

# Lesson 5: Non parametric Bayes and application to relational model

Richard Yi Da Xu

School of Computing & Communication, UTS

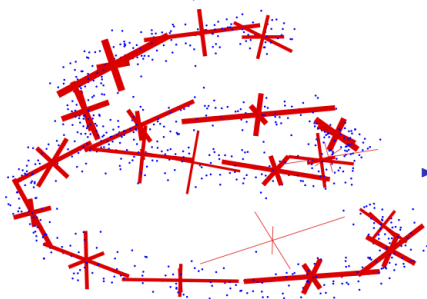
June 10, 2015

# Non parametric Bayes - Some background knowledge

- ▶ Getting Harder: introduction to Dirichlet Process - A machine learning research topic
- ▶ Also see MCMC used in practice
- ▶ Get the chance to study a real modelling application - Relational Model

# Dirichlet Process: A diagrammatic representation

Rasmussen, Infinite Gaussian Mixture Model (1999):



- ▶ For a mixture model:  
Let  $\mathbf{X} = x_1, \dots, x_N$ :

$$P(\mathbf{X}|\theta_1, \dots, \theta_K, w_1, \dots, w_K) = \sum_{l=1}^K w_l f(\mathbf{X}|\theta_l)$$

$$\text{where } \sum_{l=1}^K w_l = 1$$

- ▶ If we allow  $K$  to also vary, what happens if you want to:

$$\arg \max_{\theta_1, \dots, \theta_K, w_1, \dots, w_K, K} P(\mathbf{X}|\theta_1, \dots, \theta_K, w_1, \dots, w_K, K)?$$

- ▶  $K = N$  for Gaussian case. Of course it's not desirable!

- ▶ For data  $x_1, \dots, x_N$ , each  $x_i$  is associating with a parameter  $\theta_i$
- ▶ We need to a good prior for  $\Pr(\theta_1 \dots \theta_N)$ :
- ▶ You also want  $K$  potentially be infinite
- ▶ A “clustering” property, controllable through a single parameter  $\alpha$
- ▶ Let’s define it using Hierarchical prior, its marginal is:

$$p(\theta_1, \dots \theta_n) = \int_G \Pr(\theta_1, \dots, \theta_n | G) \mathbf{p}(\mathbf{G})$$

**So, we are interested in the property of  $\mathbf{G}$ :**

- ▶  $G$  needs to be **discrete** random distribution.
- ▶ Perhaps it should also some resemblance with some basic distribution  $H$ .

We say  $G$  is a Dirichlet process, distributed with base distribution  $H$  and concentration parameter  $\alpha$ :

$$G \sim DP(\alpha, H), \text{ if} \\ (G(A_1), \dots, G(A_r)) \sim \text{Dir}(\alpha H(A_1), \dots, \alpha H(A_r))$$

for every finite measurable partition  $A_1, \dots, A_r$  of  $\Theta$ .  
What does this all mean? Let's visualise it!

- ▶  $G$  is a Beta process, with base distribution  $H$  and concentration parameter  $\alpha$ :

$$G \sim BP(\alpha H), \text{ if}$$
$$G(A_k) \sim \text{Beta}(\alpha H(A_k), \alpha(1 - H(A_k)))$$

- ▶ Given an infinitesimal partition  $(A_1, \dots, A_K)$  with  $K \rightarrow \infty$  and  $H(A_k) \rightarrow 0$  the samples correspond to the density function:

$$G = \sum \pi_i \delta_{\theta_i}$$

where

$$\pi_i \sim \text{Beta}(0, \alpha)$$
$$\theta_i \sim H$$

- ▶ Beta process is a Completely Random Measure with Levy measure on product space  $[0, 1] \times \Omega$  with Levy measure:

$$\nu(d\pi d\theta) = \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi H(d\theta).$$

- ▶  $\Gamma$  is a Gamma process, distributed with base distribution  $H$  and concentration parameter  $\alpha$ :
- ▶ Given a partition  $(A_1, \dots, A_K), A_i \in \Omega \implies G(A_i) \sim \text{Gamma}(H(A_i), 1/\alpha)$
- ▶ Let  $\Gamma = \{(\pi_i, \theta_i)\}_{i=1}^\infty$  be a realization of a Gamma process in product space  $\mathbb{R}^+ \times \Theta$ :

$$\begin{aligned}\Gamma &\sim \text{GaP}(\alpha, H) \\ &= \sum \pi_i \delta_{\theta_i} \\ \text{where } \pi_i &\sim \text{Gamma}(0, \alpha) \\ \theta_i &\sim H\end{aligned}$$

- ▶ Gamma process is a Completely Random Measure with Levy measure:

$$\nu(d\pi d\theta) = \pi^{-1} \exp^{-\alpha\pi} d\pi H(d\theta)$$

- ▶  $G \sim DP(\alpha, H)$  and  $N$  data points, the probability of  $K$  is:

$$\Pr(K = k|N, \alpha) = \frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} |s(N, k)| \alpha^k, \quad k = 0, 1, \dots, N$$

- ▶ This means that,

$$\sum_{k=1}^N |s(N, k)| \alpha^k = \frac{\Gamma(N + \alpha)}{\Gamma(\alpha)}$$

- ▶ it can be sampled as  $k = \sum_{n=1}^N b_n$ ,  $b_n \sim \text{Bernoulli}\left(\frac{\alpha}{n-1+\alpha}\right)$



- ▶  $X$  is a Negative Binomial Process with base measure  $H$  and another measure  $\mathcal{P}$ :
- ▶ Let  $X = \{(n_i, \theta_i)\}_{i=1}^{\infty}$  be a realization of a Gamma process in the product space  $\mathbb{Z}^+ \times \Theta$ :

$$\begin{aligned} X &\sim \text{NBP}(\mathcal{P}, H) \\ &= \sum n_i \delta_{\theta_i} \end{aligned}$$

# (From probability notes) Relationship between Multinomial distribution and Poisson

$$\text{Poisson}(x|\lambda) = \frac{\lambda^x}{x!} \exp(-\lambda) \qquad \text{Mult}(n_1, \dots, n_k | p_1, \dots, p_k) = \frac{(\sum n_i)!}{n_1! \dots n_k!} \prod_{i=1}^k p_i^{n_i}$$

suppose:

- ▶  $x_1 \sim \text{Poisson}(x|\lambda_1), \dots, x_k \sim \text{Poisson}(x|\lambda_k) \implies$
- ▶ The above generated two random variables:

1st random variable:  $\left( n = \sum_{i=1}^k x_i \right) \sim \text{Poisson}(\lambda_1 + \lambda_2 + \dots + \lambda_k)$

2nd random variable:  $\mathbf{x} = (x_1, \dots, x_k) | n \sim \text{Mult}(n, p_1, \dots, p_k)$  where  $p_i = \frac{\lambda_i}{\sum_{j=1}^k \lambda_j}$

# Extend this Relationship to Process

- ▶ Grouped data  $x_1, \dots, x_J$  for any measurable disjoint partition  $A_1, \dots, A_Q$  of  $\Omega$ ,
- ▶ Jointly model the count random variables  $\{X_j(A_q)\}$ .
- ▶ Poisson process  $X_j \sim \text{PP}(G)$ , with a shared Completely Random Measure  $G$  on  $\Omega : X_j(A) \sim \text{Pois}(G(A))$
- ▶  $X_j \sim \text{PP}(G)$   
 $\equiv X_j \sim \text{MP}(X_j(\Omega), \tilde{G}), \quad X_j(\Omega) \sim \text{Pois}(G(\Omega)) \quad \text{where } \tilde{G} = \frac{G}{G(\Omega)}$

$$\begin{aligned} X_j &\sim \text{NBP}\left(G_0, \frac{1}{c+1}\right) = \int_G \text{PP}(X_j|G) \text{GaP}(c, G_0) dG \\ &\sim \text{NBP}(G_0, p) = \int_G \text{PP}(X_j|G) \text{GaP}\left(\frac{J(1-p)}{p}, G_0\right) dG \end{aligned}$$

$$X = \left(\sum_{j=1}^J X_j\right) \sim \text{NBP}(G_0, p) \quad X_j(A) \sim \text{NBP}(G_0(A), p)$$

# Negative Binomial Process

- ▶  $L \sim \text{CRTP}(X, G_0)$  as CRT process:

$$\text{for each } A \in \Omega : \quad L(A) = \sum_{\omega \in \Omega} L(\omega), \quad L(\omega) \sim \text{CRT}(X(\omega), G_0(\omega))$$

- ▶  $X(A)$  customer count and  $L(A)$  table count. Each  $A \in \Omega$ . Number of tables:

$$L(A) \sim \text{Pois}(-G_0(A) \ln(1 - p))$$

- ▶ assign  $\log(p)$  customers to each table, with  $X(A)$  total number of customers.
- ▶  $X(A) \sim \text{NB}(G_0(A), p)$  customers and assign them into  $L(A) \sim \sum_{\omega \in A} \text{CRT}(X(\omega), G_0(\omega))$  tables:

$$X \sim \sum_{t=1}^L \log(p), \quad L \sim \text{PP}(-G_0 \ln(1 - p))$$

---

$$\text{is equivalent: } L \sim \text{CRTP}(X, G_0), \quad X \sim \text{NBP}(G_0, p)$$

## Negative Binomial Process (2)

$$(\gamma_0 = G_0(\Omega)) \sim \text{Gamma}\left(e_0, \frac{1}{f_0}\right)$$

$$p \sim \text{Beta}\left(a_0, \frac{1}{b_0}\right)$$

$$G|X, p, G_0 \sim \text{GaP}\left(\frac{J}{p}, G_0 + X\right)$$

$$p|X, G \sim \text{Beta}(a_0 + X(\Omega), b_0 + \gamma_0)$$

$$L|X, G_0 \sim \text{CRTP}(X, G_0)$$

$$\gamma_0|L, p \sim \text{Gamma}\left(e_0 + L(\Omega), \frac{1}{f_0 - \ln(1p)}\right)$$

You need both the posterior and predictive distribution of Multinomial-Dirichlet:

**Posterior**

$$\begin{aligned}
 & P(p_1, \dots, p_k | n_1, \dots, n_k) \\
 & \propto \underbrace{\frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k p_i^{\alpha_i-1}}_{\text{Dir}(p_1, \dots, p_k | \alpha_1, \dots, \alpha_k)} \underbrace{\frac{n!}{n_1! \dots n_k!} \prod_{i=1}^k p_i^{n_i}}_{\text{Mult}(n_1, \dots, n_k | p_1, \dots, p_k)} \\
 & \propto \prod_{i=1}^k p_i^{\alpha_i-1} \prod_{i=1}^k p_i^{n_i} = \prod_{i=1}^k p_i^{\alpha_i-1+n_i} \\
 & = \text{Dir}(p_1, \dots, p_k | \alpha_1 + n_1, \dots, \alpha_k + n_k)
 \end{aligned}$$

**Marginal**

$$\begin{aligned}
 p(n_1, \dots, n_k) &= \int_{p_1, \dots, p_k} P(p_1, \dots, p_k, n_1, \dots, n_k) \\
 &= \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \frac{n!}{n_1! \dots n_k!} \int_{p_1, \dots, p_k} \prod_{i=1}^k p_i^{\alpha_i-1+n_i} \\
 &= \frac{N!}{n_1! \dots n_k!} \times \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \times \frac{\prod_{i=1}^k \Gamma(\alpha_i + n_i)}{\Gamma(N + \sum_{i=1}^k \alpha_i)}
 \end{aligned}$$

for any measurable set  $A_i \in \Theta$ : we have  $E[G(A_i)] = H(A_i)$ , why?

