

Fondements de l'apprentissage machine

Automne 2014

Roland Memisevic

Leçon 6

Roland Memisevic

Fondements de l'apprentissage machine

Apprentissage non-supervisé

- ▶ Etant donné un ensemble d'observations \mathbf{x} , la tâche d'apprentissage non-supervisé est de re-représenter les données, ce qui mène à trouver et comprendre des structures cachées.
- ▶ Exemples : **Clustering** - trouvez les groupes, **réduction de la dimensionnalité** - trouvez des représentations en basse dimension.
- ▶ L'apprentissage non-supervisé peut être formalisé comme l'apprentissage supervisé, où il manque les entrées ou les sorties.
- ▶ L'apprentissage non-supervisé suppose que le processus qui a généré les données contient des paramètres internes qu'on n'a pas observé, mais qui influencent les données tout de même.

Roland Memisevic

Fondements de l'apprentissage machine

Plan

- ▶ Apprentissage non-supervisé, variables latentes
- ▶ K -means clustering
- ▶ Modèles de mélange gaussien ("Gaussian mixture models")
- ▶ L'algorithme EM

Roland Memisevic

Fondements de l'apprentissage machine

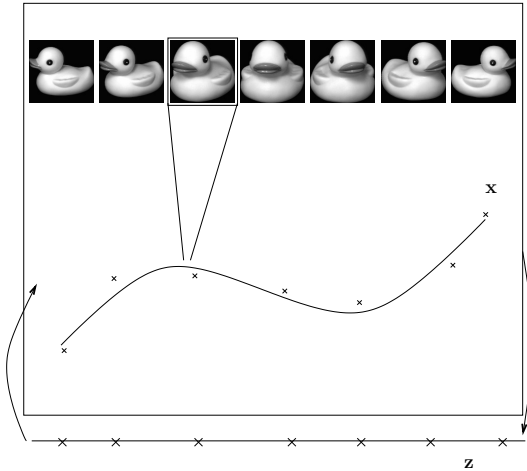
Des variables latentes

- ▶ Les variables non-observées s'appellent les **variables cachées** ou les **variables latentes**.
- ▶ Formellement, on postule qu'il y a une variable \mathbf{z}_n associée à chaque exemple d'entraînement \mathbf{x}_n .
- ▶ Le but de l'apprentissage est d'inférer les valeurs des \mathbf{z}_n en ne sachant que les données \mathbf{x}_n .

Roland Memisevic

Fondements de l'apprentissage machine

Des variables latentes

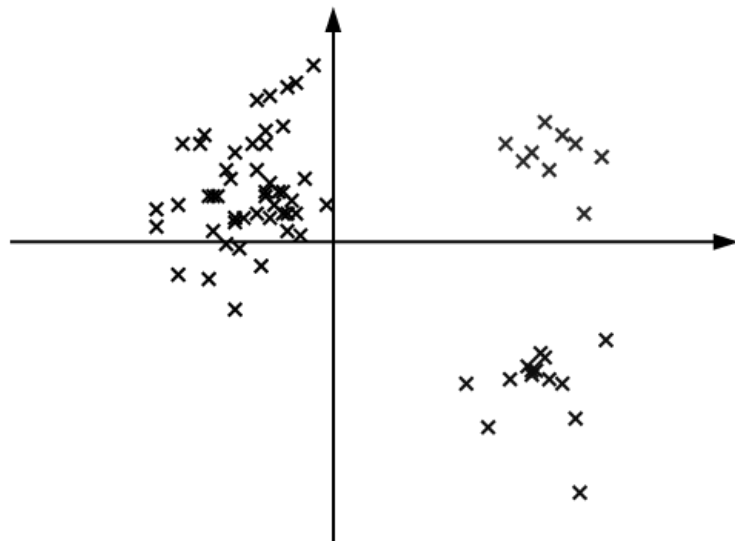


- Apprendre les “causes cachées” qui ont généré les données peut aider à compresser, comprendre ou pré-traiter les données.

