# EM Algorithm

School of Data and Computer Science
Sun Yat-sen University

# Overview

# The Gaussian Distribution

- The Gaussian distribution is a widely used model for the distribution of continuous variables.
- In the case of a single variable $x$, the Gaussian distribution can be written in the form

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\{-\frac{1}{2\sigma^2}(x - \mu)^2\} \qquad (1)$$

where $\mu$ is the mean and $\sigma^2$ is the variance.

For a $D$-dimensional vector $\mathbf{x}$, the multivariate Gaussian distribution takes the form

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}\} \qquad (2)$$

where $\boldsymbol{\mu}$ is a $D$-dimensional mean vector, $\boldsymbol{\Sigma}$ is a $D \times D$ covariance matrix, and $|\boldsymbol{\Sigma}|$ denotes the determinant of $|\boldsymbol{\Sigma}|$.

## Mixtures of Gaussians

- The Gaussian distribution suffers from significant limitations when it comes to modelling real data sets.
- The data set forms two dominant clumps, and that a simple Gaussian distribution is unable to capture this structure, whereas a linear superposition of two Gaussians gives a better characterization of the data set.
- A linear combination of Gaussians can give rise to very complex densities.
- By using a sufficient number of Gaussians, and by adjusting their means and covariances as well as the coefficients in the linear combination, almost any continuous density can be approximated to arbitrary accuracy.

# Mixtures of Gaussians (cont'd)

Example of a Gaussian mixture distribution in one dimension showing three Gaussians (each scaled by a coefficient) in blue and their sum in red.
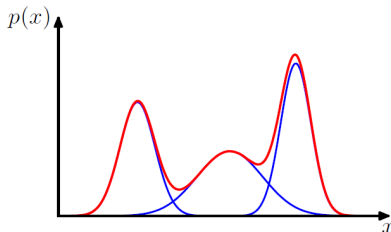


Figure 1: Example of a Gaussian mixture distribution

# Mixtures of Gaussians (cont'd)

- Consider a superposition of $K$ Gaussian densities of the form

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{3}$$

which is called a mixture of Gaussians.

- Each Gaussian density $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ is called a component of the mixture and has its own mean $\boldsymbol{\mu}_k$ and covariance $\boldsymbol{\Sigma}_k$.

- The parameters $\pi_k$ in (3) are called *mixing coefficients*.

# Mixtures of Gaussians (cont'd)

- Integrate both sides of (3) with respect to **x**, and note that both $p(\mathbf{x})$ and the individual Gaussian components are normalized, we obtain

$$\sum_{k=1}^{K} \pi_k = 1. \tag{4}$$

- $p(\mathbf{x}) \geq 0$, together with $\mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k) \geq 0$, implies $\pi_k \geq 0$ for all $k$. Combining this with condition (4) we obtain

$$0 \leq \pi_k \leq 1. \tag{5}$$

- The mixing coefficients satisfy the requirements to be probabilities.

## Mixtures of Gaussians (cont'd)

- The marginal density is given by

$$p(\mathbf{x}) = \sum_{k=1}^{K} p(k)p(\mathbf{x}|k) \qquad (6)$$

- We can view $\pi_k = p(k)$ as the prior probability of picking the $k^{th}$ component, and the density $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = p(\mathbf{x}|k)$ as the probability of $\mathbf{x}$ conditioned on $k$.

- From Bayes' theorem these are given by

$$\gamma_k(\mathbf{x}) = p(k|\mathbf{x}) = \frac{p(k)p(\mathbf{x}|k)}{\sum_l p(l)p(\mathbf{x}|l)} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_l \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}. \qquad (7)$$

## Mixtures of Gaussians (cont'd)

- The form of the Gaussian mixture distribution is governed by the parameters $\pi$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, where we have used the notation $\boldsymbol{\pi} = \{\pi_1, ..., \pi_K\}$, $\boldsymbol{\mu} = \{\boldsymbol{\mu}_1, ..., \boldsymbol{\mu}_k\}$ and $\boldsymbol{\Sigma} = \{\boldsymbol{\Sigma}_1, ..., \boldsymbol{\Sigma}_K\}$.

- One way to set the values of there parameters is to use maximum likelihood. From (3) the log of the likelihood function is given by

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln\big\{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\big\} \qquad (8)$$

where $X = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$.