For a dirichlet distribution:

$$f(x_1, \dots, x_{K-1} | \alpha_1, \dots, \alpha_K) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i)} \prod_{i=1}^K x_i^{\alpha_i-1}$$

The expectation:  $E[X_i] = \frac{\alpha_i}{\sum_k \alpha_k}$

Therefore,

$$E[G(A_i)] = \frac{\alpha H(A_i)}{\sum_j \alpha H(A_j)} = \frac{\alpha H(A_i)}{\alpha \sum_j H(A_j)} = H(A_i)$$

Variances for Dirichlet Distribution:  $Var[X_i] = \frac{\alpha_i \left( \left( \sum_i^K \alpha_{i=1} \right) - \alpha_i \right)}{\left( \sum_i^K \alpha_{i=1} \right)^2 \left( \sum_i^K \alpha_{i=1} + 1 \right)}$

Therefore:

$$\begin{aligned} Var(G(A_i)) &= \frac{\alpha H(A_i) (\alpha - \alpha H(A_i))}{\alpha^2 (\alpha + 1)} \\ &= \frac{H(A_i) (1 - H(A_i))}{(\alpha + 1)} \end{aligned}$$



From well-known multinomial-dirichlet conjugacy, we have:

$$G' = G(A_1), \dots, G(A_r) | \theta_1, \dots, \theta_n \sim \text{Dir}(\alpha H(A_1) + n_1, \dots, \alpha H(A_k) + n_k)$$

This is equivalently of saying,

$$G' \sim \text{DP} \left( \alpha + n, \frac{\alpha H + \sum_{i=1}^n \delta_{\theta_i}}{\alpha + n} \right) \text{ or,}$$
$$G' \sim \text{DP} \left( \alpha + n, \frac{\alpha}{\alpha + n} H + \frac{\sum_{i=1}^n \delta_{\theta_i}}{\alpha + n} \right)$$

DP provides a conjugate family of priors over distributions that is **closed under posterior updates given observations**

Let  $P(\theta_{n+1} \in A|G) = G(A)$ :

$$\begin{aligned}P(\theta_{n+1} \in A|\theta_1, \dots, \theta_n) &= \int_G P(\theta_{n+1} \in A|G)P(G|\theta_1, \dots, \theta_n)dG \\&= E(G(A)|\theta_1, \dots, \theta_n) \\&= E(G'(A))\end{aligned}$$

We know that  $E(G(A)) = H(A) \implies E(G'(A)) = \frac{\alpha}{\alpha+n}H(A) + \frac{\sum_{i=1}^n \delta_{\theta_i}}{\alpha+n}$

- ▶  $\beta_k \sim \text{Beta}(1, \alpha)$
- ▶  $\theta^k \sim H$
- ▶  $\pi_k = \beta_k \prod_{l=1}^{k-1} (1 - \beta_l)$
- ▶  $G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}$

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

$$\begin{aligned}
 \Pr(z_i = m | \mathbf{z}_{-1}) &= \int_{p_1, \dots, p_K} P(z_i = m | p_1, \dots, p_K) P(p_1, \dots, p_K | n_{1,-i}, \dots, n_{K,-i}) \\
 &= \frac{\int_{p_1, \dots, p_K} P(z_i = m | p_1, \dots, p_K) P(n_{1,-i}, \dots, n_{K,-i} | p_1, \dots, p_K) P(p_1, \dots, p_K)}{P(n_{1,-i}, \dots, n_{K,-i})} \\
 &= \frac{\int_{p_1, \dots, p_K} P(z_i = m | p_1, \dots, p_K) P(n_{1,-i}, \dots, n_{K,-i} | p_1, \dots, p_K) P(p_1, \dots, p_K)}{\int_{p_1, \dots, p_K} P(n_1^{-i}, \dots, n_K^{-i} | p_1, \dots, p_K) P(p_1, \dots, p_K)} \quad (1) \\
 &= \frac{\Gamma(\frac{\alpha}{k} + n_{m,-i} + 1) \prod_{l=1, l \neq m}^k \Gamma(\frac{\alpha}{k} + n_{l,-i})}{\Gamma(N + \alpha)} \times \frac{\Gamma(N - 1 + \alpha)}{\prod_{l=1}^k \Gamma(\frac{\alpha}{k} + n_{l,-1})} \\
 &= \frac{\frac{\alpha}{k} + n_{m,-i}}{N + \alpha - 1} \quad \text{Let } k \rightarrow \infty = \frac{n_{m,-i}}{N + \alpha - 1}
 \end{aligned}$$

$$\Pr(z_i = \text{new}) = \frac{\alpha}{N + \alpha - 1}.$$

- ▶ Hierarchical Dirichlet Process (HDP)
- ▶ HDP-Hidden Markov Model
- ▶ Indian Buffet Process

# Hierarchical Dirichlet Process (HDP)

## Generative model

$$G_0 \sim \text{DP}(\gamma, H)$$

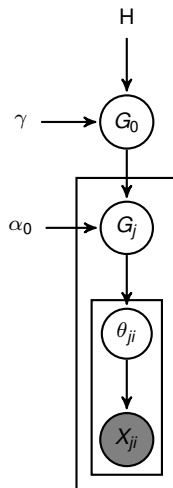
$$G_j \sim \text{DP}(\alpha_0, G_0)$$

$$\theta_{ji} \sim G_j$$

$$X_{ji} \sim F(x|\theta_{ji})$$

Drawing  $G_0 \sim \text{DP}(\cdot)$  is difficult. Therefore, we need some “construction” method:

## Graphical model



## Generative model

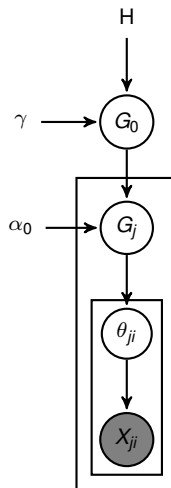
$$\beta \sim GEM(\gamma) \quad G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}$$

$$\pi_j \sim DP(\alpha_0, \beta) \quad G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\phi_k}$$

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

- ▶ See in the next slide how to sample  $\pi_j$  DIRECTLY from stick-breaking process from  $\beta$ , i.e., NO NEED to perform stick-breaking.
- ▶ Using  $\beta$  as a base, discrete distribution define on range  $\{0 \dots \infty\}$ . (takes advantage that partition of space is given)

## Graphical model



## Sample $\pi_j$ DIRECTLY from stick-breaking process from $\beta$ :

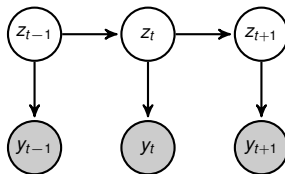
Suppose  $\beta|\gamma \sim \text{GEM}(\gamma)$  and  $\pi|\alpha, \beta \sim \text{DP}(\alpha, \beta)$ . Notice that the support is  $\{1, \dots, k, \dots, \infty\}$ :

$$\begin{aligned} & (G_j(A_1), \dots, G_j(A_r)) \sim \text{Dir}(\alpha G_0(A_1), \dots, \alpha G_0(A_r)) \\ \Rightarrow & \left( \sum_{k \in K_1} \pi_k, \dots, \sum_{k \in K_r} \pi_k \right) \sim \text{Dir} \left( \alpha \sum_{k \in K_1} \beta_k, \dots, \alpha \sum_{k \in K_r} \beta_k \right) \\ \Rightarrow & \left( \sum_{l=1}^{k-1} \pi_l, \pi_k, \sum_{l=k+1}^{\infty} \pi_l \right) \sim \text{Dir} \left( \alpha \sum_{l=1}^{k-1} \beta_l, \alpha \beta_k, \sum_{l=k+1}^{\infty} \beta_l \right) \\ \Rightarrow & \left( \frac{\pi_k}{1 - \sum_{l=1}^{k-1} \pi_l}, \frac{\sum_{l=k+1}^{\infty} \pi_l}{1 - \sum_{l=1}^{k-1} \pi_l} \right) \sim \text{Dir} \left( \alpha \beta_k, \sum_{l=k+1}^{\infty} \beta_l \right) \\ \Rightarrow & \left( \frac{\pi_k}{1 - \sum_{l=1}^{k-1} \pi_l}, \frac{1 - \sum_{l=1}^k \pi_l}{1 - \sum_{l=1}^{k-1} \pi_l} \right) \sim \text{Dir} \left( \alpha \beta_k, 1 - \sum_{l=1}^k \beta_l \right) \\ \Rightarrow & \frac{\pi_k}{1 - \sum_{l=1}^{k-1} \pi_l} \sim \text{Beta} \left( \alpha \beta_k, 1 - \sum_{l=1}^k \beta_l \right) \text{ or } \pi' \sim \text{Beta} \left( \alpha \beta_k, 1 - \sum_{l=1}^k \beta_l \right) \end{aligned}$$



Under normal HMM, you have a transition matrix  $A$ , let the  $j^{\text{th}}$  row of  $A$  to be  $\pi_j$ , 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  $z_t \sim \pi_{z_{t-1}}$ .

well, the same idea has been extended to non-parametric bayes, to allow  $\pi_j$  to have infinite many components. Therefore, 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  $A$ .

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

## Generative model

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

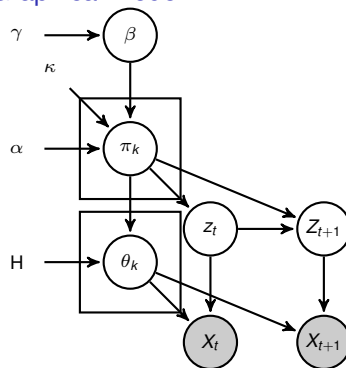
$$\pi_j \sim \text{DP} \left( \alpha + \kappa, \frac{\alpha \beta + \kappa \delta_j}{\alpha + \kappa} \right)$$

$$z_{t+1} \sim \pi_{z_t}$$

$$\theta_k \sim H$$

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

## Graphical model



# Indian Buffet Process: Its relationship with DP

## DP

- ▶  $\Pr(z_1 \dots z_N)$ , where  $z_i \in (1 \dots K)$  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  $\in (1, \dots K)$
- ▶ This is equivalent 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.

# The big $Z$ matrix

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  $u_k \sim \text{Beta}(\frac{\alpha}{k}, 1)$  is also distributed independently.

