



Bayesian Non-Parametrics

Advanced Machine Learning for NLP

Jordan Boyd-Graber

TEXT ANALYSIS

What about text?

- Gaussian distributions can't model text
- So typically use multinomial distribution as the base distribution
- Remember multinomial:

$$P(N | n, \theta) = \frac{n!}{\prod_j N_j!} \prod_j \theta_j^{N_j} \quad (1)$$

What about text?

- Gaussian distributions can't model text [or can they?]
- So typically use multinomial distribution as the base distribution
- Remember multinomial:

$$P(N | n, \theta) = \frac{n!}{\prod_j N_j!} \prod_j \theta_j^{N_j} \quad (1)$$

Dirichlet Process Multinomial Mixture Model

- Break off sticks

$$V_1, V_2, \dots \sim \text{iid Beta}(1, \alpha) \quad (2)$$

$$C_k \equiv V_k \prod_{j=1}^{k-1} (1 - V_j) \quad (3)$$

Dirichlet Process Multinomial Mixture Model

- Break off sticks

$$V_1, V_2, \dots \sim \text{iid Beta}(1, \alpha) \quad (2)$$

$$C_k \equiv V_k \prod_{j=1}^{k-1} (1 - V_j) \quad (3)$$

- Draw atoms

$$\Phi_1, \Phi_2, \dots \sim \text{iid Dir}(\beta) \quad (4)$$

Dirichlet Process Multinomial Mixture Model

- Break off sticks

$$V_1, V_2, \dots \sim \text{iid Beta}(1, \alpha) \quad (2)$$

$$C_k \equiv V_k \prod_{j=1}^{k-1} (1 - V_j) \quad (3)$$

- Draw atoms

$$\Phi_1, \Phi_2, \dots \sim \text{iid Dir}(\beta) \quad (4)$$

- Merge into complete distribution

$$\Theta = \sum_k C_k \delta_{\Phi_k} \quad (5)$$

Dirichlet Process Multinomial Mixture Model

- Break off sticks

$$V_1, V_2, \dots \sim \text{iid Beta}(1, \alpha) \quad (2)$$

$$C_k \equiv V_k \prod_{j=1}^{k-1} (1 - V_j) \quad (3)$$

- Draw atoms

$$\Phi_1, \Phi_2, \dots \sim \text{iid Dir}(\beta) \quad (4)$$

- Merge into complete distribution

$$\Theta = \sum_k C_k \delta_{\Phi_k} \quad (5)$$

- Draw document word counts

$$\phi_d \sim \Theta \quad (6)$$

$$w_d \sim \phi_d \quad (7)$$

Extending DPMM for text: HDP

- Topic models can use multiple topics per document
- Mixture model can only use one
- HDP is the non-parametric extension

Hierarchical Dirichlet Process

- Draw a global distribution over topics (e.g., $H \equiv \text{Dir}(\alpha)$)

$$G_0 \sim \text{DP}(\gamma, H) \quad (8)$$

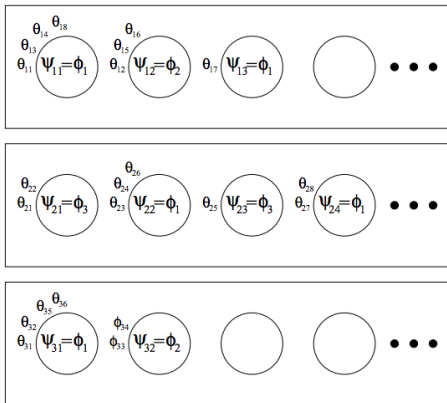
- For each document d , draw distribution over topics

$$\phi_d \sim \text{DP}(\alpha, G_0) \quad (9)$$

- For each word $w_{d,n}$ in the document, draw it from document distribution

$$w_{d,n} \sim \phi_d \quad (10)$$

Chinese Restaurant Franchise



t : Assignment at global table

z : Assignment at document table

Gibbs Sampling

$$p(z_{dn} = k, t_{dn} = j | \mathbf{z}^{-ji}, \mathbf{t}^{-ji}) \propto \begin{cases} \frac{n_{d,k}}{n_{d,\cdot} + \alpha} f(w_{dn} | \Psi_k) & k, j \text{ existing} \\ \frac{\alpha m_j}{\gamma + m_j} f(w_{dn} | \Psi_k) & k \text{ new, } j \text{ existing} \\ \alpha \gamma f(w_{dn} | H_0) & k, j \text{ new} \end{cases} \quad (11)$$

