# Hierarchical Dirichlet Processes

AMS 241, Fall 2010

Vadim von Brzeski

vvonbrze@ucsc.edu

# Reference

- *Hierarchical Dirichlet Processes*, Y. Teh, M. Jordan, M. Beal, D. Blei, Technical Report 653, Statistics, UC Berkeley, 2004.
  - Also published in NIPS 2004 : *Sharing Clusters among Related Groups: Hierarchical Dirichlet Processes*

- Some figures and equations shown here are directly taken from the above references (indicated if so)

# The HDP Prior

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

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

*group index*

*$G_0$ is discrete :*
*$G_j$ DPs necessarily*
*share atoms !*

Source: Teh, 2004.

## Stick Breaking Construction:

**1st Level :**

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

$$\beta_k' \sim \text{Beta}(1, \gamma) \qquad \beta_k = \beta_k' \prod_{l=1}^{k-1} (1 - \beta_l')$$

**2nd Level :**

$$G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k}$$

Going back to original definition of DP, we can derive relationship between β and π :

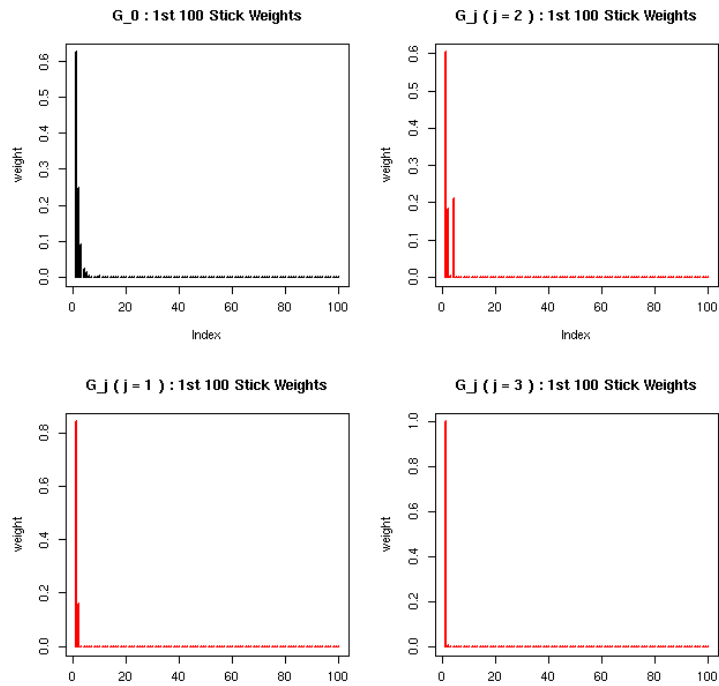$$(G_j(A_1), ..., G_j(A_r)) \sim Dirichlet\left(\alpha_0 G_0(A_1), ..., \alpha_0 G_0(A_r)\right)$$

$$\left(\sum_{k \in K_1} \pi_{jk}, ..., \sum_{k \in K_r} \pi_{jk}\right) \sim Dirichlet\left(\alpha_0 \sum_{k \in K_1} \beta_k, ..., \alpha_0 \sum_{k \in K_r} \beta_k\right) \qquad K_l = \{k : \theta_k \in A_l\}$$
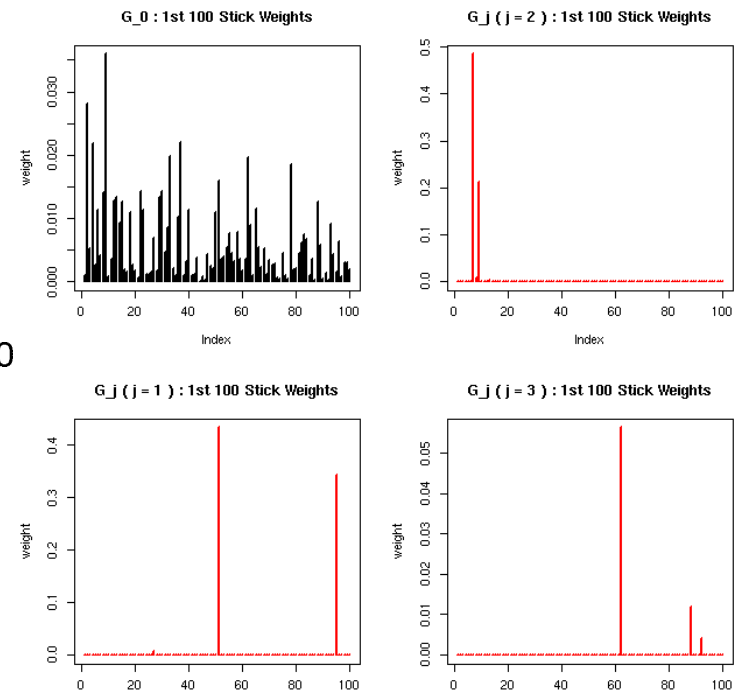
$$\pi_{jk}' \sim \text{Beta}\left(\alpha_0 \beta_k, \alpha_0 \left(1 - \sum_{l=1}^{k} \beta_l\right)\right) \qquad \pi_{jk} = \pi_{jk}' \prod_{l=1}^{k-1} (1 - \pi_{jl}')$$
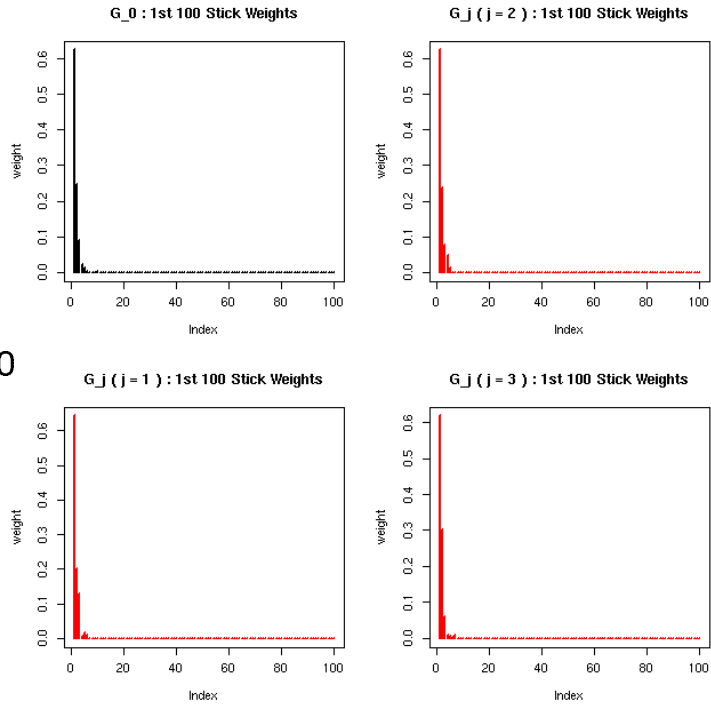
Source: Teh, 2004.

