

## EXPECTATION MAXIMIZATION FOR MULTINOMIAL MIXTURE MODEL

Setting :-

- $N$  'documents'
- $n$ th document contains  $W_n$  'words'
- $x_{nj}$  :-  $j$ th word in  $n$ th document  
:- a category out of  $V$  words.

Assumption :-  $K$  multinomial distribution

:- Each document generated from  
1 of the  $K$  distributions.

Generative model :-

Given  $K$  multinomial distributions  $\theta_1, \dots, \theta_K$   
and their rate vector  $\pi$

- for each document  $n: 1 \rightarrow N$ 
  - Sample  $z_n \sim \text{Multinomial}(\pi)$
  - for each word  $j: 1 \rightarrow W_n$ 
    - Sample  $x_{nj} | z_n \sim \text{Multinomial}(\theta_{z_n})$

Recall :-

$$\log P(x|\theta) \geq E_q \log P(x,z|\theta) - E_q \log q(z)$$

For above model,  $P(x, z | \theta)$  :-

$$= \prod_{n=1}^N P(x_{n,1}, \dots, x_{n,W_n}, z_n | \theta_1, \dots, \theta_K, \pi)$$

$$= \prod_{n=1}^N P(z_n | \pi) \prod_{j=1}^{W_n} P(x_{n,j} | \theta_{z_n})$$

So,

$$P(z_n | \pi) = \prod_{k=1}^K \pi_k^{z_{n,k}}, \quad \log P(z_n | \pi) = \sum_{k=1}^K z_{n,k} \log \pi_k$$

and,

$$P(x_{n,j} | \theta_{z_n}) = \prod_{v=1}^V \theta_{z_n, v}^{x_{n,j,v}}$$

$$= \prod_{k=1}^K \left[ \prod_{v=1}^V \theta_{k,v}^{x_{n,j,v}} \right]^{z_{n,k}}$$

$$\log P(x_{n,j} | \theta_{z_n}) = \sum_{k=1}^K z_{n,k} \left( \sum_{v=1}^V x_{n,j,v} \log \theta_{k,v} \right)$$

therefore,

$$\log P(x, z | \theta) = \sum_{n=1}^N \left\{ \sum_{k=1}^K z_{n,k} \log \pi_k + \right.$$

$w_n$

$K$

$\dots$

$\dots$

$$\sum_{j=1} \left[ \sum_{k=1} z_{n,k} \left( \sum_{v=1}^V x_{n,j,v} \log \theta_{k,v} \right) \right]$$

A

Now compute  $E_{q(z)} \log P(x, z | \theta)$

- What is  $q(z)$  ?

- Recall  $q(z)$  is chosen to 'best' capture our beliefs about latent variables

What about  $P(z|x, \theta)$  ?

→ Captures all uncertainty about  $z$   
given the data and model parameters

Can  $q(z) \stackrel{?}{=} P(z|x, \theta)$

Can  $P(z|x, \theta)$  be computed easily?

If yes, proceed

If no, further assumptions necessary  
(subject matter of variational inference)

Comments:-

also depends on what is designated as latent variables.

For above model,

$$P(z|x,\theta,\pi) = P(z_1, \dots, z_N | x, \theta, \pi)$$

$$= \frac{P(x | z_1, \dots, z_N, \theta) \cdot P(z_1, \dots, z_N | \pi)}{P(x | \theta, \pi)}$$

$$= \frac{P(x_1, \dots, x_N | z_1, \dots, z_N, \theta) \prod_{n=1}^N P(z_n | \pi)}{P(x | \theta, \pi)}$$

$$P(z_1, \dots, z_N | x, \theta) = \frac{\prod_{n=1}^N P(x_n | \theta, z_n) \cdot P(z_n | \pi)}{P(x | \theta, \pi)}$$

What is  $P(x | \theta, \pi)$

$$P(x | \theta, \pi) = P(x_1, \dots, x_N | \theta, \pi)$$

$$= \sum_{z_1, \dots, z_N} P(x_1, z_1, \dots, x_N, z_N | \theta, \pi)$$

$$= \sum_{z_1, \dots, z_N} \prod_{n=1}^N P(x_n, z_n | \theta, \pi)$$

$$= \prod_{n=1}^N \left( \sum_{z_n} P(x_n, z_n | \theta, \pi) \right)$$

$$= \prod_{n=1}^N P(x_n | \theta, \pi) \quad \left\{ \begin{array}{l} \text{our original} \\ \text{incomplete data} \\ \text{likelihood} \end{array} \right\}$$

