

Text

Yuki Shiraito

POLSCI 798 Advanced Topics in Quantitative Methodology
University of Michigan

Winter 2019

Analysis of Text

- Measurement using text (Catalinac, 2014. "From Pork to Policy.")
 - Does intraparty competition leads to particularistic campaigning?
 - Measuring the type of campaigns: **Topics** in manifestos
 - Whether mention particularistic or programmatic goods

Number	Topic Label	%	Classification
1	postal privatization	98.89	prog
2	from concrete to people	98.89	prog
3	appropriator for the district	0.379	part
4	fixer-upper for the district	0.379	part
5	political reform, protect the constitution	98.89	prog
6	foreign and national security policy	98.89	prog
7	no more unfair taxes, peace constitution	98.89	prog
8	building a society kind to women	51.00	prog
9	primary industries and tourism	0.379	part
10	transportation	0.379	part
11	no tax increases, no U.S.-Japan alliance	98.89	prog
12	economic recovery	46.21	prog
13	vision for Japan	98.89	prog
14	politics for the civilian, not for bureaucrats	98.89	prog
15	political and administrative reform	98.89	prog

Catalinac (2014) Table D4 in Online Appendix

Modeling Text Data

- Notation

- *Vocabulary*: Set of all unique words, $\{1, \dots, L\}$
- *Word*: L -dimensional standard basis vector, $(0, \dots, 0, 1, 0, \dots, 0)^\top$
- *Document*: Sequence of N_j words, $\mathbf{W}_j = (W_{1j}, \dots, W_{N_jj})$
- *Corpus*: Collection of J documents, $(\mathbf{W}_1, \dots, \mathbf{W}_J)$

- **Bag-of-words** assumption: Document as *unordered* set of words

- *Term-document matrix (TDM)*:

$$\begin{aligned}
 \text{TDM} &\equiv \left(\begin{array}{ccc} \sum_{i=1}^{N_1} W_{i1} & \cdots & \sum_{i=1}^{N_J} W_{iJ} \end{array} \right) \\
 &= \underbrace{\left(\begin{array}{ccc} \# \text{ of term 1 in doc 1} & \cdots & \# \text{ of term 1 in doc J} \\ \vdots & \ddots & \vdots \\ \# \text{ of term V in doc 1} & \cdots & \# \text{ of term V in doc J} \end{array} \right)}_{V \times J \text{ matrix of integers}}
 \end{aligned}$$

- TDM fully represents the data \Leftrightarrow Only term frequencies matter
- Simple, easy to work with
- Important dependence across words may be ignored

Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA)

- One of *topic models*
- For word i in document j , latent topic Z_{ij} is assigned
- Z_{ij} : K -dimensional standard basis vector
- Z_{ij} determines the distribution of W_{ij}
- The model:

- Topic assignment of word i, j

$$Z_{ij} \stackrel{\text{indep.}}{\sim} \text{Multinomial}(\theta_j)$$

- Observed word W_{ij}

$$W_{ij} | Z_{ijk} = 1 \stackrel{\text{indep.}}{\sim} \text{Multinomial}(\eta_k)$$

- Prior:

- Topic proportions of document $\theta_j, j = 1, \dots, J$

$$\theta_j \stackrel{\text{i.i.d.}}{\sim} \text{Dirichlet}(\alpha)$$

- Word distribution of topic $\eta_k, k = 1, \dots, K$

$$\eta_k \stackrel{\text{i.i.d.}}{\sim} \text{Dirichlet}(\beta)$$

Dirichlet and Multinomial Distributions

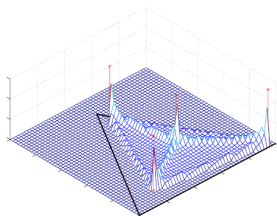
- Multinomial distribution with one trial:

$$p(x) = p_1^{x_1} \cdots p_M^{x_M}$$

- x : Vector of length M , one element is 1 and $M - 1$ elements are 0
- $M = 2$: Bernoulli distribution
- Dirichlet distribution

$$p(x) = \frac{\Gamma(\sum_{m=1}^M \xi_m)}{\prod_{m=1}^M \Gamma(\xi_m)} x_1^{\xi_1-1} \cdots x_M^{\xi_M-1}$$

- x : Vector in the $(M - 1)$ -simplex, Δ^{M-1}
- $x \in \Delta^{M-1} \Leftrightarrow x_m \geq 0$ for all m and $\sum_{m=1}^M x_m = 1$
- $M = 2$: Beta distribution



Blei, Ng, and Jordan (2003) Figure 2.

