

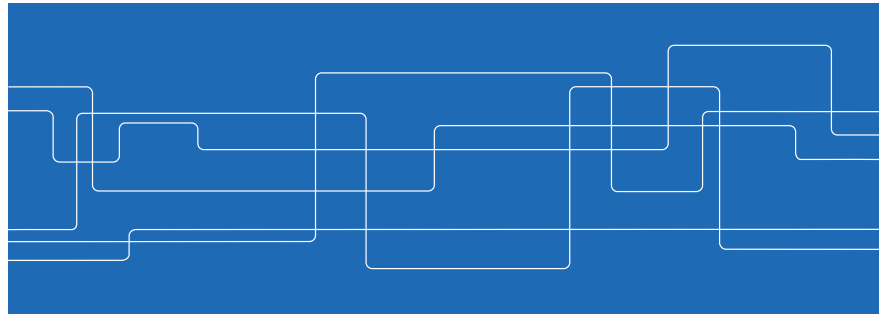


DD2434 Machine Learning, Advanced Course Lecture 11: Topic Models

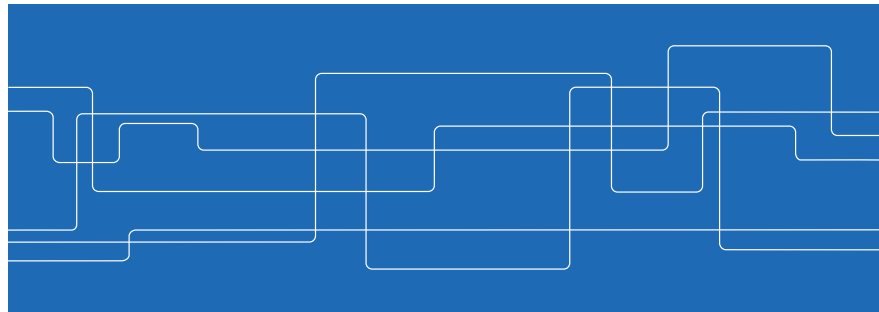
Hedvig Kjellström

hedvig@kth.se

<https://www.kth.se/social/course/DD2434/>



Text Documents and Topics



Today

The idea of modeling text documents according to topics
(Blei and Lafferty)

Text data and the bag of words model (Blei and Lafferty)

Latent Dirichlet Allocation (LDA) (Blei and Lafferty)

2



Probabilistic topic models

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

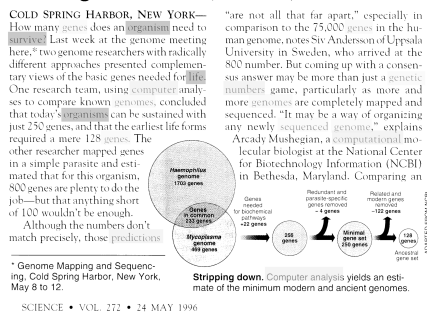
Discuss with your neighbor (2 min):
Can you see patterns in how words appear in the 4 columns?

Slide from Blei

4

Latent Dirichlet allocation (LDA)

Seeking Life's Bare (Genetic) Necessities

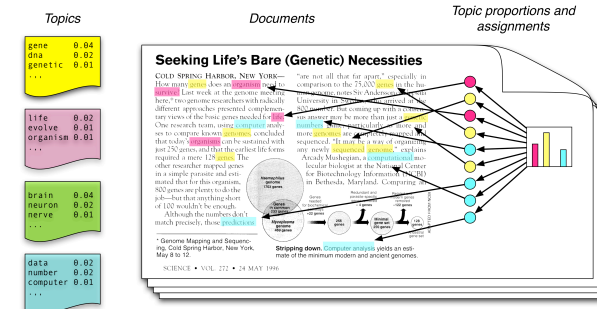


Simple intuition: Documents exhibit multiple topics.

Slide from Blei

5

Latent Dirichlet allocation (LDA)

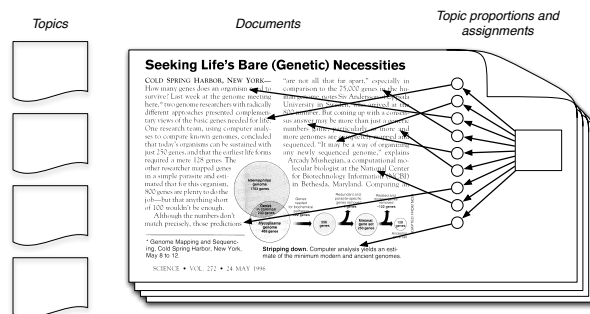


- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

Slide from Blei

6

Latent Dirichlet allocation (LDA)



