# Bayesian Non Parametrics Extensions

A/Prof Richard Yi Da Xu
Yida.Xu@uts.edu.au
https://github.com/roboticcam/machine-learning-notes

University of Technology Sydney (UTS)

May 1, 2018

- **Hierarchical Dirichlet Process (HDP)**
- HDP-Hidden Marko Model
- Indian Buffet Process

# Hierarchical Dirichlet Process (HDP)
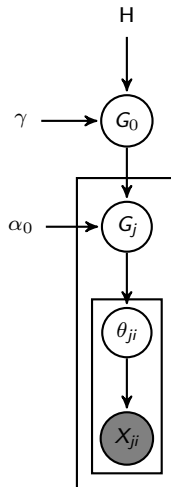
### Generative model

$$G_0 \sim \text{DP}(\gamma, H)$$
$$G_j \sim \text{DP}(\alpha_0, G_0)$$
$$\theta_{ji} \sim G_j$$
$$X_{ji} \sim F(x|\theta_{ij})$$

▶ Drawing $G_0 \sim \text{DP}(.)$ can be done using stick breaking process, i.e., $\sim \text{Beta}(1, \gamma)$.

▶ What about stick breaking construction for $G_j$?

▶ Certainly, it's NOT $\sim \text{Beta}(1, \alpha_0)$
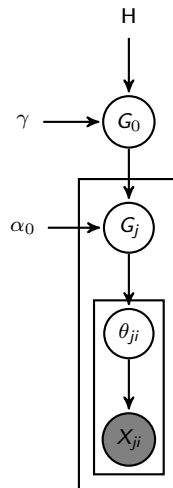
### Graphical model

## Generative model

$$\boldsymbol{\beta} \sim \text{GEM}(\gamma) \quad G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}$$

$$\boldsymbol{\pi}_j \sim \text{DP}(\alpha_0, \boldsymbol{\beta}) \qquad G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\phi_k}$$

$$z_{ji} \sim \pi_j \quad \phi_k \sim H \quad X_{ji} \sim F(x|\phi_{z_{ji}})$$

- ▶ Using $\boldsymbol{\beta}$ as a base, discrete distribution define on range $\{0 \ldots \infty\}$.

## Graphical model

# New Stick breaking for $\pi_{jk}$ using $\boldsymbol{\beta}$

▶ Dirichlet Process:

$$v_k \sim \text{Beta}(1, \alpha) \qquad \pi_k = v_k \prod_{l=1}^{k-1} (1 - v_l)$$

$$\theta_k \sim H \qquad G_0 = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}$$

▶ Hierarchical Dirichlet Process:

$$v_{jk} = \frac{\pi_k}{1 - \sum_{l=1}^{k-1} \pi_l} \sim \text{Beta}\left(\alpha\beta_k, 1 - \sum_{l=1}^{k} \beta_l\right) \qquad \pi_{jk} = v_{jk} \prod_{l=1}^{k-1} \left(1 - v_{jl}\right)$$

▶ In DP, each $v_k$ is distributed iid from $\text{Beta}(1\alpha)$
▶ In HDP, each $v_{jk}$ is distributed independently, but having different distribution

Suppose $\beta|\gamma \sim \text{GEM}(\gamma)$ and $\pi|\alpha, \beta \sim \text{DP}(\alpha, \beta)$. Notice that the support is $\{1, \ldots, k, \ldots, \infty\}$:

$$(G_j(A_1), \ldots, G_j(A_r)) \sim \text{Dir}\left(\alpha G_0(A_1), \ldots, \alpha G_0(A_r)\right)$$

$$\implies \left(\sum_{k \in K_1} u_k, \ldots, \sum_{k \in K_r} u_k\right) \sim \text{Dir}\left(\alpha \sum_{k \in K_1} \beta_k, \ldots, \alpha \sum_{k \in K_r} \beta_k\right)$$

$$\implies \left(\sum_{l=1}^{k-1} u_l, u_k, \sum_{l=k+1}^{\infty} u_l\right) \sim \text{Dir}\left(\alpha \sum_{l \in 1}^{k-1} \beta_l, \alpha\beta_k, \sum_{l=k+1}^{\infty} \beta_l\right)$$

$$\implies \left(\frac{u_k}{1 - \sum_{l=1}^{k-1} u_l}, \frac{\sum_{l=k+1}^{\infty} u_l}{1 - \sum_{l=1}^{k-1} u_l}\right) \sim \text{Dir}\left(\alpha\beta_k, \sum_{l=k+1}^{\infty} \beta_l\right) \quad \textbf{exercise} \text{ prove this}$$

$$\implies \left(\frac{u_k}{1 - \sum_{l=1}^{k-1} u_l}, \frac{\sum_{l=k+1}^{\infty} u_l}{1 - \sum_{l=1}^{k-1} u_l}\right) \sim \text{Dir}\left(\alpha\beta_k, 1 - \sum_{l=1}^{k} \beta_l\right)$$

$$\implies \left(v = \frac{u_k}{1 - \sum_{l=1}^{k-1} u_l}\right) \sim \text{Beta}\left(\alpha\beta_k, 1 - \sum_{l=1}^{k} \beta_l\right)$$

$$\left( \sum_{l=1}^{k-1} u_l, u_k, \sum_{l=k+1}^{\infty} u_l \right) \sim \text{Dir} \left( \alpha \sum_{l \in 1}^{k-1} \beta_l, \alpha \beta_k, \sum_{l=k+1}^{\infty} \beta_l \right)$$

$$\implies \left( \frac{u_k}{1 - \sum_{l=1}^{k-1} u_l}, \frac{\sum_{l=k+1}^{\infty} u_l}{1 - \sum_{l=1}^{k-1} u_l} \right) \sim \text{Dir} \left( \alpha \beta_k, \sum_{l=k+1}^{\infty} \beta_l \right)$$