Dirichlet-Multinomial Model

- Dirichlet: Conjugate prior for the multinomial
 \rightsquigarrow analytically tractable compound (mixture) distribution
- **Dirichlet-Multinomial model**

$$Y \sim \text{Dirichlet}(\xi)$$

$$X | Y \sim \text{Multinomial}(Y)$$

- Marginal probability of X (**Gamma tricks!**):

$$\begin{aligned} p(X) &= \int p(X, y) dy = \int p(X | y) p(y) dy \\ &= \int \dots \int y_1^{X_1} \dots y_M^{X_M} \frac{\Gamma(\sum_{m=1}^M \xi_m)}{\prod_{m=1}^M \Gamma(\xi_m)} y_1^{\xi_1-1} \dots y_M^{\xi_M-1} dy_1 \dots dy_M \\ &= \frac{\Gamma(\sum_{m=1}^M \xi_m)}{\prod_{m=1}^M \Gamma(\xi_m)} \int \dots \int \underbrace{y_1^{\xi_1+X_1-1} \dots y_M^{\xi_M+X_M-1}}_{\text{Unnormalized Dirichlet density}} dy_1 \dots dy_M \\ &= \frac{\Gamma(\sum_{m=1}^M \xi_m)}{\prod_{m=1}^M \Gamma(\xi_m)} \frac{\prod_{m=1}^M \Gamma(\xi_m + X_m)}{\Gamma(\sum_{m=1}^M (\xi_m + X_m))} = \frac{\xi_{m^*}}{\sum_{m=1}^M \xi_m} \end{aligned}$$

where m^* is such that $X_{m^*} = 1$

Collapsing Parameters in LDA

- Joint posterior for LDA

$$p(\mathbf{Z}, \theta, \eta \mid \mathbf{W}) \propto \left(\prod_{k=1}^K p(\eta_k) \right) \prod_{j=1}^J p(\theta_j) \prod_{i=1}^{N_j} p(Z_{ij} \mid \theta_j) p(W_{ij} \mid Z_{ij}, \eta)$$

- High dimensionality of parameters in LDA

- η_k : K vectors of length L
- θ_j : J vectors of length K

- Want to **collapse** (integrate) parameters out

$$\begin{aligned} p(\mathbf{Z} \mid \mathbf{W}) &= \int \int p(\mathbf{Z}, \theta, \eta \mid \mathbf{W}) d\theta d\eta \\ &\propto \prod_{j=1}^J \int p(\theta_j) \prod_{i=1}^{N_j} p(Z_{ij} \mid \theta_j) d\theta_j \\ &\quad \times \int \cdots \int \left(\prod_{k=1}^K p(\eta_k) \right) \prod_{j=1}^J \prod_{i=1}^{N_j} p(W_{ij} \mid Z_{ij}, \eta) d\eta_1 \cdots d\eta_K \end{aligned}$$

Collapsed Gibbs Sampler

- **Collapsed Gibbs Sampler**

$s = 0$ Initialize $\mathbf{Z}^{(0)}$

$s \geq 1$ Loop over all words in the corpus:

- Draw $Z_{ij}^{(s)}$ conditional on all the other topic assignments

- Very simple and fast steps: Random draws from $\{1, \dots, K\}$

- Posterior probability that $Z_{i^*j^*k} = 1$, i.e., the topic of $W_{i^*j^*}$ is k :

$$p(Z_{i^*j^*k} = 1 \mid \mathbf{Z}_{-(i^*,j^*)}, \mathbf{W})$$

$$\propto \underbrace{\int p(\theta_{j^*}) p(Z_{i^*j^*k} = 1 \mid \theta_{j^*}) \prod_{i \neq i^*} p(Z_{ij^*} \mid \theta_{j^*}) d\theta_{j^*}}_{\text{Dirichlet-Multinomial conditional prior: } p(Z_{i^*j^*k}=1 \mid \mathbf{Z}_{-(i^*,j^*)})}$$

