

# Probabilistic Models for Text and Graphs

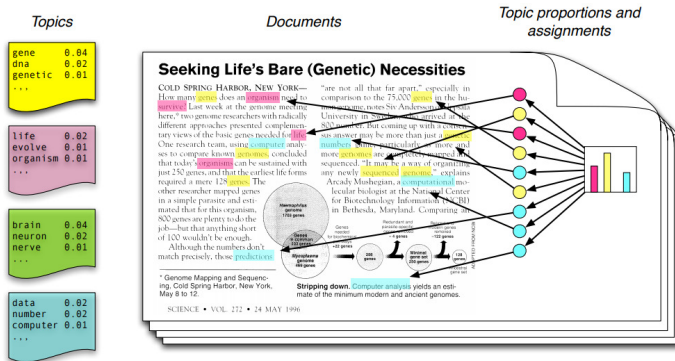
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

April 5, 2018

# Topic Models for Text

- Goal of topic modeling is to learn the topics that underlie the document collection, and represent each document in terms of these topics (fraction of various topics)



- Topic models are based on **assigning words in each document to clusters/topics** (each cluster can be thought of as representing a "topic"): Essentially word clustering with many documents!

# What's a Topic?

- Informally, a topic represents a group of similar words representative of that topic
- Suppose the vocabulary consists of  $V$  words
- Formally, each topic is a probability distribution over the  $V$  words

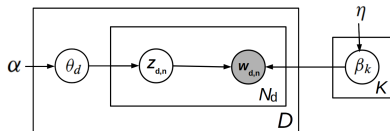
Topic 247		Topic 5		Topic 43		Topic 56	
word	prob.	word	prob.	word	prob.	word	prob.
DRUGS	.069	RED	.202	MIND	.081	DOCTOR	.074
DRUG	.060	BLUE	.099	THOUGHT	.066	DR.	.063
MEDICINE	.027	GREEN	.096	REMEMBER	.064	PATIENT	.061
EFFECTS	.026	YELLOW	.073	MEMORY	.037	HOSPITAL	.049
BODY	.023	WHITE	.048	THINKING	.030	CARE	.046
MEDICINES	.019	COLOR	.048	PROFESSOR	.028	MEDICAL	.042
PAIN	.016	BRIGHT	.030	FELT	.025	NURSE	.031
PERSON	.016	COLORS	.029	REMEMBERED	.022	PATIENTS	.029
MARIJUANA	.014	ORANGE	.027	THOUGHTS	.020	DOCTORS	.028
LABEL	.012	BROWN	.027	FORGOTTEN	.020	HEALTH	.025
ALCOHOL	.012	PINK	.017	MOMENT	.020	MEDICINE	.017
DANGEROUS	.011	LOOK	.017	THINK	.019	NURSING	.017
ABUSE	.009	BLACK	.016	THING	.016	DENTAL	.015
EFFECT	.009	PURPLE	.015	WONDER	.014	NURSES	.013
KNOWN	.008	CROSS	.011	FORGET	.012	PHYSICIAN	.012
PILLS	.008	COLORED	.009	RECALL	.012	HOSPITALS	.011

- Each topic, say topic number  $k$ , can be represented by a **probability vector**  $\beta_k$  of size  $V$ 
  - $\beta_{kv}$  denotes the probability of word  $v$  in topic  $k$  and  $\sum_{v=1}^V \beta_{kv} = 1$
  - High-probability words are representative of that topic
- $\{\beta_k\}_{k=1}^K$  denotes the set of  $K$  topics

# Latent Dirichlet Allocation (LDA) - Blei et al, 2003

Total  $D$  docs,  $w_{d,n}$  :  $n$ -th word in the  $d$ -th document. Assume  $N_d$  = no. of words in document  $d$

LDA associates each observed word  $w_{d,n} \in \{1, \dots, V\}$  with a topic  $z_{d,n} \in \{1, \dots, K\}$



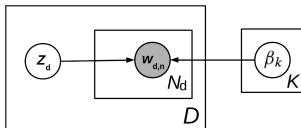
The LDA (assuming hyperparams  $\alpha, \eta$  known) generative model is like a **mixture model** for words