The marginal distribution:

# Bernoulli- Beta vs Multinomial-Dirichlet: Posterior

## Multinomial-Dirichlet

$$\begin{aligned} P(p_1, \dots, p_k | n_1, \dots, n_k) \\ &\propto \underbrace{\frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k p_i^{\alpha_i-1}}_{\text{Dir}(p_1, \dots, p_k | \alpha_1, \dots, \alpha_k)} \underbrace{\frac{n!}{n_1! \dots n_k!} \prod_{i=1}^k p_i^{n_i}}_{\text{Mult}(n_1, \dots, n_k | p_1, \dots, p_k)} \\ &\propto \prod_{i=1}^k p_i^{\alpha_i-1} \prod_{i=1}^k p_i^{n_i} = \prod_{i=1}^k p_i^{\alpha_i-1+n_i} \\ &= \text{Dir}(p_1, \dots, p_k | \alpha_i + n_i, \dots, \alpha_k + n_k) \end{aligned}$$

## Bernoulli-Binomial

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

# Bernoulli- Beta vs Multinomial-Dirichlet: Marginal

## Multinomial-Dirichlet

$$\begin{aligned} & \int_{p_1, \dots, p_K} P(p_1, \dots, p_K, n_1, \dots, n_K) \\ &= \frac{N!}{n_1! \dots n_K!} \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(\alpha_i + n_i)}{\Gamma\left(N + \sum_{i=1}^K \alpha_i\right)} \end{aligned}$$

## Bernoulli-Beta

$$\begin{aligned} & \int_p P(p, n_1, n_2) \\ &= \frac{N!}{k!(N-k)!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + k)\Gamma(\beta + N - k)}{\Gamma(N + \alpha + \beta)} \end{aligned}$$

# Bernoulli-Beta Predictive

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

$n_{k,-i}$  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  $m$ .

$$\begin{aligned}\Pr(z_{ik} = 1 | \mathbf{z}_{-i,k}) &= \int_p \Pr(z_{ik} = 1 | p) P(p | \underbrace{n_{-i,k}}_{n_1}, \underbrace{i-1-n_{-i,k}}_{n_2}) \\&= \frac{\int_p \Pr(z_{ik} = 1 | p) \Pr(n_1, n_2 | p) P(p)}{\Pr(n_1, n_2)} = \frac{\int_p \Pr(z_{ik} = 1 | p) \Pr(n_1, n_2 | p) P(p)}{\int_p \Pr(n_{-i,k}, i-1-n_{-i,k} | p) P(p)} \\&= \frac{\Gamma(\frac{\alpha}{k} + n_{-i,k} + 1) \Gamma(1 + i - 1 - n_{-i,k})}{\Gamma(i + \frac{\alpha}{k} + 1)} \frac{\Gamma(i - 1 + \frac{\alpha}{k} + 1)}{\Gamma(\frac{\alpha}{k} + n_{-i,k}) \Gamma(1 + i - 1 - n_{-i,k})} = \frac{\frac{\alpha}{k} + n_{-i,k}}{i + \frac{\alpha}{k}}\end{aligned}$$



# One more factor: relationship between Binomial and Poisson

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

Let  $\lambda = np$ :

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

$$\begin{aligned}\lim_{n \rightarrow \infty} \text{Binomial}(x|n, p) &= \lim_{n \rightarrow \infty} \binom{n}{x} p^x (1-p)^{n-x} \\&= \frac{\lambda^x}{x!} \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right) \dots \lim_{n \rightarrow \infty} \left(1 - \frac{x-1}{n}\right) \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-x} = \frac{\lambda^x}{x!} \exp(-\lambda)\end{aligned}$$

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

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

$$\text{Let } k \rightarrow \infty : \quad = \frac{n_{-i,k}}{i}$$

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:

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

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

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

# An motivational example of IBP: Factor Analysis

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

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

Or in a matrix form:  $\mathbf{Y} = \mathbf{GX} + \mathbf{E}$ .

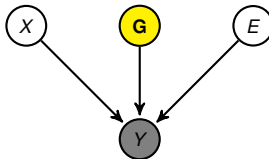
What this means is that a person's  $i$ 's raw mark is interpreted as:

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

- ▶ 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  $\mathbf{E}$ . This is not as silly as it seems, as DoF is much reduced.

## The Bayesian Treatment:

$$\begin{aligned} e_i &\sim \mathcal{N}(0, \sigma_e^2 \mathbf{I}) & \sigma_e^2 &\sim \text{IG}(a, b) \\ g_k &\sim \mathcal{N}(0, \sigma_G^2) & \sigma_G^2 &\sim \text{IG}(c, d) \\ x_{ki} &\sim \mathcal{N}(0, 1) & y_i &= \mathbf{G}x_i + e_i \end{aligned}$$



# Infinite Factor Analysis

- ▶ Knowles, d and Ghahramani, Z, Infinite Sparse Factor Analysis
- ▶  $K$  should be known beforehand. What about making  $K$  a variable?
- ▶ Although  $[x_{1,i}, \dots, x_{K,i}]^T$  has a reduced dimension, it can still cause “overfitting”.
- ▶ We need to introduce variable number of latent factors  $K$ , at the same time, have **sparsity**!

How?

$$\mathbf{e}_i \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{e}}^2 \mathbf{I})$$

$$\sigma_{\mathbf{e}}^2 \sim \text{IG}(a, b)$$

$$g_k \sim \mathcal{N}(0, \sigma_G^2)$$

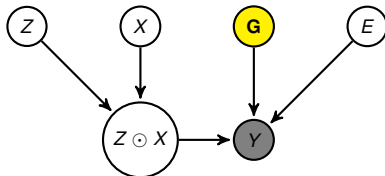
$$\sigma_G^2 \sim \text{IG}(c, d)$$

$$Z \sim \text{IBP}(\alpha)$$

$$\alpha \sim \mathcal{G}(\mathbf{e}, f)$$

$$x_{ki} \sim \mathcal{N}(0, 1)$$

$$y_i = \mathbf{G}(\mathbf{x}_i \odot \mathbf{z}_i) + \mathbf{e}_i$$



# 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**.

$$\mathbf{e}_i \sim \mathcal{N}(0, \sigma_e^2 \mathbf{I})$$

$$\sigma_e^2 \sim \mathcal{IG}(a, b)$$

$$g_k \sim \mathcal{N}(0, \sigma_G^2)$$

$$\sigma_G^2 \sim \mathcal{IG}(c, d)$$

$$Z \sim \mathcal{IBP}(\alpha)$$

$$\alpha \sim \mathcal{G}(\mathbf{e}, f)$$

$$x_{ki} \sim \mathcal{N}(0, 1)$$

$$y_i = \mathbf{G}(x_i \odot z_i) + \mathbf{e}_i$$

$$\mathbf{e}'_i \sim \mathcal{N}(0, \sigma_e'^2 \mathbf{I})$$

$$\sigma_e'^2 \sim \mathcal{IG}(a', b')$$

$$g_k \sim \mathcal{N}(0, \sigma_G^2)$$

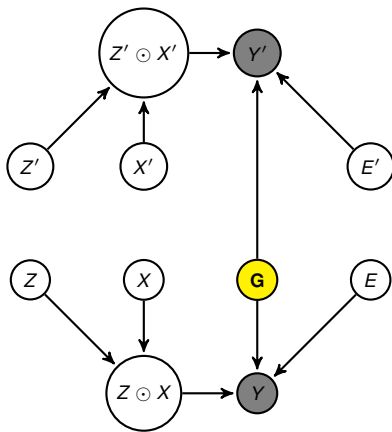
$$\sigma_G^2 \sim \mathcal{IG}(c, d)$$

$$Z' \sim \mathcal{IBP}(\alpha')$$

$$\alpha' \sim \mathcal{G}(\mathbf{e}', f')$$

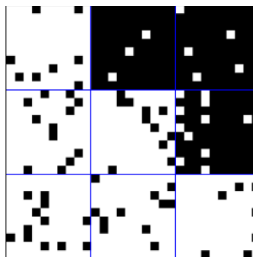
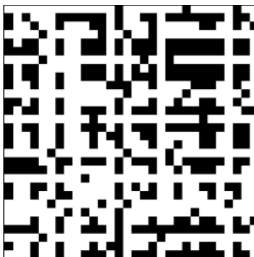
$$x'_{ki} \sim \mathcal{N}(0, 1)$$

$$y'_i = \mathbf{G}(x'_i \odot z'_i) + \mathbf{e}'_i$$



# Introduction of Relational Model

- ▶ Community learning is an emerging topic applicable to many social networking problems and “hot” in machine learning.
- ▶ Partition a network of nodes into different groups based on their pairwise and directional binary observations.
- ▶ Many models were proposed in the last few years and they become increasing sophisticated.
- ▶ Data is **directional** i.e., I like you doesn't mean you like me.





# Simple Stochastic Block Model assumption: Fixed $K$ communities

The model:

- ▶ There is a hidden “compatibility” matrix  $\mathbf{B}$ , size  $K \times K$ , each element  $\mathbf{B}_{kl} \sim \text{Beta}(\lambda_1, \lambda_2)$ , a realization example:

0.5	0.2	0.1	0.1	0	...	0.1
0.3	0.91	0.2	0.4	0.2	...	0.5
...	...	...	...	...	...	0.2
0.32	0.2	0.96	0.4	0.7	...	0.9

- ▶ Suppose that person  $i$  is in latent community 2, i.e.,  $z_i = 2$  and person  $j$  is in latent community 3, i.e.,  $z_j = 3$ .
- ▶ Then  $\mathbf{e}_{ij} \sim \text{Bernoulli}(\mathbf{B}_{z_i=2, z_j=3}) = \text{Bernoulli}(0.2)$ .
- ▶  $z_i \sim \pi$ : some weights of communities

Inference:

- ▶ Then, our task is to perform posterior inference on:  $\Pr(z_1, \dots, z_n, \pi, \mathbf{B} | \{\mathbf{e}_{ij}\})$

Early work assumes a fixed number of  $K$  communities exist a node  $i$  can potentially belong to. However, in many applications, an accurate guess of  $K$  can be impractical.

## **Infinite Relational Model** (Kemp 2006)

- ▶ Infinite Relational Model was incorporated to address this problem, where  $K$  can be inferred from the data itself, and potentially be  $\infty$ .
- ▶ That's where Non-parametric Bayes comes in!
- ▶ Still a drawback: assumes each node  $i$  must belong to only a single community  $k$  (i.e.,  $z_i = k$ ).

