

## Dirichlet Process

1. Bayesian mixture model with K components

$$\Rightarrow \begin{cases} \pi | \alpha \sim \text{Dir}(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}) \\ z_i | \pi \sim \text{Mult}(\pi) \end{cases} \quad \theta_k^* \sim H \quad \left[ z_i | z_i, \{\theta_k^*\}_{k=1}^K \sim F(\theta_{z_i}^*) \right]$$

2. Dirichlet Process ( $G \sim DP(\alpha, H)$ )

- : distribution over the probability measure on  $\Theta$  whose marginals are dirichlet
- $\Rightarrow (G(A_1), \dots, G(A_K)) \sim \text{Dir}(\alpha H(A_1), \dots, \alpha H(A_K))$  for every finite partitions of  $\Theta$
- $\Rightarrow \forall A, E[G(A)] = H(A) \rightarrow H$  is the mean of DP
- $\Rightarrow \forall A, \text{Var}[G(A)] = H(A)(1-H(A))/(d+1) \rightarrow d$  is the inverse variance of DP
- (As  $d \rightarrow \infty$ ,  $G \rightarrow H$  weakly or pointwise)

① Existence

$$\Rightarrow \text{Exchangeability: } p(\theta_1, \dots, \theta_n) = \prod_{i=1}^n p(\theta_i | \theta_1, \dots, \theta_{i-1}) = p(\theta_{(1)}, \dots, \theta_{(n)}) \rightarrow n \text{ may be } \infty$$

$$\Rightarrow \text{De Finetti's theorem: } \exists G \text{ s.t. } p(\theta_1, \dots, \theta_n) = \int p(G) \prod_{i=1}^n G(\theta_i) dG$$

: for infinitely exchangeable sequence, there is a random distribution  $G$  s.t. the sequence is composed of i.i.d. draws from it.

② Posterior distribution ( $G | \theta_1, \dots, \theta_n \sim DP(\alpha+n, \frac{\alpha}{\alpha+n}H + \frac{n}{\alpha+n} \sum_{i=1}^n \delta_{\theta_i})$ )

$$\Rightarrow \theta_1, \dots, \theta_n \sim \text{Mult}(G(A_1), \dots, G(A_K)) \text{ where } A_1, \dots, A_K \subset \Theta, n_k = |\{i : \theta_i \in A_k\}|$$

$$\Rightarrow (G(A_1), \dots, G(A_K)) | \theta_1, \dots, \theta_n \sim \text{Dir}(\alpha H(A_1) + n_1, \dots, \alpha H(A_K) + n_K)$$

$\Rightarrow$  As  $\alpha \rightarrow 0$ , prior  $H$  becomes non-informative

$\Rightarrow$  As  $n \rightarrow \infty$ , posterior  $G | \theta_1, \dots, \theta_n$  is dominated by the empirical distribution

③ Predictive distribution ( $\theta_{n+1} | \theta_1, \dots, \theta_n \sim \frac{\alpha}{\alpha+n}H + \frac{n}{\alpha+n} \sum_{i=1}^n \delta_{\theta_i}$ )

$$\begin{aligned} \Rightarrow \forall A, p(\theta_{n+1} \in A | \theta_1, \dots, \theta_n) &= \int p(G | \theta_1, \dots, \theta_n) p(\theta_{n+1} \in A | G, \theta_1, \dots, \theta_n) dG \\ &= \int p(G | \theta_1, \dots, \theta_n) G(A) dG \\ &= E[G(A) | \theta_1, \dots, \theta_n] = \frac{\alpha}{\alpha+n} H(A) + \frac{n}{\alpha+n} \cdot \frac{\sum_{i=1}^n \delta_{\theta_i}(A)}{n} \end{aligned}$$

$$\Rightarrow \theta_{n+1} | \theta_1, \dots, \theta_n \sim \frac{\alpha}{\alpha+n}H + \frac{n}{\alpha+n} \sum_{i=1}^n \delta_{\theta_i}$$