The  $N_d$  words in each document are generated via a mixture model, as follows

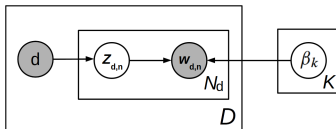
- Draw  $K$  topics  $\{\beta_k\}_{k=1}^K$  shared by all documents  $d = 1, \dots, D$   
 $\beta_k \sim \text{Dirichlet}(\eta, \dots, \eta) \quad k = 1, \dots, K$
- For each document  $d = 1, \dots, D$ 
  - Draw its  $K$ -dim topic proportion vector  $\theta_d \sim \text{Dirichlet}(\alpha, \dots, \alpha)$
  - For each word  $n = 1, \dots, N_d$  in document  $d$ 
    - First choose the topic for this word:  $z_{d,n} \sim \text{multinoulli}(\theta_d)$
    - Now generate the word from the **chosen topic**:  $w_{d,n} \sim \text{multinoulli}(\beta_{z_{d,n}})$

# Some LDA Precursors

- A naïve mixture model: Each document having a single topic, all words from that topic

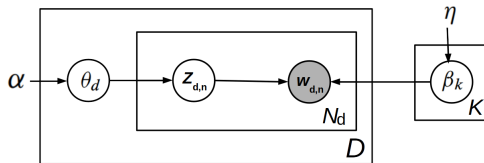


- Each word in a doc has its own topic: Probabilistic Latent Semantic Analysis (Hofmann (2001))



- PLSA has too many parameters  $p(w, d) = p(d)p(w|d) = p(d) \sum_{k=1}^K p(w|z = k)p(z = k|d)$ 
  - $VK$  params for  $p(w|z = k)$ ,  $KD$  params for  $p(z = k|d)$ . Grows in the number of documents
  - LDA however models  $p(z = k|d)$  via **random variables**  $\theta_d$ , not tied to only training data, so can extend to previously unseen documents unlike PLSA. A subtle but important difference (Blei et al, 2003)
- Also, unlike LDA, PLSA is not a Bayesian model

# Inference for LDA

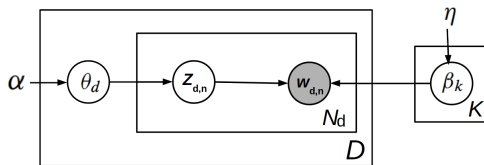


- The goal is to infer the posterior distribution over all the latent variables

$$p(\mathbf{Z}, \theta, \beta | \mathbf{W}, \alpha, \eta) = \frac{p(\mathbf{W} | \beta, \mathbf{Z}) p(\mathbf{Z} | \theta) p(\beta | \eta) p(\theta | \alpha)}{p(\mathbf{W} | \alpha, \eta)} \quad (\text{assuming hyperparams } \alpha, \eta \text{ are fixed})$$

- $\mathbf{Z}$ : Topic assignments of all words across all documents
- $\theta = \{\theta_1, \dots, \theta_D\}$ : Topic proportion vectors; each  $\theta_d$  is a  $K$ -dim vector
- $\beta = \{\beta_1, \dots, \beta_K\}$ : The  $K$  topics; each  $\beta_k$  is a  $V$ -dim vector
- $\mathbf{W} = \{\mathbf{w}_{d,n}\}, d = 1, \dots, D, n = 1, \dots, N_d$ : Collection of all words in all the documents

# Inference for LDA



- A wide variety of inference methods exist for LDA (MCMC, VB, EP, batch, online, ...)
- LDA is a **locally conjugate model** (recall that the model uses Dirichlet/multinoulli distr.)
  - The likelihood and priors are all exponential family distributions
- MCMC, in particular Gibbs sampling, is very easy to derive due to local conjugacy
- Local conjugacy also allow deriving a simple VB algorithm (the original LDA paper<sup>†</sup> used VB)
- Note: Can even **collapse** some variables and do collapsed Gibbs or collapsed VB