γ = 1
α₀ = 1

γ = 100
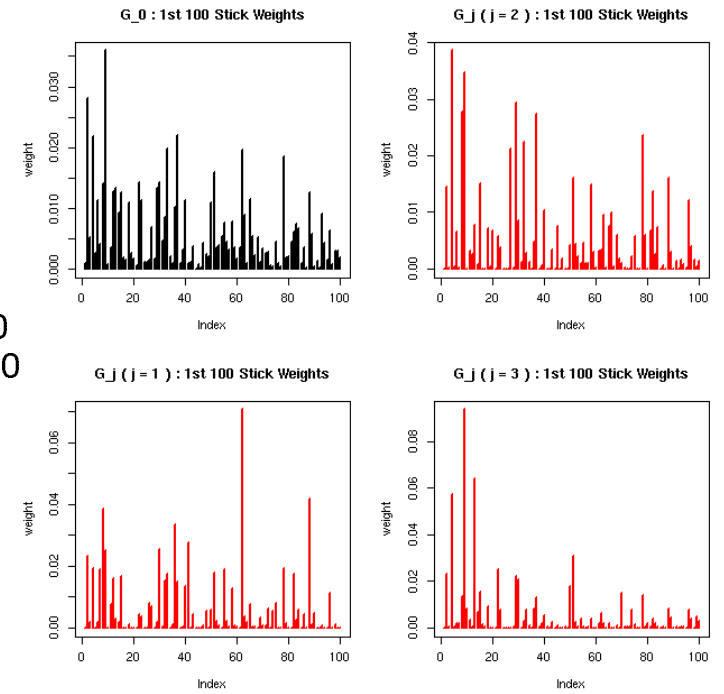α₀ = 1

γ = 1
α₀ = 100

γ = 100
α₀ = 100

# H : Normal(0,1)



G_0 and G_j draws (via stick-breaking)
alpha0 = 1  gamma = 1

G_0 and G_j draws (via stick-breaking)
alpha0 = 1  gamma = 100

G_0 and G_j draws (via stick-breaking)
alpha0 = 100  gamma = 1

G_0 and G_j draws (via stick-breaking)
alpha0 = 100  gamma = 100

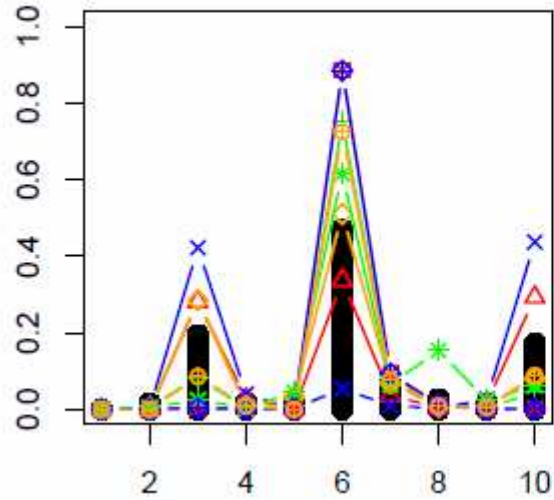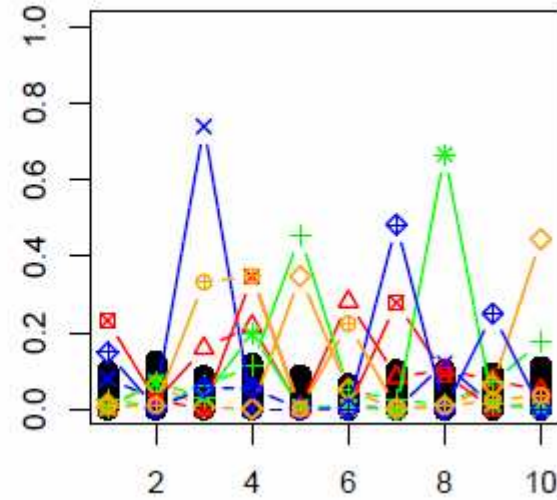$G_0$

