

Lecture 3: Topic Models

Mirella Lapata

School of Informatics
University of Edinburgh
mlap@inf.ed.ac.uk

January 24, 2016
Slides based on Blei 2011

1 Vector Space Models

2 Topic Models

- Latent Dirichlet Allocation
- Inference with Gibbs Sampling
- Representing Meaning with LDA

3 Discussion

Reading: Griffiths et al. (2007).

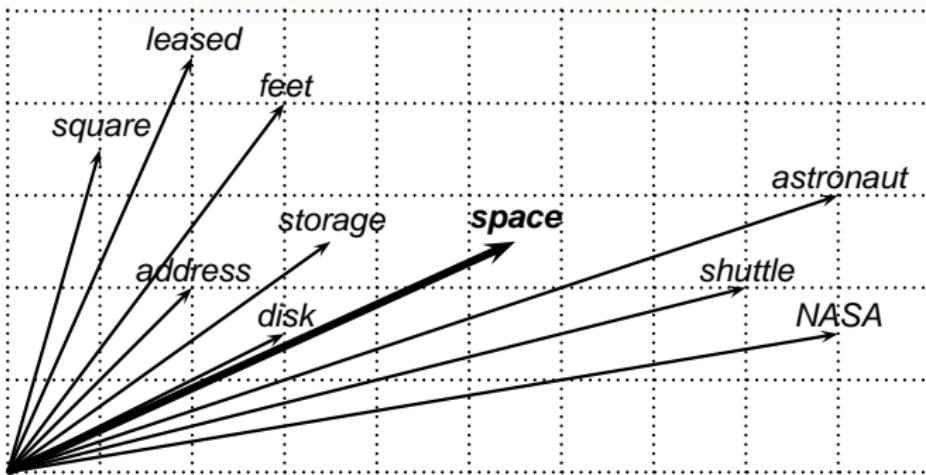
Motivation

We'd like to represent the meaning of words and documents, to help us do things like:

- Retrieve documents relevant to a query.
- Extract features for downstream applications (e.g., sentiment classification, entailment recognition, MT)
- Explain human learning and processing of words (word associations, speed of acquisition, etc).
- Analyze large document collections.

Vector-based Models

Represent the meaning of words in high-dimensional space:



Multiple word meanings are conflated, no probability model.

Topic Models

- We want to find themes (or topics) in documents which is useful for e.g., search or browsing.
- We don't want to do supervised topic classification (neither fix topics in advance nor do manual annotation).
- Approach must automatically tease out the topics.
- Essentially a **clustering** problem: both words and documents are being clustered.

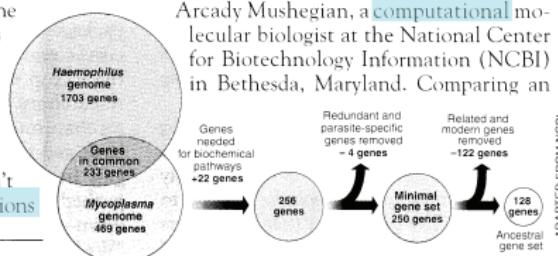
Topic Models

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 of 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



ADAPTED FROM NCBI

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

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Latent Dirichlet Allocation

Topics	
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
...	
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,⁹ 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, a biologist at Stockholm University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a question of numbers. Genes, particularly a more and more genomes are being completed 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 all

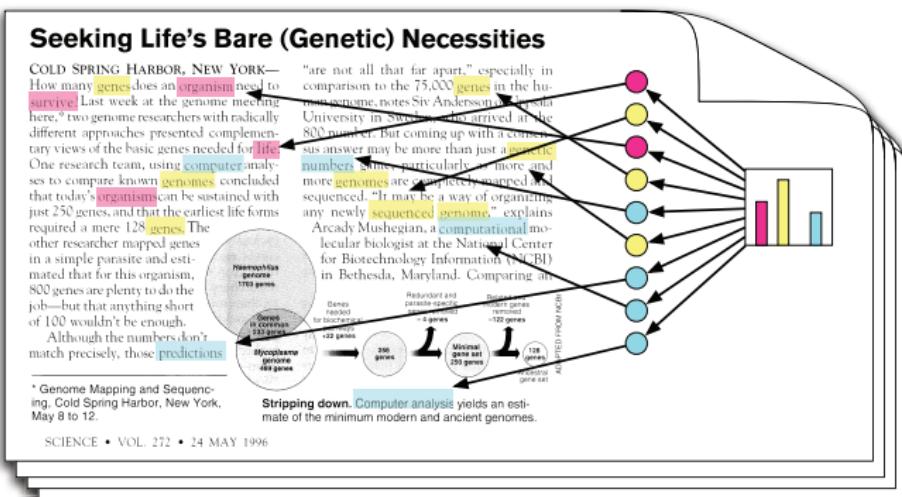
the genomes, he says, "is the only way to find the minimum set of genes that are needed for life."