Dirichlet-Multinomial conditional prior: $p(Z_{i^*j^*k}=1 \mid \mathbf{Z}_{-(i^*,j^*)})$

$$\times \underbrace{\int p(\eta_k) p(W_{i^*j^*} \mid Z_{i^*j^*k} = 1, \eta_k) \prod_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} p(W_{ij} \mid Z_{ij}, \eta_k) d\eta_k}_{\text{Dirichlet-Multinomial conditional likelihood: } p(W_{i^*j^*} \mid Z_{i^*j^*k}=1, \mathbf{Z}_{-(i^*,j^*)}, \mathbf{W}_{-(i^*,j^*)})}$$

Dirichlet-Multinomial conditional likelihood: $p(W_{i^*j^*} \mid Z_{i^*j^*k}=1, \mathbf{Z}_{-(i^*,j^*)}, \mathbf{W}_{-(i^*,j^*)})$

- Conditional prior $p(Z_{i^*j^*k} = 1 \mid \mathbf{Z}_{-(i^*,j^*)})$:

$$\begin{aligned}
 & \int \underbrace{\frac{\Gamma(\sum_{k'=1}^K a_{k'})}{\prod_{k'=1}^K \Gamma(a_{k'})} \theta_{j^*,1}^{a_{1}-1} \cdots \theta_{j^*,K}^{a_{K}-1}}_{p(\theta_{j^*}): \text{Dirichlet}} \underbrace{\theta_{j^*,k}}_{p(Z_{i^*j^*,k}=1|\theta_{j^*})} \\
 & \times \underbrace{\theta_{j^*,1}^{\sum_{i \neq i^*} Z_{ij^*,1}} \cdots \theta_{j^*,K}^{\sum_{i \neq i^*} Z_{ij^*,K}}}_{\prod_{i \neq i^*} p(Z_{ij^*}|\theta_{j^*}): \text{Multinomial}} d\theta_{j^*} \\
 & \propto \frac{\Gamma(a_k + \sum_{i \neq i^*} Z_{ij^*k} + 1) \prod_{k' \neq k} \Gamma(a_{k'} + \sum_{i \neq i^*} Z_{ij^*k'})}{\Gamma(a_k + \sum_{i \neq i^*} Z_{ij^*k} + 1 + \sum_{k' \neq k} (a_{k'} + \sum_{i \neq i^*} Z_{ij^*k'}))} \\
 & \propto \frac{(a_k + \sum_{i \neq i^*} Z_{ij^*k}) \Gamma(a_k + \sum_{i \neq i^*} Z_{ij^*k}) \prod_{k' \neq k} \Gamma(a_{k'} + \sum_{i \neq i^*} Z_{ij^*k'})}{\Gamma(\sum_{k'=1}^K a_{k'} + \underbrace{1 + \sum_{k'=1}^K \sum_{i \neq i^*} Z_{ij^*k'}}_{=N_{j^*}: \text{Does not depend on } k})} \\
 & \propto (a_k + \sum_{i \neq i^*} Z_{ij^*k}) \prod_{k'=1}^K \Gamma(a_{k'} + \sum_{i \neq i^*} Z_{ij^*k'}) \propto a_k + \sum_{i \neq i^*} Z_{ij^*k}
 \end{aligned}$$

- $\sum_{i \neq i^*} Z_{ij^*k}$: # of other words with topic k in the same document

- Conditional likelihood $p(W_{ijl} = 1 \mid Z_{i^*j^*k} = 1, \mathbf{Z}_{-(i^*,j^*)}, \mathbf{W}_{-(i^*,j^*)})$

$$\int \frac{\Gamma(\sum_{l'=1}^L \beta_{l'})}{\prod_{l'=1}^L \Gamma(\beta_{l'})} \eta_{k1}^{\beta_1-1} \cdots \eta_{kL}^{\beta_L-1} \eta_{kl}$$

$$\times \eta_{k1}^{\sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ij1}} \cdots \eta_{kL}^{\sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ijL}} d\eta_k$$

$$\propto \frac{\beta_l + \sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ijl}}{\sum_{l=1}^L \beta_l + \sum_{l=1}^L \sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ijl}}$$