- One approach to maximizing the likelihood function is to employ a powerful framework called expectation maximization (EM).

# About Latent Variables

- Let us introduce a $K$-dimensional binary random variable $\mathbf{z}$ having a 1-of-$K$ representation in which a particular element $z_k$ is equal to 1 and all other elements are equal to 0.

- The values of $z_k$ therefore satisfy $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$, and we see that there are $K$ possible states for the vector $\mathbf{z}$ according to which element is nonzero.

- We shall define the joint distribution $p(\mathbf{x}, \mathbf{z})$ in terms of a marginal distribution $p(\mathbf{z})$ and a conditional distribution $p(\mathbf{x}|\mathbf{z})$.

## About Latent Variables (cont'd)

The marginal distribution over **z** is specified in terms of the mixing coefficients $\pi_k$, such that

$$p(z_k = 1) = \pi_k \tag{9}$$

where the parameters $\{\pi_k\}$ must satisfy

$$0 \leq \pi_k \leq 1 \tag{10}$$

together with

$$\sum_{k=1}^{K} \pi_k = 1 \tag{11}$$

in order to be valid probabilities. Because **z** uses a 1-of-$K$ representation, we can also write this distribution in the form

$$p(\mathbf{z}) = \prod_{k=1}^{K} \pi_k^{z_k}. \tag{12}$$

Similarly, the conditional distribution of $\mathbf{x}$ given a particular value for $\mathbf{z}$ is a Gaussian

$$p(\mathbf{x}|z_k = 1) = (\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{13}$$

which can also be written in the form

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^{K} p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}. \tag{14}$$

## About Latent Variables (cont'd)

The joint distribution is given by $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$, and the marginal distribution of $\mathbf{x}$ is then obtained by summing the joint distribution over all possible states of $\mathbf{z}$ to give

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \qquad (15)$$

where we have made use of (13) and (14).

## About Latent Variables (cont'd)

We use $\gamma(z_k)$ to denote $p(z_k = 1|\mathbf{x})$, whose value can be found using Bayes' theorem

$$\gamma(z_k) = p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^{K} p(z_j = 1)p(\mathbf{x}|z_j = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$
(16)

We shall view $\pi_k$ as the prior probability of $z_k = 1$, and the quantity $\gamma(z_k)$ as the corresponding posterior probability once we have observed $\mathbf{x}$. As we shall see later, $\gamma(z_k)$ can also be viewed as the responsibility that component $k$ takes for 'explaining' the observation $\mathbf{x}$.

# EM for Gaussian Mixtures

- Initially, we shall motivate the EM algorithm by giving a relatively informal treatment in the context of the Gaussian mixture model.

- Let us begin by writing down the conditions that must be satisfied at a maximum of the likelihood function. Setting the derivatives of $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ with respect to the means $\boldsymbol{\mu}_k$ of the Gaussian components to zero, we obtain

$$0 = -\sum_{n=1}^{n} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{\gamma(z_{nk})} \sum_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \qquad (17)$$

## EM for Gaussian Mixtures (cont'd)

Multiplying by $\boldsymbol{\Sigma}_k^{-1}$ (which we assume to be nonsingular) and rearranging we obtain

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_n \tag{18}$$

where we have defined

$$N_k = \sum_{n=1}^{N} \gamma(z_{nk}). \tag{19}$$

# EM for Gaussian Mixtures (cont'd)

- We can interpret $N_k$ as the effective number of points assigned to cluster $k$.
- The mean $\boldsymbol{\mu}_k$ for the $k^{th}$ Gaussian component is obtained by taking a weighted mean of all of the points in the data set, in which the weighting factor for data point $\mathbf{x}_n$ is given by the posterior probability $\gamma(z_{nk})$ that component $k$ was responsible for generating $\mathbf{x}_n$.

Finally, we maximize $\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ with respect to the mixing coefficients $\pi_k$. Here we must take account of the constraint $\sum_{k=1}^{K} \pi_k = 1$. This can be achieved using a Lagrange multiplier and maximizing the following quantity

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda(\sum_{k=1}^{K} \pi_k - 1) \tag{20}$$

which gives

$$0 = \sum_{n=1}^{N} \frac{\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \tag{21}$$

where again we see the appearance of the responsibilities.