* Genome Mapping and Sequencing. Cold Spring Harbor, New York, May 8 to 12.

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

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



- 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

The Posterior Distribution

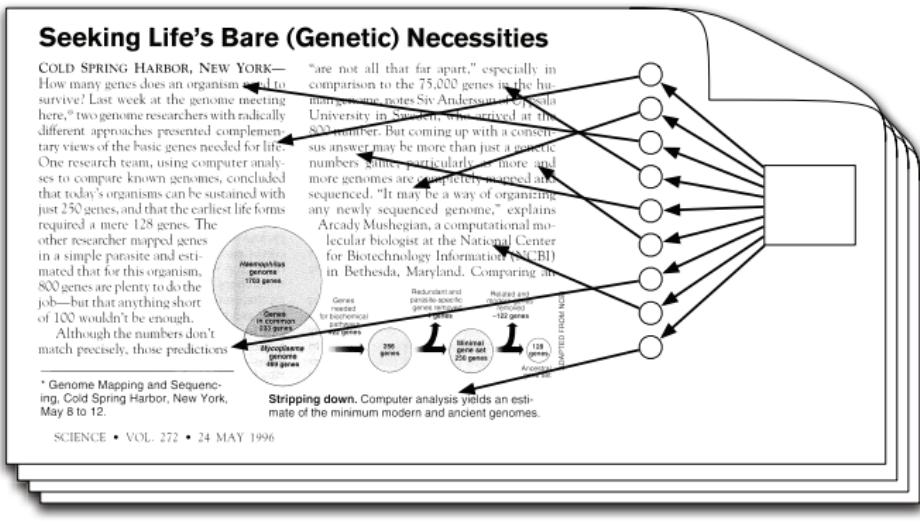
Topics



Documents



Topic proportions and assignments



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

The Posterior Distribution

Topics



Documents



Topic proportions and assignments



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,¹ 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, these predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who served at the SGD meeting. 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 all

Genes (10,000 genes, 123 genes)
 Genes (10,000 genes, 299 genes)
 Genes (10,000 genes, 122 genes)
 Genes (10,000 genes, 118 genes)

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

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

Reduced and
 parasite-specific
 genes
 123 genes

Reduced and
 modern-specific
 genes
 122 genes

Reduced and
 ancient-specific
 genes
 118 genes

Adapted from NCBI

- Our goal is to **infer** the hidden variables
- I.e., compute their distribution conditioned on the documents
 $p(\text{topics, proportions, assignments} | \text{documents})$

LDA: Key Assumptions

- Documents exhibit multiple topics (but typically not many)
- LDA is a probabilistic model with a corresponding **generative process** (each document is generated by this process)
- A **topic** is a distribution over a fixed vocabulary (topics are assumed to be generated first, before the documents)
- Only the number of topics is specified in advance

Document Generation

- ① Choose a number of words N the document will have
- ② Choose topic mixture for document according to Dirichlet distribution over a set of T topics

Dirichlet distribution $\text{Dir}(\alpha)$; α is **prior**: what kind of topic mixtures can I generally expect (determines shape of distribution)

- ③ Generate each word w_i in the document by:
 - ① randomly choosing a topic from the distribution over topics
 - ② randomly choose a word from the corresponding topic (distribution over the vocabulary)

Dirichlet distribution $\text{Dir}(\beta)$

- ⚠ Note we need a distribution over a distribution (for step 1)
- ⚠ Note that words are generated independently of other words (unigram bag-of-words model)

Document Generation: Example

- ① Pick 5 to be the number of words in D .
- ② Decide D will be 1/2 about **food** and 1/2 about **cute animals**.
- ③ Pick first word to come from the **food** topic: *broccoli*.
- ④ Pick second word to come from **cute animals** topic: *panda*.
- ⑤ Pick third word from the **cute animals** topic: *adorable*.
- ⑥ Pick fourth word from the **food** topic: *cherries*.
- ⑦ Pick fifth word to come from **food** topic: *eating*.

