

Royal Institute of Technology

## MACHINE LEARNING 2 - THE EM ALGORITHM Lecture 9

## LAST LECTURE

- ★ Backward
- ★ Smoothing
- ★ Sampling
- ★ Viterbi
- ★ K-means (inspiration)
- ★ GMM (towards EM)

# THIS LECTURE

- ★ GMM
- ★ EM
  - \* given "structure" and observations find parameters
- \* EM algorithm for GMM
- ★ EM algorithm for BMM
- \* Baum-Welch EM algorithm for training an HMM

#### ASSIGNING POINTS TO MULTIPLE MEANS



## K-MEANS AS GMM

 $\star$  Fixed variance, a Gaussian and mean per cluster, i.e.,  $\theta_c = (\mu_c, \sigma^2)$ 

 $\star$  Idea: each point can belong to several means (clusters)

 $\star$  Use responsibilities to find means

1

$$r_{nc} = p(z_n = c | \boldsymbol{x}_n, \boldsymbol{\theta}) = \frac{p(z_n = c | \boldsymbol{\theta}) p(\boldsymbol{x}_n | z_n = c, \boldsymbol{\theta})}{\sum_{c=1}^{C} p(z_n = c | \boldsymbol{\theta}) p(\boldsymbol{x}_n | z_n = c, \boldsymbol{\theta})}$$

$$\boldsymbol{\mu}_{c} = \frac{1}{N_{c}} \sum_{n} r_{nc} \boldsymbol{x}_{n}, \qquad \text{where } N_{c} = \sum_{n} r_{nc}$$

#### 1-DIM GAUSSIAN MIXTURE MODELS

 $\mathcal{D} = (oldsymbol{x}_1, \dots, oldsymbol{x}_N)$ 

#### $Z \sim \operatorname{Cat}(\pi)$

 $p(\boldsymbol{X}|Z=c) = \mathcal{N}(\boldsymbol{X}|\boldsymbol{\mu}_{c},\sigma_{c})$  $\boldsymbol{\theta}_{c} = (\boldsymbol{\mu}_{c},\sigma_{c})$ 

# Z hidden $\sim \operatorname{Cat}(\pi)$

 $oldsymbol{X} \sim \mathcal{N}(oldsymbol{\mu}_c, \sigma_c)$ 

## EXAMPLE

 $z_n$  is red with probability 1/2, green with probability 3/10, blue with probability 1/5



-

The three gaussian distributions in our mixture

 $z_n = blue$ 

 $x_n$  is generated from the Gaussian indicated by  $z_n$ 

We get  $x_1, \ldots, x_N$ 

# 

So,

 $p(x_n, z_n) = p(z_n)p(x_n|z_n)$ 

and

$$p(x_n) = \sum_{c=1}^{C} p(z_n = c) p(x_n | z_n = c) = \sum_{c=1}^{C} \pi_c p(x_n | z_n = c)$$

and

$$r_{nc} = p(z_n = c | x_n) = \frac{p(z_n = c, x_n)}{p(x_n)} = \frac{\pi_c p(x_n | z_n = c)}{\sum_{c=1}^C \pi_c p(x_n | z_n = c)}$$

$$COMPLETE LOG$$

$$I(\theta'; \mathcal{D}) = \log \prod_{n} p(\boldsymbol{x}_{n}, z_{n} | \theta')$$

$$= \sum_{n} \log \prod_{c} (\pi_{c}' p(\boldsymbol{x}_{n} | Z_{n} = c, \theta')^{I(z_{n} = c)})$$

$$= \sum_{n} \sum_{c} I(z_{n} = c) \log(\pi_{c}' p(\boldsymbol{x}_{n} | \theta_{c}'))$$

$$= \sum_{n} \sum_{c} I(z_{n} = c) \log \pi_{c}' + \sum_{c} \sum_{n:I(z_{n} = c)} \log p(\boldsymbol{x}_{n} | \theta_{c}')$$

$$= \sum_{c} \sum_{n:I(z_{n} = c)} \log \pi_{c}' + \sum_{c} \sum_{n:I(z_{n} = c)} \log p(\boldsymbol{x}_{n} | \theta_{c}')$$

$$= \sum_{c} N_{c} \log \pi_{c}' + \sum_{c} \sum_{n:I(z_{n} = c)} \log p(\boldsymbol{x}_{n} | \theta_{c}')$$

$$= \sum_{n} I(z_{n} = c)$$

 $N_{o}$ 

MLE -COMPLETE DATA FOR GAUSSIAN  $\mathcal{D} = (z_1, x_1), \dots, (z_N, x_N))$ • Maximizing the complete log likelihood

$$l(\theta'; \mathcal{D}) = \sum_{c} N_c \log \pi'_c + \sum_{c} \sum_{n: I(z_n = c)} \log p(\boldsymbol{x}_n | \theta'_c)$$

Boils down to maximizing



### EM & EXPECTED LOG LIKELIHOOD (Q-TERM)

- Iteratively maximizing the expected log likelihood (expected sufficient statistics).
- Iteratively maximizing the expected log likelihood in practice always leads to a local maxima
- The expectation is over hidden variables given data and current parameters
- We maximize the expression by choosing new parameters.