Therefore,

$$\begin{aligned} P(z_1, \dots, z_N | x_1, \dots, x_N, \theta, \pi) &= \prod_{n=1}^N \left[ \frac{P(x_n | \theta, z_n) P(z_n | \pi)}{P(x_n | \theta, \pi)} \right] \\ &= \prod_{n=1}^N P(z_n | x_n, \theta, \pi) \end{aligned}$$

— (B)

- Fully factorized posterior.
- Each  $P(z_n | x_n, \theta, \pi)$  easily computable
- Hence  $q_j(z) \doteq P(z | x, \theta)$

Recall in EM :-

- Given  $\theta^{(old)}$ , we compute  $q(z | \theta^{old})$
- Compute  $E_{q(z | \theta^{old})} \log P(x, z | \theta) - E_{q(z | \theta^{old})} \underbrace{\log q(z | \theta)}_{\text{entropy}}$
- maximize above w.r.t ' $\theta$ '
- entropy term independent of ' $\theta$ '

So, we compute  $q(z|\theta^{\text{old}}) = p(z|x, \theta^{\text{old}})$

$$p(z_n|x_n, \theta^{\text{old}}, \pi^{\text{old}}) = \frac{p(x_n|\theta_{z_n}^{\text{old}}) \cdot p(z_n|\pi^{\text{old}})}{\sum_{k=1}^K \pi_k^{\text{old}} \cdot p(x_n|\theta_k^{\text{old}})}$$

$$= [r_{n,1}, \dots, r_{n,K}]$$

$$\overbrace{p(z_n = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} | x_n, \theta^{\text{old}}, \pi^{\text{old}})}$$

and compute

$$E_{q(z|\theta^{\text{old}})} \log p(x, z|\theta)$$

$$= E_{q(z|\theta^{\text{old}})} \sum_{n=1}^N \left\{ \sum_{k=1}^{w_n} z_{n,k} \log \pi_k + \sum_{j=1}^{v_n} \left[ \sum_{k=1}^K z_{n,k} \left( \sum_{v=1}^V x_{n,j,v} \log \theta_{k,v} \right) \right] \right\}$$

- Because of fully factorized posterior  $p(z|x, \theta^{\text{old}})$ ,

Computing Expectation is trivial

$$S_0, = \sum_{n=1}^N \left\{ \sum_{k=1}^K r_{n,k} \log \pi_k + \sum_{j=1}^{W_n} \left[ \sum_{k=1}^K r_{n,k} \left( \sum_{v=1}^V x_{n,j,v} \log \theta_{k,v} \right) \right] \right\}$$

a little re-arrangement,

$$= \sum_{k=1}^K \left( \sum_{n=1}^N r_{n,k} \right) \log \pi_k +$$

$$\sum_{k=1}^K \sum_{v=1}^V \left( \sum_{n=1}^N r_{n,k} \sum_{j=1}^{W_n} x_{n,j,v} \right) \log \theta_{k,v}$$

- Recognize this lets us do ML estimates of multinomials  $\pi_1, \theta_1, \dots, \theta_K$  independently
- We can just read-off result for each multinomial

$$\implies \pi_k = \sum_{n=1}^N r_{n,k} / N$$

$\theta_{k,v}$  :-  $v$ th-component of  $k$ th distribution

$$\theta_{k,v} = \frac{\sum_{n=1}^N r_{n,k} \sum_{j=1}^{w_n} x_{n,j,v}}{\sum_{v=1}^V \sum_{n=1}^N r_{n,k} \sum_{j=1}^{w_n} x_{n,j,v}}$$

These  $\pi$  and  $\theta_k$  become  $\pi^{\text{old}}$  and  $\theta_k^{\text{old}}$  for next iteration.

### Implementation Comments:-

- Compute  $E_{q(z)} \log P(x, z | \theta) - E_{q(z)} \log q(z)$

in every iteration

- Should improve in every iteration.

- Stop when improvement < threshold

### Moral

EM decomposes the mixture model in independent ML estimates of model components in each iteration.