- $\sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ijl}$: # of the same unique word with topic k

- Gibbs step for $Z_{ij}^{(s)}$

$$p(Z_{i^*j^*k} = 1 \mid \mathbf{Z}_{-(i^*,j^*)}, \mathbf{W})$$

$$\propto \left(a_k + \sum_{i \neq i^*} Z_{ij^*k} \right) \frac{\beta_l + \sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ijl}}{\sum_{l=1}^L \beta_l + \sum_{l=1}^L \sum_{\{(i,j) \neq (i^*,j^*): Z_{ijk}=1\}} W_{ijl}}$$

- Simple multinomial sampling (`sample.int()` in **R**) for each word
- Memory efficient: Only need to store $K \times L$ integers, often sparse

EM Algorithm

- **Large** number of latent variables in LDA
 \rightsquigarrow what about the EM algorithm?
- Log joint posterior density:

$$\begin{aligned} \log p(\mathbf{Z}, \theta, \eta \mid \mathbf{W}) &= \sum_{k=1}^K \sum_{l=1}^L (\beta_l - 1) \log \eta_{kl} \\ &\quad + \sum_{j=1}^J \left\{ \sum_{k=1}^K (\alpha_k - 1) \log \theta_{jk} \right. \\ &\quad \left. + \sum_{i=1}^{N_j} \sum_{k=1}^K Z_{ijk} \left(\log \theta_{jk} + \sum_{l=1}^L W_{ijl} \log \eta_{kl} \right) \right\} \\ &\quad + \text{constant} \end{aligned}$$

- Consider the EM algorithm where:
 - 1 E-step: Compute $Q(\theta, \eta \mid \theta^{(s-1)}, \eta^{(s-1)})$
 - 2 M-step: Maximize $Q(\theta, \eta \mid \theta^{(s-1)}, \eta^{(s-1)})$ w.r.t. θ and η

- Q-function:

$$\begin{aligned}
 Q(\theta, \eta \mid \theta^{(s-1)}, \eta^{(s-1)}) &\equiv \mathbb{E}_{p(\mathbf{Z} \mid \theta^{(s-1)}, \eta^{(s-1)}, \mathbf{W})} [\log p(\mathbf{Z}, \theta, \eta \mid \mathbf{W})] \\
 &= \sum_{k=1}^K \sum_{l=1}^L (\beta_l - 1) \log \eta_{kl} + \sum_{j=1}^J \left\{ \sum_{k=1}^K (\alpha_k - 1) \log \theta_{jk} \right. \\
 &\quad \left. + \sum_{i=1}^{N_j} \sum_{k=1}^K \mathbb{E}_{p(\mathbf{Z} \mid \theta^{(s-1)}, \eta^{(s-1)}, \mathbf{W})} [Z_{ijk}] \left(\log \theta_{jk} + \sum_{l=1}^L W_{ijl} \log \eta_{kl} \right) \right\} \\
 &\quad + \text{constant}
 \end{aligned}$$

- Posterior of \mathbf{Z} conditional on θ and η :

$$p(\mathbf{Z} \mid \theta, \eta, \mathbf{W}) \propto \prod_{j=1}^J \prod_{i=1}^{N_j} \prod_{k=1}^K \left(\theta_{jk} \prod_{l=1}^L \eta_{kl}^{W_{ijl}} \right)^{Z_{ijk}}$$

$\Rightarrow Z_{ij}$'s are conditionally independent given θ, η , and \mathbf{W}

- Simple E-step:

$$\tilde{Z}_{ijk^*}^{(s-1)} \equiv \mathbb{E}_{p(\mathbf{Z} \mid \theta^{(s-1)}, \eta^{(s-1)}, \mathbf{W})} [Z_{ijk^*}] = \frac{\theta_{jk^*}^{(s-1)} \eta_{k^*l^*}^{(s-1)}}{\sum_{k=1}^K \theta_{jk}^{(s-1)} \eta_{kl^*}^{(s-1)}}$$