- In reality, we only observe the documents
- The other structure are **hidden variables**

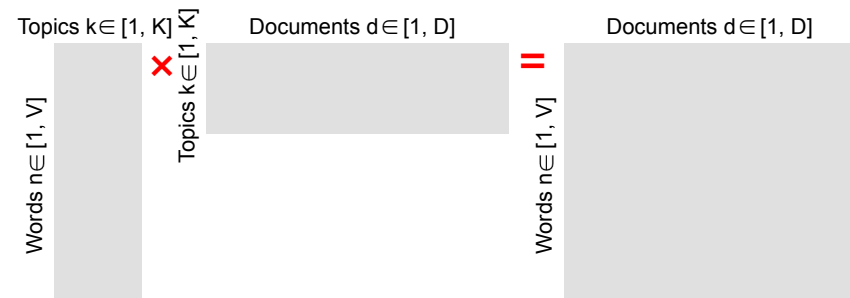
Hidden = Latent

Slide from Blei

7

Topics = Latent Low-Dimensional Representation

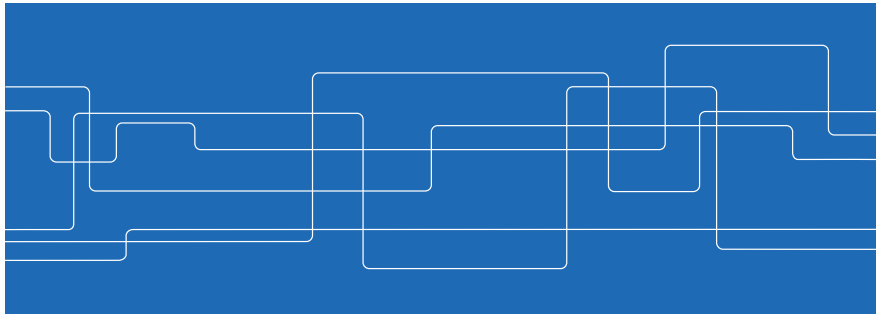
Remember from Lecture 10: Topic model is latent variable model (LVM) for discrete/categorical data!
(LVMs for continuous data: PCA, ICA, GPLVM, ...)



8



Text Data and Bag of Words

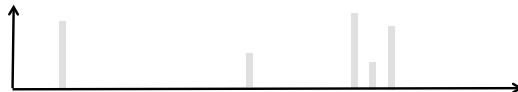


Multinomial Distribution of Text

Statespace = set of unique words in the language in which the text document is written

High-dim Sparse

Multinomial distribution (normalized histogram) of a text document is called a **bag of words**



Discuss with your neighbor (5 min):

What information have you thrown away when you represent data as a bag of words?



Multinomial Distribution of Text

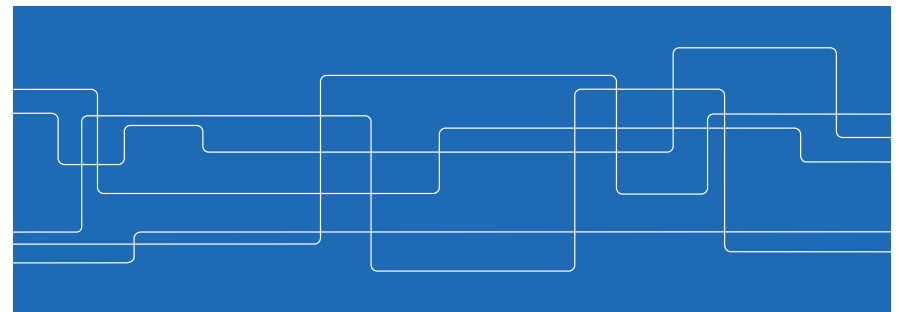
From Lecture 10: Multinomial distribution – essentially normalized histogram over a finite set of outcomes
In dice case, set of outcomes $x_i \in \{1, \dots, K\}$, $K = 6$

Discuss with your neighbor (5 min):

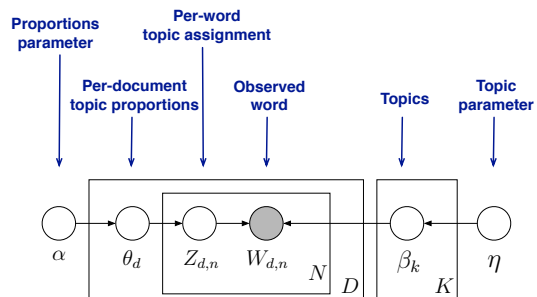
What is the set of possible outcomes if we think of a text document instead of a sequence of dice rolls?