(Base function of the posterior DP is the predictive distribution)

$$\Rightarrow \theta_{n+1} | \theta_1, \dots, \theta_n \sim \frac{\alpha}{\alpha+n}H + \frac{n}{\alpha+n} \sum_{i=1}^n \frac{\delta_{\theta_i}}{n} \text{ (rich-gets-richer phenomenon)}$$

$\rightarrow (\theta_1^*, \dots, \theta_n^*)$  are unique values in  $(\theta_1, \dots, \theta_n) \rightarrow$  clustering property

$\rightarrow n_k$  is the number of repeats of  $\theta_k^*$

$$\rightarrow \sum_{i=1}^n \frac{\alpha}{\alpha+i} \approx \alpha \log(1 + \frac{n}{\alpha}) = O(\alpha \log n) \rightarrow \text{expectation}(\# \text{ of unique values})$$

④ Chinese Restaurant Process ( $\approx$  Polya urn scheme)

: customers sit on an infinite number of round tables

: balls are colored among an infinite number of colours and put inside an urn

$\Rightarrow$  express an uncertainty about possible number of components

(pick a new table ( $\approx$  colour) with probability  $\frac{\alpha}{\alpha+n}$ )

(pick an old table ( $\approx$  colour) with probability  $\frac{n}{\alpha+n}$ )

$\rightarrow$  positive probability on picking the same value

$\rightarrow$  for a long enough sequence, any draw will be repeated by another draw

$\therefore G$  is composed of a weighted sum of point masses as discrete distribution

(support of  $G$  is countably infinite atoms drawn from  $H$ )

⑤ stick-breaking construction

$$\begin{cases} \beta_k \sim \text{Beta}(1, \alpha) \\ \pi_k \sim \beta_k \prod_{i=1}^{k-1} (1 - \beta_i) \end{cases} \rightarrow \pi \sim \text{GEM}(\alpha) \quad \theta_k^* \sim H \quad \left[ G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k^*} \right] \rightarrow \sum_{k=1}^{\infty} \pi_k = 1$$

$\rightarrow$  sample takes value  $\theta_k^*$  with probability  $\pi_k$

$\rightarrow \theta_k^*$  and  $\pi$  are mutually independent

⑥ Dirichlet Process Mixture Model

$$\begin{cases} \pi \sim \text{GEM}(\alpha) \\ z_i \sim \text{Mult}(\pi) \end{cases} \quad \theta_k^* \sim H \quad \left[ z_i | z_i, \{\theta_k^*\}_{k=1}^{\infty} \sim F(\theta_{z_i}^*) \right]$$

$\rightarrow$  infinite limit of finite mixture models

3. Hierarchical Dirichlet Process ( $\overset{\text{global}}{G_0} \sim DP(r, H)$ ,  $\overset{\text{local}}{G_j} \sim DP(\alpha, G_0)$  for  $j = 1, 2, \dots, J$ )

: clusters are shared across multiple, nested groupings of data

$\Rightarrow$  forces the base measure to be discrete as a draw from the dirichlet process

$\Rightarrow$  forces the atoms to be shared among the multiple dirichlet processes

$$(G_0 = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k^*} \text{ where } \theta_k^* \sim \text{GEM}(r))$$

$$(G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k^*} \text{ where } \pi_{jk} \sim \text{Dir}(1, \theta_k^*))$$

$$\rightarrow (G_j(A_1), \dots, G_j(A_r)) \sim \text{Dir}(\alpha_0 G_0(A_1), \dots, \alpha_0 G_0(A_r))$$

$$\rightarrow (\sum_{k=1}^{\infty} \pi_{jk}, \dots, \sum_{k=1}^{\infty} \pi_{jr}) \sim \text{Dir}(\alpha_0 \sum_{k=1}^{\infty} \pi_k, \dots, \alpha_0 \sum_{k=1}^{\infty} \pi_k)$$

① Chinese Restaurant Franchise

: customers visit a franchise with  $J$  restaurants and sit on an infinite number of round tables, each of which serves the same dishes