Let $g_i \sim \text{Gamma}(\alpha_i, 1)$ for $i = 1, \ldots, n$:

$$\left( \frac{g_1}{\sum_{i=1}^{n} g_i}, \ldots, \frac{g_n}{\sum_{i=1}^{n} g_i} \right) \sim \text{DIR}(\alpha_1, \alpha_2, \ldots \alpha_n)$$

The following is also true:

$$\left( \frac{g_2}{\sum_{i=2}^{n} g_i}, \ldots, \frac{g_n}{\sum_{i=2}^{n} g_i} \right) \sim \text{Dirichlet}(\alpha_2, \ldots \alpha_n)$$

Look at a particular term:

$$\frac{g_j}{\sum_{i=2}^{n} g_i} = \frac{\frac{g_j}{\sum_{i=1}^{n} g_i}}{\frac{\sum_{i=2}^{n} g_i}{\sum_{i=1}^{n} g_i}} = \frac{\pi_j}{\frac{\left( \sum_{i=1}^{n} g_i \right) - g_1}{\sum_{i=1}^{n} g_i}} = \frac{\pi_j}{1 - \pi_1}$$

So we can write:

$$\left( \frac{\pi_2}{1 - \pi_1}, \ldots, \frac{\pi_n}{1 - \pi_1} \right) \sim \text{Dirichlet}(\alpha_2, \ldots \alpha_n)$$

- $x_{ji}$: $i^{\text{th}}$ customer at the $j^{\text{th}}$ restaurant.
- $N$ customers at each restaurant $j$.
- each customer $x_{ji}$ associates a table index $t_{ji} \in \{1, \ldots T\}$, $T << N$.
- each table $t_{ji}$ associates with a dish number $k_{jt} \in \{1, \ldots, K\}$, $K << T$.
- a **shorthand** notation $z_{ji} = k_{jt_{ji}}$: customer $x_{ji}$ has table number $t_{ji}$ which serve dish $k_{jt}$
- $m_{\cdot}$ is the count of all dish served.

▶ the equation is:

$$p(t_{ji} = t | \mathbf{t}^{-ji}, \mathbf{k}, x_{ji}) \propto \begin{cases} n_{jt.}^{-ji} f_{k_{ji}}^{\mathbf{x}^{-ji}}(x_{ji}) & \text{IF } t \text{ is previously used} \\ \alpha_0 p(x_{ji} | \mathbf{t}^{-ji}, t_{ji} = t^{\text{new}}, \mathbf{k}) & \text{IF } t = t^{\text{new}} \end{cases}$$

▶ when $t_{ji}$ is a **new table**, $x_{ji}$ should associate a new dish $k$.

▶ just like $f(x|k^{\text{new}}) = \int_{\phi} f(x|\phi)h(\phi)d\phi$, we also need to **integrate** out possible values of $k_{jt^{\text{new}}}$:

▶ However, this dish may be an existing or a **new** one in the entire franchise.

$$p(x_{ji} | \mathbf{x}^{-ji}, t_{jt} = t^{\text{new}}, \mathbf{k}) = \underbrace{\sum_{k=1}^{K} \frac{m_{\cdot k}}{m_{\cdot \cdot} + \gamma} f_k^{\mathbf{x}^{-ji}}(x_{ji})}_{\text{part 1}} + \underbrace{\frac{\gamma}{m_{\cdot \cdot} + \gamma} f_{k^{\text{new}}}^{\mathbf{x}^{-ji}}(x_{ji})}_{\text{part 2}}$$

1. **part 1**: $k_{jt_{ji}}$ is an **existing** dish in the franchise
2. **part 2**: $k_{jt_{ji}}$ is a **new** dish in the franchise

▶ **exercise** what is **after** a customer sits in a **new** table?

- this is to decide dish for all customers of the same **table** $k_{jt}$:

$$p(k_{jt} = k | \mathbf{k}^{-jt}, \mathbf{t}, \mathbf{x}_{jt}) \propto \begin{cases} m_{\cdot k}^{-jt} f_{\mathbf{x}_{jt}}^{\mathbf{x}_{-jt}}(\mathbf{x}_{jt}) & \text{IF } k \text{ is previously used} \\ \gamma f_{k^{\text{new}}}^{\mathbf{x}_{-jt}}(\mathbf{x}_{jt}) & \text{IF } k = k^{\text{new}} \end{cases}$$

where $\mathbf{x}_{-jt}$ is every customer of the same table $t$, and $x_{ji}$ is a single customer

- there is also a single person version:

$$p(k_{jt^{\text{new}}} = k | \mathbf{k}^{-ji}, \mathbf{t}, \mathbf{x}_{jt}) \propto \begin{cases} m_{\cdot k}^{-ji} f_{\mathbf{x}_{jt}}^{\mathbf{x}_{-jt}}(\mathbf{x}_{jt}) & \text{IF } k \text{ is previously used} \\ \gamma f_{k^{\text{new}}}^{\mathbf{x}_{-jt}}(\mathbf{x}_{jt}) & \text{IF } k = k^{\text{new}} \end{cases}$$

