# Ψηφιακή Επεξεργασία Φωνής

**Ενότητα 9η:** Μοντέλα Μείξης Γκαουσιανών

Στυλιανού Ιωάννης
Τμήμα Επιστήμης Υπολογιστών

# CS578- Speech Signal Processing
## Lecture 10: Gaussian Mixture Model, GMM

Yannis Stylianou

University of Crete, Computer Science Dept., Multimedia Informatics Lab
yannis@csd.uoc.gr

Univ. of Crete

# OUTLINE

- A d-dimensional random variable follows a Gaussian, or Normal, probability law: $x \rightarrow \mathcal{N}(\mu, \Sigma)$

$$g_{(\mu, \Sigma)}(x) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} \, e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

where $\mu$ is the mean vector and $\Sigma$ is the variance-covariance matrix.

- If $x \rightarrow \mathcal{N}(0, I)$ and if $y = \sqrt{\Sigma} x + \mu$, then $y \rightarrow \mathcal{N}(\mu, \Sigma)$.
- $\sqrt{\Sigma}$ defines the *standard deviation* of the random variable $x$. Note this square root is meant in the *matrix sense*.

# FORMULAS AND DEFINITIONS

- A d-dimensional random variable follows a Gaussian, or Normal, probability law: $x \to \mathcal{N}(\mu, \Sigma)$

$$g_{(\mu, \Sigma)}(x) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} \, e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

where $\mu$ is the mean vector and $\Sigma$ is the variance-covariance matrix.

- If $x \to \mathcal{N}(0, I)$ and if $y = \sqrt{\Sigma} x + \mu$, then $y \to \mathcal{N}(\mu, \Sigma)$.

- $\sqrt{\Sigma}$ defines the *standard deviation* of the random variable $x$. Note this square root is meant in the *matrix sense*.

- A d-dimensional random variable follows a Gaussian, or Normal, probability law: $x \rightarrow \mathcal{N}(\mu, \Sigma)$

$$g_{(\mu,\Sigma)}(x) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} \, e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

where $\mu$ is the mean vector and $\Sigma$ is the variance-covariance matrix.

- If $x \rightarrow \mathcal{N}(0, I)$ and if $y = \sqrt{\Sigma}\, x + \mu$, then $y \rightarrow \mathcal{N}(\mu, \Sigma)$.
- $\sqrt{\Sigma}$ defines the *standard deviation* of the random variable $x$. Note this square root is meant in the *matrix sense*.

- **Mean estimator :**

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} x_i$$

- Unbiased covariance estimator :

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \mu)(x_i - \mu)^T$$

# Formulas and definitions

- **Mean estimator :**

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} x_i$$

- **Unbiased covariance estimator :**

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \mu)(x_i - \mu)^T$$

- **Likelihood :**

$$p(x_i|\theta) = p(x_i|\mu, \Sigma) = g_{(\mu, \Sigma)}(x_i)$$

- Joint likelihood :

$$p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta) = \prod_{i=1}^{N} p(x_i|\mu, \Sigma) = \prod_{i=1}^{N} g_{(\mu, \Sigma)}(x_i)$$

for $X = \{x_1, x_2, \cdots, x_N\}$ being a set of independent identically distributed (i.i.d.) points

- **Likelihood :**

$$p(x_i|\theta) = p(x_i|\mu, \Sigma) = g_{(\mu,\Sigma)}(x_i)$$

- **Joint likelihood :**

$$p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta) = \prod_{i=1}^{N} p(x_i|\mu, \Sigma) = \prod_{i=1}^{N} g_{(\mu,\Sigma)}(x_i)$$

for $X = \{x_1, x_2, \cdots, x_N\}$ being a set of independent identically distributed (i.i.d.) points

# FORMULAS AND DEFINITIONS

- **Log likelihood :**

$$p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta) \quad \Leftrightarrow \quad \log p(X|\theta) = \sum_{i=1}^{N} \log p(x_i|\theta)$$

- *In the Gaussian case:*

$$p(x|\theta) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$\log p(x|\theta) = -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma)) - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-$$

- *Property:*

$$p(x|\theta_1) > p(x|\theta_2) \quad \Leftrightarrow \quad \log p(x|\theta_1) > \log p(x|\theta_2)$$

