

## COLLAPSED SAMPLING

Recap :- Bayesian mixture model

## Generative model :-

Given hyper-parameters

- (i) Sample mixture components and rate parameter
  - (ii) For each data point
    - a) Sample mixture component membership
    - b) Sample data point given mixture membership and mixture components

*prior, likelihood*

*prior, likelihood*

Example :- Multinomial mixture model

Given  $\gamma = [\gamma_1, \dots, \gamma_v]$ ,  $\alpha = [\alpha_1, \dots, \alpha_k]$

$$\begin{aligned}
 & \text{Sample } \theta_k \sim \text{Dir}(\theta_k | \gamma) \quad k=1, \dots, K \\
 & \pi \sim \text{Dir}(\pi | \alpha) \\
 & z_n | \pi \sim \text{Multinomial}(z_n | \pi) \quad n=1, \dots, N \\
 & x_{n,j} | z_n, \theta \sim \text{Multinomial}(x_{n,j} | \theta_{z_n}) \quad j=1, \dots, J_n
 \end{aligned}$$

\* Different prior-likelihood give different models

\* Of particular importance are "conjugate" prior, likelihood pairs

\* In Bayesian setting, we are interested in

$$\Pr(z_{\text{new}} | x_{\text{new}}, x_{\text{Train}}, \text{hyper-parameters}) \\ \approx \frac{1}{M} \sum_m \Pr(z_{\text{new}} | x_{\text{new}}, \theta_1^{(m)}, \dots, \theta_K^{(m)}, \pi^{(m)})$$

\* Sampling based solution

\* Sampling involves drawing samples from  $P(\theta_1, \dots, \theta_K, \pi, z_1, \dots, z_n | x_{\text{Train}}, h)$

For multinomial mixture :-

$$* \theta_k | \theta_{-k}, z, x, \pi \sim \text{Dir}(\theta_k | [\sum_{n:z_n=k}^{\textcircled{A}} c_{n,1} + \gamma_1, \dots, \sum_{n:z_n=k} c_{n,V} + \gamma_V])$$

$$* \pi | z, x, \theta \sim \text{Dir}(\pi | [\sum_{n=1}^N z_{n,1} + \alpha_1, \dots, \sum_{n=1}^N z_{n,K} + \alpha_K])$$

$$* z_{n,k} | z_{-n}, \theta, \pi, x_{n,1}, \dots, x_{n,J_n} \sim \left[ \prod_{j=1}^{J_n} \text{Multinomial}(x_{n,j} | \theta_k) \right] \cdot \pi_k$$

$$\sum_{m=1}^K \left[ \prod_{j=1}^{J_n} \text{Multinomial}(x_{n,j} | \theta_m) \right] \pi_m$$

$$\textcircled{A} C_{n,v} = \sum_{j=1}^{J_n} x_{n,j,v}$$

= # word 'v' in document n

- \* Notice, Sampling of  $z_n$  only depends on  $\theta_1, \dots, \theta_K$  and  $\pi$  (and  $n$ -th-document)
- \*  $z_n$  and  $z_m$  (local states) only influence each other through  $\theta_1, \dots, \theta_K$  and  $\pi$  (global state)
- \* Intuitively, Speed of convergence depends on how frequently we can update "global state"
- \* Can we update "global state" more frequently
  - Possibly (expensive)
  - Can we simply get rid of "global state"
  - getting rid = marginalizing.

Or :-  $\Pr(z_n | z_{-n}, \chi_{\text{Train}}, \text{hyper-parameter})$  ?

Our Bayesian prediction now becomes :-

$$\Pr(z_{\text{new}} | x_{\text{new}}, \chi_{\text{Train}}, h) \approx \frac{1}{M} \sum_m \Pr(z_{\text{new}} | x_{\text{new}}, \chi_{\text{Train}}, z_{\text{Train}}^{(m)}, h)$$

$$\Pr(z_n | z_{-n}, x_n, x_{-n}, \gamma, \alpha) = \frac{\Pr(z_n, z_{-n}, x_n, x_{-n} | \gamma, \alpha)}{\sum_{z_n} \Pr(z_n, z_{-n}, x_n, x_{-n} | \gamma, \alpha)}$$

?

$$= \Pr(z_n | z_{-n}, x_{-n}, \gamma, \alpha) \cdot \Pr(x_n | z_n, z_{-n}, x_{-n}, \gamma, \alpha)$$