**exercise** think about when you may use this version?

# Likelihood function $f_{\mathbf{k}}^{\mathbf{x}^{-ji}}(x_{ji})$

▶ the likelihood function for $z_{ji} = k$, i.e., sitting on **existing** table

$$f_{\mathbf{k}}^{\mathbf{x}^{-ji}}(x_{ji}) = p(x_{ji}|\mathbf{x}_{-ji}, z_{jt} = \mathbf{k}, \mathbf{z}^{-ji})$$

$$= \int_{\phi_k} p(x_{ji}|\phi_k) p(\phi_k|\mathbf{x}_{-ji} = k) \mathrm{d}\phi_k$$

$$= \int_{\phi_k} p(x_{ji}|\phi_k) p(\mathbf{x}_{-ji} = k|\phi_k) p(\phi_k) \mathrm{d}\phi_k$$

$$\propto \int_{\phi_k} f(x_{ji}|\phi_k) \prod_{j' \neq j, i' \neq i, z_{j'i'} = k} f(x_{j'i'}|\phi_k) h(\phi_k) \mathrm{d}\phi_k$$

$$= \frac{\int_{\phi_k} f(x_{ji}|\phi_k) \prod_{j' \neq j, i' \neq i, z_{j'i'} = k} f(x_{j'i'}|\phi_k) h(\phi_k) \mathrm{d}\phi_k}{p(\mathbf{x}_{-ji}, z_{jt} = k, \mathbf{z}^{-ji})}$$

$$= \frac{\int_{\phi_k} f(x_{ji}|\phi_k) \prod_{j' \neq j, i' \neq i, z_{j'i'} = k} f(x_{j'i'}|\phi_k) h(\phi_k) d\phi_k}{\int_{\phi_k} \prod_{j' \neq j, i' \neq i, z_{j'i'} = k} f(x_{j'i'}|\phi_k) h(\phi_k) \mathrm{d}\phi_k}$$

▶ the likelihood function for $z_{ji} = $ new, i.e., sitting on **new** table:

$$f_{\mathbf{k}^{\mathrm{new}}}^{\mathbf{x}^{-ji}}(x_{ji}) = p(x_{ji}|\mathbf{x}_{-ji}, z_{jt} = \mathrm{new}, \mathbf{z}^{-ji})$$

$$= \int_{\phi} p(x_{ji}|\phi) p(\phi) \mathrm{d}\phi$$

- in previous sampling scheme, all groups are coupled since $G_0$ is integrated out.
- this is just like the DP case: $z_i|\mathbf{z}_{-1}$
- alternative sampling scheme is to have explicit $G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}$
- allow posterior conditioned on $G_0$ factorizes across groups.

- given $(\mathbf{t}, \mathbf{k})$, we can draw $G_0$ by noting:
  - $G_0 \sim \mathrm{DP}(\gamma, H)$
  - $\psi_{jt} \sim G_0$ for each table $t$
- this is just the posterior of DP we saw earlier:
  $$G' = G(A_1), \ldots, G(A_r) | \theta_1, \ldots, \theta_n \sim \mathrm{Dir}(\alpha H(A_1) + n_1, \ldots, \alpha H(A_k) + n_k)$$

  $$G_0 | \mathbf{t}, \mathbf{k}, \gamma, H, \{\psi_{jt}\} = \mathrm{DP}\left(\gamma + m_{..}, \frac{\gamma H + \sum_{k=1}^{K} m_{.k} \delta_{\phi_k}}{\gamma + m_{..}}\right)$$

- posterior of $G_0$ constructed from different elements:

  $$\boldsymbol{\beta} = (\beta_1, \ldots, \beta_K, \beta_u) \sim \mathrm{Dir}(m_{.1}, \ldots, m_{.K}, \gamma)$$

  $$p(\phi_k | \mathbf{t}, \mathbf{k}) \propto h(\phi_k) \prod_{ji:z_{ji}=k} f(x_{ji} | \phi_k)$$
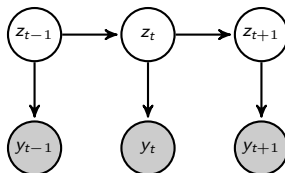
  $$G_u \sim \mathrm{DP}(\gamma, H)$$

  $$G_0 = \sum_{k=1}^{K} \beta_k \delta_{\phi_k} + \beta_u G_u$$

- when **new** component is instantiated:
  1. $b \sim \mathrm{Beta}(1, \gamma)$
  2. $K \leftarrow K + 1$
  3. $\beta_K = b\beta_u$
  4. $\beta_u \leftarrow (1 - b)\beta_u$

Under normal HMM, you have a transition matrix $A$, let the $j^{\text{th}}$ row of $A$ to be $\pi_i$, then:

$$A = \begin{bmatrix} \pi_1 \\ \pi_2 \\ \dots \\ \pi_K \end{bmatrix} = \begin{bmatrix} p(z_{t+1}=1|z_t=1) & p(z_{t+1}=2|z_t=1) & \dots & p(z_{t+1}=K|z_t=1) \\ p(z_{t+1}=1|z_t=2) & p(z_{t+1}=2|z_t=2) & \dots & p(z_{t+1}=K|z_t=2) \\ \dots & \dots & \dots & \dots \\ p(z_{t+1}=1|z_t=K) & p(z_{t+1}=2|z_t=K) & \dots & p(z_{t+1}=K|z_t=K) \end{bmatrix}$$



To obtain the current latent state, we need to sample $z_t \sim \text{Mult}(\pi_{z_{t-1}})$.

