Theta_{n-1}$  are taken as constants.
- 2  $\Theta$  is a normal variable that we wish to adjust.
- 3  $\mathcal{Y}$  is a random variable governed by distribution  $p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1})$  = marginal distribution of missing data.



## Use the Expected Value

Then, we want an iterative method to guess  $\Theta$  from  $\Theta_{n-1}$

$$Q(\Theta, \Theta_{n-1}) = E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] \quad (17)$$

Take in account that

- 1  $\mathcal{X}, \Theta_{n-1}$  are taken as constants.
- 2  $\Theta$  is a normal variable that we wish to adjust.
- 3  $\mathcal{Y}$  is a random variable governed by distribution  $p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1})$ =marginal distribution of missing data.



## Another Interpretation

Given the previous information

$$E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] = \int_{\mathcal{Y} \in \mathbb{Y}} \log p(\mathcal{X}, \mathcal{Y}|\Theta) p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1}) d\mathcal{Y}$$





## Another Interpretation

Given the previous information

$$E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] = \int_{\mathcal{Y} \in \mathbb{Y}} \log p(\mathcal{X}, \mathcal{Y}|\Theta) p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1}) d\mathcal{Y}$$

Something Notable

- 1 In the best of cases, this marginal distribution is a simple analytical expression of the assumed parameter  $\Theta_{n-1}$ .

2 In the worst of cases, this density might be very hard to obtain.



## Another Interpretation

Given the previous information

$$E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] = \int_{\mathcal{Y} \in \mathbb{Y}} \log p(\mathcal{X}, \mathcal{Y}|\Theta) p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1}) d\mathcal{Y}$$

Something Notable

- 1 In the best of cases, this marginal distribution is a simple analytical expression of the assumed parameter  $\Theta_{n-1}$ .
- 2 In the worst of cases, this density might be very hard to obtain.

$$p(\mathcal{Y}, \mathcal{X}|\Theta_{n-1}) = p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1}) p(\mathcal{X}|\Theta_{n-1}) \quad (18)$$

which is not dependent on  $\Theta$ .



## Another Interpretation

Given the previous information

$$E[\log p(\mathcal{X}, \mathcal{Y}|\Theta) | \mathcal{X}, \Theta_{n-1}] = \int_{\mathcal{Y} \in \mathbb{Y}} \log p(\mathcal{X}, \mathcal{Y}|\Theta) p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1}) d\mathcal{Y}$$

Something Notable

- 1 In the best of cases, this marginal distribution is a simple analytical expression of the assumed parameter  $\Theta_{n-1}$ .
- 2 In the worst of cases, this density might be very hard to obtain.

Actually, we use

$$p(\mathcal{Y}, \mathcal{X}|\Theta_{n-1}) = p(\mathcal{Y}|\mathcal{X}, \Theta_{n-1}) p(\mathcal{X}|\Theta_{n-1}) \quad (18)$$

which is not dependent on  $\Theta$ .



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - **Analogy**
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# Back to the $Q$ function

## The intuition

We have the following analogy:

- Consider  $h(\theta, Y)$  a function
  - ▶  $\theta$  a constant
  - ▶  $Y \sim p_Y(y)$ , a random variable with distribution  $p_Y(y)$ .



# Back to the $Q$ function

## The intuition

We have the following analogy:

- Consider  $h(\theta, \mathbf{Y})$  a function
  - ▶  $\theta$  a constant
  - ▶  $\mathbf{Y} \sim p_{\mathbf{Y}}(y)$ , a random variable with distribution  $p_{\mathbf{Y}}(y)$ .

Thus, if  $\mathbf{Y}$  is a discrete random variable

$$q(\theta) = E_{\mathbf{Y}} [h(\theta, \mathbf{Y})] = \sum_y h(\theta, y) p_{\mathbf{Y}}(y) \quad (19)$$



# Back to the $Q$ function

## The intuition

We have the following analogy:

- Consider  $h(\theta, \mathbf{Y})$  a function

- ▶  $\theta$  a constant

- ▶  $\mathbf{Y} \sim p_{\mathbf{Y}}(y)$ , a random variable with distribution  $p_{\mathbf{Y}}(y)$ .

Thus, if  $\mathbf{Y}$  is a discrete random variable

$$q(\theta) = E_{\mathbf{Y}} [h(\theta, \mathbf{Y})] = \sum_y h(\theta, y) p_{\mathbf{Y}}(y) \quad (19)$$



# Back to the $Q$ function

## The intuition

We have the following analogy:

- Consider  $h(\theta, \mathbf{Y})$  a function
  - ▶  $\theta$  a constant
  - ▶  $\mathbf{Y} \sim p_{\mathbf{Y}}(y)$ , a random variable with distribution  $p_{\mathbf{Y}}(y)$ .

Thus, if  $\mathbf{Y}$  is a discrete random variable

$$q(\theta) = E_{\mathbf{Y}} [h(\theta, \mathbf{Y})] = \sum_y h(\theta, y) p_{\mathbf{Y}}(y) \quad (19)$$





# Back to the $Q$ function

## The intuition

We have the following analogy:

- Consider  $h(\theta, \mathbf{Y})$  a function
  - ▶  $\theta$  a constant
  - ▶  $\mathbf{Y} \sim p_{\mathbf{Y}}(y)$ , a random variable with distribution  $p_{\mathbf{Y}}(y)$ .

Thus, if  $\mathbf{Y}$  is a discrete random variable

$$q(\theta) = E_{\mathbf{Y}}[h(\theta, \mathbf{Y})] = \sum_y h(\theta, y) p_{\mathbf{Y}}(y) \quad (19)$$



# Why E-step!!!

From here the name

This is basically the E-step

The second step

It tries to maximize the  $Q$  function

$$\Theta_n = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta_{n-1}) \quad (20)$$



# Why E-step!!!

From here the name

This is basically the E-step

The second step

It tries to maximize the  $Q$  function

$$\Theta_n = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta_{n-1}) \quad (20)$$



# Why E-step!!!

From here the name

This is basically the E-step

The second step

It tries to maximize the  $Q$  function

$$\Theta_n = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta_{n-1}) \quad (20)$$



# Derivation of the EM-Algorithm

The likelihood function we are going to use

Let  $\mathcal{X}$  be a random vector which results from a parametrized family:

$$\mathcal{L}(\Theta) = \ln \mathcal{P}(\mathcal{X}|\Theta) \quad (21)$$

**Note:**  $\ln(x)$  is a strictly increasing function.

We wish to compute  $\Theta$

Based on an estimate  $\Theta_n$  (After the  $n^{\text{th}}$ ) such that  $\mathcal{L}(\Theta) > \mathcal{L}(\Theta_n)$

Of this maximization of the difference

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) = \ln \mathcal{P}(\mathcal{X}|\Theta) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \quad (22)$$



# Derivation of the EM-Algorithm

The likelihood function we are going to use

Let  $\mathcal{X}$  be a random vector which results from a parametrized family:

$$\mathcal{L}(\Theta) = \ln \mathcal{P}(\mathcal{X}|\Theta) \quad (21)$$

**Note:**  $\ln(x)$  is a strictly increasing function.

We wish to compute  $\Theta$

Based on an estimate  $\Theta_n$  (After the  $n^{th}$ ) such that  $\mathcal{L}(\Theta) > \mathcal{L}(\Theta_n)$

Of this maximization of the difference

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) = \ln \mathcal{P}(\mathcal{X}|\Theta) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \quad (22)$$



# Derivation of the EM-Algorithm

The likelihood function we are going to use

Let  $\mathcal{X}$  be a random vector which results from a parametrized family:

$$\mathcal{L}(\Theta) = \ln \mathcal{P}(\mathcal{X}|\Theta) \quad (21)$$

**Note:**  $\ln(x)$  is a strictly increasing function.

We wish to compute  $\Theta$

Based on an estimate  $\Theta_n$  (After the  $n^{\text{th}}$ ) such that  $\mathcal{L}(\Theta) > \mathcal{L}(\Theta_n)$

Or the maximization of the difference

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) = \ln \mathcal{P}(\mathcal{X}|\Theta) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \quad (22)$$



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - **Hidden Features**
    - Proving Concavity
    - Using the Concave Functions for Approximation
    - From The Concave Function to the EM
    - The Final Algorithm
    - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm





# Introducing the Hidden Features

Given that the hidden random vector  $\mathcal{Y}$  exists with  $y$  values

$$\mathcal{P}(\mathcal{X}|\Theta) = \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \quad (23)$$

Thus, using our first constraint  $\mathcal{L}(\Theta) = \mathcal{L}(\Theta_n)$

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) = \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \quad (24)$$



# Introducing the Hidden Features

Given that the hidden random vector  $\mathcal{Y}$  exists with  $y$  values

$$\mathcal{P}(\mathcal{X}|\Theta) = \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \quad (23)$$

Thus, using our first constraint  $\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n)$

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) = \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \quad (24)$$



# Here, we introduce some concepts of convexity

## For Convexity

### Theorem (Jensen's inequality)

Let  $f$  be a convex function defined on an interval  $I$ . If  $x_1, x_2, \dots, x_n \in I$  and  $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$  with  $\sum_{i=1}^n \lambda_i = 1$ , then

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i) \quad (25)$$



# Proof:

For  $n = 1$

We have the trivial case

For  $n = 2$

The convexity definition.

Now, the inductive hypothesis.

We assume that the theorem is true for some  $n$ .



# Proof:

For  $n = 1$

We have the trivial case

For  $n = 2$

The convexity definition.

Now, the inductive hypothesis

We assume that the theorem is true for some  $n$ .



# Proof:

For  $n = 1$

We have the trivial case

For  $n = 2$

The convexity definition.

Now the inductive hypothesis

We assume that the theorem is true for some  $n$ .



Now, we have

The following linear combination for  $\lambda_i$

$$\begin{aligned} f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) &= f\left(\lambda_{n+1} x_{n+1} + \sum_{i=1}^n \lambda_i x_i\right) \\ &= f\left(\lambda_{n+1} x_{n+1} + \frac{(1-\lambda_{n+1})}{(1-\lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \\ &\leq \lambda_{n+1} f(x_{n+1}) + (1-\lambda_{n+1}) f\left(\frac{1}{(1-\lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \end{aligned}$$



Now, we have

The following linear combination for  $\lambda_i$

$$\begin{aligned} f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) &= f\left(\lambda_{n+1} x_{n+1} + \sum_{i=1}^n \lambda_i x_i\right) \\ &= f\left(\lambda_{n+1} x_{n+1} + \frac{(1 - \lambda_{n+1})}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \end{aligned}$$

$$\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) f\left(\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right)$$





Now, we have

The following linear combination for  $\lambda_i$

$$\begin{aligned} f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) &= f\left(\lambda_{n+1} x_{n+1} + \sum_{i=1}^n \lambda_i x_i\right) \\ &= f\left(\lambda_{n+1} x_{n+1} + \frac{(1 - \lambda_{n+1})}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \\ &\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) f\left(\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \end{aligned}$$



# Did you notice?

## Something Notable

$$\sum_{i=1}^{n+1} \lambda_i = 1$$

Thus

$$\sum_{i=1}^n \lambda_i = 1 - \lambda_{n+1}$$

Finally

$$\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i = 1$$

Did you notice?