Latent Dirichlet Allocation (LDA)



LDA as a graphical model

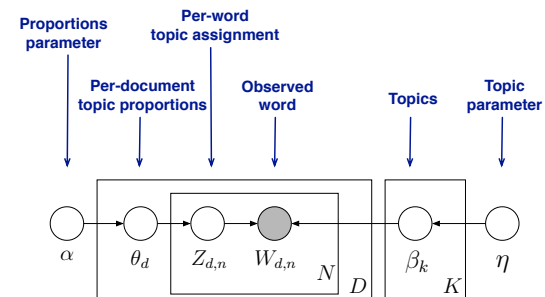


- Encodes **assumptions**
- Defines a **factorization** of the joint distribution
- Connects to **algorithms** for computing with data

Slide from Blei

13

LDA as a graphical model

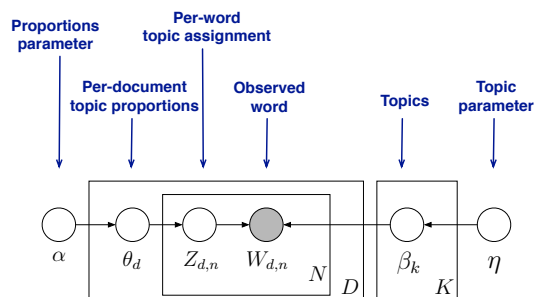


- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed.
- Plates indicate replicated variables.

Slide from Blei

14

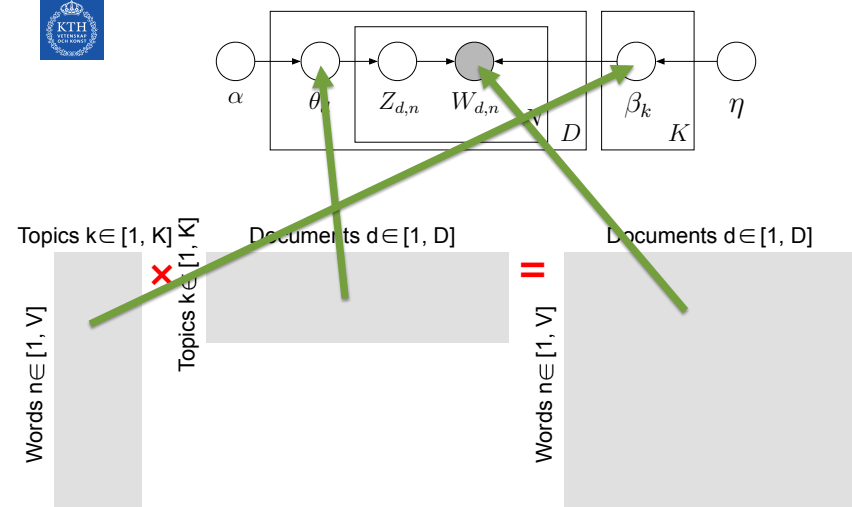
LDA as a graphical model



$$\prod_{i=1}^K p(\beta_i | \eta) \prod_{d=1}^D p(\theta_d | \alpha) \left(\prod_{n=1}^N p(z_{d,n} | \theta_d) p(w_{d,n} | \beta_{1:K}, z_{d,n}) \right)$$

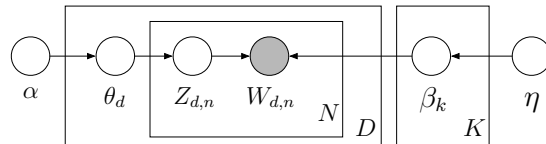
Slide from Blei

15



16

LDA as a graphical model

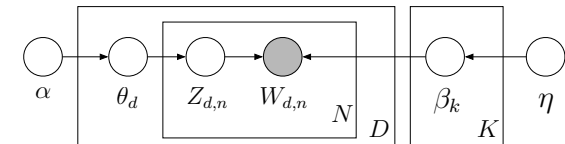


- This joint defines a posterior.
- From a collection of documents, infer
 - Per-word topic assignment $z_{d,n}$
 - Per-document topic proportions θ_d
 - Per-corpus topic distributions β_k
- Then use posterior expectations to perform the task at hand, e.g., information retrieval, document similarity, exploration, ...