- Same idea has been extended to non-parametric bayes,
- Allow $\pi_j$ to have infinite many components.
- Matrix $A$ has size $\infty \times \infty$. But the "recovered" number of states are finite, so you only "jumping around" in the upper-left corner of matrix $A$.

$$
\begin{bmatrix}
p(z_{t+1}=1|z_t=1) & p(z_{t+1}=2|z_t=1) & \ldots & p(z_{t+1}=\infty|z_t=1) \\
p(z_{t+1}=1|z_t=2) & p(z_{t+1}=2|z_t=2) & \ldots & p(z_{t+1}=\infty|z_t=2) \\
\ldots & \ldots & \ldots & \ldots \\
p(z_{t+1}=1|z_t=\infty) & p(z_{t+1}=2|z_t=\infty) & \ldots & p(z_{t+1}=\infty|z_t=\infty)
\end{bmatrix}
$$

### Generative model

$$\boldsymbol{\beta} \sim \mathsf{GEM}(\gamma)$$
$$\boldsymbol{\pi}_j \sim \mathsf{DP}\left(\alpha, \boldsymbol{\beta}\right)$$
$$z_t \sim \mathsf{Mult}(\pi_{z_{t-1}})$$
$$\theta_k \sim H$$
$$X_t \sim F(x|\theta_{z_t})$$

### Graphical model

- let $t - 1 = 3, \mathbf{t} = \mathbf{4}, t + 1 = 5$
- $n_{ij}$ is the number of transitions from state $i$ to $j$ **excluding** time steps $t - 1$ and $t$:

$$
\begin{array}{llll}
n_{1,1} = 0 & n_{1,2} = 1 & n_{1,3} = 2 & \mathbf{n}_{1,:} = 3 \\
n_{2,1} = 1 & n_{2,2} = 1 & n_{2,3} = 0 & \mathbf{n}_{2,:} = 2 \\
n_{3,1} = 1 & n_{3,2} = 2 & n_{3,3} = 0 & \mathbf{n}_{3,:} = 3 \\
\mathbf{n}_{:,1} = 2 & \mathbf{n}_{:,2} = 4 & \mathbf{n}_{:,3} = 2 &
\end{array}
$$

- $\mathbf{n}_{:,k}$ is the number of transitions **INTO** state $k$
- $\mathbf{n}_{k,:}$ is the number of transitions **FROM** state $k$

$$\Pr(z_t = k | \mathbf{z}_{-t}) \propto \Pr\left(\{z_t = k | z_{t-1} = \mathbf{2}\}_{t=2:T}\right) \Pr\left(\{z_{t+1} = \mathbf{1} | z_t = k\}_{t=1:T-1}\right)$$

$$\Pr(z_t = 1 | \mathbf{z}_{-t}) \propto \Pr\left(\{z_t = 1 | z_{t-1} = \mathbf{2}\}_{t=2:T}\right) \Pr\left(\{z_{t+1} = \mathbf{1} | z_t = 1\}_{t=1:T-1}\right)$$

$$= \frac{n_{2,1}}{\mathbf{n}_{:,1}} \frac{n_{1,1}}{\mathbf{n}_{1,:}}$$

$$\Pr(z_t = 2 | \mathbf{z}_{-t}) \propto \Pr\left(\{z_t = 2 | z_{t-1} = \mathbf{2}\}_{t=2:T}\right) \Pr\left(\{z_{t+1} = \mathbf{1} | z_t = 2\}_{t=1:T-1}\right)$$

$$= \frac{n_{2,2}}{\mathbf{n}_{:,2}} \frac{n_{2,1}}{\mathbf{n}_{2,:} + 1} \quad \text{exercise why denominator increase by 1? What happens when } z_{t+1} = z_t$$

$$\Pr(z_t = 3 | \mathbf{z}_{-t}) \propto \Pr\left(\{z_t = 3 | z_{t-1} = \mathbf{2}\}_{t=2:T}\right) \Pr\left(\{z_{t+1} = \mathbf{1} | z_t = 3\}_{t=1:T-1}\right)$$

$$= \frac{n_{2,3}}{\mathbf{n}_{:,3}} \frac{n_{3,1}}{\mathbf{n}_{3,:}}$$

$$\Pr(z_t|z_{t-1}, \boldsymbol{\beta}, \mathbf{Y}, \alpha, H) \propto p(y_t|z_t, \mathbf{z}_{-t}, \mathbf{y}_{-t}, H) \underbrace{\Pr(z_t|\mathbf{z}_{-t}, \boldsymbol{\beta}, \alpha)}$$