What is the document generated under the LDA model?

Learning

- ① Go through each document, and randomly assign each word in the document to one of T topics.
- ② For each document d , go through each word w in d and for each topic t compute:
 - ① $p(\text{topic } t | \text{document } d)$
 - ② $p(\text{word } w | \text{topic } t)$
- ③ Reassign w a new topic, where we choose topic t with probability: $p(\text{topic } t | \text{document } d) * p(\text{word } w | \text{topic } t)$
- ④ After repeating the previous step a large number of times, you'll eventually reach a roughly steady state where your assignments are pretty good.

Generative Process

For $j = 1 \dots T$ topics,

Choose $\varphi^{(j)} \sim \text{Dirichlet}(\beta)$.

$\varphi_1^{(j)} \dots \varphi_V^{(j)}$: prob. of each wd. in topic j

For $d = 1 \dots D$ documents,

Choose $\theta^{(d)} \sim \text{Dirichlet}(\alpha)$.

$\theta_1^{(d)} \dots \theta_T^{(d)}$: prob. of each topic in doc. d

For $i = 1 \dots N_d$ words in doc d ,

Choose $z_i \sim \text{Multinomial}(\theta^{(d)})$.

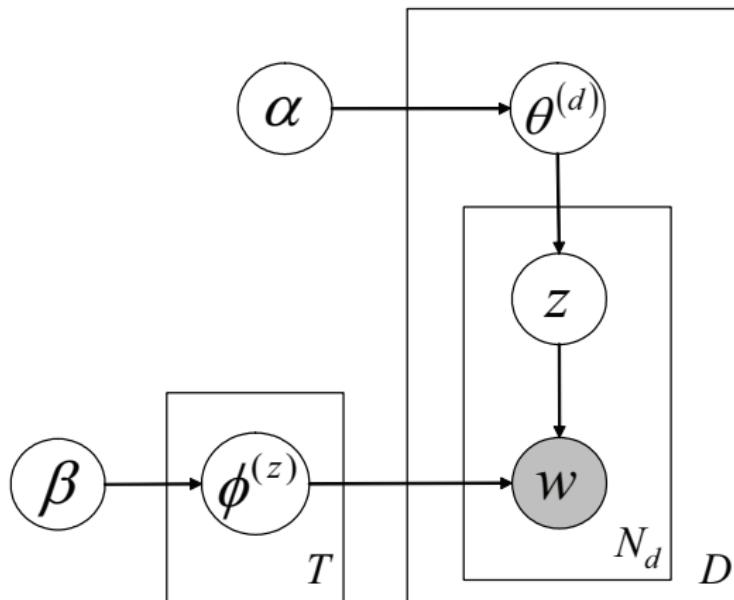
z_i : topic of word i

Choose $w_i \sim \text{Multinomial}(\varphi^{(z_i)})$.

w_i : identity of word i

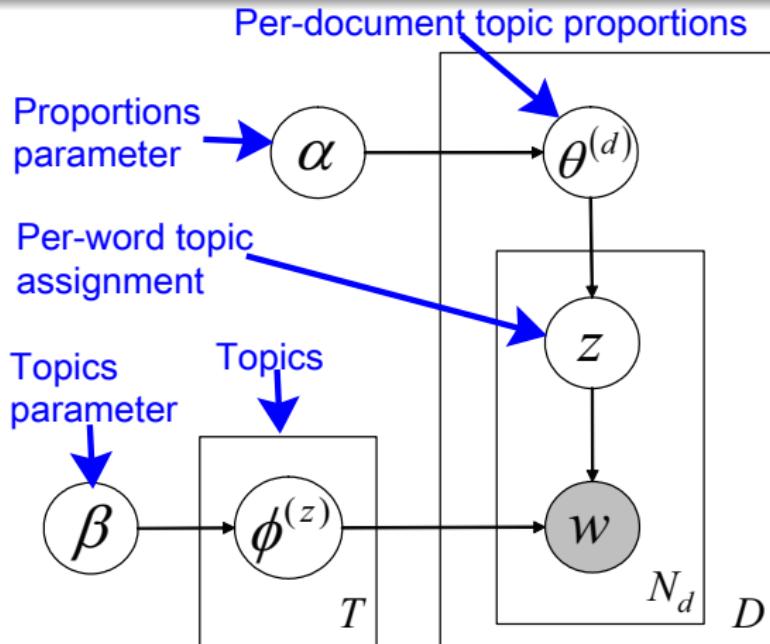
Note that, except in the appendices, Griffiths et al. (2007) use g (for gist) to refer to θ . They also use the term Discrete rather than Multinomial.