## Something Notable

$$\sum_{i=1}^{n+1} \lambda_i = 1$$

Thus

$$\sum_{i=1}^n \lambda_i = 1 - \lambda_{n+1}$$

Finally

$$\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i = 1$$

Did you notice?

Something Notable

$$\sum_{i=1}^{n+1} \lambda_i = 1$$

Thus

$$\sum_{i=1}^n \lambda_i = 1 - \lambda_{n+1}$$

Finally

$$\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i = 1$$

We have that

$$f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) \leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) f\left(\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right)$$

$$\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) \frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i f(x_i)$$

$$\leq \lambda_{n+1} f(x_{n+1}) + \sum_{i=1}^n \lambda_i f(x_i) \quad \text{Q.E.D.}$$



We have that

$$\begin{aligned}
 f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) &\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) f\left(\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \\
 &\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) \frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i f(x_i) \\
 &\leq \lambda_{n+1} f(x_{n+1}) + \sum_{i=1}^n \lambda_i f(x_i) \quad \text{Q.E.D.}
 \end{aligned}$$



We have that

$$\begin{aligned} f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) &\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) f\left(\frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i x_i\right) \\ &\leq \lambda_{n+1} f(x_{n+1}) + (1 - \lambda_{n+1}) \frac{1}{(1 - \lambda_{n+1})} \sum_{i=1}^n \lambda_i f(x_i) \\ &\leq \lambda_{n+1} f(x_{n+1}) + \sum_{i=1}^n \lambda_i f(x_i) \quad \text{Q.E.D.} \end{aligned}$$



Thus, for concave functions

It is possible to show that

Given  $\ln(x)$  a concave function:

$$\ln \left[ \sum_{i=1}^n \lambda_i x_i \right] \geq \sum_{i=1}^n \lambda_i \ln(x_i)$$

If we take in consideration

Assume that the  $\lambda_i = \mathcal{P}(y|\mathcal{X}, \Theta_n)$ . We know that

- $\mathcal{P}(y|\mathcal{X}, \Theta_n) \geq 0$
- $\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) = 1$





Thus, for concave functions

It is possible to show that

Given  $\ln(x)$  a concave function:

$$\ln \left[ \sum_{i=1}^n \lambda_i x_i \right] \geq \sum_{i=1}^n \lambda_i \ln(x_i)$$

If we take in consideration

Assume that the  $\lambda_i = \mathcal{P}(y|\mathcal{X}, \Theta_n)$ . We know that

- 1  $\mathcal{P}(y|\mathcal{X}, \Theta_n) \geq 0$
- 2  $\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) = 1$



## We have

First

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) = \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n)$$

$$= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \frac{\mathcal{P}(y|\mathcal{X}, \Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n)$$

$$= \ln \left( \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n)$$

$$\geq \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \dots$$

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \mathcal{P}(\mathcal{X}|\Theta_n) \text{ Why this?}$$

We have

First

$$\begin{aligned}\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) &= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \frac{\mathcal{P}(y|\mathcal{X}, \Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &= \ln \left( \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &\geq \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \dots \\ &\quad \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \mathcal{P}(\mathcal{X}|\Theta_n) \text{ Why this?}\end{aligned}$$

We have

First

$$\begin{aligned}\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) &= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \frac{\mathcal{P}(y|\mathcal{X}, \Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &= \ln \left( \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n)\end{aligned}$$

$$\geq \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \dots$$

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \mathcal{P}(\mathcal{X}|\Theta_n) \text{ Why this?}$$

We have

First

$$\begin{aligned}\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) &= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &= \ln \left( \sum_y \mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta) \frac{\mathcal{P}(y|\mathcal{X}, \Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &= \ln \left( \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \ln \mathcal{P}(\mathcal{X}|\Theta_n) \\ &\geq \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n)} \right) - \dots \\ &\quad \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \mathcal{P}(\mathcal{X}|\Theta_n) \text{ Why this?}\end{aligned}$$

## Next

Because

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) = 1$$

Then

$$\mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) \geq \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$



## Next

Because

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) = 1$$

Then

$$\begin{aligned} \mathcal{L}(\Theta) - \mathcal{L}(\Theta_n) &\geq \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \\ &= \Delta(\Theta|\Theta_n) \end{aligned}$$



Then, we have

Then, we have proved that

$$\mathcal{L}(\Theta) \geq \mathcal{L}(\Theta_n) + \Delta(\Theta|\Theta_n) \quad (26)$$

Then, we define a new function

$$l(\Theta|\Theta_n) = \mathcal{L}(\Theta_n) + \Delta(\Theta|\Theta_n) \quad (27)$$

Thus  $l(\Theta|\Theta_n)$

It is bounded from above by  $\mathcal{L}(\Theta)$  i.e  $l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$





Then, we have

Then, we have proved that

$$\mathcal{L}(\Theta) \geq \mathcal{L}(\Theta_n) + \Delta(\Theta|\Theta_n) \quad (26)$$

Then, we define a new function

$$l(\Theta|\Theta_n) = \mathcal{L}(\Theta_n) + \Delta(\Theta|\Theta_n) \quad (27)$$

This  $l(\Theta|\Theta_n)$

It is bounded from above by  $\mathcal{L}(\Theta)$  i.e.  $l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$



Then, we have

Then, we have proved that

$$\mathcal{L}(\Theta) \geq \mathcal{L}(\Theta_n) + \Delta(\Theta|\Theta_n) \quad (26)$$

Then, we define a new function

$$l(\Theta|\Theta_n) = \mathcal{L}(\Theta_n) + \Delta(\Theta|\Theta_n) \quad (27)$$

Thus  $l(\Theta|\Theta_n)$

It is bounded from above by  $\mathcal{L}(\Theta)$  i.e  $l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$



Now, we can do the following

We evaluate in  $\Theta_n$

$$l(\Theta_n|\Theta_n) = \mathcal{L}(\Theta_n) + \Delta(\Theta_n|\Theta_n)$$

$$= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta_n) \mathcal{P}(y|\Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$

$$= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta_n)}{\mathcal{P}(\mathcal{X}, y|\Theta_n)} \right)$$

$$= \mathcal{L}(\Theta_n)$$



Now, we can do the following

We evaluate in  $\Theta_n$

$$\begin{aligned}l(\Theta_n|\Theta_n) &= \mathcal{L}(\Theta_n) + \Delta(\Theta_n|\Theta_n) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta_n) \mathcal{P}(y|\Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta_n)}{\mathcal{P}(\mathcal{X}, y|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n)\end{aligned}$$

Intuitively:

For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal



Now, we can do the following

We evaluate in  $\Theta_n$

$$\begin{aligned}l(\Theta_n|\Theta_n) &= \mathcal{L}(\Theta_n) + \Delta(\Theta_n|\Theta_n) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta_n) \mathcal{P}(y|\Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta_n)}{\mathcal{P}(\mathcal{X}, y|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n)\end{aligned}$$

Intuitively:

For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal



Now, we can do the following

We evaluate in  $\Theta_n$

$$\begin{aligned}l(\Theta_n|\Theta_n) &= \mathcal{L}(\Theta_n) + \Delta(\Theta_n|\Theta_n) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta_n) \mathcal{P}(y|\Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta_n)}{\mathcal{P}(\mathcal{X}, y|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n)\end{aligned}$$

This means that

For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal



Now, we can do the following

We evaluate in  $\Theta_n$

$$\begin{aligned}l(\Theta_n|\Theta_n) &= \mathcal{L}(\Theta_n) + \Delta(\Theta_n|\Theta_n) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta_n) \mathcal{P}(y|\Theta_n)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta_n)}{\mathcal{P}(\mathcal{X}, y|\Theta_n)} \right) \\ &= \mathcal{L}(\Theta_n)\end{aligned}$$

This means that

For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal



# Therefore

The function  $l(\Theta|\Theta_n)$  has the following properties

- 1 It is bounded from above by  $\mathcal{L}(\Theta)$  i.e  $l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$ .
- 2 For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal.
- 3 The function  $l(\Theta|\Theta_n)$  is concave... How?





# Therefore

The function  $l(\Theta|\Theta_n)$  has the following properties

- 1 It is bounded from above by  $\mathcal{L}(\Theta)$  i.e  $l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$ .
- 2 For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal.

3 The function  $l(\Theta|\Theta_n)$  is concave... How?



# Therefore

The function  $l(\Theta|\Theta_n)$  has the following properties

- 1 It is bounded from above by  $\mathcal{L}(\Theta)$  i.e  $l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$ .
- 2 For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal.
- 3 The function  $l(\Theta|\Theta_n)$  is concave... How?



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# First

We have the value  $\mathcal{L}(\Theta_n)$

We know that  $\mathcal{L}(\Theta_n)$  is constant i.e. an offset value

What about  $\Delta \mathcal{L}(\Theta)$ ?

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$

We have that the  $\ln$  is a concave function

$$\ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$



# First

We have the value  $\mathcal{L}(\Theta_n)$

We know that  $\mathcal{L}(\Theta_n)$  is constant i.e. an offset value

What about  $\Delta(\Theta|\Theta_n)$

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$

We have that this is a concave function

$$\ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$



# First

We have the value  $\mathcal{L}(\Theta_n)$

We know that  $\mathcal{L}(\Theta_n)$  is constant i.e. an offset value

What about  $\Delta(\Theta|\Theta_n)$

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$

We have that the  $\ln$  is a concave function

$$\ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$



Therefore

Each element is concave

$$\mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$

Therefore, the sum of concave functions is a concave function

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$



Therefore

Each element is concave

$$\mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$

Therefore, the sum of concave functions is a concave function

$$\sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right)$$





# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - **Using the Concave Functions for Approximation**
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm

# Given the Concave Function

Thus, we have that

① We can select  $\Theta_n$  such that  $l(\Theta|\Theta_n)$  is maximized.

② Thus, given a  $\Theta_n$ , we can generate  $\Theta_{n+1}$ .



# Given the Concave Function

Thus, we have that

- 1 We can select  $\Theta_n$  such that  $l(\Theta|\Theta_n)$  is maximized.
- 2 Thus, given a  $\Theta_n$ , we can generate  $\Theta_{n+1}$ .

The process can be seen in the following graph

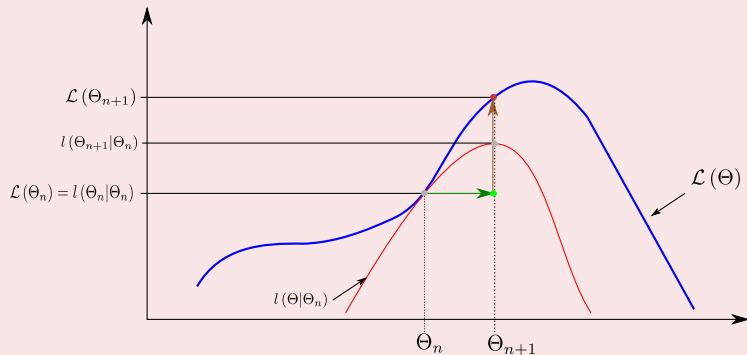


# Given the Concave Function

Thus, we have that

- 1 We can select  $\Theta_n$  such that  $l(\Theta|\Theta_n)$  is maximized.
- 2 Thus, given a  $\Theta_n$ , we can generate  $\Theta_{n+1}$ .