$$\Pr(z_t = k|\mathbf{z}_{-t}, \boldsymbol{\beta}, \alpha) \propto \begin{cases} \left(\frac{n_{z_{t-1},k} + \alpha\beta_k}{\mathbf{n}_{:,k} + \alpha}\right)\left(\frac{n_{k,z_{t+1}} + \alpha\beta_{z_{t+1}}}{\mathbf{n}_{k,:} + \alpha}\right) & \text{if} \quad k \leq K, k \neq z_{t-1} \\ \left(\frac{n_{z_{t-1},k} + \alpha\beta_k}{\mathbf{n}_{:,k} + \alpha}\right)\left(\frac{n_{k,z_{t+1}} + 1 + \alpha\beta_{z_{t+1}}}{\mathbf{n}_{k,:} + 1 + \alpha}\right) & \text{if} \quad k = z_{t-1} = z_{t+1} \\ \left(\frac{n_{z_{t-1},k} + \alpha\beta_k}{\mathbf{n}_{:,k} + \alpha}\right)\left(\frac{n_{k,z_{t+1}} + \alpha\beta_{z_{t+1}}}{\mathbf{n}_{k,:} + 1 + \alpha}\right) & \text{if} \quad k = z_{t-1} \neq z_{t+1} \\ \alpha\beta_k\beta_{z_{t+1}} & \text{if} \quad k = K+1 \end{cases}$$

▶ note that the DP sampling $\Pr(z_t = k|\mathbf{z}_{-t}, \alpha) \propto \begin{cases} \frac{n_k + \alpha}{\cancel{}} & \text{if existing} \\ \frac{\alpha}{\cancel{n+\alpha}} & \text{if new} \end{cases}$ does not apply in

HDP-HMM, as $\mathbf{n}$ is not constant.

▶ also when $k =$ new, $\mathbf{n}_{k,:} = \mathbf{n}_{:,k} = n_{z_{t-1},k} = n_{k,z_{t+1}} = 0$

▶ in DP sampling $\mathbf{n} > 0$ and remain constant.

▶ Introduce auxiliary variables $u_1, \ldots u_t$:

$$u_t \sim \mathrm{U}(0, \pi_{z_{t-1}, z_t}) \implies p(u_t | \mathbf{z}, \boldsymbol{\pi}) = p(u_t | z_{t-1}, z_t, \boldsymbol{\pi})$$

▶ Another way of writing it:

$$p(u_t | z_{t-1}, z_t, \boldsymbol{\pi}) = \frac{\mathbb{I}\left(0 < u_t < \pi_{z_{t-1}, z_t}\right)}{\pi_{z_{t-1}, z_t}}$$

$$
\begin{aligned}
p(z_t | y_{1:t}, u_{1:t}) &\propto p(z_t, u_t, y_t | y_{1:t-1}, u_{1:t-1}) \\
&= \sum_{z_{t-1}} p(z_t, u_t, y_t, z_{t-1} | y_{1:t-1}, u_{1:t-1}) \\
&= \sum_{z_{t-1}} p(y_t | z_t) \underbrace{p(u_t | z_t, z_{t-1})} p(z_t | z_{t-1}) p(z_{t-1} | y_{1:t-1}, u_{1:t-1}) \\
&= p(y_t | z_t) \sum_{z_{t-1}} \underbrace{\frac{\mathbb{I}\left(0 < u_t < \pi_{z_{t-1}, z_t}\right)}{\pi_{z_{t-1}, z_t}}} p(z_t | z_{t-1}) p(z_{t-1} | y_{1:t-1}, u_{1:t-1}) \\
&= p(y_t | z_t) \sum_{z_{t-1}} \mathbb{I}\left(u_t < \pi_{z_{t-1}, z_t}\right) p(z_{t-1} | y_{1:t-1}, u_{1:t-1})
\end{aligned}
$$

► **forward pass**:

$$\Pr(z_t|y_{1:t}, u_{1:t}) \propto \Pr(z_t, u_t, y_t|y_{1:t-1}, u_{1:t-1})$$

$$= \Pr(y_t|z_t) \sum_{z_{t-1}} \mathbb{I}\left(u_t < \pi_{z_{t-1}, z_t}\right) \Pr(z_{t-1}|y_{1:t-1}, u_{1:t-1})$$

$$= \Pr(y_t|z_t) \sum_{\{z_{t-1}\}_{u_t < \pi_{z_{t-1}, z_t}}} \Pr(z_{t-1}|y_{1:t-1}, u_{1:t-1})$$

$u_t$ truncates the above summation to **finitely many** $z_{t-1}$s that satisfy both constraints:

1. $u_t < \pi_{z_{t-1}, z_t}$
2. $\Pr(z_{t-1}|y_{1:t-1}, u_{1:t-1}) > 0$

► To sample the whole trajectory $z_{1:t}$:

1. Sample $\mathbf{z_T} \sim \Pr(z_T|y_{1:T}, u_{1:T})$ - which is used in the "likelihood" function for $z_{T-1}$:
2. then, perform a **backward pass**, where we sample:

$$z_t|z_{t+1} : \Pr(z_t|z_{t+1}, y_{1:T}, u_{1:T}) \propto \Pr(\mathbf{z_{t+1}}|z_t, u_{t+1}) \Pr(z_t|y_{1:t}, u_{1:t})$$