Slide from Blei

17

LDA as a graphical model



Approximate posterior inference algorithms

- Mean field variational methods (Blei et al., 2001, 2003)
- Expectation propagation (Minka and Lafferty, 2002)
- Collapsed Gibbs sampling (Griffiths and Steyvers, 2002)
- Collapsed variational inference (Teh et al., 2006)
- Online variational inference (Hoffman et al., 2010)

More in Lecture 12

Also see Mukherjee and Blei (2009) and Asuncion et al. (2009).

Slide from Blei

18

Example inference



- **Data:** The OCR'd collection of *Science* from 1990–2000
 - 17K documents
 - 11M words
 - 20K unique terms (stop words and rare words removed)
- **Model:** 100-topic LDA model using variational inference.

Slide from Blei

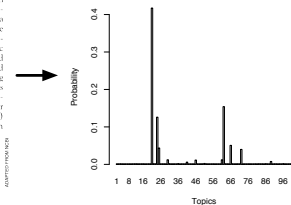
19

Example inference

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 252 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.
SCIENCE • VOL. 272 • 24 MAY 1996



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

Slide from Blei

20



Example inference

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

Slide from Blei

21



Aside: The Dirichlet distribution

- The Dirichlet distribution is an exponential family distribution over the simplex, i.e., positive vectors that sum to one

$$p(\theta | \vec{\alpha}) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i - 1}.$$

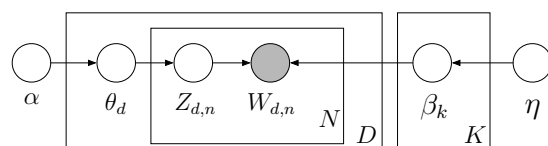
- It is **conjugate** to the multinomial. Given a multinomial observation, the posterior distribution of θ is a Dirichlet.
- The parameter α controls the mean shape and sparsity of θ .
- The topic proportions are a K dimensional Dirichlet. The topics are a V dimensional Dirichlet.

Slide from Blei

22



LDA as a graphical model



Discuss with your neighbor (5 min):

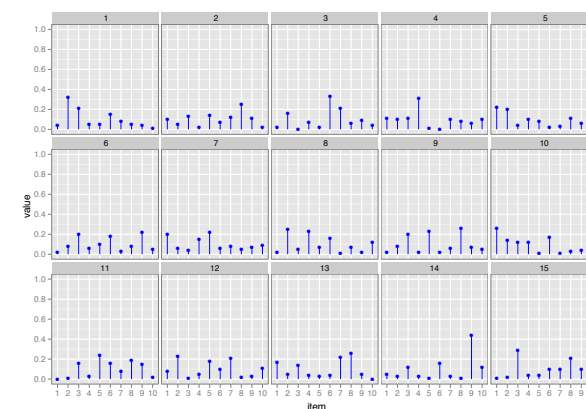
What would happen to the topic distribution if we changed the Dirichlet priors to uniform priors?