- Constrained optimization w.r.t. θ and η :

$$\max_{\theta, \eta} Q(\theta, \eta \mid \theta^{(s-1)}, \eta^{(s-1)})$$

$$- \sum_{j=1}^J \lambda_{\theta_j} \left(\sum_{k=1}^K \theta_{jk} - 1 \right) - \sum_{k=1}^K \lambda_{\eta_k} \left(\sum_{l=1}^L \eta_{kl} - 1 \right)$$

where λ_{θ_j} and λ_{η_k} are Lagrange multipliers

- First order conditions:

$$\frac{\alpha_k - 1 + \sum_{i=1}^{N_j} \tilde{z}_{ijk}^{(s-1)}}{\theta_{jk}} - \lambda_{\theta_j} = 0 \text{ for all } j \text{ and } k$$

$$\sum_{k=1}^K \theta_{jk} - 1 = 0 \text{ for all } j$$

$$\frac{\beta_l - 1 + \sum_{j=1}^J \sum_{i=1}^{N_j} \tilde{z}_{ijk}^{(s-1)} w_{ijl}}{\eta_{kl}} - \lambda_{\eta_k} = 0 \text{ for all } k \text{ and } l$$

$$\sum_{l=1}^L \eta_{kl} - 1 = 0 \text{ for all } k$$

- M-step is simple as well:

$$\theta_{jk}^{(s)} = \frac{a_k - 1 + \sum_{i=1}^{N_j} \tilde{z}_{ijk}^{(s-1)}}{\sum_{k'=1}^K a_{k'} - K + \sum_{k'=1}^K \sum_{i=1}^{N_j} \tilde{z}_{ijk'}^{(s-1)}}$$

$$\eta_{kl}^{(s)} = \frac{\beta_l - 1 + \sum_{j=1}^J \sum_{i=1}^{N_j} \tilde{z}_{ijk}^{(s-1)} w_{ijl}}{\sum_{l'=1}^L \beta_{l'} - L + \sum_{l'=1}^L \sum_{j=1}^J \sum_{i=1}^{N_j} \tilde{z}_{ijk}^{(s-1)} w_{ijl'}}$$

- Maximizers $\theta^{(s)}$ and $\eta^{(s)}$ are separately computed
 \rightsquigarrow no need iterative steps as in the two-parameter IRT model
- An EM algorithm for LDA
 - 1 Initialize $\theta^{(0)}$ and $\eta^{(0)}$
 - 2 Repeat:
 - 1 E-step: Compute $\tilde{z}_{ijk}^{(s-1)}$ for all (i, j, k)
 - 2 M-step: Compute $\theta_{jk}^{(s)}$ for all (j, k) and $\eta_{kl}^{(s)}$ for all (k, l)
 until $(\theta^{(s)}, \eta^{(s)})$ converge
- Key: Z_{ij} 's are conditionally independent given θ, η , and \mathbf{W}
- Differs from Blei et. al., which approximates $Q(\eta | \eta^{(s-1)})$

EM with the Collapsed Posterior?

- The EM algorithm updates high-dimensional θ and η
- What if collapse those parameters?
- Revisit $p(\mathbf{Z} | \mathbf{W})$:

$$\begin{aligned}
 p(\mathbf{Z} | \mathbf{W}) &\propto \prod_{j=1}^J \frac{\prod_{k=1}^K \Gamma(a_k + \sum_{i=1}^{N_j} Z_{ijk})}{\Gamma\left(\sum_{k=1}^K (a_k + \sum_{i=1}^{N_j} Z_{ijk})\right)} \\
 &\quad \times \prod_{k=1}^K \frac{\prod_{l=1}^L \Gamma(\beta_l + \sum_{j=1}^J \sum_{i=1}^{N_j} Z_{ijk} W_{ijl})}{\Gamma\left(\sum_{l=1}^L \beta_l + \sum_{l=1}^L \sum_{j=1}^J \sum_{i=1}^{N_j} Z_{ijk} W_{ijl}\right)} \\
 &\propto \prod_{j=1}^J \prod_{k=1}^K \Gamma(a_k) \prod_{m=1}^{\sum_{i=1}^{N_j} Z_{ijk}} (a_k + m) \\
 &\quad \times \prod_{k=1}^K \frac{\prod_{l=1}^L \Gamma(\beta_l) \prod_{m=1}^{\sum_{j=1}^J \sum_{i=1}^{N_j} Z_{ijk} W_{ijl}} (\beta_l + m)}{\Gamma\left(\sum_{l=1}^L \beta_l\right) \prod_{m=1}^{\sum_{l=1}^L \sum_{j=1}^J \sum_{i=1}^{N_j} Z_{ijk} W_{ijl}} \left(\sum_{l=1}^L \beta_l + m\right)}
 \end{aligned}$$