## Generative model

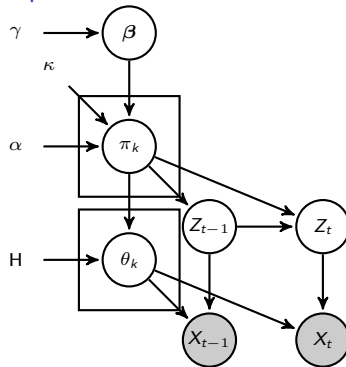$$\boldsymbol{\beta} \sim \text{GEM}(\gamma)$$

$$\boldsymbol{\pi}_j \sim \text{DP}\left(\alpha + \kappa, \frac{\alpha\boldsymbol{\beta} + \kappa\delta_j}{\alpha + \kappa}\right)$$

$$z_t \sim \text{Mult}(\pi_{z_{t-1}})$$

$$\theta_k \sim H$$

$$X_t \sim F(x|\theta_{z_t})$$

## Graphical model

- Hierarchical Dirichlet Process (HDP)
- HDP-Hidden Marko Model
- **Indian Buffet Process**

## DP

- $\Pr(z_1 \ldots z_N)$, where $z_i \in (1 \ldots K)$ indicate category.
- You also want $K$ potentially be infinite
- A "clustering" property, controllable through a single parameter $\alpha$
- Can also be thought as a special $N \times K$ $Z$ matrix, where there is only one "1" in each row.

## IBP

- More general than DP: $z_i$ can take multiple values $\in (1, \ldots K)$
- This is equivalently of saying that, $z_i$ is a binary vector of $K$ elements.
- Given $N$ such data, we have a binary matrix of size $N \times K$
- A "clustering" property, controllable through a single parameter $\alpha$, a column with more $1$, results it to have more 1s.

An example of $Z$ matrix:

| 1 | 0 | 1 | 1 | 0 | ... | 1 |
|---|---|---|---|---|-----|---|
| 0 | 1 | 0 | 0 | 0 | ... | 0 |
| ... | ... | ... | ... | ... | ... | 0 |
| 1 | 1 | 0 | 0 | 0 | ... | 0 |

For each column: $Pr(z_{ik} = 1) \sim \text{Ber}(\mu_k)$ independently.
Each $u_k \sim \text{Beta}\left(\frac{\alpha}{k}, 1\right)$ is also distributed independently.
The marginal distribution:

**Multinomial-Dirichlet**

$P(p_1, \ldots, p_k | n_1, \ldots, n_k)$

$\propto \underbrace{\frac{\Gamma\left(\sum_{i=1}^{k} \alpha_i\right)}{\prod_{i=1}^{k} \Gamma(\alpha_i)} \prod_{i=1}^{k} p_i^{\alpha_i - 1}}_{\text{Dir}(p_1, \ldots, p_k | \alpha_1, \ldots, \alpha_k)} \underbrace{\frac{n!}{n_1! \ldots n_k!} \prod_{i=1}^{k} p_i^{n_i}}_{\text{Mult}(n_1, \ldots, n_k | p_1, \ldots p_k)}$

$\propto \prod_{i=1}^{k} p_i^{\alpha_i - 1} \prod_{i=1}^{k} p_i^{n_i} = \prod_{i=1}^{k} p_i^{\alpha_i - 1 + n_i}$

$= \text{Dir}(p_1, \ldots p_k | \alpha_i + n_i, \ldots \alpha_k + n_k)$

**Bernoulli-Binomial**

$P(p | n_1 = m)$

$\propto \underbrace{\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha - 1}(1 - p)^{\beta - 1}}_{\text{Beta}(p | \alpha, \beta)} \underbrace{\frac{N!}{m!(N - m)!} p^k (1 - p)^{N-k}}_{\text{Binomial}(n_1, n_2 | p)}$

$\propto p^{\alpha - 1}(1 - p)^{\beta - 1} p^k (1 - p)^{N-k}$

$= p^{\alpha - 1 + k}(1 - p)^{\beta - 1 + N - k}$

$= \text{Beta}(p | \alpha + k, \beta + N - k)$

**Multinomial-Dirichlet**

$$\int_{p_1,\ldots,p_k} P(p_1,\ldots,p_k,n_1,\ldots,n_k)$$

$$= \frac{N!}{n_1!\ldots n_k!} \frac{\Gamma\left(\sum_{i=1}^k \alpha_i\right)}{\prod_{i=1}^k \Gamma(\alpha_i)} \frac{\prod_{i=1}^k \Gamma(\alpha_i+n_i)}{\Gamma\left(N+\sum_{i=1}^k \alpha_i\right)}$$

**Bernoulli-Beta**

$$\int_p P(p,n_1,n_2)$$

$$= \frac{N!}{k!(N-k)!} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha+k)\Gamma(\beta+N-k)}{\Gamma(N+\alpha+\beta)}$$

$\mu_k \sim \text{Beta}\left(\frac{\alpha}{k}, 1\right) \qquad \Pr(z_{ik} = 1) \sim \text{Ber}(\mu_k).$

$n_{k,-i}$ is the number of 1s of $k^{\text{th}}$ column, above row $i$.

