

# Machine Learning for Data Science (CS4786)

## Lecture 16

EM Algorithm

Mar 24, 2015

Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2015sp/>

# EXAMPLES

- Gaussian Mixture Model

- Each  $\theta$  consists of mixture distribution  $\pi = (\pi_1, \dots, \pi_K)$ , means  $\mu_1, \dots, \mu_K \in \mathbb{R}^d$  and covariance matrices  $\Sigma_1, \dots, \Sigma_K$
- At time  $t$  we generate a new tree as follows:

$$c_t \sim \pi, \quad x_t \sim N(\mu_{c_t}, \Sigma_{c_t})$$

# PROBABILISTIC MODELS

More generally:

- $\Theta$  consists of set of possible parameters
- We have a distribution  $P_\theta$  over the data induced by each  $\theta \in \Theta$
- Data is generated by one of the  $\theta \in \Theta$
- Learning: Estimate value or distribution for  $\theta^* \in \Theta$  given data

# MAXIMUM LIKELIHOOD PRINCIPAL

Pick  $\theta \in \Theta$  that maximizes probability of observation

$$\theta_{MLE} = \operatorname{argmax}_{\theta \in \Theta} \underbrace{\log P_{\theta}(x_1, \dots, x_n)}_{\text{Likelihood}}$$

# MAXIMUM A POSTERIORI

Pick  $\theta \in \Theta$  that is most likely given data

$$\begin{aligned}\theta_{MAP} &= \operatorname{argmax}_{\theta \in \Theta} \log P(\theta | x_1, \dots, x_n) \\ &= \operatorname{argmax}_{\theta \in \Theta} \underbrace{P(x_1, \dots, x_n | \theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}\end{aligned}$$

# LATENT VARIABLES

- We only observe  $x_1, \dots, x_n$ , cluster assignments  $c_1, \dots, c_n$  are not observed
- Finding  $\theta \in \Theta$  (even for 1-d GMM) that directly maximizes Likelihood or A Posteriori given  $x_1, \dots, x_n$  is hard!
- Given latent variables  $c_1, \dots, c_n$ , the problem of maximizing likelihood (or a posteriori) became easy

Can we use latent variables to device an algorithm?

# EXPECTATION MAXIMIZATION ALGORITHM

- For demonstration we shall consider the problem of finding MLE (MAP version is very similar)

# EXPECTATION MAXIMIZATION ALGORITHM

- For demonstration we shall consider the problem of finding MLE (MAP version is very similar)
- Initialize  $\theta^{(0)}$  arbitrarily, repeat until convergence:

(E step) For every  $t$ , define distribution  $Q_t$  over the latent variable  $c_t$  as:

$$Q_t^{(i)}(c_t) = P(c_t|x_t, \theta^{(i-1)})$$

(M step)

$$\theta^{(i)} = \operatorname{argmax}_{\theta \in \Theta} \sum_{t=1}^n \sum_{c_t} Q_t^{(i)}(c_t) \log P(x_t, c_t|\theta)$$



## EXAMPLE: EM FOR GMM

- E step: For every  $k \in [K]$ ,

$$\begin{aligned} Q_t^{(i)}(c_t = k) &= P(c_t = k | x_t, \theta^{(i-1)}) = P(x_t | c_t = k, \theta^{(i-1)}) \times P(c_t = k | \theta^{(i-1)}) \\ &\propto \underbrace{\phi(x_t; \mu_k^{(i-1)}, \Sigma_k^{(i-1)})}_{\text{gaussian p.d.f.}} \times \pi_k^{(i-1)} \end{aligned}$$

# EXAMPLE: EM FOR GMM

- E step: For every  $k \in [K]$ ,

$$\begin{aligned} Q_t^{(i)}(c_t = k) &= P(c_t = k | x_t, \theta^{(i-1)}) = P(x_t | c_t = k, \theta^{(i-1)}) \times P(c_t = k | \theta^{(i-1)}) \\ &\propto \underbrace{\phi(x_t; \mu_k^{(i-1)}, \Sigma_k^{(i-1)})}_{\text{gaussian p.d.f.}} \times \pi_k^{(i-1)} \end{aligned}$$