# FORMULAS AND DEFINITIONS

- **Log likelihood :**

$$p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta) \quad \Leftrightarrow \quad \log p(X|\theta) = \sum_{i=1}^{N} \log p(x_i|\theta)$$

- *In the Gaussian case:*

$$p(x|\theta) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$\log p(x|\theta) = -\frac{d}{2}\log(2\pi) - \frac{1}{2}\log(\det(\Sigma)) - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-$$

- *Property:*

$$p(x|\theta_1) > p(x|\theta_2) \quad \Leftrightarrow \quad \log p(x|\theta_1) > \log p(x|\theta_2)$$

# Formulas and definitions

- **Log likelihood :**

$$p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta) \quad \Leftrightarrow \quad \log p(X|\theta) = \sum_{i=1}^{N} \log p(x_i|\theta)$$

- *In the Gaussian case:*

$$p(x|\theta) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$\log p(x|\theta) = -\frac{d}{2}\log(2\pi) - \frac{1}{2}\log(\det(\Sigma)) - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-$$

- *Property:*

$$p(x|\theta_1) > p(x|\theta_2) \quad \Leftrightarrow \quad \log p(x|\theta_1) > \log p(x|\theta_2)$$

# OUTLINE

- A-priori class probability
- Gaussian modeling of classes

- A-priori class probability
- Gaussian modeling of classes

## Formulas and definitions

- **Bayes' decision rule :**

$$X \in q_k \text{ if } P(q_k|X, \Theta) \geq P(q_j|X, \Theta), \ \forall j \neq k$$

with $P(q_k|X, \Theta)$ being *a posteriori probability* (while $P(q_k|\Theta)$ is the *a priori probability*) for classes $q_k$. Note $\Theta$ represents the set of all $\theta$.

- A posteriori probability :

$$P(q_k|X, \Theta) = \frac{p(X|q_k, \Theta)P(q_k|\Theta)}{p(X|\Theta)}$$

(*Bayes' law*)

- For speech :

$$\forall k, \ P(q_k|X, \Theta) \propto p(X|q_k, \Theta)P(q_k|\Theta)$$

- or in log domain :

$$\log P(q_k|X, \Theta) \simeq \log p(X|q_k, \Theta) + \log P(q_k|\Theta)$$

# Formulas and definitions

- **Bayes' decision rule :**

$$X \in q_k \text{ if } P(q_k|X,\Theta) \geq P(q_j|X,\Theta), \ \forall j \neq k$$

with $P(q_k|X,\Theta)$ being *a posteriori probability* (while $P(q_k|\Theta)$ is the *a priori probability*) for classes $q_k$. Note $\Theta$ represents the set of all $\theta$.

- **A posteriori probability :**

$$P(q_k|X,\Theta) = \frac{p(X|q_k,\Theta)P(q_k|\Theta)}{p(X|\Theta)}$$

(*Bayes' law*)

- For speech :

$$\forall k, \ P(q_k|X,\Theta) \propto p(X|q_k,\Theta)P(q_k|\Theta)$$

- or in log domain :

$$\log P(q_k|X,\Theta) \simeq \log p(X|q_k,\Theta) + \log P(q_k|\Theta)$$

- **Bayes' decision rule :**

$$X \in q_k \text{ if } P(q_k|X,\Theta) \geq P(q_j|X,\Theta), \ \forall j \neq k$$

  with $P(q_k|X,\Theta)$ being *a posteriori probability* (while $P(q_k|\Theta)$ is the *a priori probability*) for classes $q_k$. Note $\Theta$ represents the set of all $\theta$.

- **A posteriori probability :**

$$P(q_k|X,\Theta) = \frac{p(X|q_k,\Theta)P(q_k|\Theta)}{p(X|\Theta)}$$

  (*Bayes' law*)

- *For speech :*

$$\forall k, \ P(q_k|X,\Theta) \propto p(X|q_k,\Theta)P(q_k|\Theta)$$

- or in log domain :

$$\log P(q_k|X,\Theta) \simeq \log p(X|q_k,\Theta) + \log P(q_k|\Theta)$$

# FORMULAS AND DEFINITIONS

- **Bayes' decision rule :**

$$X \in q_k \text{ if } P(q_k|X, \Theta) \geq P(q_j|X, \Theta), \ \forall j \neq k$$

  with $P(q_k|X, \Theta)$ being *a posteriori probability* (while $P(q_k|\Theta)$ is the *a priori probability*) for classes $q_k$. Note $\Theta$ represents the set of all $\theta$.