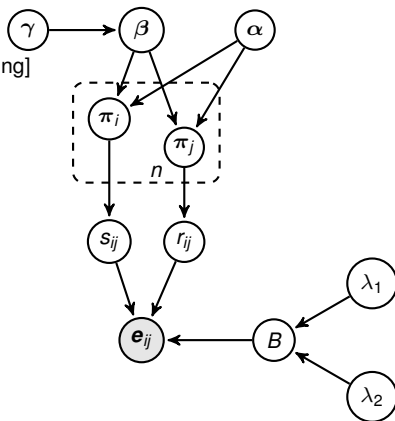
## Mixed-Membership Stochastic Blockmodel (MMSB)

[Airoldi et al.(2008)Airoldi, Blei, Fienberg, and Xing]

- ▶ mixed-membership concept: each node  $i$  may belong to multiple communities. Having individual distribution  $\pi_i$
- ▶  $e_{ij}$  no longer dependant only on each pair of community indicators  $z_i$  and  $z_j$ . Instead, they are sampled from pairs of interactions between nodes  $i$  and  $j$ .  $(s_{ij}, r_{ij})$ .

## Generative Model

1.  $\beta \sim GEM(\gamma)$
2.  $\{\pi_i\}_{i=1}^n \sim DP(\alpha \cdot \beta)$
3.  $s_{ij} = \pi_i, r_{ij} = \pi_j$
4.  $B_{k,l} \sim Beta(\lambda_1, \lambda_2), \forall k, l;$
5.  $e_{ij} \sim Bernoulli(B_{s_{ij}, r_{ij}})$ .



# Mixed Membership Stochastic Block Model

The priors:

- ▶ Each element  $\mathbf{B}_{kl} \sim \text{Beta}(\lambda_1, \lambda_2)$  still: a realization example:

0.5	0.2	0.1	0.1	0	...	0.1
0.3	0.91	0.2	0.4	0.2	...	0.5
...	...	...	...	...	...	0.2
0.32	0.2	0.96	0.4	0.7	...	0.9

- ▶ Suppose that interaction  $i$  **sent to**  $j$  is of latent community 2, i.e.,  $s_{ij} = 2$ ,
- ▶ Interaction  $j$  **received from**  $i$  is in latent community 3, i.e.,  $r_{ji} = 3$ .
- ▶ Note that  $s_{ij}$  do not generally equal  $r_{ji}$ .
- ▶ Then  $\mathbf{e}_{ij} \sim \text{Bernoulli}(\mathbf{B}_{s_{ij}=2, r_{ji}=3}) = \text{Bernoulli}(0.2)$ .
- ▶  $\{s_{i,k}, r_{i,k}\} \sim \pi_j$ : There are altogether  $N$   $\pi$ s

The posterior

- ▶ Then, our task is to perform posterior inference on:  
 $\Pr(\{s_{i,j}, r_{j,i}\}_{\forall 1 \leq i, j \leq N}, \mathbf{B}, \pi_1, \dots, \pi_N | \{\mathbf{e}_{ij}\})$

A few variants were subsequently proposed from MMSB, examples include:

- ▶ [?, Xing et al.(2010)Xing, Fu, and Song] extends the mixture-membership model with a dynamic setting;
- ▶ [Koutsourelakis and Eliassi-Rad(2008)] extends the MMSB into the infinite case; and
- ▶ [Kim et al.(2012)Kim, Hughes, and Sudderth] incorporates the node's metadata information into MMSB.

Sudderth's method:

**The model**

$$v_{:i} \sim \mathcal{N}(\eta^T \phi_{:i}, \lambda_{v_i}^{-1})$$

Instead of  $v_i \sim \beta(1, \alpha)$  :

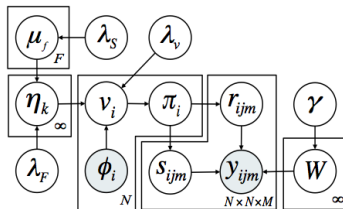
$$\pi_{ki} = \psi(v_{ki}) \prod_{l=1}^{k-1} \psi(-v_{li})$$

$$\text{where } \psi(v_{ki}) = \frac{1}{1 + \exp(-v_{ki})}$$

For each community  $k$ :

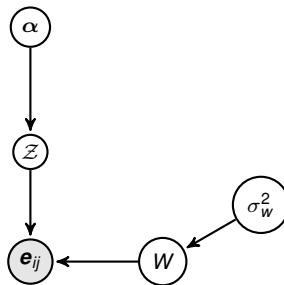
$$\eta_{:k} \sim \mathcal{N}(\mu, \lambda_F^{-1} I_F)$$

**Graphic Model**



# Literatures: Infinite Latent Feature Relational Model (LFRM)

$$\begin{aligned}Z &\sim \text{IBP}(\alpha) \\ e_{ij} &\sim Z_i W Z_j^T \\ W_{k,k'} &\sim \mathcal{N}(0, \sigma_w)\end{aligned}$$



# Our work: Copula Mixed-Membership Stochastic Blockmodel with Subgroup Correlation

- ▶ Despite MMSB's powerful representations, it assumes that the distributions of relational membership indicators between the two nodes are independent.
- ▶ Under many social network settings, possible that certain known subgroups of people may have higher correlations in terms of their membership categories towards each other
- ▶ We introduce a new framework where individual Copula function is to be employed to model jointly the membership pairs of those nodes within the subgroup of interest.
- ▶ Various Copula functions may be used to suit the scenario, while maintaining the membership's marginal distribution, as needed for modeling membership indicators with other nodes outside of the subgroup of interest.
- ▶ Experimental results shows a superior performance when comparing with the existing models on both the synthetic and real world datasets.

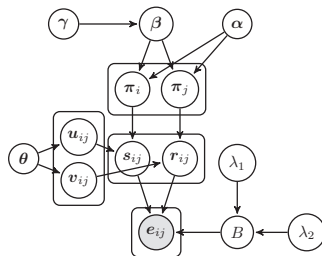


# The model

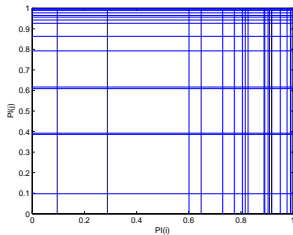
## Generative Model

1.  $\beta \sim GEM(\gamma)$
2.  $\{\pi_i\}_{i=1}^n \sim DP(\alpha \cdot \beta)$
3.  $\begin{cases} (u_{ij}, v_{ij}) \sim Copula(\theta), & g_{ij} = 1; \\ u_{ij}, v_{ij} \sim U(0, 1), & g_{ij} = 0. \end{cases}$
4.  $s_{ij} = \Pi_i^{-1}(u_{ij}), r_{ij} = \Pi_j^{-1}(v_{ij})$
5.  $B_{k,l} \sim Beta(\lambda_1, \lambda_2), \forall k, l;$
6.  $e_{ij} \sim Bernoulli(B_{s_{ij}, r_{ij}})$ .

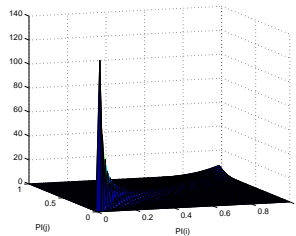
## Graphical model



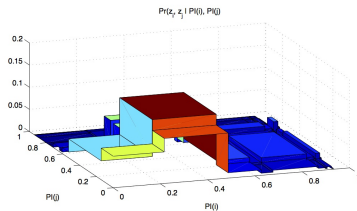
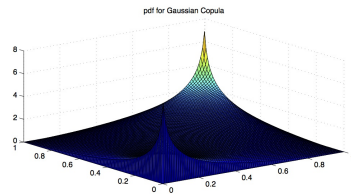
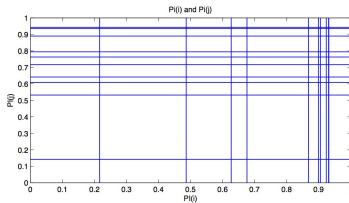
1.  $\pi_i, \pi_j \text{simDP}(\alpha \cdot \beta)$



1.  $(u_{ij}, v_{ij}) \sim \text{Copula}(\theta), g_{ij} = 1$



# Diagrammatic Representation 2



# So what is Copula

- ▶ A bivariate copula function  $C(u, v)$  is a Cumulative Distribution Function over the interval  $[0, 1] \times [0, 1]$  with uniform marginal distribution.
- ▶ *Sklar's theorem*: Let  $X$  and  $Y$  be random variables with distribution functions  $F$  and  $G$  respectively and joint distribution function  $H$ . Then there exists a Copula  $C$  such that for all  $(x, y) \in R \times R$ :

$$H(x, y) = C(F(x), G(y))$$

- ▶  $C$  is unique if  $F$  and  $G$  are continuous, then the joint probability density function is:

$$h(x, y) = c(F(x), G(y)) \cdot f(x)g(y) \quad (2)$$

Here  $c(u, v) = \frac{\partial^2 C(u, v)}{\partial u \partial v}$  is notes for copula density function.

## So what is Copula Cont.

- ▶ Sklar's theorem ensures the uniqueness of copula function  $C(F(x), G(y))$
- ▶ Change Copula function does not change the marginal distributions. This is what we want!
- ▶ Copula is popular! Many are availability: Commonly used copula functions includes, Gaussian Copula (Gaussian, t), Archimedean Copula (Clayton, Gumbel, Frank, etc.), Empirical Copula.

$$\begin{aligned} \Pr(s_{ij}, r_{ij}) = & \int_{\pi_{i,1}, \dots, \pi_{i,K+1}} \int_{\pi_{j,1}, \dots, \pi_{j,K+1}} \int_{(u_{ij}, v_{ij})} \\ & \cdot \mathbf{1}(s_{ij} = \Pi_i^{-1}(u_{ij}), r_{ij} = \Pi_j^{-1}(v_{ij})) \\ & \cdot dC(u_{ij}, v_{ij}) dF(\pi_{i1}, \dots, \pi_{iK+1}) dF(\pi_{j1}, \dots, \pi_{jK+1}) \end{aligned} \quad (3)$$

Unfortunately, we cannot get it to an analytical form without any integrals present.

$$\begin{aligned}
 p_{ij}^{kl}(\pi_i, \pi_j) &\equiv \Pr(s_{ij} = k, r_{ij} = l | \pi_i, \pi_j, \theta_d) \\
 &= \int_{\hat{\pi}_i^{k-1}}^{\hat{\pi}_i^k} \int_{\hat{\pi}_j^{l-1}}^{\hat{\pi}_j^l} dC(u, v | \theta_d) \\
 &= C(\hat{\pi}_i^k, \hat{\pi}_j^l) + C(\hat{\pi}_i^{k-1}, \hat{\pi}_j^{l-1}) - C(\hat{\pi}_i^k, \hat{\pi}_j^{l-1}) - C(\hat{\pi}_i^{k-1}, \hat{\pi}_j^l)
 \end{aligned}$$

$$\hat{\pi}_i^k = \begin{cases} 0, & k = 0; \\ \sum_{q=1}^k \pi_{iq}, & k > 0 \end{cases}$$

- ▶ Easily calculate this “rectangular” area.
- ▶ When no correlations,  $p_{ij}^{kl}(\pi_i, \pi_j) = \pi_{ik}\pi_{jl}$