Mean-Field Approximation

- $p(\mathbf{Z} | \mathbf{W}) \neq \prod_{j=1}^J \prod_{i=1}^{N_j} p(Z_{ij} | \mathbf{W}) \rightsquigarrow \mathbb{E}_{p(\mathbf{z}|\mathbf{w})}[Z_{ij}]$ is not easy
- Consider an approximate version of the EM algorithm
 - Intractable due to dependence across Z_{ij} 's
 - *Let's approximate assuming independence*
- Recall the key inequality for the EM algorithm:

$$\begin{aligned} \log p(\psi | \text{Data}) &= \log \int \frac{p(\psi, x | \text{Data})}{p(x | \underline{\psi}, \text{Data})} p(x | \underline{\psi}, \text{Data}) dz \\ &\geq \mathbb{E}_{p(X|\underline{\psi}, \text{Data})} \left[\log \frac{p(\psi, X | \text{Data})}{p(X | \underline{\psi}, \text{Data})} \right] \end{aligned}$$

- Analogously using Jensen's inequality,

$$\begin{aligned} \log p(\text{Data}) &= \log \int \int \frac{p(\psi, x, \text{Data})}{q(x | \delta) q(\psi | \zeta)} q(x | \delta) q(\psi | \zeta) dz d\psi \\ &\geq \mathbb{E}_{q(X|\delta)q(\psi|\zeta)} \left[\log \frac{p(\psi, X, \text{Data})}{q(X | \delta) q(\psi | \zeta)} \right] \end{aligned}$$

for arbitrary distributions $q(X | \delta)$ and $q(\psi | \zeta)$

- Independence between ψ and X : **Mean-field approximation**

Variational Inference

- **Variational inference:**

- *Variational lower bound (VLB)*

$$\mathbb{E}_{q(X|\delta)q(\psi|\zeta)} \left[\log \frac{p(\psi, X, \text{Data})}{q(X|\delta)q(\psi|\zeta)} \right] \leq \log p(\text{Data})$$

- Maximum at $\log p(\text{Data})$ under $q(X|\delta)q(\psi|\zeta) = p(\psi, X|\text{Data})$
- Never attained if $p(\psi, X|\text{Data}) \neq p(\psi|\text{Data})p(X|\text{Data})$
- Maximize w.r.t. δ and $\zeta \rightsquigarrow q(X|\delta)q(\psi|\zeta)$ "close to" $p(\psi, X|\text{Data})$
- EM as a special case

- Alternating maximization:

$$\begin{aligned} \text{VLB} = & \mathbb{E}_{q(X|\delta)} \left[\mathbb{E}_{q(\psi|\zeta)} [\log p(\psi, X, \text{Data})] \right] \\ & - \mathbb{E}_{q(X|\delta)} [\log q(X|\delta)] - \mathbb{E}_{q(\psi|\zeta)} [\log q(\psi|\zeta)] \end{aligned}$$

- ① Maximization w.r.t. δ :

$$\begin{aligned} & \mathbb{E}_{q(X|\delta)} \left[\mathbb{E}_{q(\psi|\zeta)} [\log p(\psi, X, \text{Data})] - \log q(X|\delta) \right] \\ = & -\text{KL} \left(q(X|\delta) \parallel \frac{\exp(\mathbb{E}_{q(\psi|\zeta)} [\log p(\psi, X, \text{Data})])}{\int \exp(\mathbb{E}_{q(\psi|\zeta)} [\log p(\psi, x, \text{Data})]) dx} \right) + \text{constant} \\ \Rightarrow & q(X|\delta^{(s)}) \propto e^{\mathbb{E}_{q(\psi|\zeta^{(s-1)})} [\log p(\psi, X, \text{Data})]} \end{aligned}$$

- ② Maximization w.r.t. ζ : $q(\psi|\zeta^{(s)}) \propto e^{\mathbb{E}_{q(X|\delta^{(s)})} [\log p(\psi, X, \text{Data})]}$