Gibbs Sampling

$$p(z_{dn} = k, t_{dn} = j | \mathbf{z}^{-ji}, \mathbf{t}^{-ji}) \propto \begin{cases} \frac{n_{d,\cdot} + \alpha}{n_{d,\cdot} + \alpha} f(w_{dn} | \Psi_k) & k, j \text{ existing} \\ \frac{\alpha m_j}{\gamma + m_j} f(w_{dn} | \Psi_k) & k \text{ new, } j \text{ existing} \\ \alpha \gamma f(w_{dn} | H_0) & k, j \text{ new} \end{cases} \quad (11)$$

Number of tokens seated in lower-level table

Gibbs Sampling

$$p(z_{dn} = k, t_{dn} = j | \mathbf{z}^{-ji}, \mathbf{t}^{-ji}) \propto \begin{cases} \frac{n_{d,k}}{n_{d,\cdot} + \alpha} f(w_{dn} | \Psi_k) & k, j \text{ existing} \\ \frac{\alpha m_j}{\gamma + m_j} f(w_{dn} | \Psi_k) & k \text{ new, } j \text{ existing} \\ \alpha \gamma f(w_{dn} | H_0) & k, j \text{ new} \end{cases} \quad (11)$$

Number of tokens seated at higher-level table

Gibbs Sampling

$$p(z_{dn} = k, t_{dn} = j | \mathbf{z}^{-ji}, \mathbf{t}^{-ji}) \propto \begin{cases} \frac{n_{d,k}}{n_{d,\cdot} + \alpha} f(w_{dn} | \Psi_k) & k, j \text{ existing} \\ \frac{\alpha m_j}{\gamma + m_j} f(w_{dn} | \Psi_k) & k \text{ new, } j \text{ existing} \\ \alpha \gamma f(w_{dn} | H_0) & k, j \text{ new} \end{cases} \quad (11)$$

Lower-level concentration

Gibbs Sampling

$$p(z_{dn} = k, t_{dn} = j | \mathbf{z}^{-ji}, \mathbf{t}^{-ji}) \propto \begin{cases} \frac{n_{d,k}}{n_{d,\cdot} + \alpha} f(w_{dn} | \Psi_k) & k, j \text{ existing} \\ \frac{\alpha m_j}{\gamma + m_j} f(w_{dn} | \Psi_k) & k \text{ new, } j \text{ existing} \\ \alpha \gamma f(w_{dn} | H_0) & k, j \text{ new} \end{cases} \quad (11)$$

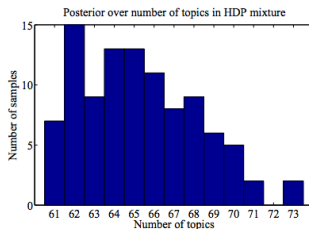
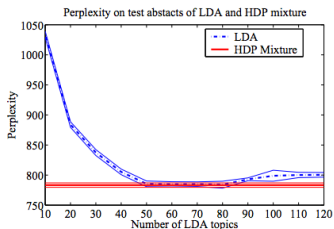
Higher-level concentration

Gibbs Sampling

$$p(z_{dn} = k, t_{dn} = j | \mathbf{z}^{-ji}, \mathbf{t}^{-ji}) \propto \begin{cases} \frac{n_{d,k}}{n_{d,\cdot} + \alpha} f(w_{dn} | \Psi_k) & k, j \text{ existing} \\ \frac{\alpha m_j}{\gamma + m_j} f(w_{dn} | \Psi_k) & k \text{ new, } j \text{ existing} \\ \alpha \gamma f(w_{dn} | H_0) & k, j \text{ new} \end{cases} \quad (11)$$

Multinomial (or whatever base distribution)

Discovers Dimensionality



- Discovers dimensionality
- Additional layers can capture different aspects of data
- But only unsupervised objective