Properties of a Copula function: marginal of  $\Pr(s_{ij} = k, r_{ij} = l | \pi_i, \pi_j, \theta_d)$  remain  $\pi_i$  and  $\pi_j$  respectively:

$$\begin{aligned}
 \sum_{l=1}^{K+1} \Pr(s_{ij} = k, r_{ij} = l | \pi_i, \pi_j, \theta_d) &= \pi_{ik}; \\
 \sum_{k=1}^{K+1} \Pr(s_{ij} = k, r_{ij} = l | \pi_i, \pi_j, \theta_d) &= \pi_{jl}.
 \end{aligned} \tag{4}$$

## Marginal conditional on $u, v$ only ( $cMMSB^{uv}$ )

- ▶ Integrate over  $\{\pi_i\}_{i=1}^n$  given  $\{(u_{ij}, v_{ij})\}_{i,j}$
- ▶ Given  $\{(u_{ij}, v_{ij})\}_{i,j}$ ,  $\Pr(s_{ij} = k)$  and  $\Pr(r_{ij} = l)$  are independent.
- ▶ Copula function leaves marginal distributions of  $s_{ij}$  and  $r_{ij}$  invariant, which remains the same as the classical posterior of  $\pi_i$  in *MMSB*:

$$\pi_i | \alpha, \beta, \{N_{ik}^{-ij}\}_{k=1}^K \sim \text{Dir}(\alpha\beta_1 + N_{i1}^{-ij}, \dots, \alpha\beta_K + N_{iK}^{-ij}, \alpha\beta_{K+1})$$

- ▶  $\Pr(s_{ij} = k)$  is equal to computing the probability of  $u_{ij}$  falling in  $\pi_i$ 's  $k^{\text{th}}$  interval:

$$\Pr\left(\sum_{d=1}^{k-1} \pi_{id} \leq u_{ij} < \sum_{d=1}^k \pi_{id}\right)$$

- ▶ The fact that the set  $\{u_{ij} \in [0, 1] | \sum_{d=1}^{k-1} \pi_{id} \leq u_{ij}\}$  can be decomposed into two **disjoint** sets:

$$\begin{aligned} & \{u_{ij} \in [0, 1] | \sum_{d=1}^{k-1} \pi_{id} \leq u_{ij}\} \\ &= \{u_{ij} \in [0, 1] | \sum_{d=1}^{k-1} \pi_{id} \leq u_{ij} < \sum_{d=1}^k \pi_{id}\} \cup \{u_{ij} \in [0, 1] | \sum_{d=1}^k \pi_{id} \leq u_{ij}\} \end{aligned} \tag{5}$$

where  $\sum_{d=1}^k \pi_{id} \sim \text{Beta}(\sum_{d=1}^k \alpha\beta_d + N_{id}, \sum_{d=k+1}^{K+1} \alpha\beta_d + N_{id})$ .



We have:

$$\begin{aligned} & \Pr\left(\sum_{d=1}^{k-1} \pi_{id} \leq u_{ij} < \sum_{d=1}^k \pi_{id}\right) \\ &= \Pr\left(\sum_{d=1}^{k-1} \pi_{id} \leq u_{ij}\right) - \Pr\left(\sum_{d=1}^k \pi_{id} \leq u_{ij}\right) \\ &= l_{u_{ij}}(h_i^{k-1}, \hat{h}_i^{k-1}) - l_{u_{ij}}(h_i^k, \hat{h}_i^k) \end{aligned}$$

$$h_i^k = \sum_{d=1}^k \alpha \beta_d + N_{id} \qquad \hat{h}_i^k = \sum_{d=k+1}^{K+1} \alpha \beta_d + N_{id}$$

$l_u(a, b)$  is Beta c.d.f. of  $u$  with parameter  $a, b$

Non-negativity is guaranteed by the fact that

$\{u_{ij} \in [0, 1] \mid \sum_{d=1}^k \pi_{id} \leq u_{ij}\} \subseteq \{u_{ij} \in [0, 1] \mid \sum_{d=1}^{k-1} \pi_{id} \leq u_{ij}\}$  on the same  $\pi_j$ .

- ▶ In  $cMMSB^\pi$ , variables of interest are  $\{\pi_i\}, \{s_{ij}, r_{ij}\}, \beta$ .
- ▶ In  $cMMSB^{uv}$ , variables of interest include  $\{u_{ij}, v_{ij}\}, \{s_{ij}, r_{ij}\}, \beta$ , and an auxiliary variable  $\mathbf{m}$ .

# Inference $cMMSB^\pi$ - Sampling $\pi_i$

- ▶ When a Copula is introduced,  $p(\pi_i)$  and  $\Pr(s_{ij}|\pi_i)$  are no longer a conjugate pair.
- ▶ Therefore, resort to the use of Metropolis-Hastings (M-H) Sampling in each  $(\tau)$ -th Gibbs iteration

For each node  $i$ , posterior distribution of  $\pi_i$  is:

$$p(\pi_i | \alpha, \beta, \{s_{ij}, r_{ij}\}_{i,j}) \\ \propto \prod_{k=1}^{K+1} \pi_{ik}^{\alpha\beta_k - 1} \cdot \prod_{j=1}^n \left[ p_{ij}^{s_{ij}r_{ij}}(\pi_i, \pi_j) p_{ji}^{s_{ji}r_{ji}}(\pi_j, \pi_i) \right]$$

Corresponding proposal of  $\pi_i$ :

$$q(\pi_i^* | \alpha, \beta, \{s_{ij}, r_{ij}\}_{i,j}) \propto \prod_{k=1}^{K+1} [\pi_{ik}^*]^{\alpha\beta_k + N_{ik} - 1}$$

Acceptance ratio becomes:

$$A(\pi_i^*, \pi_i^{(\tau)}) = \min(1, a) \quad (6)$$

$$a = \frac{\prod_{j=1}^n \left[ p_{ij}^{s_{ij}r_{ij}}(\pi_i^*, \pi_j) p_{ji}^{s_{ji}r_{ji}}(\pi_j, \pi_i^*) \right]}{\prod_{j=1}^n \left[ p_{ij}^{s_{ij}r_{ij}}(\pi_i^{(\tau)}, \pi_j) p_{ji}^{s_{ji}r_{ji}}(\pi_j, \pi_i^{(\tau)}) \right]} \cdot \frac{\prod_{k=1}^{K+1} [\pi_{ik}^{(\tau)}]^{N_{ik}}}{\prod_{k=1}^{K+1} [\pi_{ik}^*]^{N_{ik}}} \quad (7)$$

$$\begin{aligned} & \Pr(s_{ij}, r_{ij} | e_{ij}, \lambda_1, \lambda_2, \theta_d, \pi_i, \pi_j, \{(s_{ij}, r_{ij})\}_{i,j}) \\ & \propto p_{ij}^{s_{ij}, r_{ij}}(\pi_i, \pi_j) \cdot p(e_{ij} | \lambda_1, \lambda_2, \{(s_{ij}, r_{ij})\}_{i,j}) \end{aligned}$$

$$p(e_{ij} | \lambda_1, \lambda_2, \{(s_{ij}, r_{ij})\}_{i,j}) = \begin{cases} n_{s_{ij}, r_{ij}}^{1, -e_{ij}} + \lambda_1, & e_{ij} = 1; \\ n_{s_{ij}, r_{ij}}^{0, -e_{ij}} + \lambda_2, & e_{ij} = 0. \end{cases}$$

- ▶ obvious choice of M-H proposal of  $\beta$  its prior  $p(\beta|\gamma) = GEM(\gamma)$ .
- ▶ this proposal can be non-informative, which results in a low acceptance rate.
- ▶ We sample  $\beta^*$  conditioned on an auxiliary variable  $\mathbf{m}$ :  
 $(\beta_1^*, \dots, \beta_K^*, \beta_{K+1}^*) \sim Dir(\mathbf{m}_1, \dots, \mathbf{m}_K, \gamma)$ , in order to increase the M-H's acceptance rate
- ▶ instead of sampling  $\beta$  directly from  $\mathbf{m}$  as in [Teh et al.(2006)Teh, Jordan, Beal, and Blei], we only use it for our proposal distribution, as we have explicitly sampled  $\{\pi_i\}_{i=1}^n$ . The acceptance ratio is hence:

$$A(\beta^*, \beta^{(\tau)}) = \min(1, a)$$
$$a = \frac{\prod_{i=1}^n \left[ \prod_{d=1}^{K+1} \Gamma(\alpha \beta_d^{(\tau)}) \cdot \pi_{id}^{\alpha \beta_d^*} \right]}{\prod_{i=1}^n \left[ \prod_{d=1}^{K+1} \Gamma(\alpha \beta_d^*) \cdot \pi_{id}^{\alpha \beta_d^{(\tau)}} \right]} \cdot \frac{\prod_{d=1}^K [\beta_d^{(\tau)}]^{m_d - \gamma}}{\prod_{d=1}^K [\beta_d^*]^{m_d - \gamma}} \quad (8)$$

The Copula function is used as its proposal, and therefore, its corresponding acceptance ratio becomes that of:

$$A\left((u_{ij}^{(\tau)}, v_{ij}^{(\tau)}), (u_{ij}^*, v_{ij}^*)\right) = \min(1, a)$$

$$a = \frac{l_{u_{ij}^*}(h_i^{k-1}, \hat{h}_i^{k-1}) - l_{u_{ij}^*}(h_i^k, \hat{h}_i^k)}{l_{u_{ij}^{(\tau)}}(h_i^{k-1}, \hat{h}_i^{k-1}) - l_{u_{ij}^{(\tau)}}(h_i^k, \hat{h}_i^k)} \cdot \frac{l_{v_{ij}^*}(h_j^{l-1}, \hat{h}_j^{l-1}) - l_{v_{ij}^*}(h_j^l, \hat{h}_j^l)}{l_{v_{ij}^{(\tau)}}(h_j^{l-1}, \hat{h}_j^{l-1}) - l_{v_{ij}^{(\tau)}}(h_j^l, \hat{h}_j^l)}$$

Here  $h_i^k, \hat{h}_i^k$ 's definitions are the same as in Eq. (7) in the paper; assuming  $s_{ij} = k, r_{ij} = l$ .

$$\begin{aligned}
 & \Pr(s_{ij} = k, r_{ij} = l | \mathbf{e}_{ij}, \lambda_1, \lambda_2, n_{kl}, u_{ij}, v_{ij}, \{h_i^k\}_k, \{\hat{h}_i^k\}_k, \{h_j^k\}_k, \{\hat{h}_j^k\}_k) \\
 & \propto \Pr(s_{ij} = k | u_{ij}, \{h_i^k\}_k, \{\hat{h}_i^k\}_k) \cdot \Pr(r_{ij} = l | v_{ij}, \{h_j^k\}_k, \{\hat{h}_j^k\}_k) \cdot \Pr(\mathbf{e}_{ij} | \lambda_1, \lambda_2, n_{kl}) \\
 & \propto (l_{u_{ij}}(h_i^{k-1}, \hat{h}_i^{k-1}) - l_{u_{ij}}(h_i^k, \hat{h}_i^k)) \cdot (l_{v_{ij}}(h_j^{l-1}, \hat{h}_j^{l-1}) - l_{v_{ij}}(h_j^l, \hat{h}_j^l)) \cdot \Pr(\mathbf{e}_{ij} | \lambda_1, \lambda_2, n_{kl})
 \end{aligned}$$