Roland Memisevic

Fondements de l'apprentissage machine

Exemple



Roland Memisevic

Fondements de l'apprentissage machine

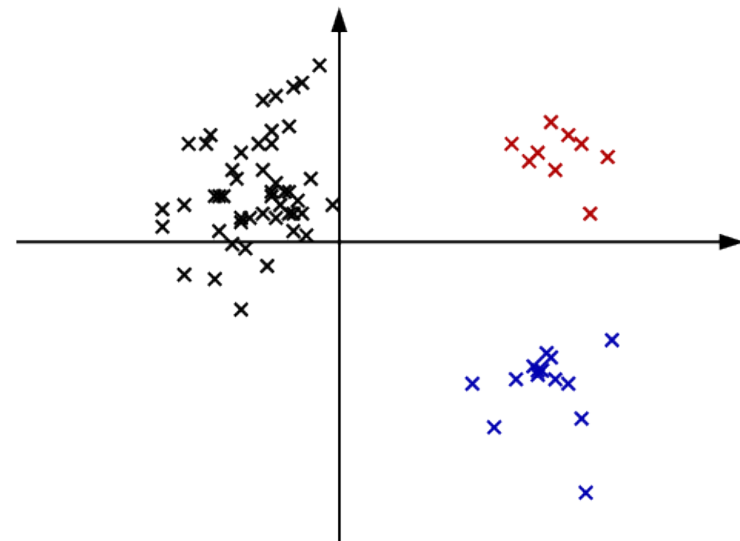
K -means clustering

- Nous considérons les variables latentes discrètes dans cette leçon.
- Pour modéliser des variables latentes discrètes, nous utilisons l'encodage 1-de- K pour \mathbf{z}_n .
- L'inférence d'une variable discrète cachée est connue comme **clustering** parce qu'elle mène à associer chaque point avec un groupe parmi K , où \mathbf{z}_n représente le groupe assigné.
- La méthode la plus connue pour le clustering est **K -means**.

Roland Memisevic

Fondements de l'apprentissage machine

Exemple



Roland Memisevic

Fondements de l'apprentissage machine

K-means clustering

- ▶ Notation : Nous empilons des \mathbf{z}_n en codage one-hot dans une matrice \mathbf{Z} .
- ▶ Nous supposons qu'il existe K prototypes μ_1, \dots, μ_K qui représentent les K groupes. La dimensionnalité des prototypes est la même que celle des observations \mathbf{x} .
- ▶ Supposez (pour le moment) que pour chaque point \mathbf{x}_n , nous connaissons le groupe \mathbf{z}_n y étant associé.
- ▶ La fonction de perte que K -means clustering essaie de minimiser est la distance moyenne entre les points \mathbf{x} et leur représentant μ :

$$J = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

Roland Memisevic

Fondements de l'apprentissage machine

Trouver les \mathbf{z}_n optimaux

- ▶ Notez que sachant les μ_k , nous pouvons optimiser tous les \mathbf{z}_n indépendamment, car l'objectif est la somme sur n .
- ▶ L'erreur au carré est minimale si $z_{nk} = 1$ pour le plus proche μ_k .
- ▶ La solution de l'optimisation est donc :

$$z_{nk} = \begin{cases} 1 & \text{si } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{sinon.} \end{cases}$$

pour chaque z_{nk} .

Roland Memisevic

Fondements de l'apprentissage machine

K-means clustering

- ▶ L'apprentissage correspond à trouver les prototypes μ_k ainsi que les associations \mathbf{z}_n qui minimisent J .
- ▶ C'est un problème d'optimisation difficile, car les \mathbf{z}_n sont discrets et les μ_k continus.
- ▶ L'apprentissage peut être simplifié si on découple l'optimisation par rapport aux \mathbf{z}_n de l'optimisation par rapport aux μ_k .
- ▶ Cela mène à une *descente par blocs de coordonnées* (block coordinate descent), qui est un cas spécial d'une approche d'optimisation dans le domaine d'apprentissage non-supervisé connue sous le nom d'algorithme EM (EM-algorithm).

