

Mixture Models and Order Selection

Lesson 6 : Lab Session

Advanced Machine Learning, CentraleSupélec

Teacher's Assistant: Omar CHEHAB

Professors : Emilie CHOUZENOUX, Frederic PASCAL



General Information

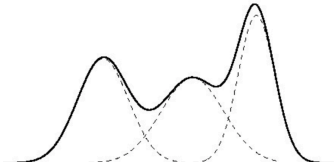
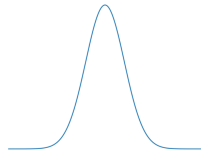
- **Assignment** : alone or in pairs, you will code the algorithms you learnt in ‘scikit-learn formalism’, and apply them to images and text.
- **Due** : the 5 lab assignments for lessons 3-7 are due a week from when they are given, at aml.centralesupelec.2020@gmail.com
- **Grading** : each assignment is worth 4 points — your 4 best labs out of the 5 will be retained and will count for half of your final grade.
- **Questions** : questions or feedback are welcome after class or by email at l-emir-omar.chehab@inria.fr

Mixture Model

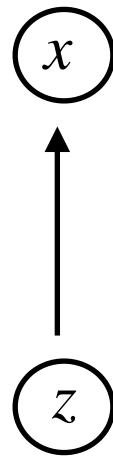
Builds a rich (e.g. multimodal) distribution
...from (a weighted sum of) simpler ones.

$$p_{\theta}(x) = \sum_{z=1}^K p_{\theta}(x | z = k) p_{\theta}(z = k)$$

complex model simple model k its weight

Gaussian Mixture Model



$$X | z = k \sim \mathcal{N}(\mu_k, \Sigma_k)$$

each cluster is gaussian
(mean and variance)

$$= \theta$$

$$Z \sim \text{Multinom}(\pi_1, \dots, \pi_K)$$

prior on K clusters

Inference

Maximum-Likelihood

X and Z are observed

data-point and
cluster assignment $(x_i, z_i)_{i=1, \dots, n}$
 $\in \mathbb{R}^p \in \mathbb{R}$

‘easy’ $\max_{p_\theta} \log p_\theta(x, z)$

direct

MLE

$$\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n 1_{z_i=k}$$

$$\hat{\mu}_k = \frac{1}{n \hat{\pi}_k} \sum_{i=1}^n 1_{z_i=k} x_i$$

$$\hat{\sigma}_k^2 = \frac{1}{n \hat{\pi}_k} \sum_{i=1}^n 1_{z_i=k} (x_i - \hat{\mu}_k) (x_i - \hat{\mu}_k)^T$$

X only is observed (Z is *hidden*)

$(x_i)_{i=1, \dots, n}$ data-point without
cluster assignment
 $\in \mathbb{R}^p$

$\max_{p_\theta} \log p_\theta(x)$ ‘hard’

indirect (via E.M.)

E-Step

update
(soft) assignment

$$p_{\theta_t}(Z = k | X = x_i) = \frac{\hat{\pi}_k^t \times p_{\theta_t}(x_i | z = k)}{\sum_{l=1}^K \hat{\pi}_l^t \times p_{\theta_t}(x_i | z = l)}$$

like in K-means but ‘soft’