The likelihood is:

$$\begin{aligned}
 & \Pr(\mathbf{e}_{ij} | s_{ij} = k, r_{ij} = l, \lambda_1, \lambda_2, n_{kl}^{-\mathbf{e}_{ij}}) \\
 & \propto P(\mathbf{e}_{ij}, \mathbf{e} \setminus \{\mathbf{e}_{ij}\}, s_{ij} = k, r_{ij} = l, n_{kl}^{-\mathbf{e}_{ij}}, \lambda_1, \lambda_2) \\
 & = \frac{\Gamma(\mathbf{e}_{ij} + n_{k,l}^1 + \lambda_1) \Gamma(1 - \mathbf{e}_{ij} + n_{k,l}^0 + \lambda_2)}{\Gamma(1 + n_{k,l} + \lambda_1 + \lambda_2)} \\
 & P(\mathbf{e}_{ij} | \mathbf{e} \setminus \{\mathbf{e}_{ij}\}, s_{ij} = k, r_{ij} = l, \mathbf{Z} \setminus \{s_{ij}, r_{ij}\}, \lambda_1, \lambda_2) \\
 & = \begin{cases} \frac{n_{k,l}^1 + \lambda_1}{n_{k,l} + \lambda_1 + \lambda_2} & \text{if } \mathbf{e}_{ij} = 1; \\ \frac{n_{k,l}^0 + \lambda_2}{n_{k,l} + \lambda_1 + \lambda_2} & \text{if } \mathbf{e}_{ij} = 0. \end{cases}
 \end{aligned}$$

Here  $n_{k,l} = \sum_{i',j'} \mathbf{1}(s_{i'j'} = k, r_{i'j'} = l)$ ,  $n_{k,l}^1 = \sum_{s_{i'j'}=k, r_{i'j'}=l} \mathbf{e}_{i'j'}$ , and  $n_{k,l}^0 = n_{k,l} - n_{k,l}^1$ .

- ▶ Selected and report the results on 3 datasets. NIPS Co-authorship dataset, the lazega-lawfirm dataset and the MIT Reality Mining dataset
- ▶ ten-folds cross-validation to complete this task, where we randomly select one out of ten for each node's link data as test data and the others as training data
- ▶ he criteria in evaluating this predict ability includes the train error ( $0 - 1$  loss), the test error ( $0 - 1$  loss), the test log likelihood and the AUC (Area Under the roc Curve) score  $cMMSB^{uv}$  obtain either comparable or better performance with other state-of-the-art.



Table: Different models' performance (Mean  $\pm$  Standard Deviation) on Real world datasets

dataset		<i>Train error</i>	<i>Test error</i>	<i>Test log likelihood</i>	<i>AUC</i>
NIPS co-author	<i>IRM</i>	0.0317 $\pm$ 0.0004	0.0423 $\pm$ 0.0014	-135.0467 $\pm$ 7.3816	0.8901 $\pm$ 0.0162
	<i>LFRM</i>	0.0482 $\pm$ 0.0794	0.0239 $\pm$ 0.0735	-105.2166 $\pm$ 179.5505	0.9348 $\pm$ 0.1667
	<i>MMSB</i>	0.0132 $\pm$ 0.0042	0.0301 $\pm$ 0.0064	-86.2134 $\pm$ 10.1258	0.9524 $\pm$ 0.0215
	<i>iMMM</i>	<b>0.0061 <math>\pm</math> 0.0019</b>	0.0253 $\pm$ 0.0035	-83.4264 $\pm$ 9.4293	0.9574 $\pm$ 0.0155
	<i>cMMSB<sup><math>\pi</math></sup></i>	0.0066 $\pm$ 0.0038	<b>0.0231 <math>\pm</math> 0.0043</b>	-83.4261 $\pm$ 9.4280	0.9569 $\pm$ 0.0159
	<i>cMMSB<sup>uv</sup></i>	0.0097 $\pm$ 0.0047	0.0240 $\pm$ 0.0065	<b>-83.4257 <math>\pm</math> 9.4292</b>	<b>0.9581 <math>\pm</math> 0.0153</b>
MIT realty	<i>IRM</i>	0.0627 $\pm$ 0.0002	0.0665 $\pm$ 0.0004	-133.8037 $\pm$ 1.1269	0.8261 $\pm$ 0.0047
	<i>LFRM</i>	0.0397 $\pm$ 0.0017	0.0629 $\pm$ 0.0037	-143.6067 $\pm$ 10.0592	0.8529 $\pm$ 0.0179
	<i>MMSB</i>	0.0243 $\pm$ 0.0105	0.0716 $\pm$ 0.0043	-129.4354 $\pm$ 7.6549	0.8561 $\pm$ 0.0176
	<i>iMMM</i>	0.0297 $\pm$ 0.0055	0.0625 $\pm$ 0.0015	-126.7876 $\pm$ 3.4774	0.8617 $\pm$ 0.0124
	<i>cMMSB<sup><math>\pi</math></sup></i>	0.0246 $\pm$ 0.0016	0.0489 $\pm$ 0.0016	-125.3876 $\pm$ 3.2689	0.8794 $\pm$ 0.0159
	<i>cMMSB<sup>uv</sup></i>	0.0283 $\pm$ 0.0035	0.0438 $\pm$ 0.0015	-123.3876 $\pm$ 3.1254	0.8738 $\pm$ 0.0364
Lazega lawfirm	<i>IRM</i>	0.0987 $\pm$ 0.0003	0.1046 $\pm$ 0.0012	-201.7912 $\pm$ 3.3500	0.7056 $\pm$ 0.0167
	<i>LFRM</i>	0.0566 $\pm$ 0.0024	0.1051 $\pm$ 0.0064	-222.5924 $\pm$ 16.1985	0.7970 $\pm$ 0.0197
	<i>MMSB</i>	0.0391 $\pm$ 0.0071	0.0913 $\pm$ 0.0030	-212.1256 $\pm$ 3.2145	0.7789 $\pm$ 0.0102
	<i>iMMM</i>	0.0487 $\pm$ 0.0068	0.1096 $\pm$ 0.0026	-202.7148 $\pm$ 5.3076	0.7874 $\pm$ 0.0141
	<i>cMMSB<sup><math>\pi</math></sup></i>	0.0246 $\pm$ 0.0050	0.1023 $\pm$ 0.0056	-201.0154 $\pm$ 5.2167	0.8273 $\pm$ 0.0148
	<i>cMMSB<sup>uv</sup></i>	0.0276 $\pm$ 0.0043	0.1143 $\pm$ 0.0019	-204.0289 $\pm$ 9.5460	0.8215 $\pm$ 0.0167

# The second work: Dynamic Infinite Mixed-Membership Stochastic Blockmodel

Summary of its advantage:

- ▶ allows the infinite number of communities;
- ▶ it allows mixed-membership for each node
- ▶ the model extends to the dynamic settings.
- ▶ it is apparent that in many social networking applications, a node's membership may become persistent over consecutive times, for example, a person's opinion of his peer is more likely to be consistent in two consecutive times.

# Continue with our model: Notations (1)

Continue, ...

- ▶  $E = \{e_{ij}^t\}_{n \times n}^{1:T}$ : entire set observations: if  $i$  has a relationship to node  $j$  at time  $t$ , it implies  $e_{ij}^t = 1$ . Otherwise,  $e_{ij}^t = 0$ .
- ▶ Each  $e_{ij}^t$  is specific to each pair of communities membership indicators  $(s_{ij}^t, r_{ij}^t)$ .
- ▶ For each node pair  $i$  and  $j$ , at  $t$ ,  $s_{ij}^t$  refers to the sender's community membership indicator.  $r_{ij}^t$  is for the receiver's community membership indicator.
- ▶  $Z$  to denote all the hidden labels  $\{s_{ij}^t, r_{ij}^t\}$ .

- ▶ Each node  $i$  at time  $t$ , has a mixed-membership distribution,  $\pi_i^t$  having infinite components, and the  $k^{\text{th}}$  component of  $\pi_i^t$ :  $\pi_{ik}^t$  represents the “significance” of community  $k$  for node  $i$ .
- ▶ Role-compatibility matrix  $W$ : its  $(k, l)^{\text{th}}$  entry, i.e.,  $W_{k,l}$  represents compatibilities between communities  $k$  and  $l$ . dimension of  $W$  can potentially be  $\infty \times \infty$ . Each  $W_{k,l}$  is i.i.d from  $Beta(\lambda_1, \lambda_2)$  which gives conjugacy to the Bernoulli distribution used to generate  $e_{ij}^t$

- ▶  $n_{k,l}^t$  denote the number of links from communities  $k$  to  $l$ , i.e., the number of times in which  $s_{ij}^t = k$  and  $r_{ij}^t = l$  simultaneously.  $n_{k,l}^t = n_{k,l}^{t,1} + n_{k,l}^{t,0}$ . **scalar**
- ▶  $n_{k,l}^{t,1}$  denotes the part of  $n_{k,l}^t$  where the corresponding  $e_{ij}^t = 1$ .
- ▶ The number of times that a node  $i$  has participated in community  $k$  (both as a sending and receiving) at time  $t$  is represented by  $N_{ik}^t$  **vector**

# Mixture Time Variant (MTV) and Mixture Time Invariant (MTI) Models

- ▶ To address the **phenomenon** that one's social community's memberships may change over times: allow each node's mixed-membership indicators to change cross times.
- ▶ Additionally, imperative that these indicators should have some persistence with its past values which reflects the reality of social behaviour.
- ▶ The persistence is achieved in two ways:
  - ▶ “**Mixture Time Variant (MTV)**” Mixed-membership distributions itself to change over times. Membership indicator of a node at time  $t$  is dependent on the “statistics” of all membership indicators of the same node at  $t - 1$  and  $t + 1$ .
  - ▶ “**Mixture Time Invariant (MTI)**” mixed-membership distributions to stay invariant over times. Membership indicator at time  $t$  is dependent and more likely to have the same value as it was in  $t - 1$ .

# Mixture Time Variant (MTV)

## Generative model

(1) Across all times  $1 : T$ :

- ▶  $\beta \sim GEM(\gamma)$
- ▶  $W_{k,l} \sim Beta(\lambda_1, \lambda_2) \forall k, l$

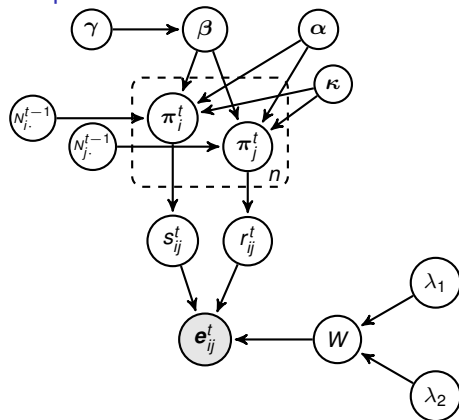
(2) Membership distributions:

- ▶  $\pi_i^t \sim DP\left(\alpha + \kappa, \frac{\alpha\beta + \frac{\kappa}{2n} \cdot \sum_k N_{ik}^{t-1} \delta_k}{\alpha + \kappa}\right)$

node  $i$ 's mixed membership distribution at  $t$ .