The Graphical Model



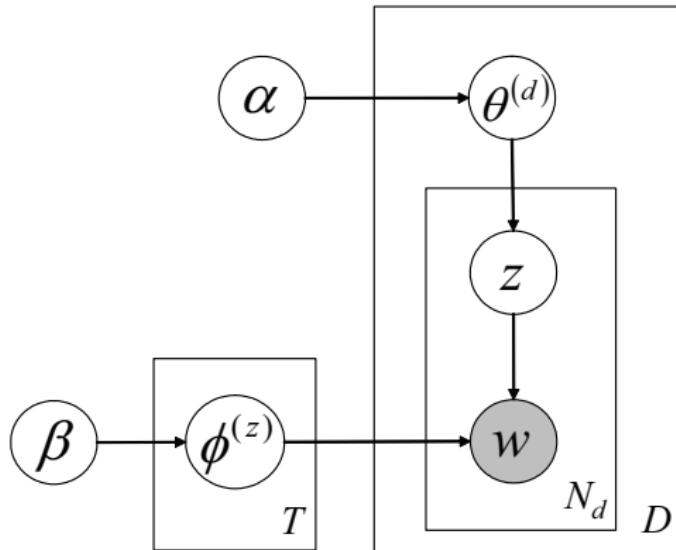
- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; plates \approx replicated variables.

The Graphical Model

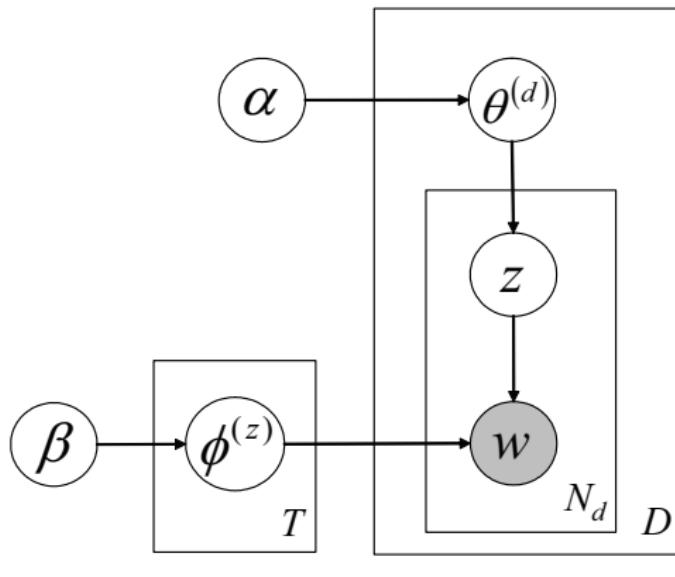


- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; plates \approx replicated variables.

The Graphical Model



The Graphical Model



$$P(w_i) = \sum_{j=1}^T P(w_i|z_i=j)P(z_i=j)$$

- $P(z_i = j)$ is probability that j th topic was sampled for i th word token
- $P(w_i|z_i = j)$ is probability of word w_i under topic j
- $\phi^{(j)} = P(z_i = j)$
- $\theta^{(d)} = P(z)$

Aside: Multinomial Distribution

For $x_i \in \{0, \dots, n\}$

$$P(\mathbf{x}|\theta) = \frac{n!}{\prod_{i=1}^d x_i!} \prod_{i=1}^d \theta_i^{x_i}, \quad n = \sum_{i=1}^d x_i, \quad \sum_{i=1}^d \theta_i = 1, \quad \theta_i > 0$$