M-Step

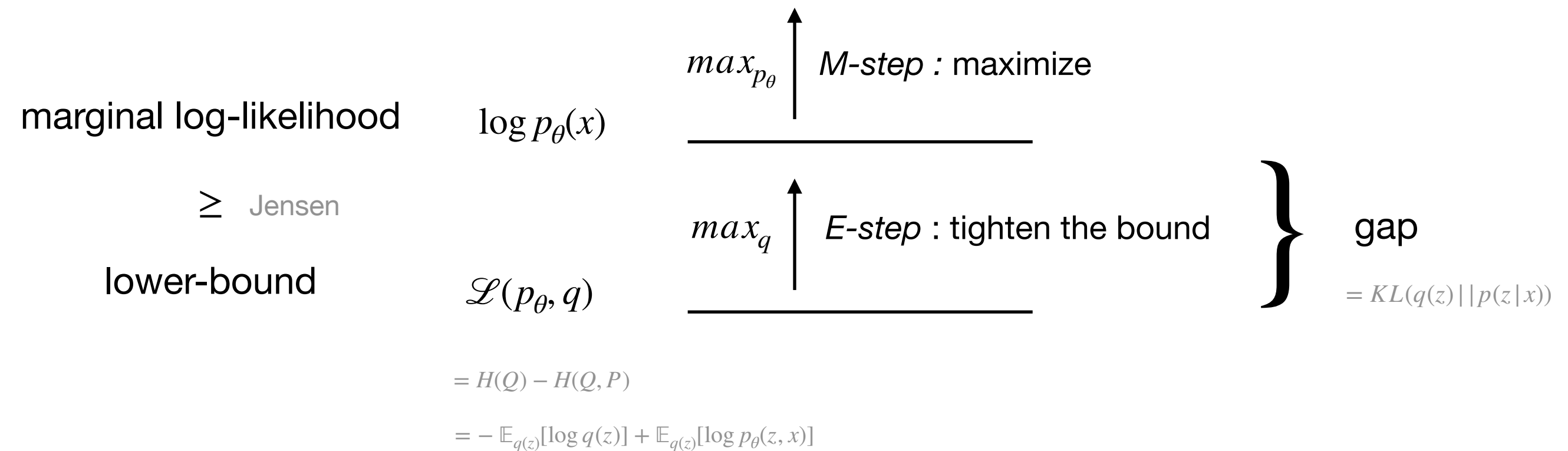
update cluster

$$\hat{\pi}_k^{t+1} = \frac{1}{n} \sum_{i=1}^n P_{\theta_t}(Z = k | X = x_i)$$

$$\hat{\mu}_k^{t+1} = \frac{1}{n \hat{\pi}_k^{t+1}} \sum_{i=1}^n P_{\theta_t}(Z = k | X = x_i) x_i$$

$$\hat{\Sigma}_k^{t+1} = \frac{1}{n \hat{\pi}_k^{t+1}} \sum_{i=1}^n P_{\theta_t}(Z = k | X = x_i) (x_i - \hat{\mu}_k^{t+1}) (x_i - \hat{\mu}_k^{t+1})^T$$

The idea behind E.M.



The marginal log-likelihood is difficult to maximize directly
...so do so via a tractable lower-bound!

E-Step $\max_q \mathcal{L}(p_\theta, q)$ or equivalently $\min_q gap(p_\theta, q) \Rightarrow q(z) = p(z|x)$

M-Step $\max_{p_\theta} \mathcal{L}(p_\theta, q) \iff \max_{p_\theta} \log \mathbb{E}_{z \sim q(z)}[p(x, z)]$

- \nearrow stochastic EM:
replace Expectation over q with a sample
- \rightarrow N-EM:
replace Expectation over q with N samples
- \searrow EM:
use Expectation

Model order

I can model my data $\mathcal{D} = (x_i)_{i \in [1, N]}$ with any model(s) from $(\mathcal{M}_1, \dots, \mathcal{M}_K)$: which is best?

Choose a model:

$$\max_j \begin{matrix} \text{model quality} \\ AIC(j) \\ BIC(j) \end{matrix} = \begin{matrix} \text{capacity to overfit} \\ \mathcal{L}(\hat{\theta}_{MLE}) \\ \mathcal{L}(\hat{\theta}_{MLE}) \end{matrix} - \begin{matrix} \text{capacity to underfit} \\ M_j \\ \frac{1}{2}M_jN \end{matrix}$$

M_j =
nb of unknown params
for model j

harsher

rewards complexity penalizes complexity

tradeoff

Reduce a model:
(specific here to Gaussian Mixture Model)

after training, merge two clusters (k, k')
by the weighted sum of their statistics:

$$\mu_{new} = \frac{n_k \mu_k + n_{k'} \mu_{k'}}{n_k + n_{k'}} \quad \Sigma_{new} = \frac{n_k \Sigma_k + n_{k'} \Sigma_{k'}}{n_k + n_{k'}}$$

$$\pi_{new} = \pi_k + \pi_{k'} \quad \text{with } n_k \text{ points in cluster } k$$

correlation criterion: if the profiles of their assigned points are 'similar enough' $PearsonR(p_{\theta^*}(z = k_1 | x = \dots), p_{\theta^*}(z = k_1 | x = \dots)) > 1 - \epsilon$

distance criterion: or if their centroids are 'close enough' $\|\mu_k - \mu_{k'}\| \leq \epsilon$

parameter criterion: or if it reduces the intra-class variance $d_{new} < d_k + d_{k'} - \epsilon$ where d_k is the variance of cluster k

Assignment: plan

1. Gaussian Mixture Model : train using E.M. algorithm (your own code)
2. Choose number of clusters : model order *selection* (AIC, BIC)
or model order *reduction* (3 criteria) (your own code)
3. Application : vision (MNIST digits dataset) (your own code)