The process can be seen in the following graph



# Given

## The Previous Constraints

- 1  $l(\Theta|\Theta_n)$  is bounded from above by  $\mathcal{L}(\Theta)$

$$l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$$

- 2 For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal

$$\mathcal{L}(\Theta_n) = l(\Theta|\Theta_n)$$

- 3 The function  $l(\Theta|\Theta_n)$  is concave



## The Previous Constraints

- 1  $l(\Theta|\Theta_n)$  is bounded from above by  $\mathcal{L}(\Theta)$

$$l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$$

- 2 For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal

$$\mathcal{L}(\Theta_n) = l(\Theta|\Theta_n)$$

- 3 The function  $l(\Theta|\Theta_n)$  is concave



## The Previous Constraints

- 1  $l(\Theta|\Theta_n)$  is bounded from above by  $\mathcal{L}(\Theta)$

$$l(\Theta|\Theta_n) \leq \mathcal{L}(\Theta)$$

- 2 For  $\Theta = \Theta_n$ , functions  $\mathcal{L}(\Theta)$  and  $l(\Theta|\Theta_n)$  are equal

$$\mathcal{L}(\Theta_n) = l(\Theta|\Theta_n)$$

- 3 The function  $l(\Theta|\Theta_n)$  is concave



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - **From The Concave Function to the EM**
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# From

The following

$$\Theta_{n+1} = \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \right\}$$

The terms with  $\Theta_n$  are constants.

$$\approx \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y|\Theta) \mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(y, \Theta)}{\frac{\mathcal{P}(y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(\Theta)} \right) \right\}$$

## From

The following

$$\begin{aligned}\Theta_{n+1} &= \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \right\}\end{aligned}$$

The terms with  $\Theta_n$  are constants.

$$\approx \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y|\Theta) \mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(y, \Theta)}{\frac{\mathcal{P}(y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(\Theta)} \right) \right\}$$

## From

The following

$$\begin{aligned}\Theta_{n+1} &= \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \right\}\end{aligned}$$

The terms with  $\Theta_n$  are constants.

$$\approx \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y|\Theta) \mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(y, \Theta)}{\frac{\mathcal{P}(y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(\Theta)} \right) \right\}$$

# From

The following

$$\begin{aligned}\Theta_{n+1} &= \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \right\}\end{aligned}$$

The terms with  $\Theta_n$  are constants.

$$\begin{aligned}&\approx \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y|\Theta) \mathcal{P}(\Theta)} \right) \right\}\end{aligned}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \frac{\mathcal{P}(y, \Theta)}{\mathcal{P}(\Theta)} \right) \right\}$$

## The following

$$\begin{aligned}\Theta_{n+1} &= \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \mathcal{L}(\Theta_n) + \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)}{\mathcal{P}(y|\mathcal{X}, \Theta_n) \mathcal{P}(\mathcal{X}|\Theta_n)} \right) \right\}\end{aligned}$$

The terms with  $\Theta_n$  are constants.

$$\begin{aligned}&\approx \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}|y, \Theta) \mathcal{P}(y|\Theta)) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y|\Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y|\Theta) \mathcal{P}(\Theta)} \right) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(y, \Theta)}{\frac{\mathcal{P}(y, \Theta)}{\mathcal{P}(\Theta)} \mathcal{P}(\Theta)} \right) \right\}\end{aligned}$$

Thus

Then

$$\theta_{n+1} = \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y, \Theta) \mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}, y|\Theta)) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$

$$\text{Then } \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \approx \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$

Thus

Then

$$\theta_{n+1} = \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y, \Theta) \mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}, y|\Theta)) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$

$$\text{Then } \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \approx \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$

Thus

Then

$$\theta_{n+1} = \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y, \Theta) \mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \right) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}, y|\Theta)) \right\}$$

$$= \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$

$$\text{Then } \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \approx \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$



Thus

Then

$$\begin{aligned}\theta_{n+1} &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y, \Theta) \mathcal{P}(\Theta)} \right) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \right) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}, y|\Theta)) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}\end{aligned}$$

Then  $\operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \approx \operatorname{argmax}_{\Theta} \{E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))]\}$

Thus

Then

$$\begin{aligned}\theta_{n+1} &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta) \mathcal{P}(y, \Theta)}{\mathcal{P}(y, \Theta) \mathcal{P}(\Theta)} \right) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln \left( \frac{\mathcal{P}(\mathcal{X}, y, \Theta)}{\mathcal{P}(\Theta)} \right) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ \sum_y \mathcal{P}(y|\mathcal{X}, \Theta_n) \ln (\mathcal{P}(\mathcal{X}, y|\Theta)) \right\} \\ &= \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}\end{aligned}$$

$$\text{Then } \operatorname{argmax}_{\Theta} \{l(\Theta|\Theta_n)\} \approx \operatorname{argmax}_{\Theta} \left\{ E_{y|\mathcal{X}, \Theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))] \right\}$$

# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - **The Final Algorithm**
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm

# The EM-Algorithm

## Steps of EM

① Expectation under hidden variables.

② Maximization of the resulting formula.



# The EM-Algorithm

## Steps of EM

- 1 Expectation under hidden variables.
- 2 Maximization of the resulting formula.

### E-Step

Determine the conditional expectation,  $E_{y|x, \theta_n} [\ln (\mathcal{P}(\mathcal{X}, y|\Theta))]$ .



# The EM-Algorithm

## Steps of EM

- 1 Expectation under hidden variables.
- 2 Maximization of the resulting formula.

## E-Step

Determine the conditional expectation,  $E_{y|\mathcal{X},\Theta_n} [\ln (\mathcal{P} (\mathcal{X}, y|\Theta))]$ .

## M-Step

Maximize this expression with respect to  $\Theta$ .



# The EM-Algorithm

## Steps of EM

- 1 Expectation under hidden variables.
- 2 Maximization of the resulting formula.

## E-Step

Determine the conditional expectation,  $E_{y|\mathcal{X},\Theta_n} [\ln (\mathcal{P} (\mathcal{X}, y|\Theta))]$ .

## M-Step

Maximize this expression with respect to  $\Theta$ .



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - **Notes and Convergence of EM**
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# Notes and Convergence of EM

## Gains between $\mathcal{L}(\Theta)$ and $l(\Theta|\Theta_n)$

Using the hidden variables it is possible to simplify the optimization of  $\mathcal{L}(\Theta)$  through  $l(\Theta|\Theta_n)$ .



# Notes and Convergence of EM

## Gains between $\mathcal{L}(\Theta)$ and $l(\Theta|\Theta_n)$

Using the hidden variables it is possible to simplify the optimization of  $\mathcal{L}(\Theta)$  through  $l(\Theta|\Theta_n)$ .

## Convergence

- Remember that  $\Theta_{n+1}$  is the estimate for  $\Theta$  which maximizes the difference  $\Delta(\Theta|\Theta_n)$ .



# Notes and Convergence of EM

## Gains between $\mathcal{L}(\Theta)$ and $l(\Theta|\Theta_n)$

Using the hidden variables it is possible to simplify the optimization of  $\mathcal{L}(\Theta)$  through  $l(\Theta|\Theta_n)$ .

## Convergence

- Remember that  $\Theta_{n+1}$  is the estimate for  $\Theta$  which maximizes the difference  $\Delta(\Theta|\Theta_n)$ .



# Therefore

Then, we have

Given the initial estimate of  $\Theta$  by  $\Theta_n$

$$\Delta(\Theta_n | \Theta_n) = 0$$

Now

If we choose  $\Theta_{n+1}$  to maximize the  $\Delta(\Theta | \Theta_n)$ , then

$$\Delta(\Theta_{n+1} | \Theta_n) \geq \Delta(\Theta_n | \Theta_n) = 0$$

We have that

The Likelihood  $\mathcal{L}(\Theta)$  is not a decreasing function with respect to  $\Theta$ .



# Therefore

Then, we have

Given the initial estimate of  $\Theta$  by  $\Theta_n$

$$\Delta(\Theta_n | \Theta_n) = 0$$

Now

If we choose  $\Theta_{n+1}$  to maximize the  $\Delta(\Theta | \Theta_n)$ , then

$$\Delta(\Theta_{n+1} | \Theta_n) \geq \Delta(\Theta_n | \Theta_n) = 0$$

We have that

The Likelihood  $\mathcal{L}(\Theta)$  is not a decreasing function with respect to  $\Theta$ .



# Therefore

Then, we have

Given the initial estimate of  $\Theta$  by  $\Theta_n$

$$\Delta(\Theta_n | \Theta_n) = 0$$

Now

If we choose  $\Theta_{n+1}$  to maximize the  $\Delta(\Theta | \Theta_n)$ , then

$$\Delta(\Theta_{n+1} | \Theta_n) \geq \Delta(\Theta_n | \Theta_n) = 0$$

We have that

The Likelihood  $\mathcal{L}(\Theta)$  is not a decreasing function with respect to  $\Theta$ .



# Notes and Convergence of EM

## Properties

When the algorithm reaches a fixed point for some  $\Theta_n$ , the value maximizes  $l(\Theta|\Theta_n)$ .



# Notes and Convergence of EM

## Properties

When the algorithm reaches a fixed point for some  $\Theta_n$ , the value maximizes  $l(\Theta|\Theta_n)$ .

## Definition

A fixed point of a function is an element on domain that is mapped to itself by the function:

$$f(\mathbf{x}) = \mathbf{x}$$

Basically the EM algorithm does the following

$$EM[\Theta^*] = \Theta^*$$



# Notes and Convergence of EM

## Properties

When the algorithm reaches a fixed point for some  $\Theta_n$ , the value maximizes  $l(\Theta|\Theta_n)$ .

## Definition

A fixed point of a function is an element on domain that is mapped to itself by the function:

$$f(\mathbf{x}) = \mathbf{x}$$

Basically the EM algorithm does the following

$$EM[\Theta^*] = \Theta^*$$

## At this moment

### We have that

The algorithm reaches a fixed point for some  $\Theta_n$ , the value  $\Theta^*$  maximizes  $l(\Theta|\Theta_n)$ .

Then, when the algorithm

- It reaches a fixed point for some  $\Theta_n$ , the value maximizes  $l(\Theta|\Theta_n)$ .
  - ▶ Basically  $\Theta_{n+1} = \Theta_n$ .



# At this moment

## We have that

The algorithm reaches a fixed point for some  $\Theta_n$ , the value  $\Theta^*$  maximizes  $l(\Theta|\Theta_n)$ .