$G_j$

6

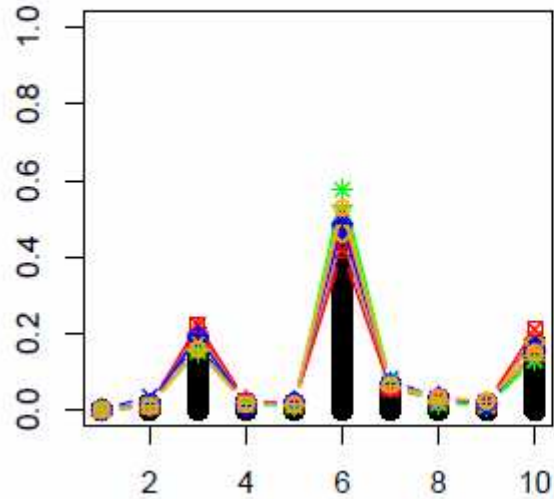# H : Dirichlet(0.1,0.1,.....,0.1), dim V=10



G_0 and G_j draws (via stick-breaking)
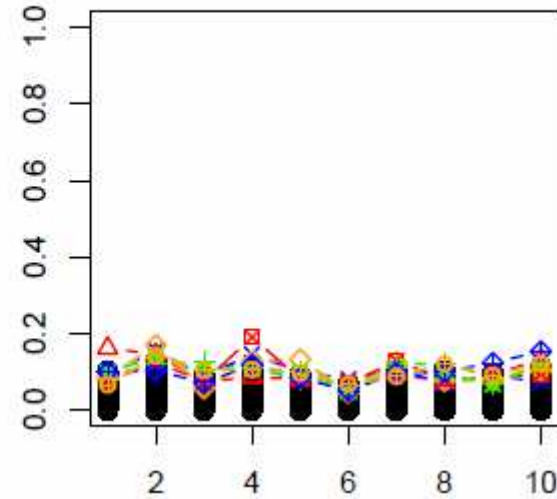alpha0 = 1  gamma = 1

G_0 and G_j draws (via stick-breaking)
alpha0 = 1  gamma = 100

G_0 and G_j draws (via stick-breaking)
alpha0 = 100  gamma = 1

G_0 and G_j draws (via stick-breaking)
alpha0 = 100  gamma = 100

$G_0$

# Prior and Data Model

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

$$G_j \mid \alpha_0, G_0 \sim \mathrm{DP}(\alpha_0, G_0)$$
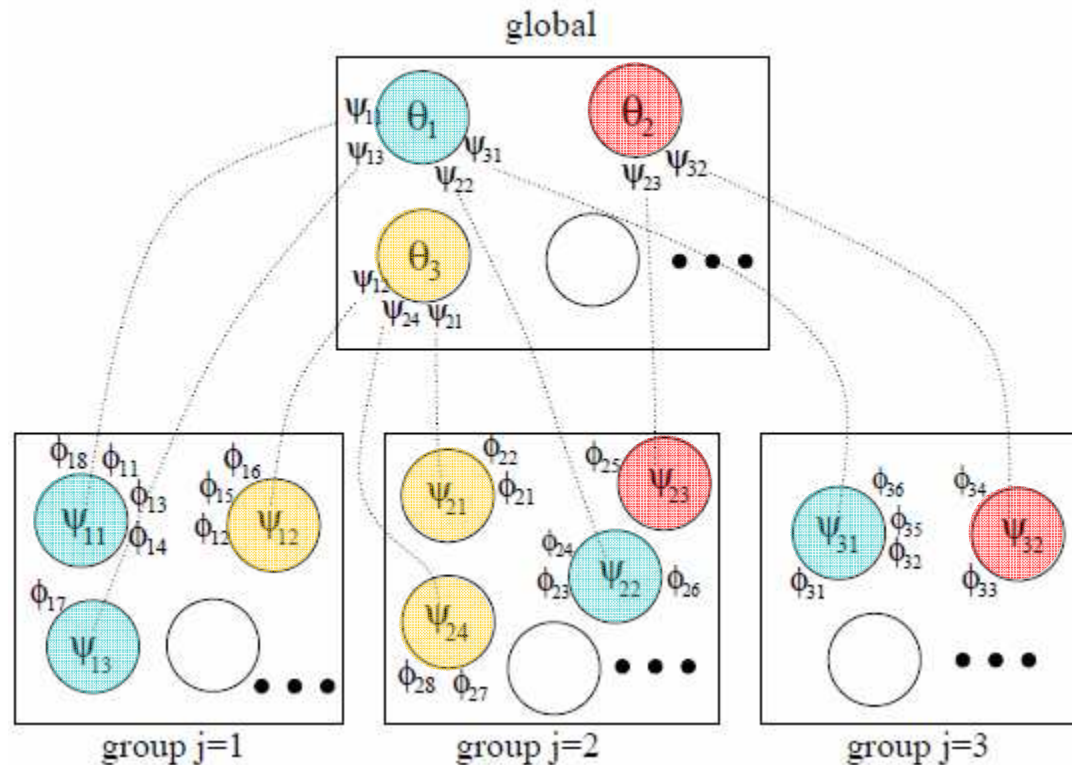
$$\phi_{ji} \mid G_j \sim G_j$$

$$x_{ji} \mid \phi_j \sim F(\phi_{ji})$$

*$i^{th}$ datum in group j*

Source: Teh, 2004.