Let $\alpha_i = \frac{\alpha}{k}$: compute the density of $i^{\text{th}}$ data belonging to existing component $m$.

$$\Pr(z_{ik} = 1 | \mathbf{z}_{-i,k}) = \int_p \Pr(z_{ik} = 1 | p) P(p | \underbrace{n_{-i,k}}_{n_1}, \underbrace{i - 1 - n_{-i,k}}_{n_2})$$

$$= \frac{\int_p \Pr(z_{ik} = 1 | p) \Pr(n_1, n_2 | p) P(p)}{\Pr(n_1, n_2)} = \frac{\int_p \Pr(z_{ik} = 1 | p) \Pr(n_1, n_2 | p) P(p)}{\int_p \Pr(n_{-i,k}, i - 1 - n_{-i,k} | p) P(p)}$$

$$= \frac{\Gamma(\frac{\alpha}{k} + n_{-i,k} + 1) \Gamma(1 + i - 1 - n_{-i,k})}{\Gamma\left(i + \frac{\alpha}{k} + 1\right)} \frac{\Gamma\left(i - 1 + \frac{\alpha}{k} + 1\right)}{\Gamma(\frac{\alpha}{k} + n_{-i,k}) \Gamma(1 + i - 1 - n_{-i,k})} = \frac{\frac{\alpha}{k} + n_{-i,k}}{i + \frac{\alpha}{k}}$$

$$\exp(x) = \lim_{n \to \infty} \left(1 + \frac{x}{n}\right)^n$$

Let $\lambda = np$:

$$
\begin{aligned}
\text{Binomial}(x|n, p) &= \binom{n}{x} p^x (1-p)^{n-x} = \frac{n!}{x!(n-x)!} \frac{\lambda^x}{n} (1 - \frac{\lambda}{n})^{n-x} \\
&= \underbrace{\frac{\lambda^x}{x!}}_{\text{constant}} \underbrace{\frac{n!}{(n-x)!} \frac{1}{n^x}}_{} \left(1 - \frac{\lambda}{n}\right)^{n-x} \\
&= \frac{\lambda^x}{x!} \frac{\overbrace{n(n-1), \dots (n-x+1)}^{n \text{ terms}}}{n^x} \left(1 - \frac{\lambda}{n}\right)^{n-x} \\
&= \frac{\lambda^x}{x!} \frac{n}{n} \frac{n-1}{n} \dots \frac{n-x+1}{n} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \\
&= \frac{\lambda^x}{x!} 1 \left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{x+1}{n}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}
\end{aligned}
$$

$$
\lim_{n \to \infty} \text{Binomial}(x|n, p) = \lim_{n \to \infty} \binom{n}{x} p^x (1-p)^{n-x}
$$
$$
= \frac{\lambda^x}{x!} \lim_{n \to \infty} \left(1 - \frac{1}{n}\right) \dots \lim_{n \to \infty} \left(1 - \frac{x+1}{n}\right) \lim_{n \to \infty} \left(1 - \frac{\lambda}{n}\right)^n \lim_{n \to \infty} \left(1 - \frac{\lambda}{n}\right)^{-x} = \frac{\lambda^x}{x!} \exp(-\lambda)
$$

$$\lim_{k \to \infty} \Pr(z_{ik}) = \lim_{k \to \infty} \frac{\frac{\alpha}{k} + n_{-i,k}}{i + \frac{\alpha}{k}} = \frac{n_{-i,k}}{i}$$

$$\lim_{n \to \infty} \text{Binomial}(\frac{\lambda}{n}, n) = \text{Poisson}(\lambda)$$

$$\text{Let } k \to \infty : \qquad\qquad = \frac{n_{-i,k}}{i}$$

For "new" dishes, i.e., $n_{-i,k} = 0$, then, $\Pr(z_{ik} = 1) = \text{Bernoulli}\left(\frac{\frac{\alpha}{K}}{i + \frac{\alpha}{K}}\right)$

i.e., how many new dishes across all columns would be: $\text{Binomial}\left(\frac{\frac{\alpha}{K}}{i + \frac{\alpha}{K}}, K\right)$

Since $\frac{\frac{\alpha}{k}}{i + \frac{\alpha}{k}} \times k = \frac{\alpha}{i + \frac{\alpha}{k}}$, we have:

$$\lim_{K \to \infty} \text{Binomial}\left(\frac{\frac{\alpha}{K}}{i + \frac{\alpha}{K}}, K\right) = \text{Poisson}\left(\frac{\alpha}{i}\right)$$

So, how many $K^+$ columns there are?

Let $n_i \sim \text{Poisson}\left(\frac{\alpha}{i}\right)$ $\qquad\qquad\qquad \left(\sum_{i=1}^{N} n_i\right) \sim \text{Poisson}\left(\sum_{i=1}^{N} \frac{\alpha}{i}\right)$

**What is Factor Analysis?** There are $N = 1000$ students, each having $(p = 10)$ scores. Therefore:

$$
\begin{bmatrix}
y_{11} & y_{12} & \cdots & y_{1N} \\
y_{21} & y_{22} & \cdots & y_{2N} \\
\cdots & \cdots & \cdots & \cdots \\
y_{p1} & y_{p2} & \cdots & y_{pN}
\end{bmatrix}
=
\begin{bmatrix}
g_{11} & \cdots & g_{1k} \\
g_{21} & \cdots & g_{2k} \\
\cdots & \cdots & \cdots \\
g_{p1} & \cdots & g_{pk}
\end{bmatrix}
\begin{bmatrix}
x_{11} & x_{12} & \cdots & x_{1N} \\
\cdots & \cdots & \cdots & \cdots \\
x_{k1} & x_{k2} & \cdots & x_{kN}
\end{bmatrix}
+ \mathbf{E}
$$

$$
\mathbf{E} =
\begin{bmatrix}
e_{11} & e_{12} & \cdots & e_{1N} \\
e_{21} & e_{22} & \cdots & e_{2N} \\
\cdots & \cdots & \cdots & \cdots \\
e_{p1} & e_{p2} & \cdots & e_{pN}
\end{bmatrix}
\text{ and } k << p
$$

