

# Inference and Representation

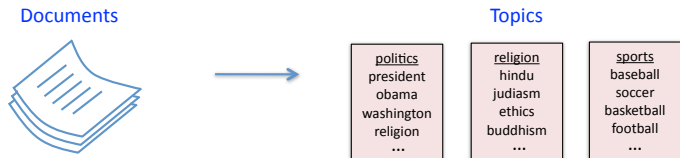
David Sontag

New York University

Lecture 8, Nov. 3, 2015

# Latent Dirichlet allocation (LDA)

- **Topic models** are powerful tools for exploring large data sets and for making inferences about the content of documents



- Many applications in information retrieval, document summarization, and classification



- LDA is one of the simplest and most widely used topic models

# Generative model for a document in LDA

- 1 Sample the document's **topic distribution**  $\theta$  (aka topic vector)

$$\theta \sim \text{Dirichlet}(\alpha_{1:T})$$

where the  $\{\alpha_t\}_{t=1}^T$  are fixed hyperparameters. Thus  $\theta$  is a distribution over  $T$  topics with mean  $\theta_t = \alpha_t / \sum_{t'} \alpha_{t'}$

- 2 For  $i = 1$  to  $N$ , sample the **topic**  $z_i$  of the  $i$ 'th word

$$z_i | \theta \sim \theta$$

- 3 ... and then sample the actual **word**  $w_i$  from the  $z_i$ 'th topic

$$w_i | z_i \sim \beta_{z_i}$$

where  $\{\beta_t\}_{t=1}^T$  are the *topics* (a fixed collection of distributions on words)

# Generative model for a document in LDA

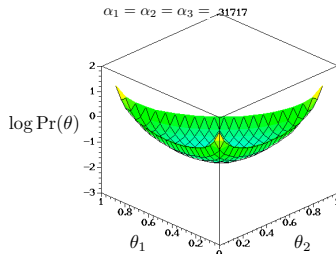
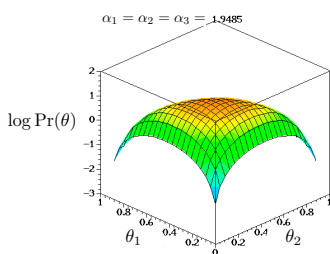
- 1 Sample the document's **topic distribution**  $\theta$  (aka topic vector)

$$\theta \sim \text{Dirichlet}(\alpha_{1:T})$$

where the  $\{\alpha_t\}_{t=1}^T$  are hyperparameters. The Dirichlet density, defined over  $\Delta = \{\vec{\theta} \in \mathbb{R}^T : \forall t \theta_t \geq 0, \sum_{t=1}^T \theta_t = 1\}$ , is:

$$p(\theta_1, \dots, \theta_T) \propto \prod_{t=1}^T \theta_t^{\alpha_t - 1}$$

For example, for  $T=3$  ( $\theta_3 = 1 - \theta_1 - \theta_2$ ):

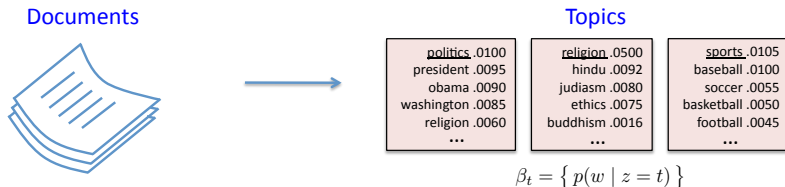


# Generative model for a document in LDA

- 3 ... and then sample the actual **word**  $w_i$  from the  $z_i$ 'th topic

$$w_i | z_i \sim \beta_{z_i}$$

where  $\{\beta_t\}_{t=1}^T$  are the *topics* (a fixed collection of distributions on words)



# Example of using LDA

## Topics

 $\beta_1$ 

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

 $\beta_T$ 

data	0.02
number	0.02
computer	0.01
...	