Roland Memisevic

Fondements de l'apprentissage machine

Trouver les μ_k optimaux

- ▶ Sachant les \mathbf{z}_n , la fonction de perte J est une fonction quadratique de μ_k . Elle peut donc être minimisée en mettant la dérivée à zéro :

$$2 \sum_{n=1}^N z_{nk} (\mathbf{x}_n - \mu_k) = 0$$

- ▶ La solution pour chaque μ_k est :

$$\mu_k = \frac{\sum_n z_{nk} \mathbf{x}_n}{\sum_n z_{nk}}$$

- ▶ Donc, le μ_k optimal est la moyenne de tous les points associés au groupe k .

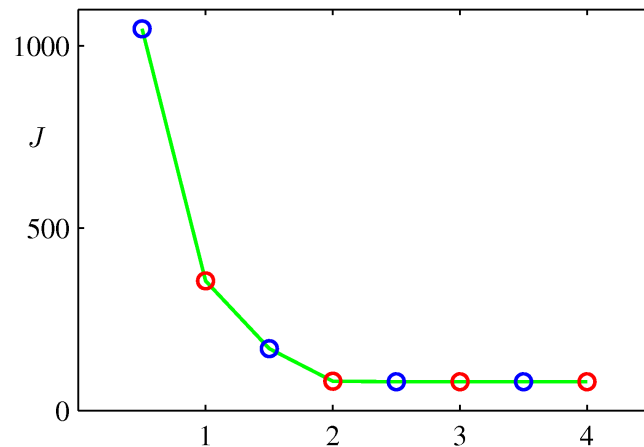
Roland Memisevic

Fondements de l'apprentissage machine

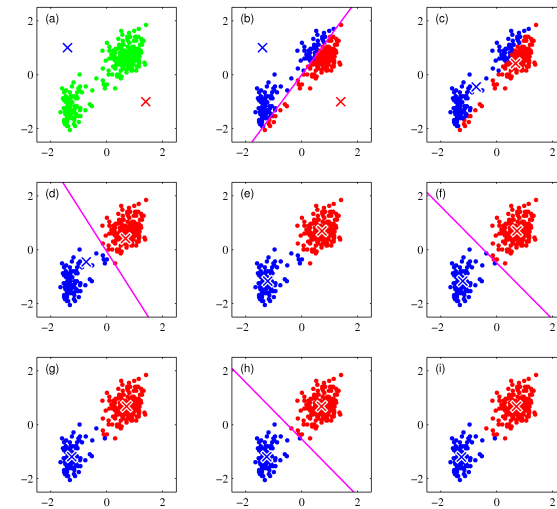
Convergence, minima local

- L'apprentissage consiste à itérer l'inférence des \mathbf{z}_n et d'adapter les paramètres μ_k jusqu'à convergence.
- Cette procédure est garantie de converger parce que J est positif et à chaque itération, soit J diminue, soit J ne change pas.
- Mais la convergence peut être à un minimum local.
- Une manière d'échapper à de mauvais minima est de faire l'optimisation à plusieurs reprises, en utilisant à chaque fois une initialisation différente des paramètres μ_k .

La valeur de J pendant l'apprentissage



Exemple



Compression

- Étant donné un modèle entraîné, nous pouvons inférer les associations de nouveaux points de test \mathbf{x} aux groupes comme ceci :

$$z_k = \begin{cases} 1 & \text{si } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{sinon.} \end{cases}$$

- Comme \mathbf{z} représente les vecteurs en haute dimension \mathbf{x} en utilisant un de K entiers, K -means est une manière de faire de la *compression* (avec perte).
- Dans ce contexte, l'ensemble des K prototypes μ_k est aussi appelé *codebook*, et K -means *vector quantization*.