When $n = 1$ the multinomial distribution simplifies to:

$$P(\mathbf{x}|\theta) = \prod_{i=1}^d \theta_i^{x_i}, \quad \sum_{i=1}^d \theta_i = 1, \quad \theta_i > 0$$

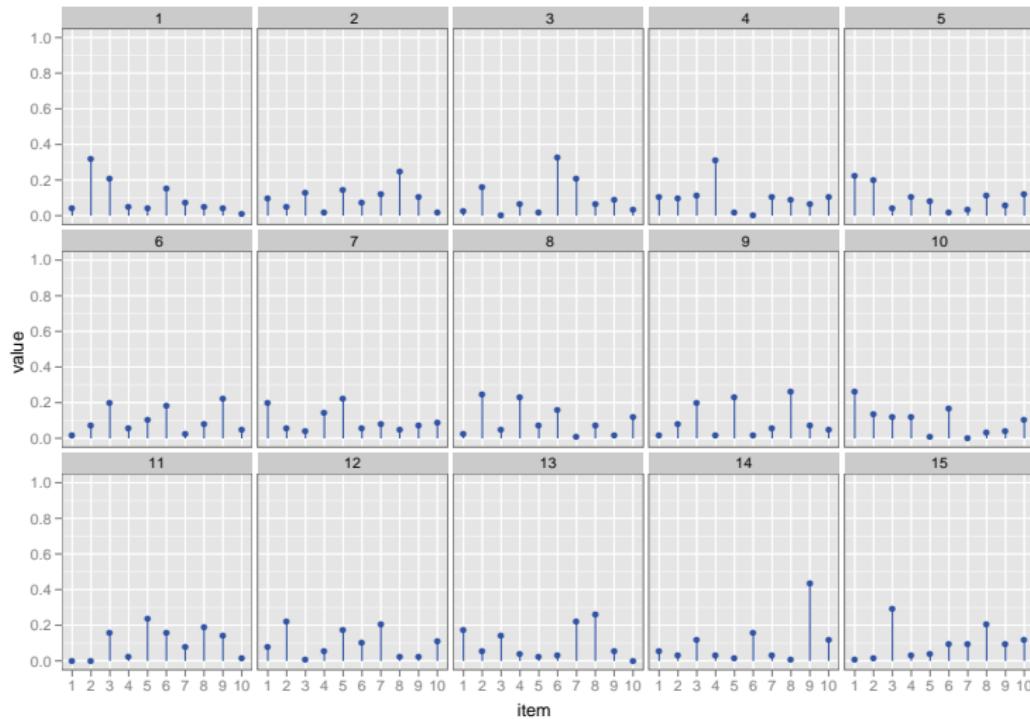
- unigram language model, **1-of-V coding** ($d = V$ vocab size)
- x_i indicates word i of the vocabulary observed ($x_i = 1$ if word i is observed and 0 otherwise)
- $\theta_i = P(w_i)$ the probability that word i is seen