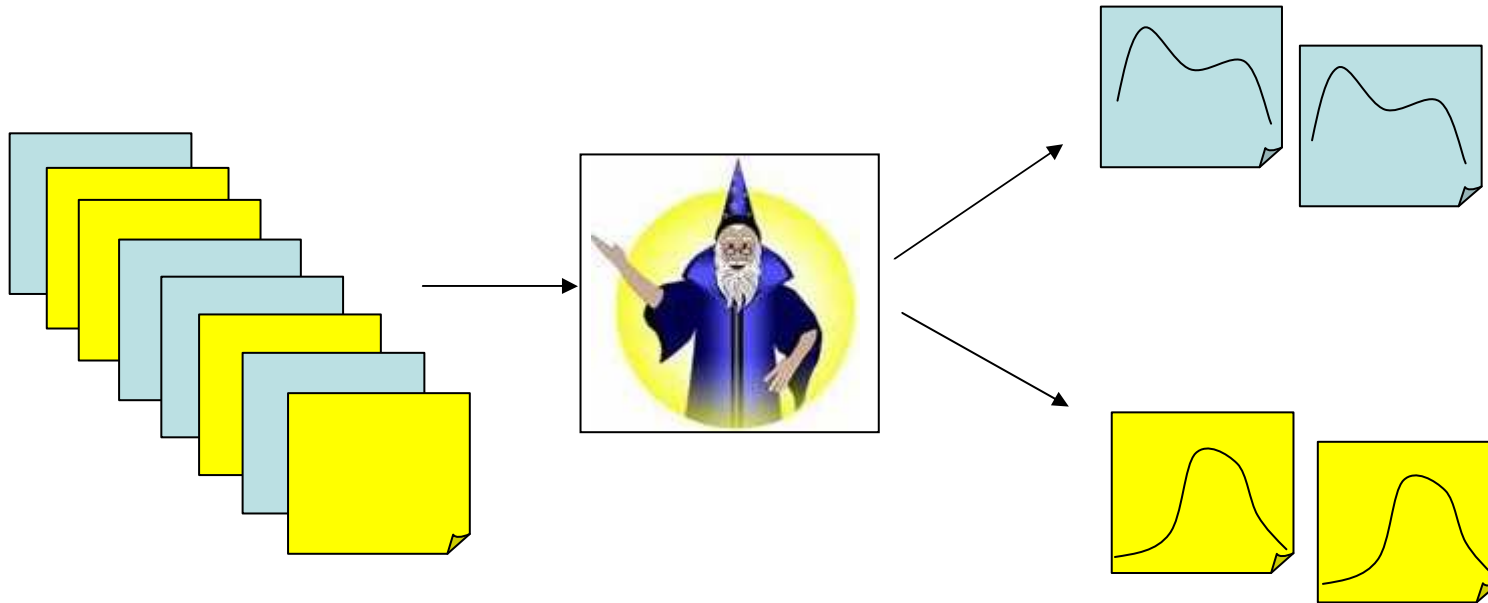# Polya Urn Sampling via Chinese Restaurant (Process) Franchise :

existing table         new table

$$\phi_{ji} \mid \phi_{j1}, \ldots, \phi_{j\,i-1}, \alpha_0, G_0 \sim \sum_{t=1}^{T_j} \frac{n_{jt}}{i-1+\alpha_0} \delta_{\psi_{jt}} + \frac{\alpha_0}{i-1+\alpha_0} G_0$$

