

# Collapsed Gibbs sampling for latent Dirichlet allocation

正田 備也

[masada@rikkyo.ac.jp](mailto:masada@rikkyo.ac.jp)

# Contents

# Joint distribution

$$\begin{aligned} & p(\mathbf{X}, \mathbf{Z}, \Theta, \Phi; \alpha, \beta) \\ &= \prod_{d=1}^D p(\boldsymbol{\theta}_d; \alpha) \times \prod_{k=1}^K p(\boldsymbol{\phi}_k; \beta) \times \prod_{d=1}^D \prod_{i=1}^{n_d} p(z_{d,i} | \boldsymbol{\theta}_d) p(x_{d,i} | \boldsymbol{\phi}_{z_{d,i}}) \quad (1) \end{aligned}$$

# Marginalization

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z}; \alpha, \beta) &= \int p(\mathbf{X}, \mathbf{Z}, \Theta, \Phi; \alpha, \beta) d\Theta d\Phi \\ &= \int \prod_{d=1}^D \left[ p(\boldsymbol{\theta}_d; \alpha) \prod_{i=1}^{n_d} p(z_{d,i} | \boldsymbol{\theta}_d) \right] d\Theta \\ &\quad \times \int \left[ \prod_{k=1}^K p(\boldsymbol{\phi}_k; \beta) \prod_{d=1}^D \prod_{i=1}^{n_d} p(x_{d,i} | \boldsymbol{\phi}_{z_{d,i}}) \right] d\Phi \\ &= \prod_d \int p(\mathbf{z}_d, \boldsymbol{\theta}_d; \alpha) d\boldsymbol{\theta}_d \times \int p(\mathbf{X}, \Phi | \mathbf{Z}; \alpha, \beta) d\Phi \\ &= \prod_d p(\mathbf{z}_d; \alpha) \times p(\mathbf{X} | \mathbf{Z}; \beta) \end{aligned} \tag{2}$$

# Topic assignments probability

$$\begin{aligned} p(\mathbf{z}_d, \boldsymbol{\theta}_d; \alpha) &= p(\boldsymbol{\theta}_d; \alpha) \prod_{i=1}^{n_d} p(z_{d,i} | \boldsymbol{\theta}_d) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_{d,k}^{\alpha_k - 1} \times \frac{(\sum_k n_{d,k})!}{\prod_k n_{d,k}!} \prod_k \theta_{d,k}^{n_{d,k}} \\ &= \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{n_d!}{\prod_k n_{d,k}!} \prod_k \theta_{d,k}^{n_{d,k} + \alpha_k - 1} \end{aligned} \quad (3)$$

where  $n_{d,k} \equiv \sum_{i=1}^{n_d} \delta(z_{d,i} = k)$ .

$$\begin{aligned} \int \left[ p(\boldsymbol{\theta}_d; \alpha) \prod_{i=1}^{n_d} p(z_{d,i} | \boldsymbol{\theta}_d) \right] d\boldsymbol{\theta}_d &= \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{n_d!}{\prod_k n_{d,k}!} \int \prod_k \theta_{d,k}^{n_{d,k} + \alpha_k - 1} d\boldsymbol{\theta}_d \\ &= \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{n_d!}{\prod_k n_{d,k}!} \frac{\prod_k \Gamma(n_{d,k} + \alpha_k)}{\Gamma(n_d + \sum_k \alpha_k)} \end{aligned} \quad (4)$$

$$\therefore p(\mathbf{z}_d; \alpha) = \int p(\mathbf{z}_d, \boldsymbol{\theta}_d; \alpha) d\boldsymbol{\theta}_d = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{n_d!}{\prod_k n_{d,k}!} \frac{\prod_k \Gamma(n_{d,k} + \alpha_k)}{\Gamma(n_d + \sum_k \alpha_k)} \quad (5)$$

# Topic posterior

$$\begin{aligned} p(\boldsymbol{\theta}_d | \mathbf{z}_d; \alpha) &= \frac{p(\mathbf{z}_d, \boldsymbol{\theta}_d; \alpha)}{p(\mathbf{z}_d; \alpha)} \\ &= \frac{\frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{n_d!}{\prod_k n_{d,k}!} \prod_k \theta_{d,k}^{n_{d,k} + \alpha_k - 1}}{\frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{n_d!}{\prod_k n_{d,k}!} \frac{\prod_k \Gamma(n_{d,k} + \alpha_k)}{\Gamma(n_d + \sum_k \alpha_k)}} \\ &= \frac{\Gamma(n_d + \sum_k \alpha_k)}{\prod_k \Gamma(n_{d,k} + \alpha_k)} \prod_k \theta_{d,k}^{n_{d,k} + \alpha_k - 1} \end{aligned} \tag{6}$$