$N_{ik}^{t-1} = \sum_{l=1}^N \mathbf{1}(s_{il}^{t-1} = k) + \sum_{l=1}^N \mathbf{1}(r_{li}^{t-1} = k)$ , count number of nodes associated with a community  $k$  at time  $t - 1$ .

## Graphical model



## Generative model

(3) Relationship Sampling: For each pair of  $i, j \in \{1, \dots, n\}, t \in \{1, \dots, T\}$

- ▶  $s_{ij}^t \sim \text{Multi}(\pi_i^t)$ : sending community's indicator;
- ▶  $r_{ij}^t \sim \text{Multi}(\pi_j^t)$ : receiving community's indicator;
- ▶  $e_{ij}^t \sim \text{Bernoulli}(W_{s_{ij}^t, r_{ij}^t})$

Graphical model only shows for  $t$ , and omit the other times, where the structure is identical.

## Graphical model

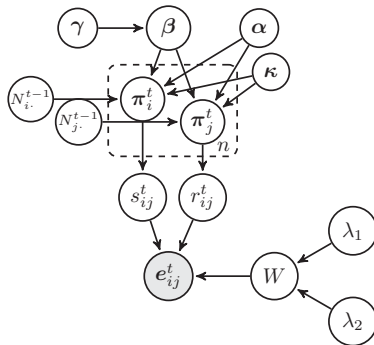


Figure: The MTV-DIM3 Model



# Mixture Time Variant (MTV): more explanation

- ▶  $\beta$  “global” representing the “significance” of all existing communities at all times, while  $W$  is the communities’ compatibility matrix.
- ▶ The prior  $P(W)$  is element-wise *Beta* distributed, which is conjugate to the Bernoulli distribution  $P(e_{i,j}^t|\cdot)$ : can obtain a marginal distribution of  $P(e_{i,j}^t) = \int_W p(e_{i,j}^t|W)p(W)d(W)$  analytically. Do not explicitly sample values of  $W$ .
- ▶ Mixed-membership distribution  $\{\pi_i^t\}_{1:n}^{1:T}$  is sampled from the Dirichlet Process with a concentration parameter  $(\alpha + \kappa)$  and a base measure  $\frac{\alpha\beta + \frac{\kappa}{2n} \sum_k N_{jk}^{t-1} \delta_k}{\alpha + \kappa}$ . There will be  $N \times T$  of these distributions. They jointly describe each node’s activities. It should be noted that each  $\pi_i^t$  is responsible to generate both the senders’ label  $\{s_{ij}^t\}_{j=1}^n$  from node  $i$  and receivers’ label  $\{r_{ji}^t\}_{j=1}^n$  to node  $i$ .

# Mixture Time Variant (MTV): persistency

- ▶ Sticky parameter  $\kappa$  stands for each node's time influence on its mixed-membership distribution.
- ▶ In another words, we assume that each node's mixed-membership distribution at time  $t$  will be largely influenced by its activities at time  $t - 1$ .
- ▶ This is reflected in the hidden label's multinomial distribution that the previous explicit activities will occupy a fixed proportion  $\frac{\kappa}{\alpha + \kappa}$  to the current distribution. The larger the value of  $\kappa$ , the more weight that the activities at  $t - 1$  is going to play at time  $t$ .

## Generative model

1. Across all times  $1 : T$ .
  - ▶  $\beta \sim GEM(\gamma)$
  - ▶  $W_{k,l} \sim Beta(\lambda_1, \lambda_2), \forall k, l$
2. Membership distribution
  - ▶  $\pi_i^{(k)} \sim DP\left(\alpha_i + \kappa, \frac{\alpha_i \beta + \kappa \delta_k}{\alpha_i + \kappa}\right)$
3. Relationship Sampling, For  $i, j \in \{1, \dots, n\}, t \in \{1, \dots, T\}$ 
  - ▶  $s_{ij}^t \sim Multi(\pi_i^{(s_{ij}^{t-1})})$
  - ▶  $r_{ij}^t \sim Multi(\pi_j^{(r_{ij}^{t-1})})$
  - ▶  $e_{ij}^t \sim Bernoulli(W_{s_{ij}^t, r_{ij}^t})$  relation from nodes  $i$  to  $j$  at time  $t$ .

## Graphical model

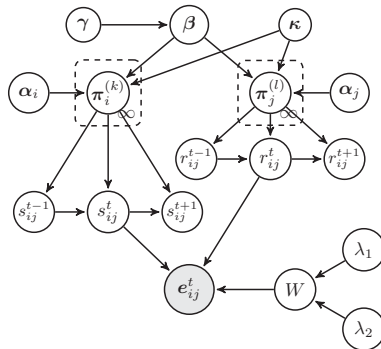


Figure: The MTI-DIM3 Model

# Mixture Time Invariant (MTI) Model: More explanation

- ▶ Each node has a variable number of membership distributions associated with, potentially infinite. At each time  $t$ , its membership indicator  $s_{ij}^t$  is generated from  $\pi_{s_{ij}^{t-1}}$ . In order to encourage persistence, each  $\pi_{ik}$  was generated from a corresponding  $\beta$ , where  $\kappa$  was added to  $\beta$ 's  $k^{\text{th}}$  component [?].
- ▶  $\beta$  and  $W$ 's generation is the same MTV The set of membership indicators  $\{s_{ij}^t, r_{ji}^t | j = 1, \dots, n, t = 1, \dots, T\}$  will be sampled from the time-invariant mixed-membership distribution set,  $\{\pi_i^{(k)}\}_{k=1}^{\infty}$ , where each member is independently distributed from a Dirichlet Process with a concentration parameter  $(\alpha + \kappa)$  and a base measure  $\frac{\alpha\beta + \kappa\delta_k}{\alpha + \kappa}$ .
- ▶ At time  $t$ , a membership indicator  $s_{ij}^t$  (or  $r_{ji}^t$ ) is sampled from the distribution  $\pi_i^{(s_{ij}^{t-1})}$  (or  $\pi_i^{(r_{ji}^{t-1})}$ )  $\forall i \in \{1, \dots, n\}$ .

Two sampling schemes are implemented for *MTV-DIM3*:

- ▶ Standard Gibbs sampling
- ▶ Slice-Efficient sampling

*MTI-DIM3* sampling is very similar, so will not duplicate in these slides.

Two sampling schemes are implemented for *MTV-DIM3*:

- ▶ Standard Gibbs sampling
- ▶ Slice-Efficient sampling

- ▶ Largely based on [Teh et al.(2006)Teh, Jordan, Beal, and Blei].
- ▶ Variables of interest are:  $\beta$ ,  $Z$  and auxiliary variables  $\hat{m}$ , where  $\hat{m}$  refers to the number of tables eating dish  $k$  without counting the tables that generated from the sticky portion, i.e.,  $\kappa N_{ik}^{t-1}$ .
- ▶ We do not sample  $\{\pi_{ij}^t\}_{1:n}^{1:T}$ , as it gets integrated out.

Sampling  $\beta$ :

$\beta$  is the prior for all  $\{\pi_i^t\}$ s, can be thought as the ratios between the community components for all communities. Posterior distribution is obtained through auxiliary variable  $\hat{\mathbf{m}}$ :

$$(\beta_1, \dots, \beta_K, \beta_\mu) \sim \text{Dir}(\hat{\mathbf{m}}_{\cdot 1}, \dots, \hat{\mathbf{m}}_{\cdot K}, \gamma)$$



# Gibbs Sampling details

Posterior  $P(s_{ij}^t = k, r_{ij}^t = l | Z \setminus \{s_{ij}^t, r_{ij}^t\}, \mathbf{e}, \boldsymbol{\beta}, \alpha, \lambda_1, \lambda_2, \kappa) \propto T_1 \times T_2 \times T_3$ :

$$T_1 = \frac{\Gamma(\alpha\beta_k + N_{ik}^{t+1} + \kappa N_{ik}^{t, -s_{ij}^t} + \kappa)}{\Gamma(\alpha\beta_k + N_{ik}^{t+1} + \kappa N_{ik}^{t, -s_{ij}^t})} \frac{\Gamma(\alpha\beta_k + \kappa N_{ik}^{t, -s_{ij}^t})}{\Gamma(\alpha\beta_k + \kappa N_{ik}^{t, -s_{ij}^t} + \kappa)} \times$$

$$\begin{cases} \alpha\beta_k + \kappa N_{ik}^{t-1} + N_{ik}^{t, -s_{ij}^t} & k \in \{1, \dots, K\} \\ \alpha\beta_u & k = K + 1 \end{cases}$$

$$T_2 = \frac{\Gamma(\alpha\beta_l + N_{jl}^{t+1} + \kappa N_{jl}^{t, -r_{ij}^t} + \kappa)}{\Gamma(\alpha\beta_l + N_{jl}^{t+1} + \kappa N_{jl}^{t, -r_{ij}^t})} \frac{\Gamma(\alpha\beta_l + \kappa N_{jl}^{t, -r_{ij}^t})}{\Gamma(\alpha\beta_l + \kappa N_{jl}^{t, -r_{ij}^t} + \kappa)} \times$$

$$\begin{cases} \alpha\beta_l + \kappa N_{jl}^{t-1} + N_{jl}^{t, -r_{ij}^t}, & l \in \{1, \dots, l\} \\ \alpha\beta_u, & l = l + 1 \end{cases}$$

$$T_3 = \begin{cases} \frac{n^{t, 1, -e_{ij}^t + \lambda_1}}{n^{t, -e_{ij}^t + \lambda_1 + \lambda_2}}, & \text{if } e_{ij}^t = 1 \\ \frac{n^{t, 0, -e_{ij}^t + \lambda_2}}{n^{t, -e_{ij}^t + \lambda_1 + \lambda_2}}, & \text{if } e_{ij}^t = 0 \end{cases}$$

Using the restaurant-table-dish analogy, we denote  $\mathbf{m}_{ik}^t$  as the number of tables eating dish  $k \forall i, k, t$ . This is related to the variable  $\hat{m}$  used in sampling  $\beta$ , but also including the counts of the “un-sticky” portion, i.e.,  $\alpha\beta_k$ .

The sampling of  $\mathbf{m}_{ik}^t$  is to incorporate a similar strategy as [Teh et al.(2006)Teh, Jordan, Beal, and Blei, ?], which is independently distributed from:

$$\begin{aligned} \Pr(\mathbf{m}_{ik}^t = m | \alpha, \beta_k, N_{ik}^{t-1}, \kappa) \\ \propto S(N_{ik}^t, m) (\alpha\beta_k + \kappa N_{ik}^{t-1})^m \end{aligned}$$

Here  $S(\cdot, \cdot)$  is the Stirling number of first kind.