Une application simple

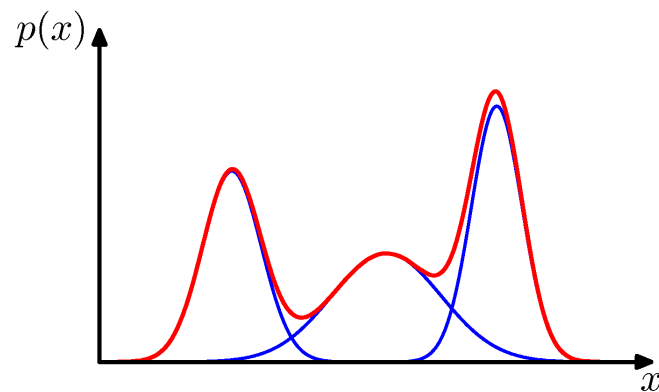
- Remplacer la valeur RGB (un vecteur 3D) de chaque pixel par un de K prototypes :



Roland Memisevic

Fondements de l'apprentissage machine

Exemple



Roland Memisevic

Fondements de l'apprentissage machine

Modèles de mélange gaussien ("Mixture of Gaussians")

- K -means est lié au

Modèle de mélange gaussien

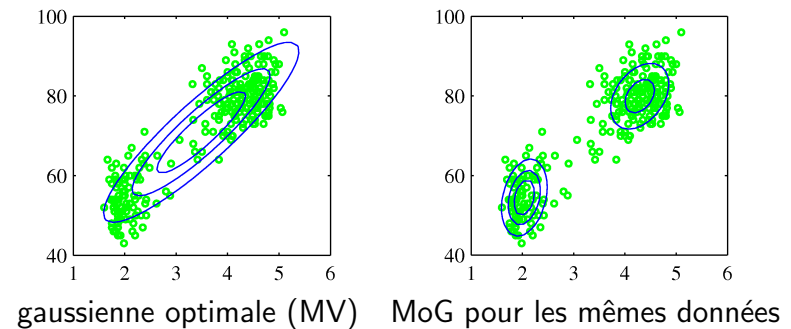
$$p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ sont des paramètres.
- π_k s'appelle *mixing proportion*, et chaque distribution gaussienne s'appelle un *mixture component*.
- Le modèle est simplement une somme pondérée de gaussiennes.
- Il est beaucoup plus puissant qu'une seule gaussienne parce qu'il peut modéliser des distributions *multimodales*.

Roland Memisevic

Fondements de l'apprentissage machine

Exemple



Roland Memisevic

Fondements de l'apprentissage machine

Modèles de mélange gaussien

$$p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ Pour que $p(\mathbf{x})$ soit une distribution probabiliste, il faut que $\sum_k \pi_k = 1$ et que $\pi_k > 0 \quad \forall k$
- ▶ Donc, il faut que les π_k représentent aussi des probabilités.
- ▶ Cela suggère d'introduire des variables latentes discrètes \mathbf{z} et de redéfinir le modèle comme suit :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z})$$

où $p(\mathbf{x} | \mathbf{z})$ est une gaussienne (conditionnelle).

Roland Memisevic

Fondements de l'apprentissage machine

L'algorithme EM

- ▶ On pourrait utiliser la descente de gradient pour l'optimisation, mais il existe une procédure plus courante semblable à celle utilisée pour K -means, l'algorithme EM.
- ▶ Pour un ensemble d'entraînement, $\{\mathbf{x}_n\}$, nous avons :

$$\begin{aligned} L := \sum_n \log p(\mathbf{x}_n) &= \sum_n \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) \\ &= \sum_n \log \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \frac{p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n)}{q(\mathbf{z}_n)} \\ &\geq \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log \frac{p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n)}{q(\mathbf{z}_n)} \\ &:= \mathcal{L} \end{aligned}$$

où nous avons utilisé l'**inégalité de Jensen** :

$$\log \sum_i a_i b_i \geq \sum_i a_i \log b_i \quad \text{if } \forall i : a_i > 0 \text{ et } \sum_i a_i = 1$$

Roland Memisevic

Fondements de l'apprentissage machine

Modèles de mélange gaussien

- ▶ Cela nous permet de penser au modèle comme le processus génératif suivant pour générer des observations : Nous tirons un composant d'une distribution discrète, puis nous tirons une observation d'une gaussienne dont les paramètres dépendent de ce premier choix.
- ▶ Pour inférer le composant sachant l'observation, \mathbf{x}_n , nous pouvons utiliser la règle de Bayes :

$$p(\mathbf{z}_n | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n)}{\sum_{\mathbf{z}_n} p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n)}$$

ce qui représente la vraisemblance pour le point d'être tiré de chaque gaussienne.