## Documents

### 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 250 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** are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson at Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



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

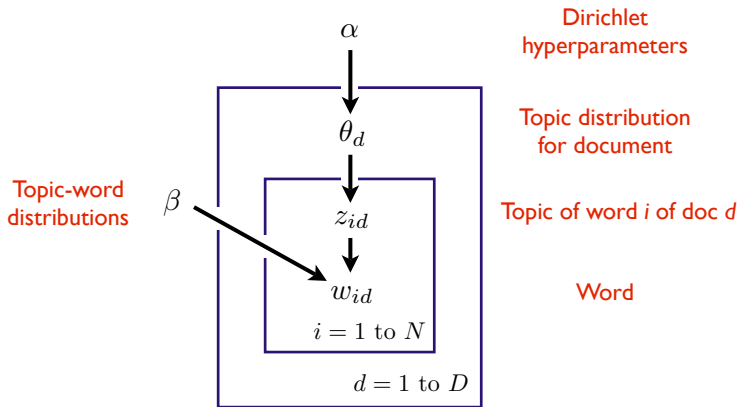
SCIENCE • VOL. 272 • 24 MAY 1996

## Topic proportions and assignments

 $z_{1d}$ 
 $\theta_d$ 
 $z_{Nd}$ 

(Blei, *Introduction to Probabilistic Topic Models*, 2011)

# “Plate” notation for LDA model



Variables within a plate are replicated in a conditionally independent manner

- How to learn topic models?
  - Importance of hyperparameters
  - Choosing number of topics
  - Evaluating topic models
- Examples of extending LDA
  - Polylingual topic models
  - Author-topic model



# Learning algorithm: Gibbs Sampling

By putting a prior distribution on the parameters, they become random variables which can be sampled within the Gibbs Sampling algorithm:

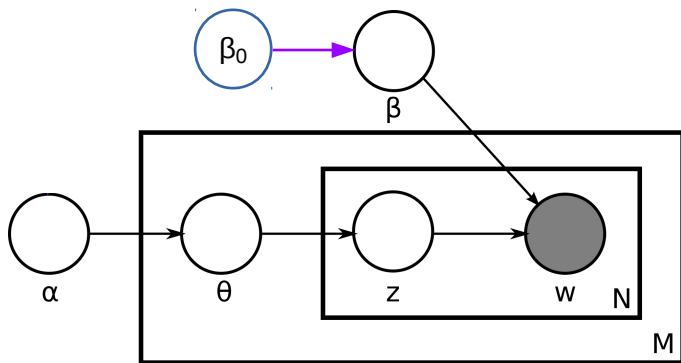


Figure: Putting a Bayesian prior on the parameters:  $\beta \sim \text{Dirichlet}(\cdot; \beta_0)$

# Collapsed Gibbs sampler (Griffiths and Steyvers '04)

- Learn using a *collapsed* Gibbs sampler
- After marginalizing out  $\theta_d$  for all documents  $d$  and  $\beta$ , we get:

$$P(z_i = t \mid \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,t}^{(w_i)} + \beta_0}{n_{-i,t}^{(\cdot)} + W\beta_0} \frac{n_{-i,t}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

$n$  derived from  $z$ , the assignments of words to topics  
( $W$  words,  $T$  topics, and uniform hyperparameters  $\alpha$  and  $\beta_0$ )

- First ratio is probability of  $w_i$  under topic  $t$ , second ratio is probability of topic  $t$  in document  $d_i$
- Given a sample, can get an estimate for  $\beta$  and  $\theta_d$  by:

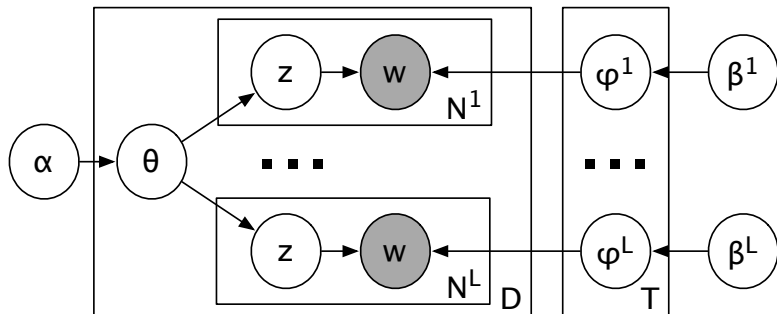
$$\hat{\beta}_{w,t} = \frac{n_t^{(w)} + \beta_0}{n_t^{(\cdot)} + W\beta_0}$$

$$\hat{\theta}_t^{(d)} = \frac{n_t^{(d)} + \alpha}{n_{\cdot}^{(d)} + T\alpha}$$

# Polylingual topic models (Mimno et al., EMNLP '09)

- Goal: topic models that are aligned across languages
- Training data: corpora with multiple documents in each language
  - EuroParl corpus of parliamentary proceedings (11 western languages; exact translations)
  - Wikipedia articles (12 languages; not exact translations)
- How to do this?

# Polylingual topic models (Mimno et al., EMNLP '09)



- DA centralbank europæiske ecb s lån centralbanks
- DE zentralbank ezb bank europäischen investitionsbank darlehen
- EL τράπεζα τράπεζας κεντρική εκτ κεντρικής τράπεζες
- EN **bank central ecb banks european monetary**
- ES banco central europeo bce bancos centrales
- FI keskuspankin eksp n euroopan keskuspankki eip
- FR banque centrale bce européenne banques monétaire
- IT banca centrale bce europea banche prestiti
- NL bank centrale ecb europese banken leningen
- PT banco central europeu bce bancos empréstimos
- SV centralbanken europeiska ecb centralbankens s lån

- DA børn familie udnyttelse børns børnene seksuel
- DE kinder kindern familie ausbeutung familien eltern
- EL παιδιά παιδιών οικογένεια οικογένειας γονείς παιδικής
- EN **children family child sexual families exploitation**
- ES niños familia hijos sexual infantil menores
- FI lasten lapsia lapset perheen lapsen lapsiin
- FR enfants famille enfant parents exploitation familles
- IT bambini famiglia figli minori sessuale sfruttamento
- NL kinderen kind gezin seksuele ouders familie
- PT crianças família filhos sexual criança infantil
- SV barn barnen familjen sexuellt familj utnyttjande

- How would you use this?
- How could you extend this?

# Author-topic model (Rosen-Zvi et al., UAI '04)

- Goal: topic models that take into consideration author *interests*
- Training data: corpora with label for who wrote each document
  - Papers from NIPS conference from 1987 to 1999
  - Twitter posts from US politicians
- Why do this?
- How to do this?



# Author-topic model (Rosen-Zvi et al., UAI '04)

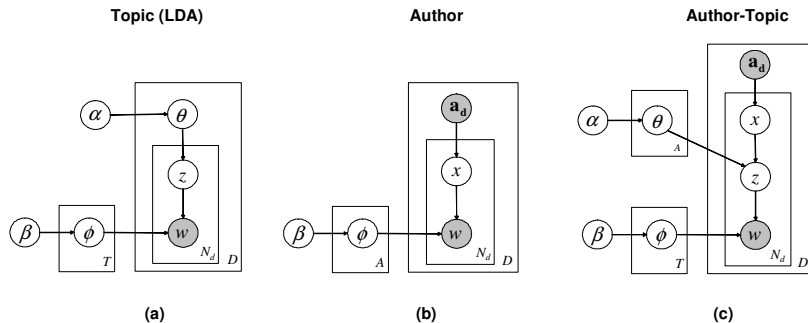


Figure 1: Generative models for documents. (a) Latent Dirichlet Allocation (LDA; Blei et al., 2003), a topic model. (b) An author model. (c) The author-topic model.