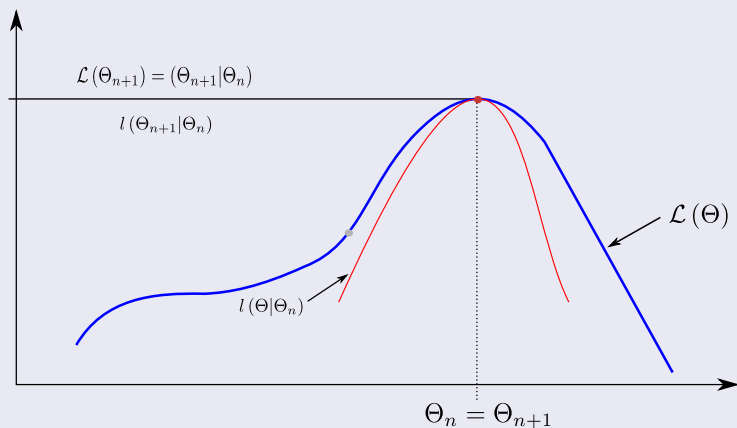
## Then, when the algorithm

- It reaches a fixed point for some  $\Theta_n$  the value maximizes  $l(\Theta|\Theta_n)$ .
  - ▶ Basically  $\Theta_{n+1} = \Theta_n$ .



Therefore

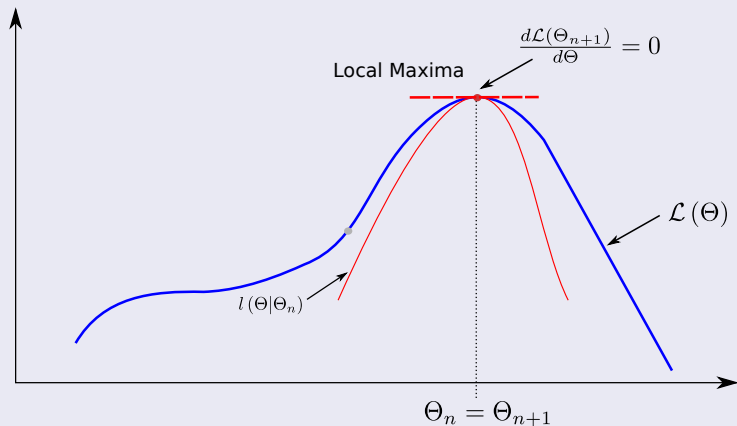
We have



# Then

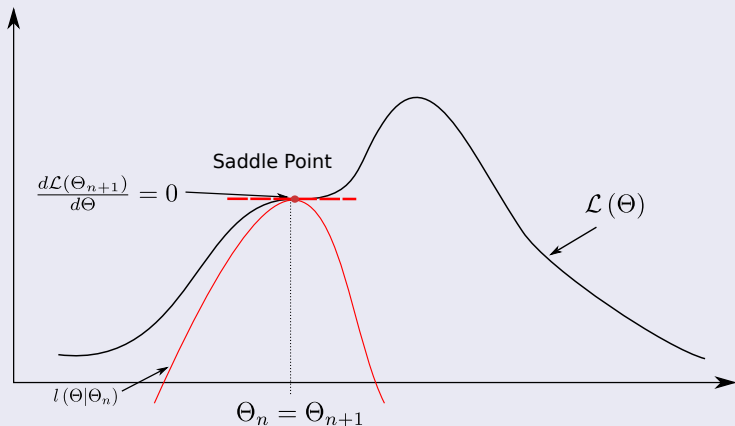
If  $\mathcal{L}$  and  $l$  are differentiable at  $\Theta_n$

- Since  $\mathcal{L}$  and  $l$  are equal at  $\Theta_n$ 
  - ▶ Then,  $\Theta_n$  is a stationary point of  $\mathcal{L}$  i.e. the derivative of  $\mathcal{L}$  vanishes at that point.



However

You could finish with the following case, no local maxima



For more on the subject

Please take a look to

Geoffrey McLachlan and Thriyambakam Krishnan, "*The EM Algorithm and Extensions*," John Wiley & Sons, New York, 1996.



Cinvestav

# Finding Maximum Likelihood Mixture Densities Parameters via EM

## Something Notable

The mixture-density parameter estimation problem is probably one of the most widely used applications of the EM algorithm in the computational pattern recognition community.

We have

$$p(\mathbf{x}|\Theta) = \sum_{i=1}^M \alpha_i p_i(\mathbf{x}|\theta_i) \quad (28)$$



# Finding Maximum Likelihood Mixture Densities Parameters via EM

## Something Notable

The mixture-density parameter estimation problem is probably one of the most widely used applications of the EM algorithm in the computational pattern recognition community.

## We have

$$p(\mathbf{x}|\Theta) = \sum_{i=1}^M \alpha_i p_i(\mathbf{x}|\theta_i) \quad (28)$$

where

①  $\Theta = (\alpha_1, \dots, \alpha_M, \theta_1, \dots, \theta_M)$

②  $\sum_{i=1}^M \alpha_i = 1$

③ Each  $p_i$  is a density function parametrized by  $\theta_i$ .

# Finding Maximum Likelihood Mixture Densities Parameters via EM

## Something Notable

The mixture-density parameter estimation problem is probably one of the most widely used applications of the EM algorithm in the computational pattern recognition community.

## We have

$$p(\mathbf{x}|\Theta) = \sum_{i=1}^M \alpha_i p_i(\mathbf{x}|\theta_i) \quad (28)$$

where

- 1  $\Theta = (\alpha_1, \dots, \alpha_M, \theta_1, \dots, \theta_M)$
- 2  $\sum_{i=1}^M \alpha_i = 1$

Each  $p_i$  is a density function parametrized by  $\theta_i$ .

# Finding Maximum Likelihood Mixture Densities Parameters via EM

## Something Notable

The mixture-density parameter estimation problem is probably one of the most widely used applications of the EM algorithm in the computational pattern recognition community.

## We have

$$p(\mathbf{x}|\Theta) = \sum_{i=1}^M \alpha_i p_i(\mathbf{x}|\theta_i) \quad (28)$$

where

- 1  $\Theta = (\alpha_1, \dots, \alpha_M, \theta_1, \dots, \theta_M)$
- 2  $\sum_{i=1}^M \alpha_i = 1$
- 3 Each  $p_i$  is a density function parametrized by  $\theta_i$ .

## A log-likelihood for this function

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(x_i | \Theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^M \alpha_j p_j(x_i | \theta_j) \right) \quad (29)$$

Note: This is too difficult to optimize due to the log function.



## A log-likelihood for this function

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(x_i | \Theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^M \alpha_j p_j(x_i | \theta_j) \right) \quad (29)$$

**Note:** This is too difficult to optimize due to the log function.

However:

We can simplify this assuming the following:

- We assume that each unobserved data  $\mathcal{Y} = \{y_i\}_{i=1}^N$  has the following range  $y_i \in \{1, \dots, M\}$
- $y_i = k$  if the  $i^{\text{th}}$  samples was generated by the  $k^{\text{th}}$  mixture.



## A log-likelihood for this function

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(x_i | \Theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^M \alpha_j p_j(x_i | \theta_j) \right) \quad (29)$$

**Note:** This is too difficult to optimize due to the log function.

However

We can simplify this assuming the following:

- We assume that each unobserved data  $\mathcal{Y} = \{y_i\}_{i=1}^N$  has the following range  $y_i \in \{1, \dots, M\}$
- $y_i = k$  if the  $i^{\text{th}}$  samples was generated by the  $k^{\text{th}}$  mixture.



## A log-likelihood for this function

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(x_i | \Theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^M \alpha_j p_j(x_i | \theta_j) \right) \quad (29)$$

**Note:** This is too difficult to optimize due to the log function.

However

We can simplify this assuming the following:

- 1 We assume that each unobserved data  $\mathcal{Y} = \{y_i\}_{i=1}^N$  has the following range  $y_i \in \{1, \dots, M\}$

$y_i = k$  if the  $i^{\text{th}}$  samples was generated by the  $k^{\text{th}}$  mixture.



## A log-likelihood for this function

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}) = \log \prod_{i=1}^N p(x_i | \Theta) = \sum_{i=1}^N \log \left( \sum_{j=1}^M \alpha_j p_j(x_i | \theta_j) \right) \quad (29)$$

**Note:** This is too difficult to optimize due to the log function.

However

We can simplify this assuming the following:

- 1 We assume that each unobserved data  $\mathcal{Y} = \{y_i\}_{i=1}^N$  has the following range  $y_i \in \{1, \dots, M\}$
- 2  $y_i = k$  if the  $i^{\text{th}}$  samples was generated by the  $k^{\text{th}}$  mixture.





Now

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}, \mathcal{Y}) = \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] \quad (30)$$

## Now

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}, \mathcal{Y}) = \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] \quad (30)$$

Remember that  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$  with  $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$  and assuming independence

$$\begin{aligned} \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] &= \log [P(x_1, x_2, \dots, x_N, y_1, y_2, \dots, y_N | \Theta)] \\ &= \log [P(x_1, y_1, \dots, x_i, y_i, \dots, x_N, y_N | \Theta)] \\ &= \log \prod_{i=1}^N P(x_i, y_i | \Theta) \\ &= \sum_{i=1}^N \log P(x_i, y_i | \Theta) \end{aligned}$$

## Now

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}, \mathcal{Y}) = \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] \quad (30)$$

Remember that  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$  with  $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$  and assuming independence

$$\begin{aligned} \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] &= \log [P(x_1, x_2, \dots, x_N, y_1, y_2, \dots, y_N | \Theta)] \\ &= \log [P(x_1, y_1, \dots, x_i, y_i, \dots, x_N, y_N | \Theta)] \\ &= \log \prod_{i=1}^N P(x_i, y_i | \Theta) \\ &= \sum_{i=1}^N \log P(x_i, y_i | \Theta) \end{aligned}$$

## Now

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}, \mathcal{Y}) = \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] \quad (30)$$

Remember that  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$  with  $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$  and assuming independence

$$\begin{aligned} \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] &= \log [P(x_1, x_2, \dots, x_N, y_1, y_2, \dots, y_N | \Theta)] \\ &= \log [P(x_1, y_1, \dots, x_i, y_i, \dots, x_N, y_N | \Theta)] \\ &= \log \prod_{i=1}^N P(x_i, y_i | \Theta) \end{aligned}$$

$$= \sum_{i=1}^N \log P(x_i, y_i | \Theta)$$

## Now

We have

$$\log \mathcal{L}(\Theta | \mathcal{X}, \mathcal{Y}) = \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] \quad (30)$$

Remember that  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$  with  $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$  and assuming independence

$$\begin{aligned} \log [P(\mathcal{X}, \mathcal{Y} | \Theta)] &= \log [P(x_1, x_2, \dots, x_N, y_1, y_2, \dots, y_N | \Theta)] \\ &= \log [P(x_1, y_1, \dots, x_i, y_i, \dots, x_N, y_N | \Theta)] \\ &= \log \prod_{i=1}^N P(x_i, y_i | \Theta) \\ &= \sum_{i=1}^N \log P(x_i, y_i | \Theta) \end{aligned}$$

Then

Thus, by the chain Rule

$$\sum_{i=1}^N \log P(x_i, y_i | \Theta) = \sum_{i=1}^N \log [P(x_i | y_i, \theta_{y_i}) P(y_i | \theta_{y_i})] \quad (31)$$

**Question** Do you need  $y_i$  if you know  $\theta_{y_i}$  or the other way around?



Then

Thus, by the chain Rule

$$\sum_{i=1}^N \log P(x_i, y_i | \Theta) = \sum_{i=1}^N \log [P(x_i | y_i, \theta_{y_i}) P(y_i | \theta_{y_i})] \quad (31)$$

**Question** Do you need  $y_i$  if you know  $\theta_{y_i}$  or the other way around?

Finally

$$\sum_{i=1}^N \log [P(x_i | y_i, \theta_{y_i}) P(y_i | \theta_{y_i})] = \sum_{i=1}^N \log [P(y_i) p_{y_i}(x_i | \theta_{y_i})] \quad (32)$$

**NOPE:** You do not need  $y_i$  if you know  $\theta_{y_i}$  or the other way around.



Finally, we have

Making  $\alpha_{y_i} = P(y_i)$

$$\log \mathcal{L}(\Theta | \mathcal{X}, \mathcal{Y}) = \sum_{i=1}^N \log [\alpha_{y_i} P(x_i | y_i, \theta_{y_i})] \quad (33)$$





# Problem

## Which Labels?

We do not know the values of  $\mathcal{Y}$ .