Collapsed Variational Inference for LDA

- Log posterior density:

$$\log p(\mathbf{Z}, \mathbf{W}) = \log p(\mathbf{Z} | \mathbf{W}) - \underbrace{\log p(\mathbf{W})}_{\text{constant}}$$

$$= \sum_{j=1}^J \sum_{k=1}^K \sum_{m=1}^{N_j} \log(a_k + m)$$

$$+ \sum_{k=1}^K \left(\sum_{l=1}^L \sum_{j=1}^J \sum_{i=1}^{N_j} Z_{ijk} W_{ijl} \log(\beta_l + m) \right)$$

$$- \sum_{m=1}^{\sum_{l=1}^L \sum_{j=1}^J \sum_{i=1}^{N_j} Z_{ijk} W_{ijl}} \log \left(\sum_{l=1}^L \beta_l + m \right)$$

+ constant

- Independent multinomial variational distributions:

$$q(\mathbf{Z} | \kappa) = \prod_{j=1}^J \prod_{i=1}^{N_j} q(Z_{ij} | \kappa_{ij}) = \prod_{j=1}^J \prod_{i=1}^{N_j} \prod_{k=1}^K \kappa_{ijk}^{Z_{ijk}}$$

- Variational update of $q(\mathbf{Z}_{i^*j^*} | \kappa_{i^*j^*})$ given all κ_{ij} s.t. $(i, j) \neq (i^*, j^*)$

$$\hat{\kappa}_{i^*j^*k^*} = \frac{\exp(\mathbb{E}_{\Pi_{(i,j) \neq (i^*, j^*)}} q(Z_{ij} | \hat{\kappa}_{ij}) [\log p(\mathbf{Z}_{-(i^*j^*)}, Z_{i^*j^*k^*} = 1, \mathbf{W})])}{\sum_{k=1}^K \exp(\mathbb{E}_{\Pi_{(i,j) \neq (i^*, j^*)}} q(Z_{ij} | \hat{\kappa}_{ij}) [\log p(\mathbf{Z}_{-(i^*j^*)}, Z_{i^*j^*k} = 1, \mathbf{W})])}$$

$$\propto e^{\mathbb{E}_{\Pi_{(i,j) \neq (i^*, j^*)}} q(Z_{ij} | \gamma_{ij}) [\log(\alpha_{k^*} + \sum_{i \neq i^*} Z_{ijk^*}) + \log(\beta_{l^*} + \sum_{(i,j) \neq (i^*, j^*)} Z_{ijk^*} W_{ijl^*})]}$$

$$\times e^{\mathbb{E}_{\Pi_{(i,j) \neq (i^*, j^*)}} q(Z_{ij} | \gamma_{ij}) [-\log(\sum_{l=1}^L \beta_l + \sum_{l=1}^L \sum_{(i,j) \neq (i^*, j^*)} Z_{ijk^*} W_{ijl^*})]}$$

where l^* is such that $W_{i^*j^*l^*} = 1$

- 0th order approximation (“CVB0”)

$$\mathbb{E}_q \left[\log \left(\alpha_{k^*} + \sum_{i \neq i^*} Z_{ijk^*} \right) \right] \approx \log \left(\alpha_{k^*} + \underbrace{\mathbb{E}_q \left[\sum_{i \neq i^*} Z_{ijk^*} \right]}_{= \sum_{i \neq i^*} \hat{\kappa}_{ijk^*}} \right)$$

Summary

- Latent Dirichlet Allocation
 - Workhorse model for statistical text analysis
 - Measurement of “topics”
 - Bag-of-words assumption
- Computational challenges and various estimation algorithms
 - Collapsed Gibbs sampler
 - EM algorithm
 - Variational inference
- Readings for review
 - ① LDA
 - Blei et. al. (2003) “Latent Dirichlet Allocation”
 - ② Collapsed Gibbs sampler:
 - Griffiths and Steyvers (2004) “Finding Scientific Topics”
 - ③ Variational inference
 - **BDA3** Sections 13.7
 - **Bishop** Ch. 10
 - Teh et. al. (2007) “A Collapsed Variational Bayesian Inference Algorithm for Latent Dirichlet Allocation”