# Most likely author for a topic

TOPIC 31	
WORD	PROB.
SPEECH	0.0823
RECOGNITION	0.0497
HMM	0.0234
SPEAKER	0.0226
CONTEXT	0.0224
WORD	0.0166
SYSTEM	0.0151
ACOUSTIC	0.0134
PHONEME	0.0131
CONTINUOUS	0.0129

AUTHOR	PROB.
Waibel_A	0.0936
Makhoul_J	0.0238
De-Mori_R	0.0225
Bourlard_H	0.0216
Cole_R	0.0200
Rigoll_G	0.0191
Hochberg_M	0.0176
Franco_H	0.0163
Abrash_V	0.0157
Movellan_J	0.0149

TOPIC 61	
WORD	PROB.
BAYESIAN	0.0450
GAUSSIAN	0.0364
POSTERIOR	0.0355
PRIOR	0.0345
DISTRIBUTION	0.0259
PARAMETERS	0.0199
EVIDENCE	0.0127
SAMPLING	0.0117
COVARIANCE	0.0117
LOG	0.0112

AUTHOR	PROB.
Bishop_C	0.0563
Williams_C	0.0497
Barber_D	0.0368
MacKay_D	0.0323
Tipping_M	0.0216
Rasmussen_C	0.0215
Opper_M	0.0204
Attias_H	0.0155
Sollich_P	0.0143
Schottky_B	0.0128

TOPIC 71	
WORD	PROB.
MODEL	0.4963
MODELS	0.1445
MODELING	0.0218
PARAMETERS	0.0205
BASED	0.0116
PROPOSED	0.0103
OBSERVED	0.0100
SIMILAR	0.0083
ACCOUNT	0.0069
PARAMETER	0.0068

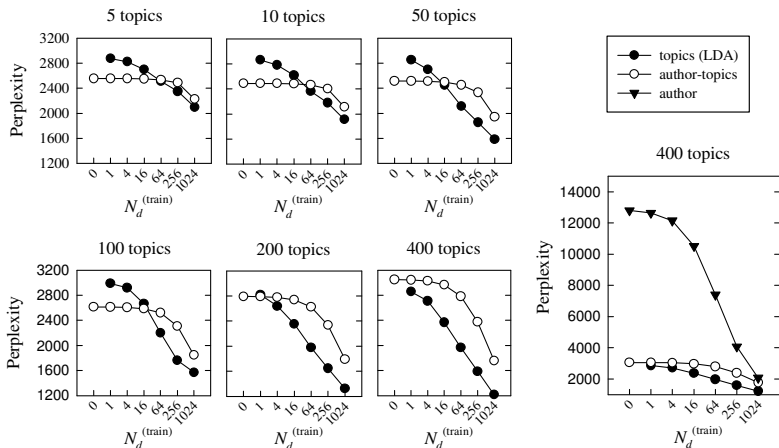
AUTHOR	PROB.
Omohundro_S	0.0088
Zemel_R	0.0084
Ghahramani_Z	0.0076
Jordan_M	0.0075
Sejnowski_T	0.0071
Atkeson_C	0.0070
Bower_J	0.0066
Bengio_Y	0.0062
Revow_M	0.0059
Williams_C	0.0054

TOPIC 100	
WORD	PROB.
HINTON	0.0329
VISIBLE	0.0124
PROCEDURE	0.0120
DAYAN	0.0114
UNIVERSITY	0.0114
SINGLE	0.0111
GENERATIVE	0.0109
COST	0.0106
WEIGHTS	0.0105
PARAMETERS	0.0096

AUTHOR	PROB.
Hinton_G	0.2202
Zemel_R	0.0545
Dayan_P	0.0340
Becker_S	0.0266
Jordan_M	0.0190
Mozer_M	0.0150
Williams_C	0.0099
de-Sa_V	0.0087
Schraudolph_N	0.0078
Schmidhuber_J	0.0056

# Perplexity as a function of number of observed words



$$\text{perplexity}(\mathbf{w}_{test,d} \mid \mathbf{w}_{train,d}, \mathbf{a}_d) = \exp \left[ -\frac{\ln p(\mathbf{w}_{test,d} \mid \mathbf{w}_{train,d}, \mathbf{a}_d)}{N_{test,d}} \right]$$