existing component        new component

$$\psi_{jt} \mid \psi_{11}, \psi_{12}, \ldots, \psi_{21}, \ldots, \psi_{j\,t-1}, \gamma, H \sim \sum_{k=1}^{K} \frac{m_k}{\sum_k m_k + \gamma} \delta_{\theta_k} + \frac{\gamma}{\sum_k m_k + \gamma} H$$
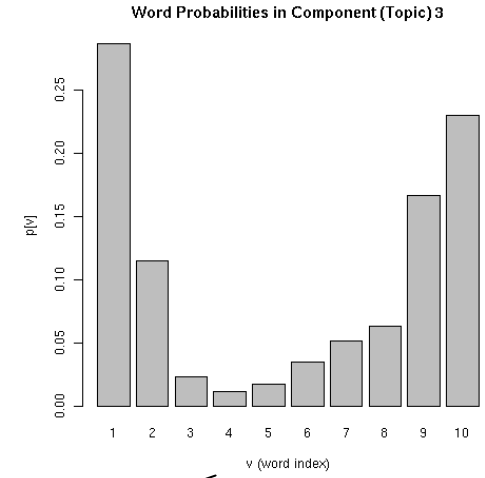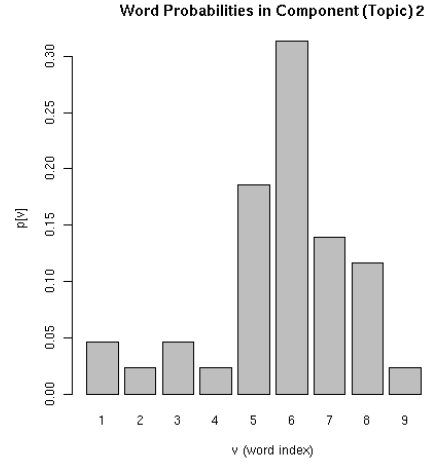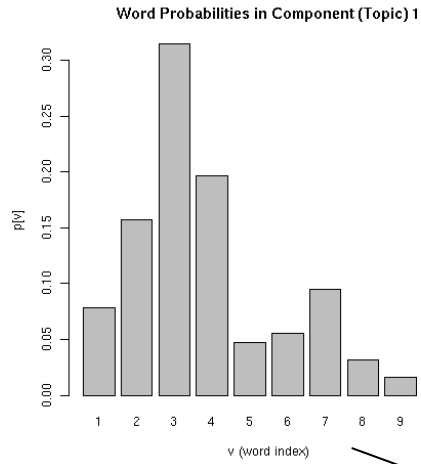


Source: Teh, 2004.

# Application : Topic Modeling

- Topic = (multinomial) distribution over words
  - Fixed size vocabulary; p(word | topic)
  - F : Multinomial kernel, H : Dirichlet()
- Document = mixture of one or more topics
- Goal = recover latent topics; use topics for clustering, finding related documents, etc.
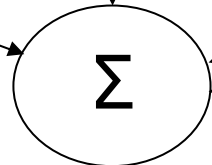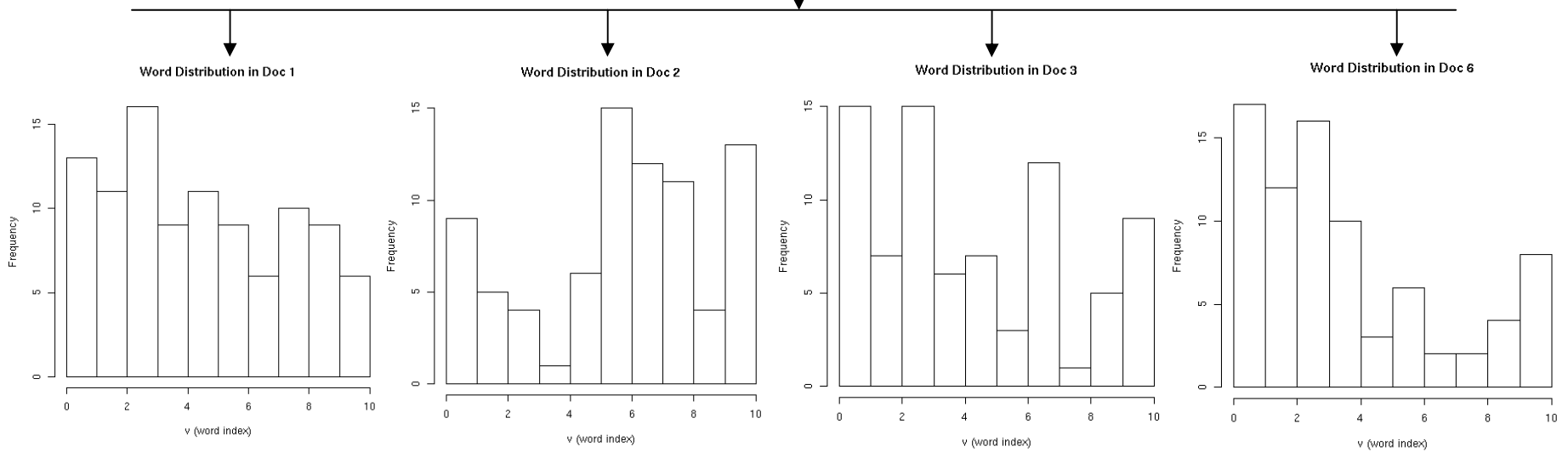
# Study Model Inference Using Simulated Data



**3 TRUE TOPICS**

Word Probabilities in Component (Topic) 1

Word Probabilities in Component (Topic) 2

Word Probabilities in Component (Topic) 3

J = 6 docs (80 – 100 words / doc)
2 – 3 mix components / doc
V (vocabulary size) = 10