- **A posteriori probability :**

$$P(q_k|X, \Theta) = \frac{p(X|q_k, \Theta)P(q_k|\Theta)}{p(X|\Theta)}$$

  (*Bayes' law*)

- *For speech :*

$$\forall k, \ P(q_k|X, \Theta) \propto p(X|q_k, \Theta)P(q_k|\Theta)$$

- or in log domain :

$$\log P(q_k|X, \Theta) \simeq \log p(X|q_k, \Theta) + \log P(q_k|\Theta)$$

# OUTLINE

# PRELIMINARIES

All unsupervised training algorithm assume:

- a set of models $q_k$ (not necessarily Gaussian), defined by some parameters $\Theta$ (means, variances, priors,...);

- a measure of membership, telling to which extent a data point "belongs" to a model;

- the above implicitly defines global criterion of "goodness of fit" of the models to the data, e.g. :

    - in the case of a distance, the models that are globally closer from the data characterize it better;

    - in the case of a probability measure, the models bringing a better likelihood for the data explain it better.

- a "recipe" to update the model parameters in function of the membership information.

# Preliminaries

All unsupervised training algorithm assume:

- a set of models $q_k$ (not necessarily Gaussian), defined by some parameters $\Theta$ (means, variances, priors,...);
- a measure of membership, telling to which extent a data point "belongs" to a model;
- the above implicitly defines global criterion of "goodness of fit" of the models to the data, e.g. :
    - in the case of a distance, the models that are globally closer from the data characterize it better;
    - in the case of a probability measure, the models bringing a better likelihood for the data explain it better.
- a "recipe" to update the model parameters in function of the membership information.

All unsupervised training algorithm assume:

- a set of models $q_k$ (not necessarily Gaussian), defined by some parameters $\Theta$ (means, variances, priors,...);
- a measure of membership, telling to which extent a data point "belongs" to a model;
- the above implicitly defines global criterion of "goodness of fit" of the models to the data, e.g. :
  - in the case of a distance, the models that are globally closer from the data characterize it better;
  - in the case of a probability measure, the models bringing a better likelihood for the data explain it better.
- a "recipe" to update the model parameters in function of the membership information.

All unsupervised training algorithm assume:

- a set of models $q_k$ (not necessarily Gaussian), defined by some parameters $\Theta$ (means, variances, priors,...);
- a measure of membership, telling to which extent a data point "belongs" to a model;
- the above implicitly defines global criterion of "goodness of fit" of the models to the data, e.g. :
    - in the case of a distance, the models that are globally closer from the data characterize it better;
    - in the case of a probability measure, the models bringing a better likelihood for the data explain it better.
- a "recipe" to update the model parameters in function of the membership information.

All unsupervised training algorithm assume:

- a set of models $q_k$ (not necessarily Gaussian), defined by some parameters $\Theta$ (means, variances, priors,...);
- a measure of membership, telling to which extent a data point "belongs" to a model;
- the above implicitly defines global criterion of "goodness of fit" of the models to the data, e.g. :
    - in the case of a distance, the models that are globally closer from the data characterize it better;
    - in the case of a probability measure, the models bringing a better likelihood for the data explain it better.
- a "recipe" to update the model parameters in function of the membership information.

All unsupervised training algorithm assume:

- a set of models $q_k$ (not necessarily Gaussian), defined by some parameters $\Theta$ (means, variances, priors,...);
- a measure of membership, telling to which extent a data point "belongs" to a model;
- the above implicitly defines global criterion of "goodness of fit" of the models to the data, e.g. :
    - in the case of a distance, the models that are globally closer from the data characterize it better;
    - in the case of a probability measure, the models bringing a better likelihood for the data explain it better.
- a "recipe" to update the model parameters in function of the membership information.

# K-MEANS ALGORITHM

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do** :

  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:

     $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$

  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if:

     $$d_k(x_n) \leq d_l(x_n), \ \forall l \neq k$$

  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;

  4. Go to 1.

- **Until** : no further change occurs.

# K-means algorithm

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do**:
  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:

  $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$

  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if:

  $$d_k(x_n) \leq d_l(x_n), \ \ \forall l \neq k$$

  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;
  4. Go to 1.
- **Until**: no further change occurs.