$\left( \psi_{ji} \sim G_j : \text{the dish for the } i\text{th customer heading to the } j\text{th restaurant} \right)$

$\left( \psi_{jt} \sim G_0 : \text{the dish for the } t\text{th table in the } j\text{th restaurant} \right)$

$\theta_k^* \sim H : k\text{th dish}$

$$(m_{jk} : \# \text{ of } \psi_{it}\text{'s associated with } \theta_k^* \rightarrow m_k = \sum_j m_{jk})$$

$$(n_{jt} : \# \text{ of } \psi_{it}\text{'s associated with } \psi_{jt})$$

$\Rightarrow$  the dish is picked according to its popularity in the whole restaurants

$$\rightarrow \psi_{j+1} | \psi_1, \psi_2, \dots, \psi_j, r, H \sim \frac{r}{r + \sum_k m_k} H + \sum_{k=1}^K \frac{m_k}{r + \sum_k m_k} \delta_{\theta_k^*}$$

$\Rightarrow$  the table is picked according to its popularity in the restaurant

$$\rightarrow \phi_{j+1} | \phi_1, \phi_2, \dots, \phi_j, \alpha, G_0 \sim \frac{\alpha}{\alpha + j} G_0 + \sum_{t=1}^T \frac{n_{jt}}{\alpha + j} \delta_{\psi_{jt}}$$

② stick-breaking construction

$$\begin{cases} \beta \sim \text{GEM}(r) \\ \pi_{jk} \sim \text{Beta}(\alpha_0, \alpha_0(1 - \sum_{l=1}^{k-1} \beta_l)) \\ \pi_{jk} = \pi_{jk}' \prod_{l=1}^{k-1} (1 - \beta_l) \end{cases} \quad \theta_k^* \sim H \quad \left[ G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k^*} \right]$$

$$\rightarrow (\sum_{k=1}^{\infty} \pi_{jk}', \pi_{jk}, \sum_{k=1}^{\infty} \pi_{jk}') \sim \text{Dir}(\alpha_0 \sum_{k=1}^{\infty} \pi_k, \alpha_0 \sum_{k=1}^{\infty} \pi_k)$$

$$\rightarrow \frac{1}{1 - \sum_{k=1}^{\infty} \pi_{jk}} (\pi_{jk}, \sum_{k=1}^{\infty} \pi_{jk}) \sim \text{Dir}(\alpha_0 \beta, \alpha_0 \sum_{k=1}^{\infty} \pi_k)$$

$$\therefore \pi_{jk} = \pi_{jk}' \prod_{l=1}^{k-1} (1 - \beta_l) \text{ where } \pi_{jk}' = \frac{\pi_{jk}}{1 - \sum_{k=1}^{\infty} \pi_{jk}} \sim \text{Beta}(\alpha_0, \alpha_0(1 - \sum_{k=1}^{k-1} \beta_l))$$

③ Hierarchical Dirichlet Process Mixture Model

$$\begin{cases} \beta \sim \text{GEM}(r) \\ \pi_{jj} \sim DP(\alpha, \beta) \\ z_{ji} | z_{ji}, \{\theta_k^*\}_{k=1}^{\infty} \sim F(\theta_{z_{ji}}^*) \end{cases} \quad \theta_k^* \sim H \quad \left[ z_{ji} | z_{ji}, \{\theta_k^*\}_{k=1}^{\infty} \sim F(\theta_{z_{ji}}^*) \right]$$

4. Nested Chinese Restaurant Process

: customers visit  $L$  distinct restaurants among an infinite number of them following the cards on the chosen tables

$\Rightarrow$  prior to be used when modeling topic hierarchies

$\Rightarrow$  express an uncertainty about possible  $L$ -level trees

## Topic Modeling

: the words in a document are generated according to a mixture model where the mixing proportions are random and document-specific.  
 (corpus → documents → topics → words)