$\Sigma$

$p = [0.4, 0.3, 0.3]$

Word Distribution in Doc 1

Word Distribution in Doc 2

Word Distribution in Doc 3

Word Distribution in Doc 6

# Inference via Gibbs Sampling

1. $$p(t_{ji} = t | \boldsymbol{t}^{-ji}, \boldsymbol{k}, \boldsymbol{\theta}, \boldsymbol{x}) \propto \begin{cases} \alpha_0 f(x_{ji} | \theta_{k_{jt}}) & \text{if } t = t^{\text{new}}, \\ n_{jt}^{-i} f(x_{ji} | \theta_{k_{jt}}) & \text{if } t \text{ previously used.} \end{cases}$$

$$k_{jt^{\text{new}}} \mid \boldsymbol{k} \sim \sum_{k=1}^{K} \frac{m_k}{\sum_k m_k + \gamma} \delta_k + \frac{\gamma}{\sum_k m_k + \gamma} \delta_{k^{\text{new}}}$$
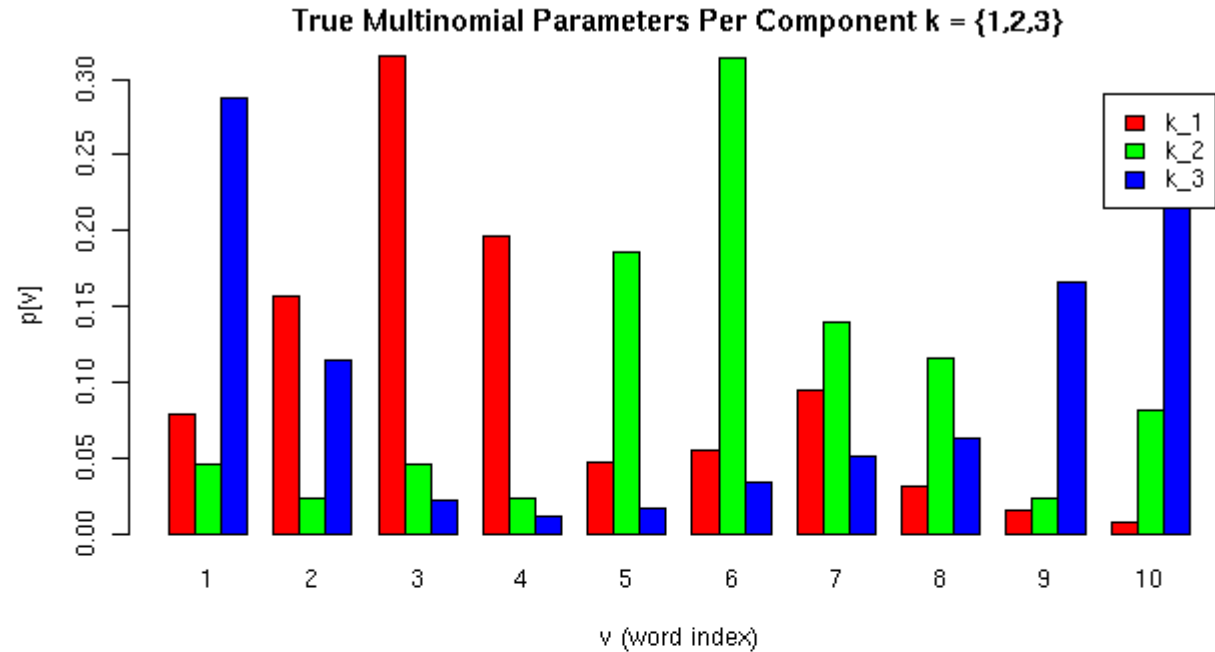
$$\theta_{k^{\text{new}}} \sim H$$

2. $$p(k_{jt} = k | \boldsymbol{t}, \boldsymbol{k}^{-jt}, \boldsymbol{\theta}, \boldsymbol{x}) \propto \begin{cases} \gamma \prod_{i:t_{ji}=t} f(x_{ji} | \theta_k) & \text{if } k = k^{\text{new}}, \\ m_k^{-t} \prod_{i:t_{ji}=t} f(x_{ji} | \theta_k) & \text{if } k \text{ is previously used.} \end{cases}$$

$$\theta_{k^{\text{new}}} \sim H$$

3. $$p(\theta_k | \boldsymbol{t}, \boldsymbol{k}, \boldsymbol{\theta}^{-k}, \boldsymbol{x}) \propto h(\theta_k) \prod_{ji:k_{jt_{ji}}=k} f(x_{ji} | \theta_k)$$

Source: Teh, 2004.

TRUTH :