---


$$\sum_{z_n} \Pr(z_n | z_{-n}, x_{-n}, \gamma, \alpha) \cdot \Pr(x_n | z_n, z_{-n}, x_{-n}, \gamma, \alpha)$$

from Session 4c

$$\Pr(z_n=l | z_{-n}, x_{-n}, \gamma, \alpha) = \underbrace{\Pr(z_n=l | \pi, z_{-n}, \gamma, \alpha)}_{d\pi} \Pr(\pi | z_{-n}, \gamma, \alpha)$$

$$= \frac{1}{N-1 + (\sum_k \alpha_k)} \left[ \left( \sum_{m \neq n} z_{m,l} \right) + \alpha_l \right]$$

number of assignments  
to component 'l'  
(excluding current  
data point)

$$\Pr(x_n | z_n=l, z_{-n}, x_{-n}, \gamma, \alpha)$$

$$= \int \Pr(x_n | z_n=l, z_{-n}, x_{-n}, \gamma, \alpha, \theta_l) \cdot \Pr(\theta_l | x_{-n}, z_{-n}, \gamma, \alpha) d\theta_l$$

↑ Why just  $\theta_l$ ?

$$= \int_{\Theta_L} d\theta_L \left[ \prod_{j=1}^{J_n} \text{Multinomial}(x_{n,j} | \theta_L) \right] \text{Dir}(\theta_L | [\sum_{\substack{m: m \neq n \\ m: m \neq n}} c_{m,1} + \gamma_1, \dots, \sum_{\substack{m: m \neq n \\ m: m \neq n}} c_{m,V} + \gamma_V])$$

and  
 $Z_m = 1$   
and  
 $Z_m = 1$

$$\left[ \int d\theta_L \prod_{v=1}^V (\theta_{L,v})^{c_{n,v}} \cdot \prod_{v=1}^V (\theta_{L,v})^{\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_{v-1}} \right] \cdot \frac{\Gamma(\sum_{v=1}^V (\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_v))}{\prod_{v=1}^V \Gamma(\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_v)}$$

$$= \prod_{v=1}^V \frac{\Gamma(\sum_{\substack{d: z_d=1}} c_{d,v} + \gamma_v)}{\Gamma(\sum_{v=1}^V (\sum_{\substack{d: z_d=1}} c_{d,v} + \gamma_v))}$$

number of words assigned to component 'l' +  $c_{n,v}$

number of words assigned to component 'l':

number of words assigned to Component l + number of words in document 'n' ( $= J_n$ ) to component 'l':

$$= \frac{\Gamma(\sum_{v=1}^V (\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_v))}{\Gamma(\sum_{v=1}^V (\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_v) + J_n)} \cdot \prod_{v=1}^V \frac{\Gamma(\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_v + c_{n,v})}{\Gamma(\sum_{\substack{m: m \neq n \\ z_m=1}} c_{m,v} + \gamma_v)}$$

Using Gamma function Recurrence

$$\therefore \Gamma(x+1) = x\Gamma(x)$$

$$= \left[ \frac{1}{\sum_{v=1}^V \left( \sum_{m: m \neq n} c_{m,v} + \gamma_v \right)} \right]_n \cdot \prod_{v=1}^V \left( \sum_{\substack{m: m \neq n \\ z_m = l}} c_{m,v} + \gamma_v \right)^{c_{n,v}}$$

In Summary :-

$$P(z_n=l | z_{-n}, x_n, x_{-n}, \alpha, \gamma) \propto \frac{1}{N-1 + (\sum \alpha_k)} \left( \sum_{m \neq n} z_{m,l} + \alpha_l \right) \times \\ \left[ \frac{1}{\sum_{v=1}^V \left( \sum_{m: m \neq n} c_{m,v} + \gamma_v \right)} \right] \prod_{v=1}^V \left( \sum_{\substack{m: m \neq n \\ z_m = l}} c_{m,v} + \gamma_v \right)^{c_{n,v}}$$

Implementation:- everything is expressed as counts

i-e \* count number of assignments of documents to mixture components.

\* count number of assignments of words to mixture components.

\* Can be very efficiently implemented by maintaining conditional frequency tables.

\* entire inference is Coupled (not easy to parallelize)

\* Collapsed inference reduces dimensionality of the problem a lot  
→ resulting in faster 'mixing'

means MCMC procedures reaches stationary distribution (joint distribution much faster)

\* Other benefits :- application to non-Parametric models.