RELATIONS BETWEEN LOG-LIKELIHOODS AND Q-TERMS

Q-term or expected complete log-likelihood

$$Q(\boldsymbol{\theta}', \boldsymbol{\theta}) = \sum_{n} E_{p(Z_n | \boldsymbol{x}_n, \boldsymbol{\theta})} \left[ l(\boldsymbol{\theta}'; Z_n, \boldsymbol{x}_n) \right]$$

Theorem: for  $\boldsymbol{\theta}' = \operatorname{argmax}_{\boldsymbol{\theta}'} Q(\boldsymbol{\theta}', \boldsymbol{\theta})$ 

 $\log p(\mathcal{D}|\boldsymbol{\theta}') \ge Q(\boldsymbol{\theta}',\boldsymbol{\theta}) - R(\boldsymbol{\theta},\boldsymbol{\theta}) \ge Q(\boldsymbol{\theta},\boldsymbol{\theta}) - R(\boldsymbol{\theta},\boldsymbol{\theta}) = \log p(\mathcal{D}|\boldsymbol{\theta})$ 

log-likelihood

So, by maximizing Q-term (through ESS), we monotonically increase the likelihood.

The Q-term may not increase in every step!

#### LOG LIKELIHOOD & EXPECTED LOG LIKELIHOOD (Q-TERM)

Complete log likelihood

$$l(\theta'; \mathcal{D}) = \sum_{n} \sum_{c} I(z_n = c) \log \pi'_c + \sum_{n} \sum_{c} I(z_n = c) \log p(\boldsymbol{x}_n | \theta'_c)$$

Expected complete log likelihood a.k.a the Q term

$$\sum_{n} E_{p(Z_n | \boldsymbol{x}_n, \boldsymbol{\theta})} \left[ l(\boldsymbol{\theta}'; Z_n) \right] = \sum_{n} \sum_{c} r_{nc} \log \pi'_c + \sum_{n} \sum_{c} r_{nc} \log p(\boldsymbol{x}_n | \boldsymbol{\theta}'_c)$$

 $\bigstar \text{ We want } \operatorname{argmax}_{\theta'} E_{p(\boldsymbol{z}_n | \boldsymbol{x}_n, \boldsymbol{\theta})} \left[ l(\theta'; \mathcal{D}) \right]$ 

**★** The 2 sums 
$$\sum_{c} \left( \sum_{n} r_{nc} \right) \log \pi'_{c}$$
 &  $\sum_{c} \sum_{n} r_{nc} \log p(\boldsymbol{x}_{n} | \theta'_{c})$ 

are independent

$$\bigstar$$
 So,  $\pi_c' = \sum_n r_{nc}/N = r_c/N$ 

 $\star$  In the second, different c indices are independent

 $\star$  So, we want to maximize each

$$\sum_{n} r_{nc} \log \frac{1}{\sqrt{2\pi\sigma_c^{\prime 2}}} \exp\left(-\frac{1}{2\sigma_c^{\prime 2}}(x_n - \mu_c)^2\right)$$

### GMM EM-ALGORITHM

- E-step: compute  $r_{nc} = p(Z_n = c | x_n, \theta)$
- M-Step: maximize (1) mixture coefficients and (2) each

$$\sum_{n} r_{nc} \log \frac{1}{\sqrt{2\pi\sigma_c'^2}} \exp\left(-\frac{1}{2\sigma_c'^2}(x_n - \mu_c)^2\right)$$

by setting

٠

$$\mu'_c = \frac{\sum_n r_{nc} x_n}{r_c} \quad \text{and} \quad \sigma'^2_c = \frac{1}{\alpha'^2_c} = \sum_n r_{nc} (x_n - \mu'_c)^2 / r_c$$
  
set  $\theta = \theta'$ 

• Stop when solution or likelihood hardly change otherwise repeat

## E AND M STEPS



#### EM-ALGORITHM IN GENERAL

- E-step: compute  $E_{p(Z_n|\boldsymbol{x}_n,\boldsymbol{\theta})}\left[l(\boldsymbol{\theta}';Z_n,\boldsymbol{x}_n)\right]$
- M-Step:

$$\theta' = \operatorname{argmax}_{\theta'} \sum_{n} E_{p(Z_n | \boldsymbol{x}_n, \boldsymbol{\theta})} \left[ l(\theta'; Z_n, \boldsymbol{x}_n) \right]$$

• set  $\theta = \theta'$ 

• Stop when solution or likelihood hardly change otherwise repeat

#### ★ Starting points

- ★ Number of starting points
- ★ Sieving starting points
- ★ The competition
  - The first iterations of EM show huge improvement in the likelihood. These are then followed by many iterations that slowly increase the likelihood. Conjugate gradient shows the opposite behaviour.