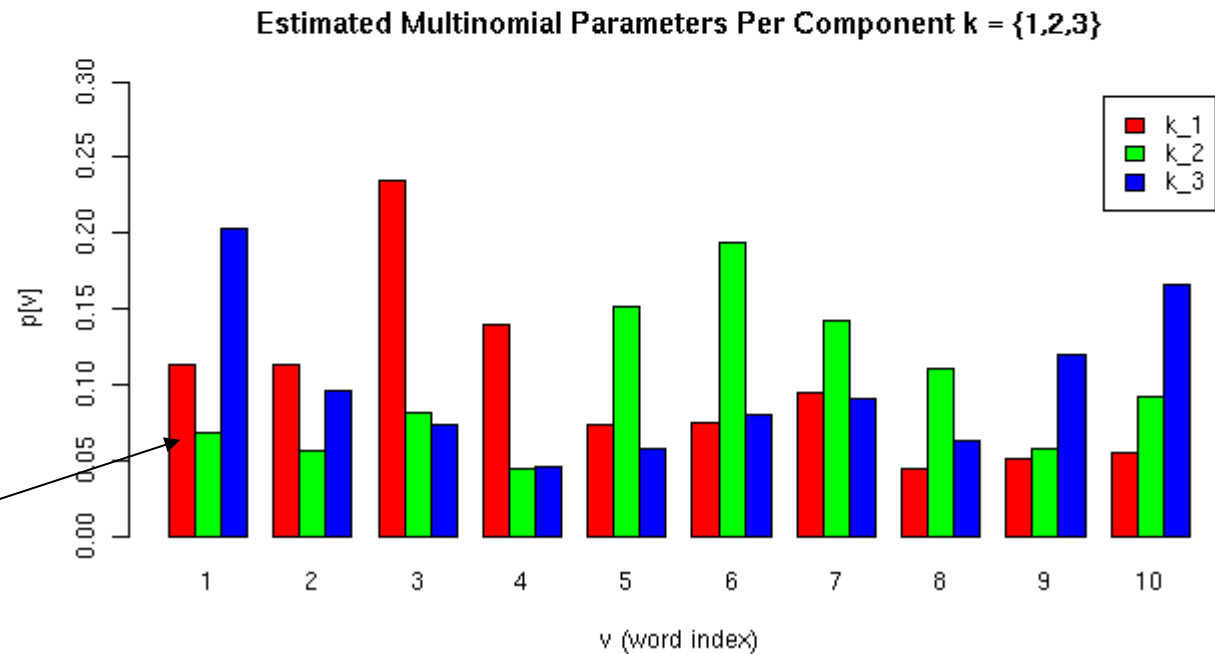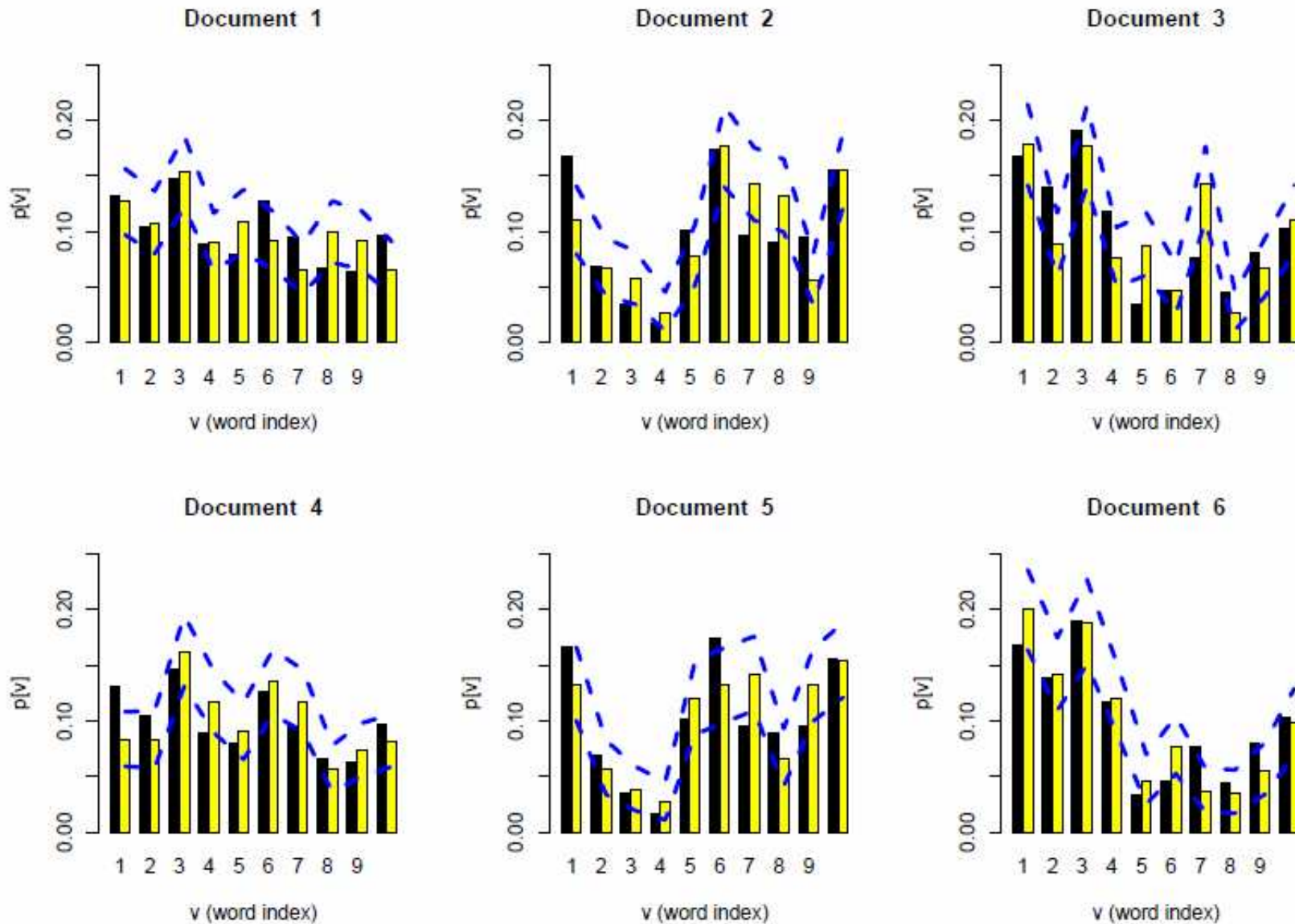True Multinomial Parameters Per Component k = {1,2,3}