We can get away by using the following idea

Assume the  $\mathcal{Y}$  is a random variable.



# Problem

## Which Labels?

We do not know the values of  $\mathcal{Y}$ .

We can get away by using the following idea

Assume the  $\mathcal{Y}$  is a random variable.



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - **The Beginning of The Process**
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



Thus

You do a first guess for the parameters at the beginning of EM

$$\Theta^g = (\alpha_1^g, \dots, \alpha_M^g, \theta_1^g, \dots, \theta_M^g) \quad (34)$$

Then, it is possible to calculate given the parametric probability

$$p_j(x_i | \theta_j^g)$$

Therefore

The mixing parameters  $\alpha_j$  can be thought of as a prior probabilities of each mixture:

$$\alpha_j = p(\text{component } j) \quad (35)$$



Thus

You do a first guess for the parameters at the beginning of EM

$$\Theta^g = (\alpha_1^g, \dots, \alpha_M^g, \theta_1^g, \dots, \theta_M^g) \quad (34)$$

Then, it is possible to calculate given the parametric probability

$$p_j(x_i | \theta_j^g)$$

Iterators

The mixing parameters  $\alpha_j$  can be thought of as a prior probabilities of each mixture:

$$\alpha_j = p(\text{component } j) \quad (35)$$



Thus

You do a first guess for the parameters at the beginning of EM

$$\Theta^g = (\alpha_1^g, \dots, \alpha_M^g, \theta_1^g, \dots, \theta_M^g) \quad (34)$$

Then, it is possible to calculate given the parametric probability

$$p_j(x_i | \theta_j^g)$$

Therefore

The mixing parameters  $\alpha_j$  can be thought of as a prior probabilities of each mixture:

$$\alpha_j = p(\text{component } j) \quad (35)$$



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - **Bayes' Rule for the components**
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm

# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - **Mixing Parameters**
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm





We want to calculate the following probability

We want to calculate

$$p(y_i|x_i, \Theta^g)$$



# We want to calculate the following probability

## We want to calculate

$$p(y_i|x_i, \Theta^g)$$

## Basically

We want a Bayesian formulation of this probability.

- Assuming that the  $y = (y_1, y_2, \dots, y_N)$  are samples identically independent samples from a distribution.



# We want to calculate the following probability

## We want to calculate

$$p(y_i | x_i, \Theta^g)$$

## Basically

We want a Bayesian formulation of this probability.

- Assuming that the  $\mathbf{y} = (y_1, y_2, \dots, y_N)$  are samples identically independent samples from a distribution.



# Using Bayes' Rule

## Compute

$$\begin{aligned} p(y_i|x_i, \Theta^g) &= \frac{p(y_i, x_i|\Theta^g)}{p(x_i|\Theta^g)} \\ &= \frac{p(x_i|\Theta^g) p(y_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)} \quad \text{We know } \theta_{y_i}^g \Rightarrow \text{Drop it} \\ &= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)} \\ &= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)} \end{aligned}$$



# Using Bayes' Rule

## Compute

$$p(y_i|x_i, \Theta^g) = \frac{p(y_i, x_i|\Theta^g)}{p(x_i|\Theta^g)}$$
$$= \frac{p(x_i|\Theta^g) p(y_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)}$$

We know  $\theta_{y_i}^g \Rightarrow$  Drop it

$$= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)}$$

$$= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)}$$



# Using Bayes' Rule

## Compute

$$\begin{aligned} p(y_i|x_i, \Theta^g) &= \frac{p(y_i, x_i|\Theta^g)}{p(x_i|\Theta^g)} \\ &= \frac{p(x_i|\Theta^g) p(y_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)} \quad \text{We know } \theta_{y_i}^g \Rightarrow \text{Drop it} \\ &= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)} \end{aligned}$$

$$\begin{aligned} &= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)} \end{aligned}$$



# Using Bayes' Rule

## Compute

$$\begin{aligned} p(y_i|x_i, \Theta^g) &= \frac{p(y_i, x_i|\Theta^g)}{p(x_i|\Theta^g)} \\ &= \frac{p(x_i|\Theta^g) p(y_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)} \quad \text{We know } \theta_{y_i}^g \Rightarrow \text{Drop it} \\ &= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{p(x_i|\Theta^g)} \\ &= \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)} \end{aligned}$$



## As in Naive Bayes

We have the fact that there is a probability per probability at the mixture and sample

$$p(y_i|x_i, \Theta^g) = \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)} \quad \forall x_i, y_i \text{ and } k \in \{1, \dots, M\}$$

This is going to be updated at each iteration of the EM algorithm.  
After the initial Guess!!! Until convergence!!!





## As in Naive Bayes

We have the fact that there is a probability per probability at the mixture and sample

$$p(y_i|x_i, \Theta^g) = \frac{\alpha_{y_i}^g p_{y_i}(x_i|\theta_{y_i}^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)} \quad \forall x_i, y_i \text{ and } k \in \{1, \dots, M\}$$

This is going to be updated at each iteration of the EM algorithm  
After the initial Guess!!! Until convergence!!!



## Additionally

We assume again that the samples  $y_i$ 's are identically and independent samples

$$p(\mathbf{y}|\mathcal{X}, \Theta^g) = \prod_{i=1}^N p(y_i|x_i, \Theta^g) \quad (36)$$

Where  $\mathbf{y} = (y_1, y_2, \dots, y_N)$



Now, using equation 17

Then

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

$$= \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{i=1}^N \log[\alpha_{y_i} p_{y_i}(x_i|\theta_{y_i})] \prod_{j=1}^N p(y_j|x_j, \Theta^g)$$



Now, using equation 17

Then

$$\begin{aligned} Q(\Theta|\Theta^g) &= \sum_{\mathbf{y} \in \mathcal{Y}} \log(\mathcal{L}(\Theta|\mathcal{X}, \mathbf{y})) p(\mathbf{y}|\mathcal{X}, \Theta^g) \\ &= \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{i=1}^N \log[\alpha_{y_i} p_{y_i}(x_i|\theta_{y_i})] \prod_{j=1}^N p(y_j|x_j, \Theta^g) \end{aligned}$$



## Here, a small stop

What is the meaning of  $\sum_{y \in \mathcal{Y}}$

It is actually a summation of all possible states of the random vector  $y$ .

Then, we can rewrite the previous summation as

$$\sum_{y \in \mathcal{Y}} = \underbrace{\sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M}_N$$

Running over all the samples  $\{x_1, x_2, \dots, x_N\}$ .



## Here, a small stop

What is the meaning of  $\sum_{\mathbf{y} \in \mathcal{Y}}$

It is actually a summation of all possible states of the random vector  $\mathbf{y}$ .

Then, we can rewrite the previous summation as

$$\sum_{\mathbf{y} \in \mathcal{Y}} = \underbrace{\sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M}_N$$

Running over all the samples  $\{x_1, x_2, \dots, x_N\}$ .



Then

We have

$$Q(\Theta|\Theta^g) = \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \left[ \log [\alpha_{y_i} p_{y_i}(x_i|\theta_{y_i})] \prod_{j=1}^N p(y_j|x_j, \Theta^g) \right]$$



We introduce the following

We have the following function

$$\delta_{l,y_i} = \begin{cases} 1 & I = y_i \\ 0 & I \neq y_i \end{cases}$$

Therefore, we can do the following

$$\alpha_i = \sum_{j=1}^M \delta_{i,j} \alpha_j$$

Then

$$\log [\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})] \prod_{j=1}^N p(y_j | x_j, \Theta^g) = \sum_{l=1}^M \delta_{l,y_i} \log [\alpha_l p_l(x_i | \theta_l)] \prod_{j=1}^N p(y_j | x_j, \Theta^g)$$



We introduce the following

We have the following function

$$\delta_{l,y_i} = \begin{cases} 1 & I = y_i \\ 0 & I \neq y_i \end{cases}$$

Therefore, we can do the following

$$\alpha_i = \sum_{j=1}^M \delta_{i,j} \alpha_j$$

Then

$$\log [\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})] \prod_{j=1}^N p(y_j | x_j, \Theta^g) = \sum_{l=1}^M \delta_{l,y_i} \log [\alpha_l p_l(x_i | \theta_l)] \prod_{j=1}^N p(y_j | x_j, \Theta^g)$$

We introduce the following

We have the following function

$$\delta_{l,y_i} = \begin{cases} 1 & I = y_i \\ 0 & I \neq y_i \end{cases}$$

Therefore, we can do the following

$$\alpha_i = \sum_{j=1}^M \delta_{i,j} \alpha_j$$

Then

$$\log [\alpha_{y_i} p_{y_i} (x_i | \theta_{y_i})] \prod_{j=1}^N p (y_j | x_j, \Theta^g) = \sum_{l=1}^M \delta_{l,y_i} \log [\alpha_l p_l (x_i | \theta_l)] \prod_{j=1}^N p (y_j | x_j, \Theta^g)$$

Thus

We have that for

$$\sum_{y_1=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \log [\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})] \prod_{j=1}^N p(y_j | x_j, \Theta^g) = *$$

$$* = \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \sum_{l=1}^M \delta_{l,y_i} \log [\alpha_l p_l(x_i | \theta_l)] \prod_{j=1}^N p(y_j | x_j, \Theta^g)$$

$$= \sum_{l=1}^M \sum_{i=1}^N \log [\alpha_l p_l(x_i | \theta_l)] \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \left[ \delta_{l,y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) \right]$$



# Thus

We have that for

$$\sum_{y_1=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \log [\alpha_{y_i} p_{y_i}(x_i | \theta_{y_i})] \prod_{j=1}^N p(y_j | x_j, \Theta^g) = *$$

$$\begin{aligned} * &= \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \sum_{l=1}^M \delta_{l,y_i} \log [\alpha_l p_l(x_i | \theta_l)] \prod_{j=1}^N p(y_j | x_j, \Theta^g) \\ &= \sum_{i=1}^N \sum_{l=1}^M \log [\alpha_l p_l(x_i | \theta_l)] \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \left[ \delta_{l,y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) \right] \end{aligned}$$

Because

$\sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M$  applies only to  $\delta_{l,y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g)$