For each node, the ratio of generating new tables can result from two factors: (1) Dirichlet prior with parameter  $\{\alpha, \beta\}$  and (2) the sticky configuration from membership indicators at  $t - 1$ , i.e.,  $\kappa N_{ik}^{t-1}$ .

To sample  $\beta$ , we need to only include tables generated from the “un-sticky” portion, i.e.,  $\hat{\mathbf{m}}$ , where each  $\hat{\mathbf{m}}_{ik}^t$  can be obtained from a single Binomial draw:

$$\hat{\mathbf{m}}_{ik}^t \sim \text{Binomial}(\mathbf{m}_{ik}^t, \frac{\alpha \beta_k}{\frac{\kappa}{2n} N_{ik}^{t-1} + \alpha \beta_k}).$$

- ▶ Incorporate the slice-efficient sampling [?][?]. The original sampling scheme was designed to sample the Dirichlet Process Mixture model.
- ▶ We use auxiliary variables  $U = \{u_{ij,s}^t, u_{ij,r}^t\}$  for each of the latent membership pair  $\{s_{ij}^t, r_{ij}^t\}$ . Having the  $U$ s, we are able to limit the number of components in which  $\pi_i$  needs to be considered, which is infinite otherwise.
- ▶ Under the slice-efficient sampling framework, the variables of interest are now extended to:  $\pi_i^t, \{u_{ij,r}^t, u_{ij,s}^t\}, \{s_{ij}^t, r_{ij}^t\}, \beta, \mathbf{m}$ :

- ▶ Incorporate the slice-efficient sampling [?][?]. The original sampling scheme was designed to sample the Dirichlet Process Mixture model.
- ▶ We use auxiliary variables  $U = \{u_{ij,s}^t, u_{ij,r}^t\}$  for each of the latent membership pair  $\{s_{ij}^t, r_{ij}^t\}$ . Having the  $U$ s, we are able to limit the number of components in which  $\pi_i$  needs to be considered, which is infinite otherwise.
- ▶ Under the slice-efficient sampling framework, the variables of interest are now extended to:  $\pi_i^t, \{u_{ij,r}^t, u_{ij,s}^t\}, \{s_{ij}^t, r_{ij}^t\}, \beta, \mathbf{m}$ :

For each node  $i = 1, \dots, N$ : we generate  $\pi_i'^t$  using sticky-breaking process [?], where each  $k^{\text{th}}$  component is generated using:

$\pi_{ik}'^t \sim \text{beta}(\pi_{ik}'^t; a_{ik}^t, b_{ik}^t)$ , where

$$a_{ik}^t = \alpha \beta_k + N_{ik}^t + \kappa N_{ik}^{t-1}$$

$$b_{ik}^t = \alpha \left(1 - \sum_{l=1}^k \beta_l\right) + N_{i, k_0 > k}^t + \kappa N_{i, k_0 > k}^{t-1}$$

Here  $\pi_k = \pi_k' \prod_{i=1}^{k-1} (1 - \pi_i')$ .

We use  $u_{ij,s}^t \sim U(0, \pi_{is_{ij}^t}^t)$ ,  $u_{ij,r}^t \sim U(0, \pi_{jr_{ij}^t}^t)$ . Then the obtained hidden label is independently sampled from the finite candidates:

$$\begin{aligned}
 & P(s_{ij}^t = k, r_{ij}^t = l | \mathbf{Z}, \mathbf{e}_{ij}^t, \beta, \alpha, \kappa, \mathbf{N}, \boldsymbol{\pi}, u_{ij,s}^t, u_{ij,r}^t) \\
 & \propto \mathbf{1}(k : \pi_{ik}^t > u_{ij,s}^t) \cdot \mathbf{1}(l : \pi_{jl}^t > u_{ij,r}^t) \\
 & \cdot \prod_{l=1}^{2n} P(z_{ij}^{t+1} | z_i^t, s_{ij}^t, s_{ij}^t = k, \beta, \alpha, \kappa, N_i^{t+1}) \\
 & \cdot \prod_{l=1}^{2n} P(z_{ij}^{t+1} | z_j^t, r_{ij}^t, r_{ij}^t = l, \beta, \alpha, \kappa, N_j^{t+1}) \\
 & \cdot P(\mathbf{e}_{ij}^t | E \setminus \{\mathbf{e}_{ij}^t\}, s_{ij}^t = k, r_{ij}^t = l, \mathbf{Z} \setminus \{s_{ij}^t, r_{ij}^t\}, \lambda_1, \lambda_2)
 \end{aligned}$$

This is the same as the Gibbs sampling.



- ▶ Hyper-parameters are  $\gamma, \alpha, \kappa$ . It is impossible to compute the posterior individually. Therefore, three prior are used:
  - ▶  $\mathcal{G}(1, 1)$  is placed on both  $\gamma$  and  $(\alpha + \kappa)$
  - ▶ A beta prior is placed on the ratio  $\frac{\kappa}{\alpha + \kappa}$ .
- ▶ Sample  $\gamma$  value, since  $\log(\gamma)$ 's posterior is log-concave (not concave!), use Adaptive Rejection Sampling (ARS) method [Rasmussen(2000)].
- ▶ Sample  $(\alpha + \kappa)$ , use auxiliary variable  $\mathbf{m}$  as proposed in [Teh et al.(2006)Teh, Jordan, Beal, and Blei].
- ▶ Sample  $\frac{\kappa}{\alpha + \kappa}$ , place  $\mathcal{B}(1, 1)$  on it, with a likelihood of  $\{\mathbf{m}_{ik}^t - \hat{\mathbf{m}}_{ik}^t, \forall i, k, t > 1\}$ , the posterior is in an analytical and samplable form, due to conjugacy.

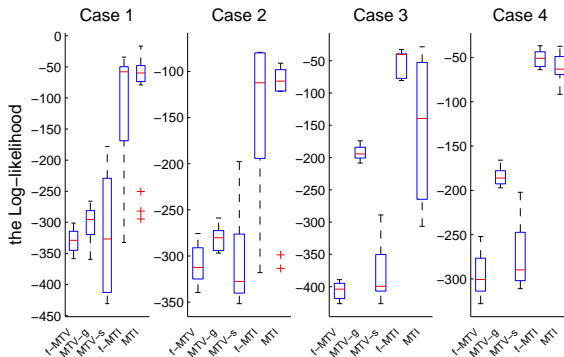
# Pros and Cons of Gibbs Sampling and Slice-Efficient Sampling

:

- ▶ Gibbs Sampling integrates out the mixed-membership distribution  $\{\pi_i^t\}$ . It is the “marginal approach” [Papaspiliopoulos and Roberts(2008)]. The property of community exchangeability makes it simple to implement. However, theoretically, it mix slowly as the sampling of each label is dependent on other labels.
- ▶ The Slice-Efficient Sampling is one “conditional approach” [?] while the membership indicators are independently sampled from  $\{\pi_i^t\}$ . In each iteration, given  $\{\pi_i^t\}$ , we can parallelize the process of sampling membership indicators, which may help to improve the computation, especially when the number of nodes, i.e.,  $N$  becomes larger, and the number of communities, i.e.,  $k$  becomes smaller.

- ▶ *DIM3* model run on synthetic datasets. Test against finite-communities case as baseline, namely *f-MTV* & *f-MTI*.
- ▶ For synthetic data generation, variables generated following [?]:  
 $N = 20, T = 3 \implies E$  is a  $20 \times 20 \times 3$  asymmetric and binary matrix. The parameters are set up such that, 20 nodes are equally partitioned into 4 groups. The ground-truth of the mixed-membership distribution for each of the groups are:  $[0.8, 0.2, 0; 0, 0.8, 0.2; 0.1, 0.05, 0.85; 0.4, 0.4, 0.2]$ .
- ▶ Consider 4 different test case role-compatibility matrix:
  - ▶ **Case 1:** large diagonal values and small non-diagonal values
  - ▶ **Case 2:** large diagonal values and mediate non-diagonal values
  - ▶ **Case 3:** large non-diagonal values and small diagonal values
  - ▶ **Case 4:** small diagonal values and mediate non-diagonal values

# Results: Log-likelihood Performance



From the log-likelihood comparison, we can see that the *MTI* model performs better than the *MTV* model.








The average  $l_2$  distance between the mixed-membership distributions and its ground-truth; and the one between the posterior role-compatibility matrix and its ground-truth:

Table: Average  $l_2$  Distance to the Ground-truth

Cases	Role-Compatibility Matrix					Mixed-Memberships				
	<i>f</i> -MTV	MTV- <i>g</i>	MTV- <i>s</i>	<i>f</i> -MTI	MTI	<i>f</i> -MTV	MTV- <i>g</i>	MTV- <i>s</i>	<i>f</i> -MTI	MTI
1	0.529	0.625	0.848	0.114	<b>0.086</b>	0.366	0.384	0.403	0.199	<b>0.191</b>
2	0.439	0.225	0.339	<b>0.195</b>	0.204	0.355	0.355	0.319	<b>0.207</b>	0.227
3	0.134	0.201	0.513	0.117	<b>0.087</b>	0.278	0.289	0.589	0.208	<b>0.187</b>
4	<b>0.195</b>	0.214	0.267	0.220	0.219	0.258	0.285	0.277	0.192	<b>0.182</b>

On the average  $l_2$  distance to the ground-truth performance, the *MTI* model also performs better.

- ▶ MCMC diagnostics
- ▶ Test real relational data
- ▶ Vision applications?

-  Airoldi, E., D. Blei, S. Fienberg, and E. Xing, 2008: Mixed membership stochastic blockmodels. *The Journal of Machine Learning Research*, **9**, 1981–2014.
-  Kim, D., M. Hughes, and E. Sudderth, 2012: The nonparametric metadata dependent relational model. *Proceedings of the 29th Annual International Conference on Machine Learning*, ACM.
-  Koutsourelakis, P. and T. Eliassi-Rad, 2008: Finding mixed-memberships in social networks. *Proceedings of the 2008 AAAI spring symposium on social information processing*.
-  Papaspiliopoulos, O. and G. Roberts, 2008: Retrospective markov chain monte carlo methods for dirichlet process hierarchical models. *Biometrika*, **95** (1), 169–186.
-  Rasmussen, C., 2000: The infinite gaussian mixture model. *Advances in neural information processing systems*, **12** (5.2), 2.
-  Teh, Y., M. Jordan, M. Beal, and D. Blei, 2006: Hierarchical dirichlet processes. *Journal of the American Statistical Association*, **101** (476), 1566–1581.
-  Xing, E., W. Fu, and L. Song, 2010: A state-space mixed membership blockmodel for dynamic network tomography. *The Annals of Applied Statistics*, **4** (2), 535–566.