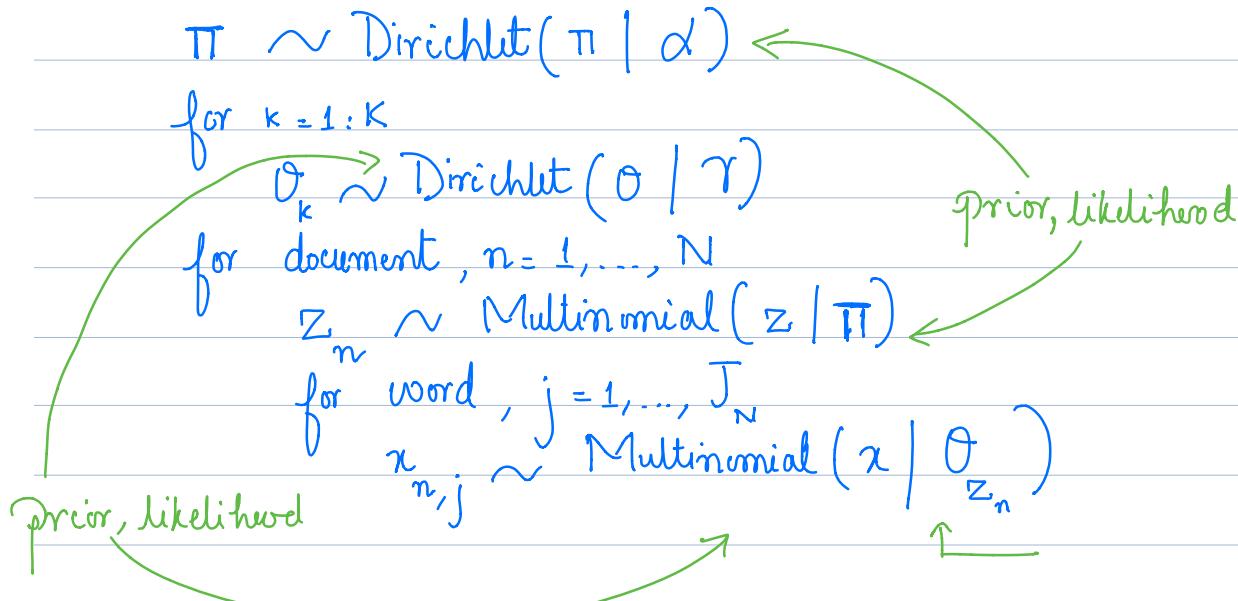


# Bayesian treatment of Multinomial Mixture

Generative model :-

Given parameters,  $K$ ,  $\alpha = [\alpha_1, \dots, \alpha_K]$ ,  $\gamma = [\gamma_1, \dots, \gamma_V]$



$$z_n := 1 - \text{out of } K, \quad z_n = k = [0, 0, \dots, \underset{k}{\frac{1}{\uparrow}}, \dots, 0]$$
$$x_{n,j} := 1 - \text{out of } V$$

Note:- • We still have parameters without prior  
i.e  $K, \alpha, \gamma$

- We can put priors on them
- Those priors will have their own parameters
- We can put priors on them
- and so forth...

Goal:-

Predict  $P(z_{\text{new}} | x_{\text{new}}, \chi_{\text{Train}}, \underbrace{\theta, \pi}_{\text{hyperparameters}}, h)$

$$\int P(z_{\text{new}} | x_{\text{new}}, \theta_1, \dots, \theta_k, \pi, \chi_{\text{Train}}, h) \cdot P(\theta_1, \dots, \theta_k, \pi | \chi_{\text{Train}}, h) \underset{\text{do}, \dots, \text{do}}{\text{d}\theta} \rightarrow A$$

$$P(z_{\text{new}} | x_{\text{new}}, \theta_1, \dots, \theta_k, \pi, \chi_{\text{Train}}, h) = \frac{P(x_{\text{new}} | z_{\text{new}}, \theta_1, \dots, \theta_k, \pi, \chi_{\text{Train}}, h)}{P(z_{\text{new}} | \theta_1, \dots, \theta_k, \pi, \chi_{\text{Train}}, h)}$$

$$P(x_{\text{new}}, \theta_1, \dots, \theta_k, \pi, \chi_{\text{Train}}, h)$$

$$= \frac{P(x_{\text{new}} | \theta_{z_{\text{new}}}) \cdot P(z_{\text{new}} | \pi)}{\sum_{z_{\text{new}}} P(x_{\text{new}} | \theta_{z_{\text{new}}}) \cdot P(z_{\text{new}} | \pi)} \rightarrow \begin{array}{l} \text{We have seen} \\ \text{this many times.} \\ (q(z) \text{ in EM}) \end{array}$$

$$= q_r(z_{\text{new}}; \theta_{z_{\text{new}}}, \pi)$$

So in  $A$ ,

$$= \int q(z_{\text{new}}; \theta_{z_{\text{new}}}, \pi) \underbrace{P(\theta_{z_{\text{new}}}, \pi | \chi_{\text{Train}}, h)}_{\text{We don't need all } \theta's} \text{d}\theta \cdot \text{d}\pi$$