# Thus

We have that for

$$\sum_{y_1=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \log [\alpha_{y_i} p_{y_i} (x_i | \theta_{y_i})] \prod_{j=1}^N p (y_j | x_j, \Theta^g) = *$$

$$\begin{aligned} * &= \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \sum_{i=1}^N \sum_{l=1}^M \delta_{l,y_i} \log [\alpha_l p_l (x_i | \theta_l)] \prod_{j=1}^N p (y_j | x_j, \Theta^g) \\ &= \sum_{i=1}^N \sum_{l=1}^M \log [\alpha_l p_l (x_i | \theta_l)] \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \left[ \delta_{l,y_i} \prod_{j=1}^N p (y_j | x_j, \Theta^g) \right] \end{aligned}$$

Because

$$\sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \text{ applies only to } \delta_{l,y_i} \prod_{j=1}^N p (y_j | x_j, \Theta^g)$$



Then, we have that

First notice the following

$$\sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \left[ \delta_{l, y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) \right] =$$

$$= \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \left\{ \left[ \sum_{y_i=1}^M \delta_{l, y_i} p(y_i | x_i, \Theta^g) \right] \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right\} \right)$$



Then, we have that

First notice the following

$$\begin{aligned} & \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \left[ \delta_{l, y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) \right] = \\ & = \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \left\{ \left[ \sum_{y_i=1}^M \delta_{l, y_i} p(y_i | x_i, \Theta^g) \right] \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right\} \right) \end{aligned}$$

Then, we have

$$\sum_{y_i=1}^M \delta_{l, y_i} p(y_i | x_i, \Theta^g) = p(l | x_i, \Theta^g)$$



Then, we have that

First notice the following

$$\begin{aligned} & \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \left[ \delta_{l,y_i} \prod_{j=1}^N p(y_j|x_j, \Theta^g) \right] = \\ & = \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \left\{ \left[ \sum_{y_i=1}^M \delta_{l,y_i} p(y_i|x_i, \Theta^g) \right] \prod_{j=1, j \neq i}^N p(y_j|x_j, \Theta^g) \right\} \right) \end{aligned}$$

Then, we have

$$\sum_{y_i=1}^M \delta_{l,y_i} p(y_i|x_i, \Theta^g) = p(l|x_i, \Theta^g)$$





In this way

Plugging back the previous equation

$$\sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \delta_{l, y_i} \prod_{j=1}^N p(y_j | x_j, \Theta^g) =$$

$$= \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M p(l | x_i, \Theta^g) \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right)$$
$$= \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right) p(l | x_i, \Theta^g)$$



In this way

Plugging back the previous equation

$$\begin{aligned} & \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \delta_{l,y_i} \prod_{j=1}^N p(y_j|x_j, \Theta^g) = \\ & = \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M p(l|x_i, \Theta^g) \prod_{j=1, j \neq i}^N p(y_j|x_j, \Theta^g) \right) \\ & = \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j|x_j, \Theta^g) \right) p(l|x_i, \Theta^g) \end{aligned}$$



In this way

Plugging back the previous equation

$$\begin{aligned} & \sum_{y_1=1}^M \sum_{y_2=1}^M \cdots \sum_{y_N=1}^M \delta_{l,y_i} \prod_{j=1}^N p(y_j|x_j, \Theta^g) = \\ & = \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M p(l|x_i, \Theta^g) \prod_{j=1, j \neq i}^N p(y_j|x_j, \Theta^g) \right) \\ & = \left( \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j|x_j, \Theta^g) \right) p(l|x_i, \Theta^g) \end{aligned}$$



Now, what about...?

The left part of the equation

$$\begin{aligned} & \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) = \\ & \left[ \sum_{y_1=1}^M p(y_1 | x_1, \Theta^g) \right] \cdots \left[ \sum_{y_{i-1}=1}^M p(y_{i-1} | x_{i-1}, \Theta^g) \right] \times \cdots \\ & \left[ \sum_{y_{i+1}=1}^M p(y_{i+1} | x_{i+1}, \Theta^g) \right] \cdots \left[ \sum_{y_N=1}^M p(y_N | x_N, \Theta^g) \right] \\ & = \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j | x_j, \Theta^g) \right] \end{aligned}$$

Now, what about...?

The left part of the equation

$$\begin{aligned} & \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) = \\ & = \left[ \sum_{y_1=1}^M p(y_1 | x_1, \Theta^g) \right] \cdots \left[ \sum_{y_{i-1}=1}^M p(y_{i-1} | x_{i-1}, \Theta^g) \right] \times \cdots \\ & \quad \left[ \sum_{y_{i+1}=1}^M p(y_{i+1} | x_{i+1}, \Theta^g) \right] \cdots \left[ \sum_{y_N=1}^M p(y_N | x_N, \Theta^g) \right] \end{aligned}$$

$$= \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j | x_j, \Theta^g) \right]$$

## Now, what about...?

The left part of the equation

$$\begin{aligned} & \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) = \\ &= \left[ \sum_{y_1=1}^M p(y_1 | x_1, \Theta^g) \right] \cdots \left[ \sum_{y_{i-1}=1}^M p(y_{i-1} | x_{i-1}, \Theta^g) \right] \times \cdots \\ & \left[ \sum_{y_{i+1}=1}^M p(y_{i+1} | x_{i+1}, \Theta^g) \right] \cdots \left[ \sum_{y_N=1}^M p(y_N | x_N, \Theta^g) \right] \\ &= \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j | x_j, \Theta^g) \right] \end{aligned}$$

Then, we have that

Plugging back to the original equation

$$\left\{ \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right\} p(l | x_i, \Theta^g) =$$

$$= \left\{ \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j | x_j, \Theta^g) \right] \right\} p(l | x_i, \Theta^g)$$



Then, we have that

Plugging back to the original equation

$$\left\{ \sum_{y_1=1}^M \cdots \sum_{y_{i-1}=1}^M \sum_{y_{i+1}=1}^M \cdots \sum_{y_N=1}^M \prod_{j=1, j \neq i}^N p(y_j | x_j, \Theta^g) \right\} p(l | x_i, \Theta^g) =$$
$$= \left\{ \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j | x_j, \Theta^g) \right] \right\} p(l | x_i, \Theta^g)$$





## We can use properties of probability

We know that

$$\sum_{y_i=1}^M p(y_i|x_i, \Theta^g) = 1 \quad (37)$$

Then

$$\left\{ \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j|x_j, \Theta^g) \right] \right\} p(l|x_i, \Theta^g) =$$

## We can use properties of probability

We know that

$$\sum_{y_i=1}^M p(y_i|x_i, \Theta^g) = 1 \quad (37)$$

Then

$$\begin{aligned} & \left\{ \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j|x_j, \Theta^g) \right] \right\} p(l|x_i, \Theta^g) = \\ & = \left\{ \prod_{j=1, j \neq i}^N 1 \right\} p(l|x_i, \Theta^g) \end{aligned}$$

$$= p(l|x_i, \Theta^g)$$

$$= \frac{\alpha_i^g p_{y_i}(x_i|\theta_i^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)}$$

## We can use properties of probability

We know that

$$\sum_{y_i=1}^M p(y_i|x_i, \Theta^g) = 1 \quad (37)$$

Then

$$\begin{aligned} & \left\{ \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j|x_j, \Theta^g) \right] \right\} p(l|x_i, \Theta^g) = \\ & = \left\{ \prod_{j=1, j \neq i}^N 1 \right\} p(l|x_i, \Theta^g) \\ & = p(l|x_i, \Theta^g) \end{aligned}$$

$$\frac{\alpha_i^g p_{y_i}(x_i|\theta_i^g)}{\sum_{k=1}^N \alpha_k^g p_k(x_i|\theta_k^g)}$$

## We can use properties of probability

We know that

$$\sum_{y_i=1}^M p(y_i|x_i, \Theta^g) = 1 \quad (37)$$

Then

$$\begin{aligned} & \left\{ \prod_{j=1, j \neq i}^N \left[ \sum_{y_j=1}^M p(y_j|x_j, \Theta^g) \right] \right\} p(l|x_i, \Theta^g) = \\ &= \left\{ \prod_{j=1, j \neq i}^N 1 \right\} p(l|x_i, \Theta^g) \\ &= p(l|x_i, \Theta^g) \\ &= \frac{\alpha_l^g p_{y_i}(x_i|\theta_l^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)} \end{aligned}$$

Thus

We can write  $Q$  in the following way

$$\begin{aligned} Q(\Theta, \Theta^g) &= \sum_{i=1}^N \sum_{l=1}^M \log [\alpha_l p_l(x_i | \theta_l)] p(l | x_i, \Theta^g) \\ &= \sum_{i=1}^N \sum_{l=1}^M \log(\alpha_l) p(l | x_i, \Theta^g) + \dots \\ &\quad \sum_{i=1}^N \sum_{l=1}^M \log(p_l(x_i | \theta_l)) p(l | x_i, \Theta^g) \end{aligned} \quad (38)$$



Thus

We can write  $Q$  in the following way

$$\begin{aligned} Q(\Theta, \Theta^g) &= \sum_{i=1}^N \sum_{l=1}^M \log [\alpha_l p_l(x_i | \theta_l)] p(l | x_i, \Theta^g) \\ &= \sum_{i=1}^N \sum_{l=1}^M \log(\alpha_l) p(l | x_i, \Theta^g) + \dots \end{aligned}$$

$$\sum_{i=1}^N \sum_{l=1}^M \log(p_l(x_i | \theta_l)) p(l | x_i, \Theta^g) \quad (38)$$



Thus

We can write  $Q$  in the following way

$$\begin{aligned} Q(\Theta, \Theta^g) &= \sum_{i=1}^N \sum_{l=1}^M \log [\alpha_l p_l(x_i | \theta_l)] p(l | x_i, \Theta^g) \\ &= \sum_{i=1}^N \sum_{l=1}^M \log(\alpha_l) p(l | x_i, \Theta^g) + \dots \\ &\quad \sum_{i=1}^N \sum_{l=1}^M \log(p_l(x_i | \theta_l)) p(l | x_i, \Theta^g) \end{aligned} \quad (38)$$



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - **Maximizing  $Q$  using Lagrange Multipliers**
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# A Method

That could be used as a general framework

To solve problems set as EM problem.

First, we will look at the Lagrange Multipliers setup

Then, we will look at a specific case using the mixture of Gaussian's

Note:

Not all the mixture of distributions will get you an analytical solution.



# A Method

That could be used as a general framework

To solve problems set as EM problem.

First, we will look at the Lagrange Multipliers setup

Then, we will look at a specific case using the mixture of Gaussian's

Not all the mixture of distributions will get you an analytical solution.



# A Method

That could be used as a general framework

To solve problems set as EM problem.

First, we will look at the Lagrange Multipliers setup

Then, we will look at a specific case using the mixture of Gaussian's

## Note

Not all the mixture of distributions will get you an analytical solution.