- ▶ $p(z_{nk} = 1 | \mathbf{x}_n)$ s'appelle *responsabilité* du composant k , que nous abrévions $\gamma(z_{kn})$.
- ▶ C'est comme un K -means "doux".

Roland Memisevic

Fondements de l'apprentissage machine

L'algorithme EM

- ▶ À la place d'optimiser L , nous allons optimiser une borne inférieure, \mathcal{L} , par rapport aux paramètres *ainsi qu'* aux variables auxiliaires $q(\mathbf{z})$.
- ▶ Nous pouvons écrire :

$$\mathcal{L} = \sum_{nk} q_{nk} \log p(\mathbf{x}_n | \mathbf{z}_n = k) p(\mathbf{z}_n = k) - \sum_{nk} q_{nk} \log q_{nk}$$

où $q_{nk} := q(\mathbf{z}_n = k)$

- ▶ Notez que le premier terme de \mathcal{L} est une espérance de $p(\mathbf{x}_n, \mathbf{z}_n)$ par rapport à $q(\mathbf{z}_n)$. Il est nommé "**expected complete log-likelihood**" en anglais.
- ▶ C'est le seul terme qui dépend des paramètres du modèle.

Roland Memisevic

Fondements de l'apprentissage machine

Optimisation par rapport aux paramètres

- Pour optimiser la borne par rapport aux paramètres, nous mettons la dérivée à zéro :

$$\frac{\partial \mathcal{L}}{\partial \mu_k} = \sum_n q_{nk} \Sigma_k (\mathbf{x}_n - \mu_k) = 0$$

$$\Leftrightarrow \mu_k = \frac{\sum_n q_{nk} \mathbf{x}_n}{\sum_n q_{nk}}$$

- Similairement on peut trouver :

$$\Sigma_k = \frac{\sum_n q_{nk} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T}{\sum_n q_{nk}}$$

et

$$\pi_k = p(\mathbf{z}_k) = \frac{\sum_n q_{nk}}{N}$$

Roland Memisevic

Fondements de l'apprentissage machine

L'algorithme EM

- La

divergence Kullback-Leibler (divergence KL)

$$\text{KL} (p_1(\mathbf{z}) \parallel p_2(\mathbf{z})) = \sum_{\mathbf{z}} p_1(\mathbf{z}) \log \frac{p_1(\mathbf{z})}{p_2(\mathbf{z})}$$

mesure la similarité entre deux distributions probabilistes

p_1, p_2 .

- La divergence KL est toujours non-négative et est zéro seulement si $p_1 = p_2$!
- Pour cette raison, $\mathcal{L} = L$ si nous mettons $q(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_n)$!
- En d'autres termes, le choix optimal pour $q(\mathbf{z}_n)$ est $q(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_n)$, ce qui rendra la borne \mathcal{L} pour L exacte !

Roland Memisevic

Fondements de l'apprentissage machine

Optimisation par rapport aux variables auxiliaires

- Pour optimiser \mathcal{L} par rapport aux variables auxiliaires, nous récrivons \mathcal{L} comme suit :

$$\begin{aligned} \mathcal{L} &= \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log \frac{p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n)}{q(\mathbf{z}_n)} \\ &= \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log \frac{p(\mathbf{z}_n | \mathbf{x}_n) p(\mathbf{x}_n)}{q(\mathbf{z}_n)} \\ &= \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log \frac{p(\mathbf{z}_n | \mathbf{x}_n)}{q(\mathbf{z}_n)} + \sum_n \sum_{\mathbf{z}_n} q(\mathbf{z}_n) \log p(\mathbf{x}_n) \\ &= - \sum_n \text{KL} (q(\mathbf{z}_n) \parallel p(\mathbf{z}_n | \mathbf{x}_n)) + L \end{aligned}$$