### 1. Latent Dirichlet Allocation (LDA)

① For each document in a corpus

① sample a  $k$ -dimensional topic proportions :  $\theta \sim \text{Dir}(\alpha)$

② For each word in a document

① sample a topic :  $z \sim \text{Mult}(\theta)$

estimate using variational EM

② sample a word from the corresponding topic :  $w \sim p(w|z|\beta)$

### 2. Hierarchical LDA

→ corpus : tree, document : path, topic : restaurant, word : table

X Finite Depth

① For each document in a corpus

① Let  $c_1$  be the root restaurant

② For each level  $l = \{2, \dots, L\}$

① Pick a table in restaurant  $c_{l-1}$  as a way of CRP( $\gamma_{l-1}$ )

② Let  $c_l$  be the restaurant referred to by the table

$$\Rightarrow c \sim nCRP(\eta)$$

③ Sample an  $L$ -dimensional topic proportions :  $\theta \sim \text{Dir}(\alpha)$

④ For each word in a document

① sample a topic :  $z \sim \text{Mult}(\theta)$

② sample a word from the corresponding topic :  $w \sim \text{CRP}(\gamma_z)$

X Infinite Depth

① For each node in the tree

① sample a topic :  $\beta \sim \text{Dir}(\eta)$

② For each document in a corpus

① sample a path :  $c \sim nCRP(r)$

② Sample a level proportions :  $\theta \sim \text{GEM}(m, \pi)$

③ For each word in a document

① sample a level :  $z \sim \text{Mult}(\theta)$

② sample a word from the corresponding level :  $w \sim \text{Mult}(\beta_z)$

## Variational Inference for Dirichlet Process Mixtures

### 1. Assumptions

→ the observed data are drawn from the exponential family distribution

→ the base distribution for the DP is the corresponding conjugate prior

$$\rightarrow p(x_n | z_n, \theta^*, \theta_1^*, \dots) = \prod_{i=1}^n \left( h(x_i) \exp \{ \theta_i^* z_i - a(\theta_i^*) \} \right)^{\mathbb{I}[z_i=i]}$$

$$\rightarrow p(\theta^* | \lambda) = h(\theta^*) \exp \{ \lambda_1^* \theta^* + \lambda_2 (-a(\theta^*)) - a(\lambda) \}$$

### 2. Mean-field Variational inference

$$\Rightarrow \log p(x) \geq \mathbb{E}_q [\log p(w, x) - \log q_v(w)]$$

$$\rightarrow \text{latent distribution} : p(w_i | w_{-i}, x) = h(w_i) \exp \{ g_i(w_i, x)^T \omega_i - a(g_i(w_i, x)) \}$$

$$\rightarrow \text{variational distribution} : q_v(w) = \prod_{i=1}^n \exp \{ v_i^T \omega_i - a(v_i) \}$$

$$\star \text{ Solution} : v_i = \mathbb{E}_q [g_i(w_i, x)]$$

$$\Rightarrow \log p(x | \alpha, \lambda) \geq \mathbb{E}_q [\log p(V | x)] + \mathbb{E}_q [\log p(\theta^* | \lambda)]$$

$$+ \sum_{n=1}^N (\mathbb{E}_q [\log p(z_n | V)] + \mathbb{E}_q [\log p(x_n | z_n)])$$

$$- \mathbb{E}_q [\log q(V, \theta^*, z)]$$

→ latent distribution

→ stick length :  $v_i \sim \text{Beta}(1, \alpha)$

→ atoms :  $\theta^* \sim H$

→ cluster assignment :  $z_n \sim \text{Mult}(\pi(V))$

→ variational distribution → "Truncated" to the level  $T$

$$\Rightarrow q(V, \theta^*, z) = \prod_{t=1}^{T-1} q_{v_t}(v_t) \prod_{t=1}^T q_{\theta^*}(\theta_t^*) \prod_{n=1}^N q_{z_n}(z_n)$$

↓                    ↓                    ↓

Beta                    exponential                    multinomial