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



## Lagrange Multipliers for $Q$

We can use the following constraint for that

$$\sum_l \alpha_l = 1 \quad (39)$$

We have the following cost function

$$Q(\theta, \theta^y) + \lambda \left( \sum_l \alpha_l - 1 \right) \quad (40)$$

Deriving by  $\alpha_l$

$$\frac{\partial}{\partial \alpha_l} \left[ Q(\theta, \theta^y) + \lambda \left( \sum_l \alpha_l - 1 \right) \right] = 0 \quad (41)$$



## Lagrange Multipliers for $Q$

We can use the following constraint for that

$$\sum_l \alpha_l = 1 \quad (39)$$

We have the following cost function

$$Q(\Theta, \Theta^g) + \lambda \left( \sum_l \alpha_l - 1 \right) \quad (40)$$

Deriving by  $\alpha_l$

$$\frac{\partial}{\partial \alpha_l} \left[ Q(\Theta, \Theta^g) + \lambda \left( \sum_l \alpha_l - 1 \right) \right] = 0 \quad (41)$$

## Lagrange Multipliers for $Q$

We can use the following constraint for that

$$\sum_l \alpha_l = 1 \quad (39)$$

We have the following cost function

$$Q(\Theta, \Theta^g) + \lambda \left( \sum_l \alpha_l - 1 \right) \quad (40)$$

Deriving by  $\alpha_l$

$$\frac{\partial}{\partial \alpha_l} \left[ Q(\Theta, \Theta^g) + \lambda \left( \sum_l \alpha_l - 1 \right) \right] = 0 \quad (41)$$



Thus

## The $Q$ function

$$Q(\Theta, \Theta^g) = \sum_{i=1}^N \sum_{l=1}^M \log(\alpha_l) p(l|x_i, \Theta^g) + \dots$$
$$\sum_{i=1}^N \sum_{l=1}^M \log(p_l(x_i|\theta_l)) p(l|x_i, \Theta^g)$$





# Deriving

We have

$$\frac{\partial}{\partial \alpha_l} \left[ Q(\Theta, \Theta^g) + \lambda \left( \sum_l \alpha_l - 1 \right) \right] = \sum_{i=1}^N \frac{1}{\alpha_l} p(l|x_i, \Theta^g) + \lambda$$



## Finally

We have making the previous equation equal to 0

$$\sum_{i=1}^N \frac{1}{\alpha_l} p(l|x_i, \Theta^g) + \lambda = 0 \quad (42)$$

Thus

$$\sum_{i=1}^N p(l|x_i, \Theta^g) = -\lambda \alpha_l \quad (43)$$

Summing over  $l$ , we get

$$\lambda = -N \quad (44)$$



## Finally

We have making the previous equation equal to 0

$$\sum_{i=1}^N \frac{1}{\alpha_l} p(l|x_i, \Theta^g) + \lambda = 0 \quad (42)$$

Thus

$$\sum_{i=1}^N p(l|x_i, \Theta^g) = -\lambda \alpha_l \quad (43)$$

Summing over  $l$ , we get

$$\lambda = -N \quad (44)$$



## Finally

We have making the previous equation equal to 0

$$\sum_{i=1}^N \frac{1}{\alpha_l} p(l|x_i, \Theta^g) + \lambda = 0 \quad (42)$$

Thus

$$\sum_{i=1}^N p(l|x_i, \Theta^g) = -\lambda \alpha_l \quad (43)$$

Summing over  $l$ , we get

$$\lambda = -N \quad (44)$$



# Lagrange Multipliers

Thus

$$\alpha_l = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g) \quad (45)$$

About  $\theta_l$

It is possible to get an analytical expressions for  $\theta_l$  as functions of everything else.

- This is for you to try!!!

For more, please look at

“Geometric Idea of Lagrange Multipliers” by John Wyatt.



# Lagrange Multipliers

Thus

$$\alpha_l = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g) \quad (45)$$

About  $\theta_l$

It is possible to get an analytical expressions for  $\theta_l$  as functions of everything else.

- This is for you to try!!!

For more, please read:

“Geometric Idea of Lagrange Multipliers” by John Wyatt.



Cinvestav

# Lagrange Multipliers

Thus

$$\alpha_l = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g) \quad (45)$$

About  $\theta_l$

It is possible to get an analytical expressions for  $\theta_l$  as functions of everything else.

- This is for you to try!!!

For more, please look at

“Geometric Idea of Lagrange Multipliers” by John Wyatt.



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - **Example on Mixture of Gaussian Distributions**
  - The EM Algorithm





## Remember?

### Gaussian Distribution

$$p_l(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_l|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_l)^T \boldsymbol{\Sigma}_l^{-1}(\mathbf{x} - \boldsymbol{\mu}_l)\right\} \quad (46)$$



# How to use this for Gaussian Distributions

For this, we need to refresh some linear algebra

①  $tr(A + B) = tr(A) + tr(B)$

②  $tr(AB) = tr(BA)$

③  $\sum_i x_i^T A x_i = tr(AB)$  where  $B = \sum_i x_i x_i^T$ .

④  $|A^{-1}| = \frac{1}{|A|}$



# How to use this for Gaussian Distributions

For this, we need to refresh some linear algebra

①  $tr(A + B) = tr(A) + tr(B)$

②  $tr(AB) = tr(BA)$

③  $\sum_i x_i^T A x_i = tr(AB)$  where  $B = \sum_i x_i x_i^T$ .

④  $|A^{-1}| = \frac{1}{|A|}$

Now, we need the derivative of a matrix function

Thus,  $\frac{\partial f(A)}{\partial A}$  is going to be the matrix with  $i, j^{th}$  entry  $\left[ \frac{\partial f(A)}{\partial a_{i,j}} \right]$  where  $a_{i,j}$  is the  $i, j^{th}$  entry of  $A$ .



# How to use this for Gaussian Distributions

For this, we need to refresh some linear algebra

- 1  $tr(A + B) = tr(A) + tr(B)$
- 2  $tr(AB) = tr(BA)$
- 3  $\sum_i x_i^T A x_i = tr(AB)$  where  $B = \sum_i x_i x_i^T$ .

4  $|A^{-1}| = \frac{1}{|A|}$

Now, we need the derivative of a matrix function

Thus,  $\frac{\partial f(A)}{\partial A}$  is going to be the matrix with  $i, j^{th}$  entry  $\left[ \frac{\partial f(A)}{\partial a_{i,j}} \right]$  where  $a_{i,j}$  is the  $i, j^{th}$  entry of  $A$ .



# How to use this for Gaussian Distributions

For this, we need to refresh some linear algebra

- 1  $tr(A + B) = tr(A) + tr(B)$
- 2  $tr(AB) = tr(BA)$
- 3  $\sum_i x_i^T A x_i = tr(AB)$  where  $B = \sum_i x_i x_i^T$ .
- 4  $|A^{-1}| = \frac{1}{|A|}$

Now, we need the derivative of a matrix function

Thus,  $\frac{\partial f(A)}{\partial A}$  is going to be the matrix with  $i, j^{th}$  entry  $\left[ \frac{\partial f(A)}{\partial a_{i,j}} \right]$  where  $a_{i,j}$  is the  $i, j^{th}$  entry of  $A$ .



# How to use this for Gaussian Distributions

For this, we need to refresh some linear algebra

- 1  $tr(A + B) = tr(A) + tr(B)$
- 2  $tr(AB) = tr(BA)$
- 3  $\sum_i x_i^T A x_i = tr(AB)$  where  $B = \sum_i x_i x_i^T$ .
- 4  $|A^{-1}| = \frac{1}{|A|}$

Now, we need the derivative of a matrix function  $f(A)$

Thus,  $\frac{\partial f(A)}{\partial A}$  is going to be the matrix with  $i, j^{th}$  entry  $\left[ \frac{\partial f(A)}{\partial a_{i,j}} \right]$  where  $a_{i,j}$  is the  $i, j^{th}$  entry of  $A$ .



## In addition

If  $A$  is symmetric

$$\frac{\partial |A|}{\partial A} = \begin{cases} \mathcal{A}_{i,j} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} \quad (47)$$

Where  $\mathcal{A}_{i,j}$  is the  $i, j^{\text{th}}$  cofactor of  $A$ .

Note: The determinant obtained by deleting the row and column of a given element of a matrix or determinant. The cofactor is preceded by a + or - sign depending whether the element is in a + or - position.

## In addition

If  $A$  is symmetric

$$\frac{\partial |A|}{\partial A} = \begin{cases} \mathcal{A}_{i,j} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} \quad (47)$$

Where  $\mathcal{A}_{i,j}$  is the  $i, j^{\text{th}}$  cofactor of  $A$ .

Note: The determinant obtained by deleting the row and column of a given element of a matrix or determinant. The cofactor is preceded by a + or - sign depending whether the element is in a + or - position.

Thus

$$\frac{\partial \log |A|}{\partial A} = \begin{cases} \frac{\mathcal{A}_{i,j}}{|A|} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} = 2A^{-1} - \text{diag}(A^{-1}) \quad (48)$$



## In addition

If  $A$  is symmetric

$$\frac{\partial |A|}{\partial A} = \begin{cases} \mathcal{A}_{i,j} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} \quad (47)$$

Where  $\mathcal{A}_{i,j}$  is the  $i, j^{\text{th}}$  cofactor of  $A$ .

**Note:** The determinant obtained by deleting the row and column of a given element of a matrix or determinant. The **cofactor** is preceded by a + or - sign depending whether the element is in a + or - position.

Thus

$$\frac{\partial \log |A|}{\partial A} = \begin{cases} \frac{\mathcal{A}_{i,i}}{|A|} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} = 2A^{-1} - \text{diag}(A^{-1}) \quad (48)$$

## In addition

If  $A$  is symmetric

$$\frac{\partial |A|}{\partial A} = \begin{cases} \mathcal{A}_{i,j} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} \quad (47)$$

Where  $\mathcal{A}_{i,j}$  is the  $i, j^{\text{th}}$  cofactor of  $A$ .

**Note:** The determinant obtained by deleting the row and column of a given element of a matrix or determinant. The **cofactor** is preceded by a + or - sign depending whether the element is in a + or - position.

Thus

$$\frac{\partial \log |A|}{\partial A} = \begin{cases} \frac{\mathcal{A}_{i,j}}{|A|} & \text{if } i = j \\ 2\mathcal{A}_{i,j} & \text{if } i \neq j \end{cases} = 2A^{-1} - \text{diag}(A^{-1}) \quad (48)$$

# Finally

The last equation we need

$$\frac{\partial \text{tr}(AB)}{\partial A} = B + B^T - \text{diag}(B) \quad (49)$$

in addition

$$\frac{\partial x^T A x}{\partial x} \quad (50)$$



# Finally

The last equation we need

$$\frac{\partial \text{tr}(AB)}{\partial A} = B + B^T - \text{diag}(B) \quad (49)$$

In addition

$$\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial \mathbf{x}} \quad (50)$$



Cinvestav

Thus, using last part of equation 38

We get, after ignoring constant terms

Remember they disappear after derivatives

$$\sum_{i=1}^N \sum_{l=1}^M \log (p_l (\mathbf{x}_i | \mu_l, \Sigma_l)) p(l | \mathbf{x}_i, \Theta^g)$$

$$= \sum_{i=1}^N \sum_{l=1}^M \left[ -\frac{1}{2} \log (|\Sigma_l|) - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right] p(l | \mathbf{x}_i, \Theta^g) \quad (51)$$



Thus, using last part of equation 38

We get, after ignoring constant terms

Remember they disappear after derivatives

$$\begin{aligned} & \sum_{i=1}^N \sum_{l=1}^M \log (p_l (\mathbf{x}_i | \mu_l, \Sigma_l)) p (l | \mathbf{x}_i, \Theta^g) \\ &= \sum_{i=1}^N \sum_{l=1}^M \left[ -\frac{1}{2} \log (|\Sigma_l|) - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right] p (l | \mathbf{x}_i, \Theta^g) \quad (51) \end{aligned}$$



# Finally

Thus, when taking the derivative with respect to  $\mu_l$

$$\sum_{i=1}^N \left[ \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) p(l|\mathbf{x}_i, \Theta^g) \right] = 0 \quad (52)$$

Then

$$\mu_l = \frac{\sum_{i=1}^N \mathbf{x}_i p(l|\mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)} \quad (53)$$



# Finally

Thus, when taking the derivative with respect to  $\mu_l$

$$\sum_{i=1}^N \left[ \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) p(l|\mathbf{x}_i, \Theta^g) \right] = 0 \quad (52)$$

Then

$$\mu_l = \frac{\sum_{i=1}^N \mathbf{x}_i p(l|\mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)} \quad (53)$$





Now, if we derive with respect to  $\Sigma_l$

First, we rewrite equation 51

$$\sum_{i=1}^N \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right] p(l|\mathbf{x}_i, \Theta^g)$$

$$= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \{ \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T \} \right]$$

$$= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \{ \Sigma_l^{-1} N_{l,i} \} \right]$$

Where  $N_{l,i} = (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T$ .