- M step: Given  $Q_1, \dots, Q_n$ , we need to find

$$\begin{aligned} \theta^{(i)} &= \operatorname{argmax}_{\theta \in \Theta} \sum_{t=1}^n \sum_{k=1}^K Q_t^{(i)}(k) \log P(x_t, c_t = k | \theta) \\ &= \operatorname{argmax}_{\theta} \sum_{t=1}^n \sum_{k=1}^K Q_t^{(i)}(k) (\log P(x_t | c_t = k, \theta) + \log P(c_t = k | \theta)) \\ &= \operatorname{argmax}_{\pi, \mu_{1,\dots,K}, \Sigma_{1,\dots,K}} \sum_{t=1}^n \sum_{c_t=1}^K Q_t^{(i)}(k) (\log \phi(x_t; \mu_k, \Sigma_k) + \log \pi_k) \end{aligned}$$

## EXAMPLE: EM FOR GMM

For every  $k \in [K]$ , the maximization step yields,

$$\mu_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k) x_t}{\sum_{t=1}^n Q_t(k)}, \quad \Sigma_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k) (x_t - \mu_k^{(i)}) (x_t - \mu_k^{(i)})^\top}{\sum_{t=1}^n Q_t(k)}$$

$$\pi_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k)}{n}$$

## EXAMPLE: EM FOR GMM

For every  $k \in [K]$ , the maximization step yields,

$$\mu_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k) x_t}{\sum_{t=1}^n Q_t(k)}, \quad \Sigma_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k) (x_t - \mu_k^{(i)}) (x_t - \mu_k^{(i)})^\top}{\sum_{t=1}^n Q_t(k)}$$

$$\pi_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k)}{n}$$

# HARD AND SOFT ASSIGNMENTS

- Hard GMM: make E-step a hard assignments with  $Q_t^{(i)} = e_{k_t^*}$  where

$$k_t^* = \operatorname{argmax}_{k \in [K]} \phi \left( x_t; \mu_k^{(i-1)}, \Sigma_k^{(i-1)} \right) \times \pi_k^{(i-1)}$$

- Soft k-means:

$$Q_t^{(i)}(k) \propto \exp(-\|x_t - \mu_k^{(i-1)}\|_2^2 / \sigma^2) \text{ and } \mu_k^{(i)} = \frac{\sum_{t=1}^n Q_t^{(i)}(k) x_t}{\sum_{t=1}^n Q_t^{(i)}(k)}$$

k-means can be seen as hard GMM with spherical covariance and

# WHY SHOULD EM WORK?

A very high level view:

- Performing E-step will never decrease log-likelihood (or log a posteriori)

# WHY SHOULD EM WORK?

A very high level view:

- Performing E-step will never decrease log-likelihood (or log a posteriori)
- Performing M-step will never decrease log-likelihood (or log a posteriori)

# WHY SHOULD EM WORK?

Steps to show that  $\log \text{Lik}(\theta^{(i+1)}) \geq \log \text{Lik}(\theta^{(i)})$  :

- Insert latent variables
- Use Jensen's inequality (log is a concave function)
- Massage the terms!



# WHY SHOULD EM WORK?

- Likelihood never decreases
- So whenever we converge we converge to a local optima
- However problem is non-convex and can have many local optimal
- In general no guarantee on rate of convergence
- In practice, do multiple random initializations and pick the best one!

# EM IN GENERAL

- There was nothing special about GMM or clustering problems
- EM can be used as a general strategy for any problem with latent/missing/unobserved variables
- The MAP version only involves an extra prior term over  $\theta$  multiplied to the likelihood
- In general probabilistic models with observed and latent variables can be represented succinctly as graphical models.  
**Next time ...**