Or in a matrix form: $\mathbf{Y} = \mathbf{G}\mathbf{X} + \mathbf{E}$.
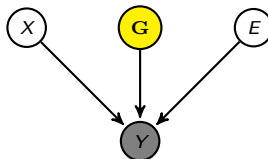
What this means is that a person's $i$'s raw mark is interpretted as:

$$\begin{bmatrix} y_{1i} \\ y_{2i} \\ \ldots \\ y_{pi} \end{bmatrix} = x_{1i} \begin{bmatrix} g_{11} \\ g_{21} \\ \ldots \\ g_{p1} \end{bmatrix} + x_{2i} \begin{bmatrix} g_{11} \\ g_{21} \\ \ldots \\ g_{p1} \end{bmatrix} + \ldots x_{ki} \begin{bmatrix} g_{1k} \\ g_{2k} \\ \ldots \\ g_{pk} \end{bmatrix} + \begin{bmatrix} e_{1i} \\ e_{2i} \\ \ldots \\ e_{pi} \end{bmatrix}$$

▶ Given a set of $k$ loading factors (vectors) each with dimension $p$: $\{\mathbf{g}_{:,i}\}_{i=1}^k$, the $x_{:,i}$ can be thought as the latent linear weights.

▶ Of course, you are only given data matrix $Y$, one has to infer the latent structure. $\mathbf{G}$, $\mathbf{X}$ and $\mathbf{E}$. Ths is not as silly as it seems, as DoF is much reduced.

**The Bayesian Treatment:**

$$e_i \sim \mathcal{N}(0, \sigma_e^2 \mathbf{I}) \qquad \sigma_e^2 \sim \mathcal{IG}(a, b)$$
$$g_k \sim \mathcal{N}(0, \sigma_G^2) \qquad \sigma_G^2 \sim \mathcal{IG}(c, d)$$
$$x_{ki} \sim \mathcal{N}(0, 1) \qquad y_i = \mathbf{G} x_i + e_i$$
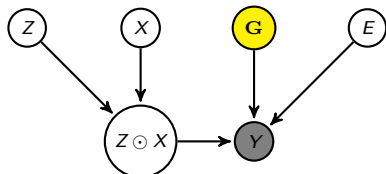
- ▶ Knowles, d and Ghahramani, Z, Infinite Sparse Factor Analysis
- ▶ $K$ should known beforehand. What about making $K$ a variable?
- ▶ Although $[x_{1,i}, \ldots x_{k,i}]^T$ has a reduced dimension, it can still cause "overfitting".
- ▶ We need to introcuce variable number of latent factors $K$, at the same time, have **sparsity**!

How?

$$e_i \sim \mathcal{N}(0, \sigma_e^2 \mathbf{I}) \qquad \sigma_e^2 \sim \mathcal{IG}(a, b)$$
$$g_k \sim \mathcal{N}(0, \sigma_G^2) \qquad \sigma_G^2 \sim \mathcal{IG}(c, d)$$
$$Z \sim \mathcal{IBP}(\alpha) \qquad \alpha \sim \mathcal{G}(e, f)$$
$$x_{ki} \sim \mathcal{N}(0, 1) \qquad y_i = \mathbf{G}(x_i \odot z_i) + e_i$$

# A proposed work

▶ What about if there are two sets of data matrix $\mathbf{Y}$ and $\mathbf{Y}'$, each having different number of entries. They share the same loading vectors $\mathbf{G}$, but with different level of **sparsities**.

$$e_i \sim \mathcal{N}(0, \sigma_e^2 \mathbf{I}) \qquad \sigma_e^2 \sim \mathcal{IG}(a, b)$$
$$g_k \sim \mathcal{N}(0, \sigma_G^2) \qquad \sigma_G^2 \sim \mathcal{IG}(c, d)$$
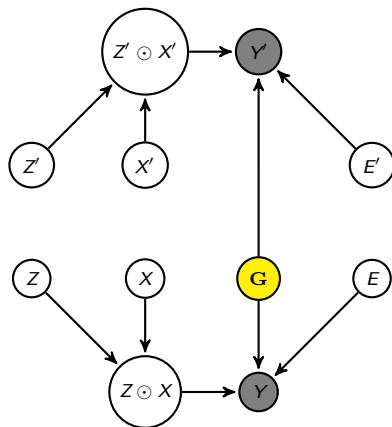$$Z \sim \mathcal{IBP}(\alpha) \qquad \alpha \sim \mathcal{G}(e, f)$$
$$x_{ki} \sim \mathcal{N}(0, 1) \qquad y_i = \mathbf{G}(x_i \odot z_i) + e_i$$

$$e_i' \sim \mathcal{N}(0, \sigma_e^2 \mathbf{I}) \qquad \sigma_e^2 \sim \mathcal{IG}(a', b')$$
$$g_k \sim \mathcal{N}(0, \sigma_G^2) \qquad \sigma_G^2 \sim \mathcal{IG}(c, d)$$
$$Z' \sim \mathcal{IBP}(\alpha') \qquad \alpha' \sim \mathcal{G}(e', f')$$
$$x_{ki}' \sim \mathcal{N}(0, 1) \qquad y_i' = \mathbf{G}(x_i' \odot z_i') + e_i'$$