Now, if we derive with respect to  $\Sigma_l$

First, we rewrite equation 51

$$\begin{aligned} & \sum_{i=1}^N \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right] p(l|\mathbf{x}_i, \Theta^g) \\ &= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \left\{ \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T \right\} \right] \\ &= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \left\{ \Sigma_l^{-1} N_{l,i} \right\} \right] \end{aligned}$$

Where  $N_{l,i} = (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T$ .



Now, if we derive with respect to  $\Sigma_l$

First, we rewrite equation 51

$$\begin{aligned} & \sum_{i=1}^N \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right] p(l|\mathbf{x}_i, \Theta^g) \\ &= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \left\{ \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T \right\} \right] \\ &= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \left\{ \Sigma_l^{-1} N_{l,i} \right\} \right] \end{aligned}$$

Where  $N_{l,i} = (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T$ .



Now, if we derive with respect to  $\Sigma_l$

First, we rewrite equation 51

$$\begin{aligned} & \sum_{i=1}^N \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) - \frac{1}{2} (\mathbf{x}_i - \mu_l)^T \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) \right] p(l|\mathbf{x}_i, \Theta^g) \\ &= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \left\{ \Sigma_l^{-1} (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T \right\} \right] \\ &= \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \left\{ \Sigma_l^{-1} N_{l,i} \right\} \right] \end{aligned}$$

Where  $N_{l,i} = (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T$ .



# Deriving with respect to $\Sigma_l^{-1}$

We have that

$$\begin{aligned} \frac{\partial}{\partial \Sigma_l^{-1}} \sum_{l=1}^M & \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \{ \Sigma_l^{-1} N_{l,i} \} \right] \\ &= -\frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2\Sigma_l - \text{diag}(\Sigma_l)) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2N_{l,i} - \text{diag}(N_{l,i})) \\ &= -\frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2M_{l,i} - \text{diag}(M_{l,i})) \\ &= -2S - \text{diag}(S) \end{aligned}$$

Where  $M_{l,i} = \Sigma_l - N_{l,i}$  and  $S = \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) M_{l,i}$



# Deriving with respect to $\Sigma_l^{-1}$

We have that

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_l^{-1}} \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \{ \Sigma_l^{-1} N_{l,i} \} \right] \\ &= \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2\Sigma_l - \text{diag}(\Sigma_l)) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2N_{l,i} - \text{diag}(N_{l,i})) \\ &= \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2M_{l,i} - \text{diag}(M_{l,i})) \\ &= 2S - \text{diag}(S) \end{aligned}$$

Where  $M_{l,i} = \Sigma_l - N_{l,i}$  and  $S = \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) M_{l,i}$



# Deriving with respect to $\Sigma_l^{-1}$

We have that

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_l^{-1}} \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \{ \Sigma_l^{-1} N_{l,i} \} \right] \\ &= \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2\Sigma_l - \text{diag}(\Sigma_l)) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2N_{l,i} - \text{diag}(N_{l,i})) \\ &= \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2M_{l,i} - \text{diag}(M_{l,i})) \\ &= 2S - \text{diag}(S) \end{aligned}$$

Where  $M_{l,i} = \Sigma_l - N_{l,i}$  and  $S = \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) M_{l,i}$



# Deriving with respect to $\Sigma_l^{-1}$

We have that

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_l^{-1}} \sum_{l=1}^M \left[ -\frac{1}{2} \log(|\Sigma_l|) \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) \text{tr} \{ \Sigma_l^{-1} N_{l,i} \} \right] \\ &= \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2\Sigma_l - \text{diag}(\Sigma_l)) - \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2N_{l,i} - \text{diag}(N_{l,i})) \\ &= \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (2M_{l,i} - \text{diag}(M_{l,i})) \\ &= 2S - \text{diag}(S) \end{aligned}$$

Where  $M_{l,i} = \Sigma_l - N_{l,i}$  and  $S = \frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) M_{l,i}$





Thus, we have

Thus

$$\text{If } 2S - \text{diag}(S) = 0 \implies S = 0$$

implying

$$\frac{1}{2} \sum_{i=1}^N p(l|x_i, \Theta^g) [\Sigma_l - N_{l,i}] = 0 \quad (54)$$

Or

$$\Sigma_l = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) N_{l,i}}{\sum_{i=1}^N p(l|x_i, \Theta^g)} = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) (x_i - \mu_l) (x_i - \mu_l)^T}{\sum_{i=1}^N p(l|x_i, \Theta^g)} \quad (55)$$

Thus, we have

Thus

$$\text{If } 2S - \text{diag}(S) = 0 \implies S = 0$$

Implying

$$\frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) [\Sigma_l - N_{l,i}] = 0 \quad (54)$$

$$\Sigma_l = \frac{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) N_{l,i}}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)} = \frac{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)} \quad (55)$$

Thus, we have

Thus

$$\text{If } 2S - \text{diag}(S) = 0 \implies S = 0$$

Implying

$$\frac{1}{2} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) [\Sigma_l - N_{l,i}] = 0 \quad (54)$$

Or

$$\Sigma_l = \frac{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) N_{l,i}}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)} = \frac{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)} \quad (55)$$

Thus, we have the iterative updates

They are

$$\alpha_l^{New} = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g)$$

$$\mu_l^{New} = \frac{\sum_{i=1}^N x_i p(l|x_i, \Theta^g)}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$

$$\Sigma_l^{New} = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) (x_i - \mu_l) (x_i - \mu_l)^T}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$



Thus, we have the iterative updates

They are

$$\alpha_l^{New} = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g)$$

$$\mu_l^{New} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l|x_i, \Theta^g)}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$

$$\Sigma_l^{New} = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$



Thus, we have the iterative updates

They are

$$\alpha_l^{New} = \frac{1}{N} \sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)$$

$$\mu_l^{New} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l|\mathbf{x}_i, \Theta^g)}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)}$$

$$\Sigma_l^{New} = \frac{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l|\mathbf{x}_i, \Theta^g)}$$



# Outline

- 1 Introduction
  - Maximum-Likelihood
  - Expectation Maximization
  - Examples of Applications of EM
- 2 Incomplete Data
  - Introduction
  - Using the Expected Value
  - Analogy
- 3 Derivation of the EM-Algorithm
  - Hidden Features
    - Proving Concavity
  - Using the Concave Functions for Approximation
  - From The Concave Function to the EM
  - The Final Algorithm
  - Notes and Convergence of EM
- 4 Finding Maximum Likelihood Mixture Densities
  - The Beginning of The Process
  - Bayes' Rule for the components
    - Mixing Parameters
  - Maximizing  $Q$  using Lagrange Multipliers
    - In Our Case
  - Example on Mixture of Gaussian Distributions
  - The EM Algorithm



# EM Algorithm for Gaussian Mixtures

## Step 1

Initialize:

- The means  $\mu_l$
- Covariances  $\Sigma_l$
- Mixing coefficients  $\alpha_l$





## Step 2 - E-Step

- Evaluate the the probabilities of component  $l$  given  $x_i$  using the current parameter values:

$$p(l|x_i, \Theta^g) = \frac{\alpha_l^g p_{y_i}(x_i|\theta_l^g)}{\sum_{k=1}^M \alpha_k^g p_k(x_i|\theta_k^g)}$$



## Step 3 - M-Step

- Re-estimate the parameters using the current iteration values:

$$\alpha_l^{New} = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g)$$

$$\mu_l^{New} = \frac{\sum_{i=1}^N x_i p(l|x_i, \Theta^g)}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$

$$\Sigma_l^{New} = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) (x_i - \mu_l) (x_i - \mu_l)^T}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$



## Step 3 - M-Step

- Re-estimate the parameters using the current iteration values:

$$\alpha_l^{New} = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g)$$

$$\mu_l^{New} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l|x_i, \Theta^g)}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$

$$\Sigma_l^{New} = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$



## Step 3 - M-Step

- Re-estimate the parameters using the current iteration values:

$$\alpha_l^{New} = \frac{1}{N} \sum_{i=1}^N p(l|x_i, \Theta^g)$$

$$\mu_l^{New} = \frac{\sum_{i=1}^N \mathbf{x}_i p(l|x_i, \Theta^g)}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$

$$\Sigma_l^{New} = \frac{\sum_{i=1}^N p(l|x_i, \Theta^g) (\mathbf{x}_i - \mu_l) (\mathbf{x}_i - \mu_l)^T}{\sum_{i=1}^N p(l|x_i, \Theta^g)}$$



# Evaluate

## Step 4

Evaluate the log likelihood:

$$\log p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\alpha}) = \sum_{i=1}^N \log \left\{ \sum_{l=1}^M \alpha_l^{New} p_l(\mathbf{x}_i | \boldsymbol{\mu}_l^{New}, \boldsymbol{\Sigma}_l^{New}) \right\}$$

## Step 5

- Check for convergence of either the parameters or the log likelihood.
- If the convergence criterion is not satisfied return to step 2.



# Evaluate

## Step 4

Evaluate the log likelihood:





$$\log p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\alpha}) = \sum_{i=1}^N \log \left\{ \sum_{l=1}^M \alpha_l^{New} p_l(\mathbf{x}_i|\boldsymbol{\mu}_l^{New}, \boldsymbol{\Sigma}_l^{New}) \right\}$$

## Step 6

- Check for convergence of either the parameters or the log likelihood.
- If the convergence criterion is not satisfied return to step 2.



# References I

-  S. Borman, “The expectation maximization algorithm—a short tutorial,” *Submitted for publication*, pp. 1–9, 2004.
-  J. Bilmes, “A gentle tutorial of the EM algorithm and its application to parameter estimation for Gaussian mixture and hidden Markov models,” *International Computer Science Institute*, vol. 4, 1998.
-  F. Dellaert, “The expectation maximization algorithm,” tech. rep., Georgia Institute of Technology, 2002.
-  G. McLachlan and T. Krishnan, *The EM algorithm and extensions*, vol. 382.  
John Wiley & Sons, 2007.