Slide from Blei

23



$\alpha = 1$

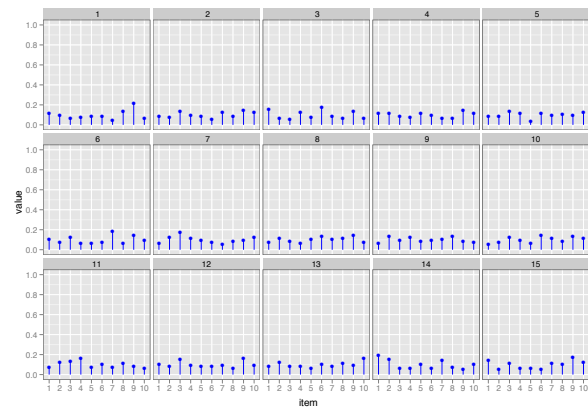


Slide from Blei

24



$\alpha = 10$

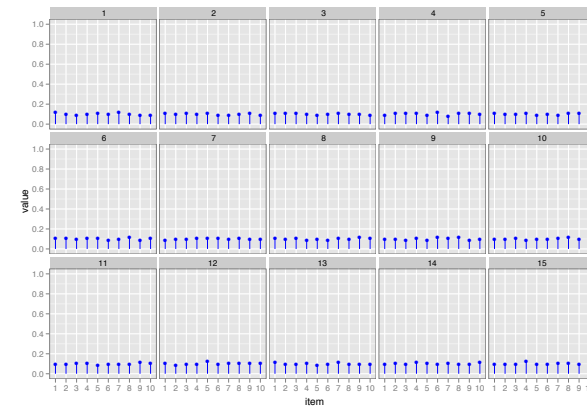


Slide from Blei

25



$\alpha = 100$

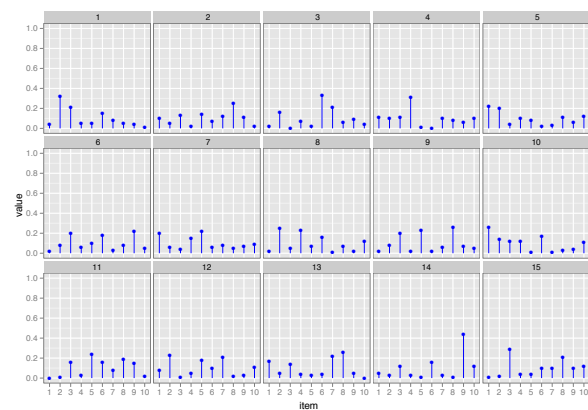


Slide from Blei

26



$\alpha = 1$

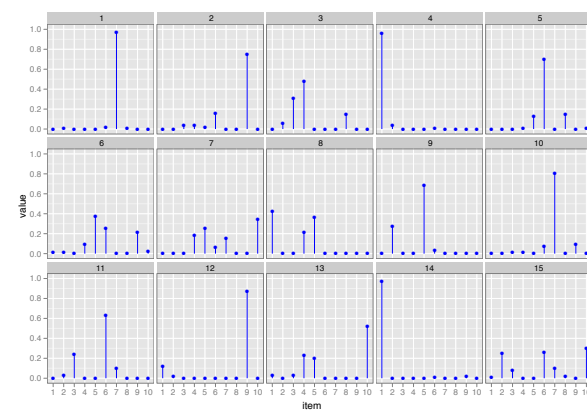


Slide from Blei

27



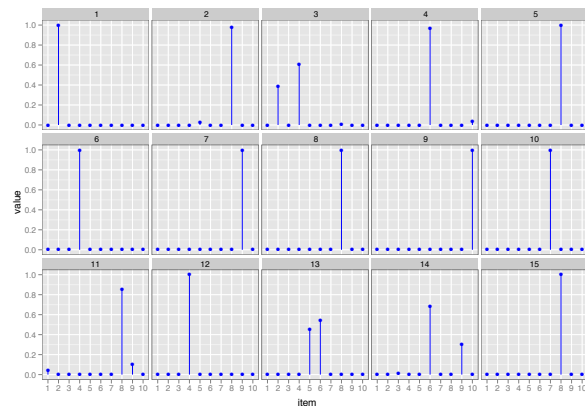
$\alpha = 0.1$



Slide from Blei

28

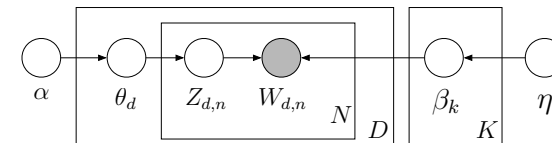
$\alpha = 0.01$



Slide from Blei

29

LDA summary

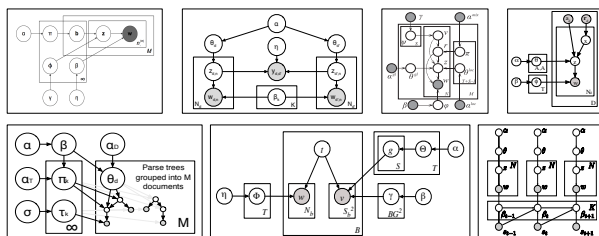


- LDA is a probabilistic model of text. It casts the problem of discovering themes in large document collections as a posterior inference problem.
- It lets us visualize the hidden thematic structure in large collections, and generalize new data to fit into that structure.
- Builds on latent semantic analysis (Deerwester et al., 1990; Hofmann, 1999)
It is mixed membership model (Erosheva, 2004).
It relates to PCA and matrix factorization (Jakulin and Buntine, 2002)
Was independently invented for genetics (Pritchard et al., 2000)

Slide from Blei

30

LDA summary



- Organizing and finding patterns in data has become important in the sciences, humanities, industry, and culture.
- LDA can be embedded in more complicated models that capture richer assumptions about the data.
- Algorithmic improvements let us fit models to massive data.

Slide from Blei

31

What is next?

Final project groups are published on the home page. Talk to your project group and select papers before **December 3**.

Continue with Assignment 2, deadline **December 16**.

Next on the schedule

Wed 2 Dec 10:15-12:00 K2

Exercise 5: Probabilistic Independent Component Analysis

Hedvig Kjellström

Optional readings: *Beckmann and Smith*

Thu 3 Dec 13:15-15:00 L51

Lecture 12: Sampling

Hedvig Kjellström

Readings: Bishop 11.1-11.4

32