# K-means algorithm

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do** :
  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:

  $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$

  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if :

  $$d_k(x_n) \leq d_l(x_n), \ \ \forall l \neq k$$

  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;
  4. Go to 1.
- **Until** : no further change occurs.

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do** :
  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:

  $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$

  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if :

  $$d_k(x_n) \leq d_l(x_n), \ \ \forall l \neq k$$

  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;
  4. Go to 1.
- **Until** : no further change occurs.

# K-means algorithm

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do** :
  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:

  $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$

  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if:

  $$d_k(x_n) \leq d_l(x_n), \ \ \forall l \neq k$$

  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;
  4. Go to 1.
- **Until** : no further change occurs.

# K-means algorithm

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do** :
  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:
  $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$
  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if :
  $$d_k(x_n) \leq d_l(x_n), \ \ \forall l \neq k$$
  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;
  4. Go to 1.
- **Until** : no further change occurs.

# K-means algorithm

- Start with $K$ initial prototypes $\mu_k$, $k = 1, \cdots, K$.
- **Do** :
  1. For each data-point $x_n$, $n = 1, \cdots, N$, compute:

  $$d_k(x_n) = (x_n - \mu_k)^T (x_n - \mu_k)$$

  2. Assign each data-point $x_n$ to its **closest** prototype $\mu_k$, i.e. assign $x_n$ to the class $q_k$ if :

  $$d_k(x_n) \leq d_l(x_n), \ \forall l \neq k$$

  3. Replace each prototype with the mean of the data-points assigned to the corresponding class;
  4. Go to 1.
- **Until** : no further change occurs.

# K-means algorithm

Global criterion :

$$J = \sum_{k=1}^{K} \sum_{x_n \in q_k} d_k(x_n)$$

# VITERBI-EM ALGORITHM FOR GAUSSIANS

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.

- **Do** :

  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.

  2. Update the parameters :

     - Update the means :

       $\mu_k^{(new)} = $ mean of the points belonging to $q_k^{(old)}$

     - Update the priors :

       $P^{(new)}(q_k) = $ fraction of the points belonging to $q_k^{(old)}$

     - Update the priors :

       $P^{(new)}(q_k) = \dfrac{\text{number of training points belonging to } q_k^{(old)}}{\text{Total number of training points}}$

  3. Go to 1.

- **Until** : no further change occurs.

# Viterbi-EM algorithm for Gaussians

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do** :
    1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
    2. Update the parameters :
        - update the means :

            $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

        - update the variances:

            $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

        - update the priors

            $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$

    3. Go to 1.
- **Until** : no further change occurs.

# Viterbi-EM algorithm for Gaussians

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do** :
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters :
     - update the means :

       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

     - update the variances :

       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

     - update the priors

       $$P(q_k^{(new)} | \Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$

  3. Go to 1.
- **Until** : no further change occurs.

# Viterbi-EM algorithm for Gaussians

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do** :
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters :
     - update the means :

       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

     - update the variances :

       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

     - update the priors :

       $$P(q_k^{(new)} | \Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$

  3. Go to 1.
- **Until** : no further change occurs.

# Viterbi-EM algorithm for Gaussians

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do**:
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters:
     - update the means:

       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

     - update the variances:

       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

     - update the priors:

       $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$

  3. Go to 1.
- **Until**: no further change occurs.

# Viterbi-EM algorithm for Gaussians

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do** :
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters :
     - update the means :

       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

     - update the variances :

       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

     - update the priors :

       $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$

  3. Go to 1.
- **Until** : no further change occurs.

# VITERBI-EM ALGORITHM FOR GAUSSIANS

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do**:
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters:
     - update the means:
       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$
     - update the variances:
       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$
     - update the priors:
       $$P(q_k^{(new)} | \Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$
  3. Go to 1.
- **Until**: no further change occurs.