ESTIMATE :

For each $\mathbf{x_{ji}}$ whose true
component was k, we have
B MCMC draws:
$\{\theta_{ji}^{(1)}, \theta_{ji}^{(2)}, \ldots, \theta_{ji}^{(B)}\}$

$$\overline{\theta_{ji}^{(B)}} = \frac{1}{B} \sum_b \theta_{ji}^{(b)}$$

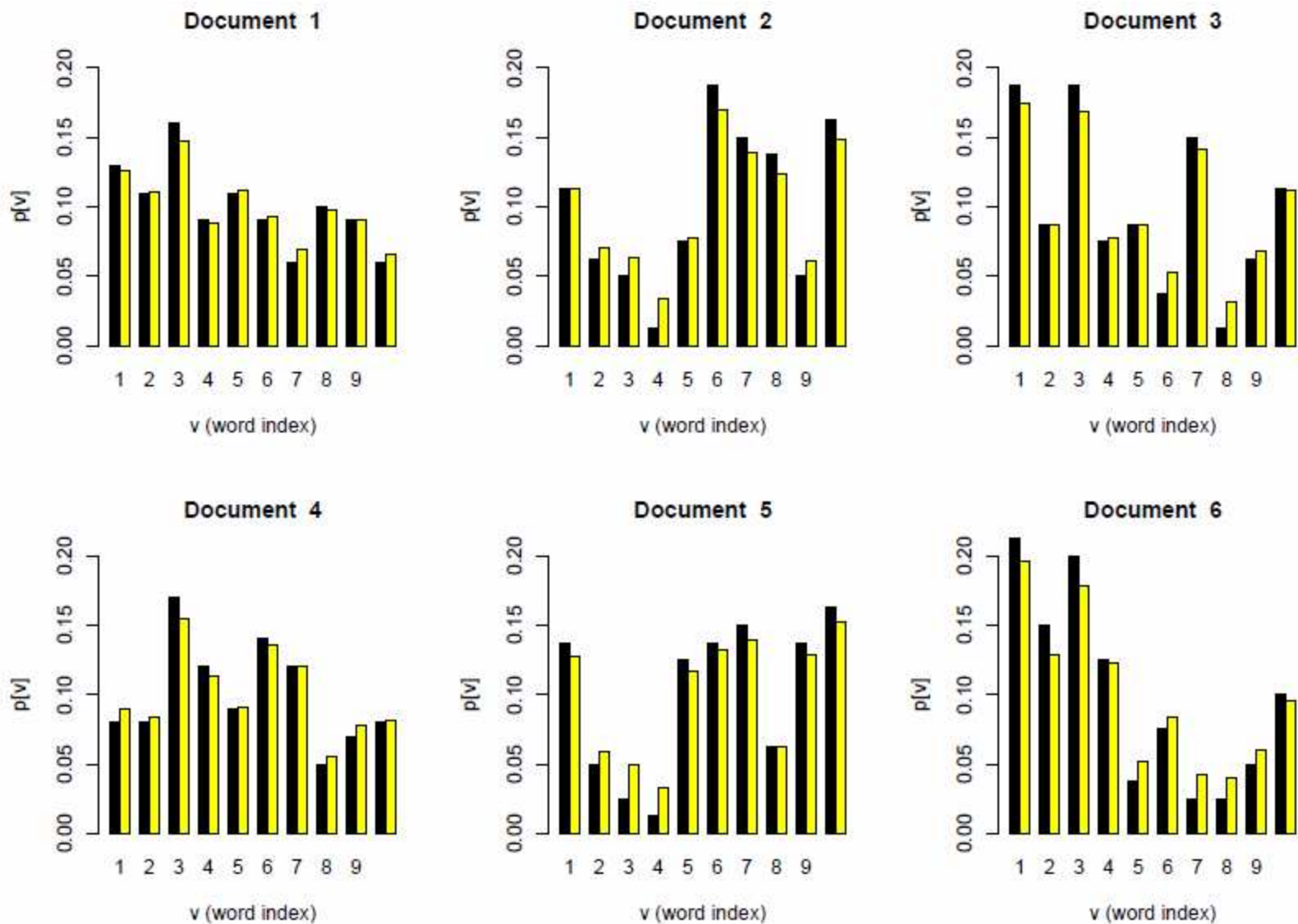$$\overline{\theta_k} = \frac{1}{n_k} \sum \overline{\theta_{ji}^{(B)}}$$



Estimated Multinomial Parameters Per Component k = {1,2,3}

# Truth vs. Posterior Point and 10/90 Interval Estimates for E[ $\theta_j$ | data ]

■ True $\theta_j$   ☐ Estimate

# Simulated Data Histograms vs. Est. Posterior Predictive : E[ $\theta_{j0}$ | data ]
## For each doc j : avg (over states b = 1..B) draws of $\theta_{j0}^{(b)}$ via CRP config @ state b.

# Simulated Data Distributions vs. Est. Posterior Predictive for New Observation $x_{j0}$

# R Code Available

- Works, but SLOOOOOOOOOW....

http://www.numberjack.net/download/classes/ams241/project/R