Roland Memisevic

Fondements de l'apprentissage machine

L'algorithme EM

L'algorithme EM

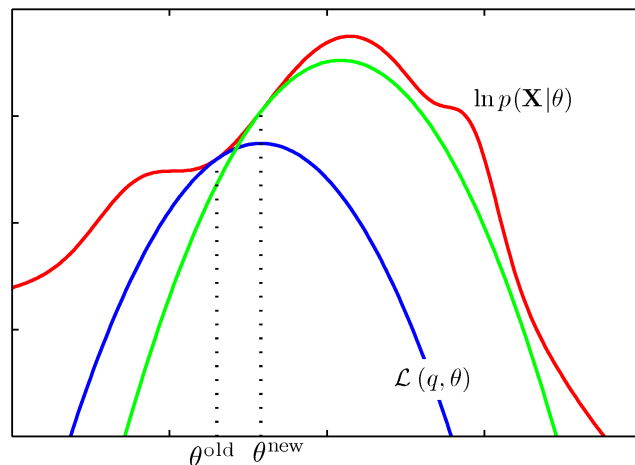
1. E-step : Évaluez les distributions postérieures $p(\mathbf{z}_n | \mathbf{x}_n)$
2. M-step : Optimisez \mathcal{L} par rapport aux paramètres, gardant $q(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_n)$.

- Le E-step calcule la "expected complete log-likelihood". Il correspond à évaluer les responsabilités $\gamma(z_{nk})$ pour chaque point.
- Le M-step maximise la "expected complete log-likelihood". Dans un modèle de mélange gaussien, cela correspond à mettre les paramètres à des sommes pondérées par les responsabilités.

Roland Memisevic

Fondements de l'apprentissage machine

EM comme optimisation d'une séquence de bornes



Roland Memisevic

Fondements de l'apprentissage machine

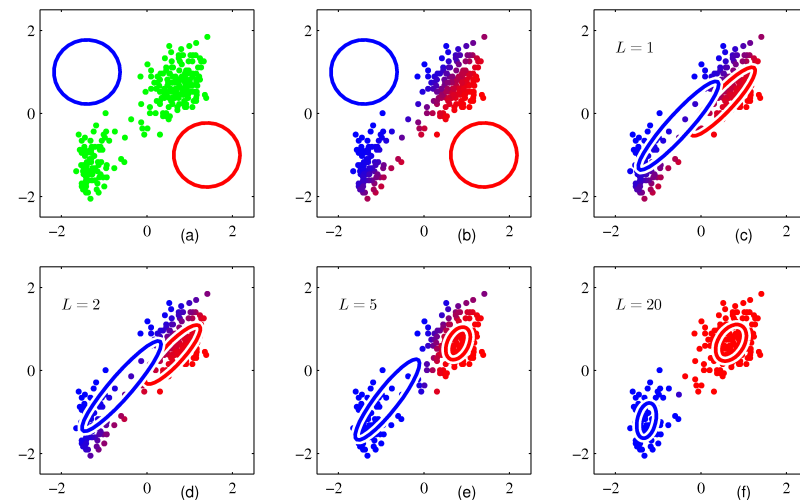
Des commentaires

- ▶ L'algorithme EM peut être appliqué à beaucoup de modèles à variables latentes (pas seulement les mélanges de gaussiennes).
- ▶ Le E-step et le M-step doivent être dérivés individuellement pour chaque modèle, mais la perspective de la borne inférieure \mathcal{L} de la log-vraisemblance est toujours la même.
- ▶ Un des premiers modèles optimisés par l'algorithme EM était le Hidden Markov Model (HMM).
- ▶ Il existe des modèles pour lesquels calculer $p(\mathbf{z}|\mathbf{x})$ n'est pas traitable ("intractable"). Dans ces cas, on peut utiliser une variation de EM où on améliore la KL-divergence dans le E-step à la place de trouver la distribution a posteriori exacte.

Roland Memisevic

Fondements de l'apprentissage machine

Exemple



Roland Memisevic

Fondements de l'apprentissage machine

Modèle de mélange gaussien et K -means

- ▶ Le modèle de mélange gaussien est comme K -means, où nous utilisons une association "douce" aux groupes.
- ▶ On peut dériver K -means formellement comme limite d'un modèle de mélange gaussien où les variances vont à zéro.

Roland Memisevic

Fondements de l'apprentissage machine