## PRACTICAL ISSUES

MIXTURE OF BERNOULLI

 $\mathcal{D} = (oldsymbol{x}_1, \dots, oldsymbol{x}_N)$ 

 $Z \sim \operatorname{Cat}(\pi)$ 

 $p(X_d|Z=c) = Ber(X|\theta_{cd})$ 

Z hidden  $\sim \operatorname{Cat}(\boldsymbol{\pi})$ 



 $X_1 \sim \operatorname{Ber}(\mu_{Z1}) \quad X_D \sim \operatorname{Ber}(\mu_{ZD})$ 

#### MIXTURE OF BERNOULLI - BASE

- $\star$  Basic model D binary variables  $x_1,\ldots,x_D$  and  $oldsymbol{x}=x_1,\ldots,x_D$
- \* Parameters  $\boldsymbol{\mu} = \mu_1, \ldots, \mu_D$  where  $p(x_i = 1 | \boldsymbol{\mu}) = \mu_i$
- \* So  $E[\boldsymbol{x}|\boldsymbol{\mu}] = \boldsymbol{\mu}$  and  $\operatorname{cov}[\boldsymbol{x}|\boldsymbol{\mu}] = \operatorname{diag}\{\mu_i(1-\mu_i)\}$

#### FULL MIXTURE MODEL -MEAN AND COVARIANCE

$$E[\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\pi}] = \sum_{k=1}^{K} \pi_k \boldsymbol{\mu}_k \qquad E[x_i|\boldsymbol{\mu},\boldsymbol{\pi}] = \sum_{k=1}^{K} \pi_k \mu_{ki}$$

$$\operatorname{cov}[x_i, x_j | \boldsymbol{\mu}, \boldsymbol{\pi}] = \sum_{k=1}^K \pi_k \mu_{ki} \mu_{kj} - \sum_{k=1}^K \sum_{k=1}^K \pi_k \pi_{k'} \mu_{kj} \mu_{k'j} \neq 0$$

#### MIXTURE OF BERNOULLI

- $\star$  Class variable  $z \in [K]$  and D binary variables  $oldsymbol{x} = x_1, \ldots, x_D$
- \* Parameters  $\boldsymbol{\mu}_k = \mu_{k1}, \ldots, \mu_{kD}$  and  $\boldsymbol{\pi} = \pi_1, \ldots, \pi_K$
- ★ Mixture with likelihood

$$p(\pmb{x}|\pmb{\mu},\pmb{\pi}) = \sum_{k=1}^{K} \pi_k p(\pmb{x}|\pmb{\mu}_k)$$
 where

$$p(x_i|\boldsymbol{\mu}_k) = \mu_{ki}^{x_i} (1 - \mu_{ki})^{1 - x_i}$$

and

$$p(\boldsymbol{x}|\boldsymbol{\mu}_k) = \prod_{i=1}^{D} \mu_{ki}^{x_i} (1 - \mu_{ki})^{1 - x_i}$$



#### BAUM-WELCH: LEARNING HMM PARAMETERS

- $\star$  Starts in the state z<sub>1</sub>
- $\star$  When in state  $z_t$ 
  - outputs p(xt|zt)  $B_{x_t,z_t}$
  - moves to  $p(z_{t+1}|z_t)$

$$A_{z_{t+1},z_t}$$

 Stops after a fixed number of steps or when reaching a stop step

The parameters

we want to learn

#### LEARNING TRANSITION AND EMISSION PARAMETERS - FULLY OBSERVED DATA

- \* Parameters
  - transition  $A_{lk} = p(Z_t = l | Z_{t-1} = k)$
  - emission  $B_{sk} = p(X_t = s | Z_t = k)$

★ Data

$$\mathcal{D} = \{ (x_{1:T}^1, z_{1:T+1}^1), \dots, (x_{1:T}^N, z_{1:T+1}^N) \}$$

\* Likelihood

$$L(\boldsymbol{\theta}; \mathcal{D}) = \prod_{n=1}^{N} \prod_{t=1}^{T} \left[ \prod_{k,s} B_{sk}^{I(x_t^n = s, z_t^n = k)} \prod_{k,l} A_{lk}^{I(z_t^n = k, z_{t+1}^n = l)} \right]$$

# - COMPLETE DATA

$$l(\boldsymbol{\theta}; \mathcal{D}) = \sum_{k,s} M_{k,s} \log B_{ks} + \sum_{k,l} N_{k,l} \log A_{kl}$$

where

$$M_{s,k} = |\{(n,t) : x_t^n = s, z_t^n = k\}|$$
$$N_{k,l} = |\{(n,t) : z_t^n = l, z_{t+1}^n = k\}|$$

Maximized by

$$B_{sk} = M_{s,k} / \sum_{s} M_{s,k} = M_{s,k} / N_k \quad \& \quad A_{kl} = N_{k,l} / \sum_{l} N_{k,l} = N_{k,l} / N_k$$
  
where  $N_k = |\{(n,t) : z_t^n = k\}|$ 

#### EN EOR HIDDEN DATA



Maximized by

$$B_{sk} = \overline{M}_{s,k} / \sum_{s} \overline{M}_{s,k} = \overline{M}_{s,k} / \overline{N}_{k} \quad \& \quad A'_{kl} = \overline{N}_{k,l} / \sum_{l} \overline{N}_{k,l} = \overline{N}_{k,l} / \overline{N}_{k}$$
  
Both obtainable from forward and backward (smoothing like)

#