



Generative Models

Donald E. Brown

School of Data Science
University of Virginia
Charlottesville, VA 22904



Causal or Generative Models

Generative
Models
2 / 7

D.E. Brown

- Generative models often contain latent variables or parameters
- In most practical problems, the parameters or latent variables will be the lower numbered indices and the higher numbered indices represent observations
- Can interpret the graphical model as producing or generating the data
- These generative graphical models can be viewed as causal (see Pearl, 1989)
- Causality flows from the relationships shown in the conditional distributions
- This process is sometimes called *ancestral sampling*



A Generative Model for Text

Generative
Models
3/7

D.E. Brown

- Goal: Find topics in text
- Approach: Latent Dirichlet Allocation (LDA) (Blei, et al., 2003)
- Notation
 - $w \in \{1, \dots, V\}$ Word in a vocabulary of length V
 - $\mathbf{w} = (w_1, \dots, w_N)$ Document with N words
 - $D = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$ Corpus of M documents
- "We wish to find a probabilistic model of a corpus that not only assigns high probability to members of the corpus, but also assigns high probability to other 'similar' documents."



LDA Processes

Generative
Models
4/7

D.E. Brown

- LDA - generative probabilistic model of a corpus, where for each $\mathbf{w} \in D$:
 - 1 Choose $N \sim \text{Poisson}(\xi)$
 - 2 Choose $\theta \sim \text{Dir}(\alpha)$
 - 3 For each of the N words w_n choose
 - 1 Choose a topic $z_n \sim \text{Multinom}(\theta)$
 - 2 Choose a word, w_n , from $p(w_n|z_n, \beta)$, a multinomial conditioned on z_n
- Assumptions
 - Dimensionality, k , of the Dirichlet is known (i.e., topics)
 - Word probabilities parameterized by β a $k \times V$ matrix, where $\beta_{ij} = p(w = i|z = j)$



LDA Distributions

Generative
Models
5/7

D.E. Brown

The joint distribution of a topic mixture θ , a set of N topics \mathbf{z} , and a set of N words \mathbf{w} is

$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n \beta)$$

The marginal distribution of a document is

$$p(\mathbf{w} | \alpha, \beta) = \int p(\theta | \alpha) \left(\prod_{n=1}^N \sum_{z_n} p(z_n | \theta) p(w_n | z_n \beta) \right) d\theta$$

Probability of a corpus is

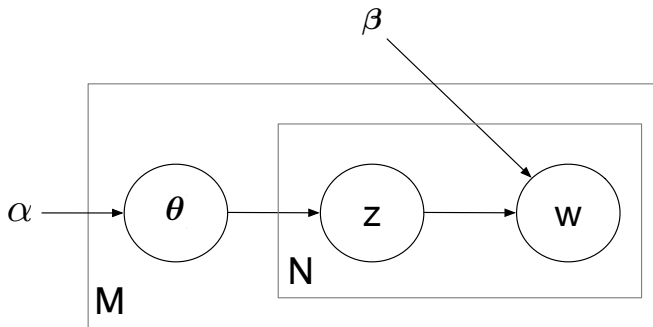
$$p(D | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left(\prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn} \beta) \right) d\theta_d$$



LDA Graphical Model

Generative
Models
6/7

D.E. Brown





LDA Inference

Generative
Models
7/7

D.E. Brown

Compute the posterior distribution of the latent variables:

$$p(\boldsymbol{\theta}, \mathbf{z} | \mathbf{w}, \alpha, \beta) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} | \alpha, \beta)}{p(\mathbf{w} | \alpha, \beta)}$$

This is intractable; need to use sampling or variational inference.