## EM for Gaussian Mixtures (cont'd)

If we now multiply both sides by $\pi_k$ and sum over $k$ making use of the constraint $\sum_{k=1}^{K} \pi_k = 1$, we find $\lambda = -N$. Using this to eliminate $\lambda$ and rearranging we obtain

$$\pi_k = \frac{N_k}{N} \qquad (22)$$

so that the mixing coefficient for the $k^{th}$ component is given by the average responsibility which that component takes for explaining the data points.

# EM Algorithm

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means $\boldsymbol{\mu}_k$, covariances $\boldsymbol{\Sigma}_k$ and mixing coefficients $\pi_k$, and evaluate the initial value of the log likelihood.

2. **E step**. Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \tag{23}$$

3. **M step**. Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_n \qquad (24)$$

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{new})(\mathbf{x}_n - \boldsymbol{\mu}_k^{new})^{\mathrm{T}} \qquad (25)$$

$$\pi_k^{new} = \frac{N_k}{N} \qquad (26)$$

where $N_k = \sum_{n=1}^{N} \gamma(z_{nk})$.

4. Evaluate the log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^{N} \ln\left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k \right\} \qquad (27)$$

and check for convergence of either the parameters or the log likelihood. If not satisfied return to step 2.

The following Chinese Football Dataset has recored the performance of 16 AFC football teams between 2005 and 2018.

| Country | 2006WorldCup | 2010WorldCup | 2014WorldCup | 2018WorldCup | 2007AsianCup | 2011AsianCup | 2015AsianCup |
|---|---|---|---|---|---|---|---|
| China | 50 | 50 | 50 | 40 | 9 | 9 | 5 |
| Japan | 28 | 9 | 29 | 15 | 4 | 1 | 5 |
| South_Korea | 17 | 15 | 27 | 19 | 3 | 3 | 2 |
| Iran | 25 | 40 | 28 | 18 | 5 | 5 | 5 |
| Saudi_Arabia | 28 | 40 | 50 | 26 | 2 | 9 | 9 |
| Iraq | 50 | 50 | 40 | 40 | 1 | 5 | 4 |
| Qatar | 50 | 40 | 40 | 40 | 9 | 5 | 9 |
| United_Arab_Emirates | 50 | 40 | 50 | 40 | 9 | 9 | 3 |
| Uzbekistan | 40 | 40 | 40 | 40 | 5 | 4 | 9 |
| Thailand | 50 | 50 | 50 | 40 | 9 | 17 | 17 |
| Vietnam | 50 | 50 | 50 | 50 | 5 | 17 | 17 |
| Oman | 50 | 50 | 40 | 50 | 9 | 17 | 9 |
| Bahrain | 40 | 40 | 50 | 50 | 9 | 9 | 9 |
| North_Korea | 40 | 32 | 50 | 50 | 17 | 9 | 9 |
| Indonesia | 50 | 50 | 50 | 50 | 9 | 17 | 17 |
| Australia | 16 | 21 | 30 | 30 | 9 | 2 | 1 |

The scoring rules are below:

- For the FIFA World Cup, teams score the same with their rankings if they enter the World Cup; teams score 50 for failing to entering the Asia Top Ten; teams score 40 for entering the Asia Top Ten but not entering the World Cup.

- For the AFC Asian Cup, teams score the same with their rankings if they finally enter the top four; teams score 5 for entering the top eight but not the top four, and 9 for entering the top sixteen but not top eight; teams score 17 for not passing the group stages.

We aim at classifying the above 16 teams into 3 classes according to their performance: the first-class, the second-class and the third-class. In our opinion, teams of Australia, Iran, South Korea and Japan belong to the first-class, while the Chinese football team belongs to the third-class.

- Assume that score vectors of teams in the same class are normally distributed, we can thus adopt the Gaussian mixture model. Please classify the teams into 3 classes by using EM algorithm.

- You should show the values of these parameters: $\gamma$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. If necessary, you can plot the clustering results. Note that $\gamma$ is essential for classifying.

- You can use MATLAB as the programming language.

# Reference

- Christopher M. Bishop, Pattern Recognition and Machine Learning.
- Stuart Russell and Peter Norvig, Artificial Intelligence: A Modern Approach.

# The End