Looks like

$$\mathbb{E}_{P(\theta_{z_{\text{new}}}, \pi | \chi_{\text{Train}}, h)} q(z_{\text{new}}; \theta_{z_{\text{new}}}, \pi)$$

$$\approx \frac{1}{M} \sum_{m=1}^M q(z_{\text{new}}; \theta_{z_{\text{new}}}^{(m)}, \pi^{(m)}) \longrightarrow \tilde{A}$$

Where  $\theta_{z_{\text{new}}}^{(m)}$  and  $\pi^{(m)}$  are mth-samples from

$$P(\theta_{z_{\text{new}}}, \pi | X_{\text{Train}}, h)$$

\* Hence, instead of finding  $P(\theta_{z_{\text{new}}}, \pi | X_{\text{Train}}, h)$

(i.e posterior distribution of parameters of interest)

We can just "draw some samples" from  $P(\theta_{z_{\text{new}}}, \pi | X_{\text{Train}}, h)$

\* We'll call it Approach ① to approximate  $\tilde{A}$ .

\* Approach ② would be to come up with an answer for  $P(\theta_{z_{\text{new}}}, \pi | X_{\text{Train}}, h)$

. If not possible to find exact form, then possible solutions are

\* Expectation propagation

\* Variational approximation.

Approach (1) :- Sampling from the posteriors.

General recipe:-

$X$  :- data,  $L = [l_1, \dots, l_N]$  all latent variables (including  $\theta_1, \dots, \theta_K, \pi$  (model parameters))  
 $h$  :- hyper-parameters

(i) Step 1 :- Write  $P(X, L | h)$

(ii), Step 2 :-  $l_n \sim P(l_n | X, L_{-n}, h) = \frac{P(X, L | h)}{P(X, L_{-n} | h)}$

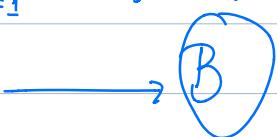
\* Should be easy!

\* true for models using conjugacy.

for multinomial mixture model,

$$= P(\theta_1, \dots, \theta_K, \pi, z_1, \dots, z_N, x_1, \dots, x_N | \gamma, \alpha)$$

$$= \left[ \prod_{k=1}^K P(\theta_k | \gamma) \right] P(\pi | \alpha) \cdot \prod_{n=1}^N P(z_n | \pi) \cdot \prod_{j=1}^{J_n} P(x_{n,j} | \theta_{z_n})$$



$$\begin{aligned}
 P(z_n | x, z_{-n}, h) &= \frac{\left[ \prod_{k=1}^K P(\theta_k | \gamma) \right] P(\pi | \alpha) \prod_{m=1}^N P(z_m | \pi) \prod_{j=1}^{J_m} P(x_{m,j} | \theta_{z_m})}{\left\{ \left[ \prod_{k=1}^K P(\theta_k | \gamma) \right] P(\pi | \alpha) \prod_{m \neq n} P(z_m | \pi) \prod_{j=1}^{J_m} P(x_{m,j} | \theta_{z_m}) \right\}} \\
 &\quad \underbrace{P(x_n | \theta_{z_1}, \dots, \theta_{z_K}, \pi)}_{\sum_{z_1} P(x_n | \theta_{z_1}) P(z_1 | \pi)} \\
 &= \frac{P(z_n | \pi) \cdot \prod_{j=1}^{J_n} P(x_{n,j} | \theta_{z_n})}{\sum_{z_1} P(x_n | \theta_{z_1}) \cdot P(z_1 | \pi)} \quad (\text{i}) \\
 P(\theta_k | \Theta_{-k}, \pi, z_N, x, h) &= P(\theta_k | \gamma) \prod_{n: z_n=k} \prod_{j=1}^{J_n} P(x_{n,j} | \theta_k) \\
 &\quad \underbrace{P(\theta_k | \gamma) \prod_{n: z_n=k} \prod_{j=1}^{J_n} P(x_{n,j} | \theta_k)}_{\text{(ii)}}
 \end{aligned}$$

We solved it in Session (4-a)

$$= \text{Dir}(\theta \mid [\sum_{n: z_n=k} x_{n,1} + \gamma_1, \dots, \sum_{n: z_n=k} x_{n,V} + \gamma_V]) \quad (\text{iii})$$

$$P(\pi | \theta_1, \dots, \theta_k, z_n, x, h) = P(\pi | \alpha) \prod_{n=1}^N P(z_n | \pi)$$

$$\int_{\Pi} \frac{P(\pi | \alpha) \prod_{n=1}^N P(z_n | \pi)}{\int_{\Pi} P(\pi | \alpha) \prod_{n=1}^N P(z_n | \pi) d\pi} d\pi$$

Same as above

$$= \text{Dir}(\pi | [\sum_{n=1}^N z_{n,1} + \alpha_1, \dots, \sum_{n=1}^N z_{n,k} + \alpha_k])$$

—(iii)

Algorithm:

Initialize  $\theta_1^{(0)}, \dots, \theta_k^{(0)}, \pi^{(0)}$   
 draw samples  $z_N^{(m)}, \theta_1^{(m)}, \dots, \theta_k^{(m)}, \pi^{(m)}$

(i), (ii), (iii) would be  
 very easy to write  
 once we cover  
 graph representations)

Use  $\tilde{M}$  samples to compute  $\tilde{A}$

From the samples,

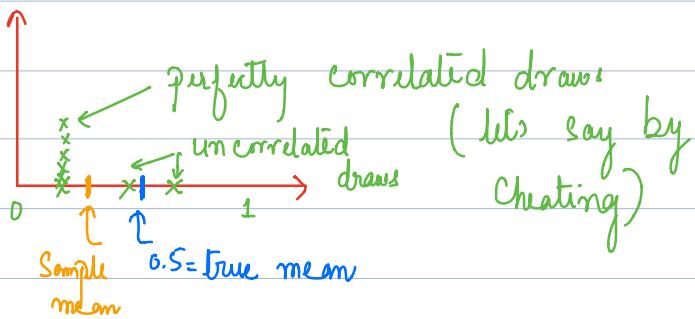
- \* average value for any distribution can be computed
- \* Samples from any marginal distribution  $P(l_n)$  can be obtained by only keeping samples of  $l_n$

\* Sampling in this fashion approximates draws

from the full-joint distribution.

## Practical Challenges :-

- \* (A) only works with independent draws from a distribution.



Solution:- \* measure empirical correlation among samples across various intervals.

\* choose an interval among samples with low empirical correlation.

\* Initial draws may be coming from a poor approximation of joint-distribution.

\* Solution:- ignore initial draws.

## Collapsing

$$P(z_n | z_{-n}, \theta_1, \dots, \theta_k, \pi, h, x)$$

$$= \frac{P(x_n | \theta_{z_n}) P(z_n | \pi)}{\sum_{z_n} P(x_n | \theta_{z_n}). P(z_n | \pi)} \quad \left\{ \text{from (ii)} \right\}$$

If we know  $z_{-n}$  (cluster-assignments) of all but nth point), do we need  $\theta_1, \dots, \theta_k$ , and  $\pi$ ?

Or

$$P(z_n | z_{-n}, x, h) = \frac{P(x | z_n, z_{-n}, h). P(z_n | z_{-n}, h)}{\sum_{z_n} P(x | z_n, z_{-n}, h). P(z_n | z_{-n}, h)}$$

$$\begin{aligned} P(z_n | z_{-n}, h) &= \int P(z_n | z_{-n}, \pi, \alpha). P(\pi | z_{-n}, \alpha) d\pi \\ &= \int \prod_{k=1}^K (\pi_k)^{z_{n,k}} P(\pi | z_{-n}, \alpha) d\pi \end{aligned}$$

$$P(\pi | z_{-n}, \alpha) = \frac{\Gamma\left(\sum_{k=1}^K \left(\sum_{m \neq n} z_{m,k} + \alpha_k\right)\right)}{\prod_k \Gamma\left(\sum_{m \neq n} z_{m,k} + \alpha_k\right)} \prod_k \pi_k^{\sum_{m \neq n} z_{m,k} + \alpha_k - 1}$$

$$= \frac{\Gamma(N-1 + \sum_k \alpha_k)}{\prod_k \Gamma(\sum_{m \neq n} z_{m,k} + \alpha_k)} \cdot \prod_k \frac{\sum_{m \neq n} z_{m,k} + \alpha_k - 1}{\prod_k (\prod_k)^{\sum_{m \neq n} z_{m,k} + \alpha_k}}$$

$$\frac{\Gamma(N-1 + \sum_k \alpha_k)}{\prod_k \Gamma(\sum_{m \neq n} z_{m,k} + \alpha_k)} \int \prod_k \prod_k^{\left( \sum_m z_{m,k} + \alpha_k - 1 \right)} d\pi$$

$$= \frac{\Gamma(N-1 + \sum_k \alpha_k)}{\prod_k \Gamma(\sum_{m \neq n} z_{m,k} + \alpha_k)} \cdot \frac{\prod_k \Gamma(\sum_m z_{m,k} + \alpha_k)}{\Gamma(N + \sum_k \alpha_k)}$$

using  $\Gamma(n+1) = n\Gamma(n)$

$$= \frac{\Gamma(N-1 + \sum \alpha_n)}{(N-1 + \sum_k \alpha_k) \Gamma(N-1 + \sum_k \alpha_k)} \prod_k \left[ \frac{\Gamma(\sum_m z_{m,k} + \alpha_k)}{\Gamma(\sum_{m \neq n} z_{m,k} + \alpha_k)} \right]$$

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

if  $z_n = [0 \cdots \underset{l}{1} \cdots 0]$

Similarly  $P(X | z_n, z_{-n}, h)$  (remaining).

Comment :-

(i) Results in very easy, numerically stable  
Simpler

(ii) One of the major ways, non-Parametric  
models are computed.