

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

- 1 GAUSSIAN STATISTICS
- 2 STATISTICAL PATTERN RECOGNITION
 - Bayesian classification
- 3 UNSUPERVISED TRAINING
- 4 ACKNOWLEDGMENTS

FORMULAS AND DEFINITIONS

- 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 μ is the mean vector and Σ 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: $x \rightarrow \mathcal{N}(\mu, \Sigma)$

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

where μ is the mean vector and Σ 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: $x \rightarrow \mathcal{N}(\mu, \Sigma)$

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

where μ is the mean vector and Σ 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 :

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

$(d \times 1) (1 \times d) = d \times d$

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

FORMULAS AND DEFINITIONS

- **Likelihood:**

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

- **Joint likelihood:**

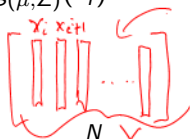
$$\rightarrow p(X|\theta) = \prod_{i=1}^N p(x_i|\theta) = \prod_{i=1}^N p(x_i|\mu, \Sigma) = \prod_{i=1}^N g_{(\mu, \Sigma)}(x_i)$$

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

$$P(x_i|\theta_1) = \alpha$$

$$P(x_i|\theta_2) = \beta$$

$$A, \alpha > \beta \Rightarrow x_i \in \theta_1$$



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:*

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

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

- *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:*

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

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

- *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:**

$$\left. \begin{array}{l} P(X | \theta_1) \\ P(X | \theta_2) \\ P(X | \theta_3) \end{array} \right\} < \epsilon$$

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

- In the Gaussian case:

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

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

- Property:

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

OUTLINE

1 GAUSSIAN STATISTICS

2 STATISTICAL PATTERN RECOGNITION

- Bayesian classification

3 UNSUPERVISED TRAINING

4 ACKNOWLEDGMENTS

DISCUSSION

- 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 Θ represents the set of all θ .

- **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 Θ represents the set of all θ .

- **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 Θ represents the set of all θ .

- **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 Θ represents the set of all θ .

- **A posteriori probability :**

$$P(q_k|X, \Theta) = \frac{p(X|q_k, \Theta)P(q_k|\Theta)}{p(X|\Theta)} = \sum_{k=1}^N P(q_k|\Theta)$$

(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

- 1 GAUSSIAN STATISTICS
- 2 STATISTICAL PATTERN RECOGNITION
 - Bayesian classification
- 3 UNSUPERVISED TRAINING
- 4 ACKNOWLEDGMENTS

All unsupervised training algorithm assume:

- a set of models q_k (not necessarily Gaussian), defined by some parameters Θ (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 Θ (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 Θ (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 Θ (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 Θ (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 Θ (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 μ_k , $k = 1, \dots, K$.

- **Do**:

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

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

- Assign each data-point x_n to its closest prototype μ_k , i.e. assign x_n to the class q_k if:

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

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

- Go to 1.

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

K-MEANS ALGORITHM

- Start with K initial prototypes μ_k , $k = 1, \dots, K$.

- **Do**:

- 1 For each data-point x_n , $n = 1, \dots, N$, 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 μ_k , i.e. assign x_n to the class q_k if:

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

- 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 μ_k , $k = 1, \dots, K$.

- **Do**:

- 1 For each data-point x_n , $n = 1, \dots, N$, 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 μ_k , i.e. assign x_n to the class q_k if:

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

- 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 μ_k , $k = 1, \dots, K$.

- **Do** :

- 1 For each data-point x_n , $n = 1, \dots, N$, 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 μ_k , i.e. assign x_n to the class q_k if:

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

- 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 μ_k , $k = 1, \dots, K$.

- **Do** :

- 1 For each data-point x_n , $n = 1, \dots, N$, 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 μ_k , i.e. assign x_n to the class q_k if:

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

- 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 μ_k , $k = 1, \dots, K$.

- **Do** :

- 1 For each data-point x_n , $n = 1, \dots, N$, 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 μ_k , i.e. assign x_n to the class q_k if:

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

- 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 $\mu_k, k = 1, \dots, K$.

- **Do**:

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

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

- 2 Assign each data-point x_n to its **closest** prototype μ_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 $\mathcal{N}(\mu_k, \Sigma_k)$, $k = 1 \dots K$, and initial prior probabilities $P(q_k) = 1/K$.

- **Do**:

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

- Update the parameters:

• Update the means:

$$\mu_k^{(new)} = \frac{1}{N_k} \sum_{i: q_i = k} x_i$$

• Update the covariances:

$$\Sigma_k^{(new)} = \frac{1}{N_k} \sum_{i: q_i = k} (x_i - \mu_k^{(new)})(x_i - \mu_k^{(new)})^T$$

• Update the priors:

$$P(q_k) = \frac{N_k}{N} \quad \text{number of training points belonging to } q_k$$

$$N = \sum_{k=1}^K N_k \quad \text{total number of training points}$$

- Go to 1.

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

VITERBI-EM ALGORITHM FOR GAUSSIANS

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

- Do:**

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

- Update the parameters:

- update the means:

$$\textcircled{1} \mu_k^{(new)} = \text{mean of the points belonging to } q_k^{(old)}$$

- update the variances:

$$\textcircled{2} \Sigma_k^{(new)} = \text{variance of the points belonging to } q_k^{(old)}$$

- update the priors:

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

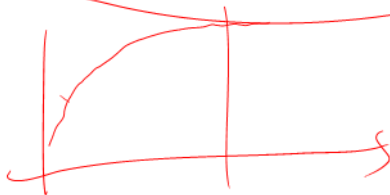
- Go to 1.

- Until:** no further change occurs.

VITERBI-EM ALGORITHM FOR GAUSSIANS

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:
 - Estimation step:

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

- Maximization step:

we replace the values:

$$\mu_k \leftarrow \frac{\sum_n P(q_k | x_n) x_n}{\sum_n P(q_k | x_n)}$$

and also the covariances:

$$\Sigma_k \leftarrow \frac{\sum_n P(q_k | x_n) (x_n - \mu_k)(x_n - \mu_k)^T}{\sum_n P(q_k | x_n)}$$

until:

$$|L(\Theta^{(old)}) - L(\Theta^{(new)})| < \epsilon$$

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:

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

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:

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

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:

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

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:

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

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:

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

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:

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

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:

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

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_k

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

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

Global criterion :

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

OUTLINE

- 1 GAUSSIAN STATISTICS
- 2 STATISTICAL PATTERN RECOGNITION
 - Bayesian classification
- 3 UNSUPERVISED TRAINING
- 4 ACKNOWLEDGMENTS

ACKNOWLEDGMENTS

Most of this material is from Hervé Bourlard's course on Speech processing and speech recognition