<sup>†</sup> Latent Dirichlet Allocation (Blei et al, 2003)

# LDA Evaluation

- Many ways to test how well LDA performs on some data
- Extrinsic measures: Perform LDA and measure performance on some external task (e.g., accuracy of a document classification model using the topic based document representation)
- Perplexity is another **intrinsic** measure

$$\text{perplexity}(D_{\text{test}}) = \exp \left\{ -\frac{\sum_{d=1}^M \log p(\mathbf{w}_d)}{\sum_{d=1}^M N_d} \right\}$$

.. where  $\mathbf{w}_d$  collectively denotes the set of words in document  $d$

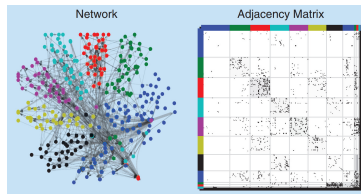
- We set aside some documents as a test set and compute the average probability the model assigns to the words in these held-out documents
- A high average probability (low perplexity) implies a good fit to the data





# Probabilistic Models for Graphs

- Graphs/networks are ubiquitous in many domains (social networks, computer networks, protein-protein networks, citation networks, etc.)
- A graph between  $N$  entities/nodes can be represented by an  $N \times N$  adjacency matrix  $\mathbf{A}$



- Given a (partially observed) graph, we may be interested in tasks such as
  - Identifying the latent structure (e.g., clusters of similar nodes)
  - Predicting possibility of a link between two node
- We can design probabilistic models for such tasks

# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$
- We can model each link/non-link  $A_{nm} \in \{0, 1\}$  via a probability model, e.g.

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \theta) = f(\mathbf{z}_n, \mathbf{z}_m, \theta)$$

where  $f$  is some function of  $\mathbf{z}_n, \mathbf{z}_m$  and params  $\theta$ , and returns a probability

- The overall probability of the observed graph

$$p(\mathbf{A} | \mathbf{Z}, \theta) = \prod_{n,m} p(A_{nm} | \mathbf{z}_n, \mathbf{z}_m, \theta)$$

- Various models differ in terms of what the embeddings  $\mathbf{z}_n$  look like, and what  $f$  is defined as
- Some representative models
  - Latent Space Model
  - Stochastic Blockmodel
  - Mixed-Membership Blockmodel (MMSB)

# Latent Space Model for Graphs

- LSM assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space
- One such possible model would be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m) = \sigma(-\|\mathbf{z}_n - \mathbf{z}_m\|)$$

- A reasonable model for link prediction
- However, the real-valued embeddings in LSM don't provide a good interpretability
  - Therefore not very ideal for discovering clusters etc.
- Blockmodels and its variants has such a properties as we will see next

# Stochastic Blockmodel

- SBM assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ ’s cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$
- Therefore the probability of a link between nodes  $n$  and  $m$  will be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like
  - For  $n = 1, \dots, N$ 
$$\mathbf{z}_n \sim \text{multinoulli}(\pi)$$
  - For  $n = 1, \dots, N$ 
    - For  $m = 1, \dots, n - 1$ 
$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_n, \mathbf{z}_m})$$
  - Note:  $\pi$  and  $\eta$  can also be given priors

# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB assumes each node  $n$  to have a probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- The full generative model looks like
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha)$$

- For  $n = 1, \dots, N$ 
  - For  $m = 1, \dots, n$

$$\mathbf{z}_{n \rightarrow m} \sim \text{multinoulli}(\pi_n)$$

$$\mathbf{z}_{m \rightarrow n} \sim \text{multinoulli}(\pi_m)$$

$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_{n \rightarrow m}, \mathbf{z}_{m \rightarrow n}})$$

- Unlike SBM in which node  $n$  has a unique one-hot  $\mathbf{z}_n$  vector, in MMSB, for each link/non-link  $A_{nm}$ , node  $n$  displays a different community membership  $\mathbf{z}_{n \rightarrow m}$  based on the node  $m$  it is interacting with (likewise for node  $m$ ).
- In contrast to SBM, in MMSB, the unique vector is  $\pi_n$ .