# Word tokens probability

$$\begin{aligned}
 p(\mathbf{X}, \Phi | \mathbf{Z}; \beta) &= \prod_{k=1}^K p(\phi_k; \beta) \prod_{d=1}^D \prod_{i=1}^{n_d} p(x_{d,i} | \phi_{z_{d,i}}) = \prod_{k=1}^K \left[ p(\phi_k; \beta) \prod_{d=1}^D \prod_{i=1}^{n_d} p(x_{d,i} | \phi_k)^{\delta(z_{d,i}=k)} \right] \\
 &= \prod_{k=1}^K \left[ \frac{\Gamma(\sum_w \beta_w)}{\prod_w \Gamma(\beta_w)} \prod_{w=1}^W \phi_{k,w}^{\beta_w - 1} \times \frac{n_k!}{\prod_w n_{k,w}!} \prod_w \phi_{k,w}^{n_{k,w}} \right] \\
 &= \prod_{k=1}^K \left[ \frac{\Gamma(\sum_w \beta_w)}{\prod_w \Gamma(\beta_w)} \frac{n_k!}{\prod_w n_{k,w}!} \prod_w \phi_{k,w}^{n_{k,w} + \beta_w - 1} \right] \tag{7}
 \end{aligned}$$

$$\begin{aligned}
 \int p(\mathbf{X}, \Phi | \mathbf{Z}; \beta) d\Phi &= \prod_{k=1}^K \left[ \frac{\Gamma(\sum_w \beta_w)}{\prod_w \Gamma(\beta_w)} \frac{n_k!}{\prod_w n_{k,w}!} \int \prod_w \phi_{k,w}^{n_{k,w} + \beta_w - 1} d\theta_k \right] \\
 &= \prod_{k=1}^K \left[ \frac{\Gamma(\sum_w \beta_w)}{\prod_w \Gamma(\beta_w)} \frac{n_k!}{\prod_w n_{k,w}!} \frac{\prod_w \Gamma(n_{k,w} + \beta_w)}{\Gamma(n_k + \sum_w \beta_w)} \right] \tag{8}
 \end{aligned}$$

## Word posterior

$$\begin{aligned} p(\Phi|\mathbf{X}, \mathbf{Z}, \beta) &= \frac{p(\mathbf{X}, \Phi|\mathbf{Z}; \beta)}{p(\mathbf{X}|\mathbf{Z}; \beta)} = \frac{p(\mathbf{X}, \Phi|\mathbf{Z}; \beta)}{\int p(\mathbf{X}, \Phi|\mathbf{Z}; \beta) d\Phi} \\ &= \frac{\prod_{k=1}^K \left[ \frac{\Gamma(\sum_w \beta_w)}{\prod_w \Gamma(\beta_w)} \frac{n_k!}{\prod_w n_{k,w}!} \prod_w \phi_{k,w}^{n_{k,w} + \beta_w - 1} \right]}{\prod_{k=1}^K \left[ \frac{\Gamma(\sum_w \beta_w)}{\prod_w \Gamma(\beta_w)} \frac{n_k!}{\prod_w n_{k,w}!} \frac{\prod_w \Gamma(n_{k,w} + \beta_w)}{\Gamma(n_k + \sum_w \beta_w)} \right]} \\ &= \prod_{k=1}^K \left[ \frac{\Gamma(n_k + \sum_w \beta_w)}{\prod_w \Gamma(n_{k,w} + \beta_w)} \prod_w \phi_{k,w}^{n_{k,w} + \beta_w - 1} \right] \end{aligned} \quad (9)$$



# Per-token topic assignment posterior (1/3)

$$\begin{aligned} p(z_{d,i} = k | x_{d,i} = v, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}, \alpha, \beta) &= \frac{p(z_{d,i} = k, x_{d,i} = v, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta)}{p(x_{d,i} = v, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta)} \\ &= \frac{p(z_{d,i} = k, x_{d,i} = v, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta)}{\sum_{k'=1}^K p(z_{d,i} = k', x_{d,i} = v, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta)} \\ &\propto p(z_{d,i} = k, x_{d,i} = v, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) \\ &= p(x_{d,i} = v, z_{d,i} = k | \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) p(\mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) \\ &= p(x_{d,i} = v | z_{d,i} = k, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) p(z_{d,i} = k | \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) p(\mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) \\ &\propto p(x_{d,i} = v | z_{d,i} = k, \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) p(z_{d,i} = k | \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) \end{aligned} \tag{10}$$

## Per-token topic assignment posterior (2/3)

$$p(z_{d,i} = k | \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) = \int p(z_{d,i} = k | \boldsymbol{\theta}_i) p(\boldsymbol{\theta}_i | \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) d\boldsymbol{\theta}_i \quad (11)$$

$$\begin{aligned} p(\boldsymbol{\theta}_i | \mathbf{X}^{\setminus d,i}, \mathbf{Z}^{\setminus d,i}; \alpha, \beta) &= p(\boldsymbol{\theta}_d | \mathbf{z}_d^{\setminus d,i}; \alpha) \\ &= \frac{\Gamma(n_d^{\setminus d,i} + \sum_k \alpha_k)}{\prod_k \Gamma(n_{d,k}^{\setminus d,i} + \alpha_k)} \prod_k \theta_{d,k}^{n_{d,k}^{\setminus d,i} + \alpha_k - 1} \end{aligned} \quad (12)$$