# VITERBI-EM ALGORITHM FOR GAUSSIANS

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do** :
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters :
     - update the means :

       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

     - update the variances :

       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

     - update the priors :

       $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$

  3. Go to 1.
- **Until** : no further change occurs.

# VITERBI-EM ALGORITHM FOR GAUSSIANS

- Assume $K$ *initial* Gaussian models $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \cdots K$, and initial prior probabilities $P(q_k) = 1/K$.
- **Do** :
  1. Classify each data-point to its most probable cluster $q_k^{(old)}$ using Bayes' rule.
  2. Update the parameters :
     - update the means :
       $$\mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$
     - update the variances :
       $$\Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$
     - update the priors :
       $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{\text{number of training points belonging to } q_k^{(old)}}{\text{total number of training points}}$$
  3. Go to 1.
- **Until** : no further change occurs.

Global criterion :

$$\mathcal{L}(\Theta) = \sum_{k=1}^{K} \sum_{x_n \in q_k} \log p(x_n | \Theta_k)$$

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do** :
  1. **Estimation step** :

$$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  2. **Maximization step** :

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do**:
  1. **Estimation step**:

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  2. **Maximization step**:
     - update the means:

     $$\mu_k^{(new)} = \frac{\sum_{n=1}^{N} x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the variances:

     $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^{\top}}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the priors:

     $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  3. Go to 1

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do** :
  1. **Estimation step** :

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  2. **Maximization step** :
     - update the means :

     $$\mu_k^{(new)} = \frac{\sum_{n=1}^{N} x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the variances :

     $$\Sigma_j^{(new)} = \frac{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^\top}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the priors :

     $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  3. Go to 1

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do** :
    1. **Estimation step** :

    $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

    2. **Maximization step** :
        - update the means :

        $$\mu_k^{(new)} = \frac{\sum_{n=1}^{N} x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

        - update the variances :

        $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^T}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

        - update the priors :

        $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})$$

    3. Go to 1

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do** :
  ### 1 Estimation step :

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  ### 2 Maximization step :
  - update the means :

  $$\mu_k^{(new)} = \frac{\sum_{n=1}^{N} x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

  - update the variances :

  $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^T}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

  - update the priors :

  $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  ### 3 Go to 1

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do** :
  1. **Estimation step** :

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  2. **Maximization step** :
     - update the means :

     $$\mu_k^{(new)} = \frac{\sum_{n=1}^N x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the variances :

     $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^T}{\sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the priors :

     $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  3. Go to 1.

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do**:
  1. **Estimation step**:

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  2. **Maximization step**:
     - update the means:

     $$\mu_k^{(new)} = \frac{\sum_{n=1}^{N} x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the variances:

     $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^T}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the priors:

     $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  3. Go to 1.

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do**:

  **1** **Estimation step**:

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  **2** **Maximization step**:
  - update the means:

  $$\mu_k^{(new)} = \frac{\sum_{n=1}^{N} x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

  - update the variances:

  $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^T}{\sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

  - update the priors:

  $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^{N} P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  **3** Go to 1

# EM algorithm for Gaussian clustering

- Assume K *initial* models $\mathcal{N}(\mu_k, \Sigma_k)$, with $P(q_k) = 1/K$.
- **Do** :
  1. **Estimation step** :

  $$P(q_k^{(old)}|x_n, \Theta^{(old)}) = \frac{P(q_k^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_j P(q_j^{(old)}|\Theta^{(old)}) \cdot p(x_n|\mu_j^{(old)}, \Sigma_j^{(old)})}$$

  2. **Maximization step** :
     - update the means :

     $$\mu_k^{(new)} = \frac{\sum_{n=1}^N x_n P(q_k^{(old)}|x_n, \Theta^{(old)})}{\sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the variances :

     $$\Sigma_k^{(new)} = \frac{\sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})(x_n - \mu_k^{(new)})(x_n - \mu_k^{(new)})^T}{\sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})}$$

     - update the priors :

     $$P(q_k^{(new)}|\Theta^{(new)}) = \frac{1}{N} \sum_{n=1}^N P(q_k^{(old)}|x_n, \Theta^{(old)})$$

  3. Go to 1

Global criterion :

$$\mathcal{L}(\Theta) = \log \sum_{k=1}^{K} P(q_k|X, \Theta) p(X|\Theta)$$

# Outline

## ACKNOWLEDGMENTS

# Τέλος Ενότητας

# Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.

- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα στο Πανεπιστήμιο Κρήτης**» έχει χρηματοδοτήσει μόνο τη αναδιαμόρφωση του εκπαιδευτικού υλικού.

- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.

**Σημειώματα**

# Σημείωμα αδειοδότησης

# Σημείωμα Αναφοράς