Dirichlet Distribution

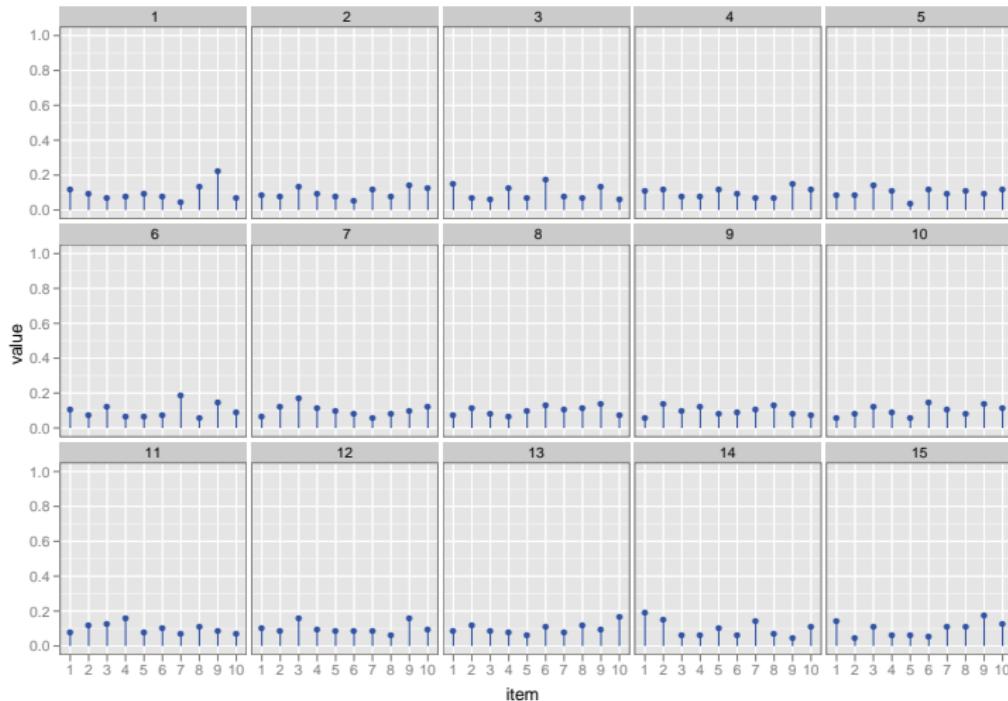
$$\text{Dir}(\alpha_1, \dots, \alpha_T) \frac{\Gamma(\sum_j \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^T p_j^{\alpha_j - 1}$$

- The Dirichlet distribution is an exponential family distribution over the simplex, i.e., positive vectors that sum to one
- It is **conjugate to the multinomial**. Given a multinomial observation, the posterior distribution of θ is a Dirichlet.
- Parameter α smoothes **topic distribution** in the document.
- Parameter β smoothes **word distribution** in every topic.

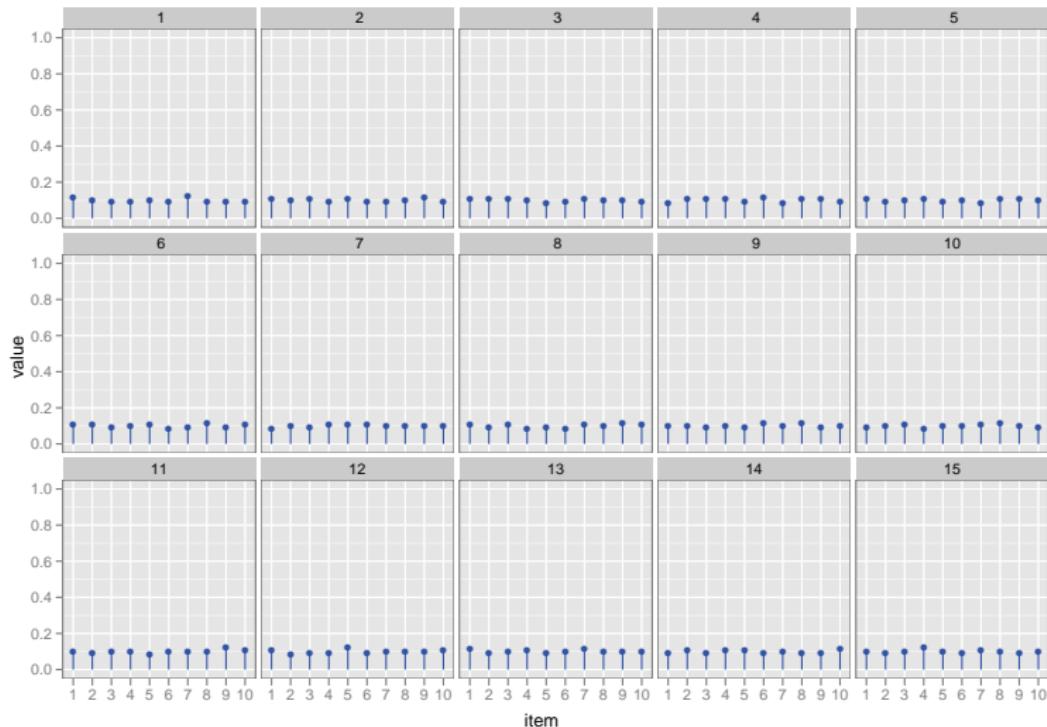
$$\alpha = 1$$



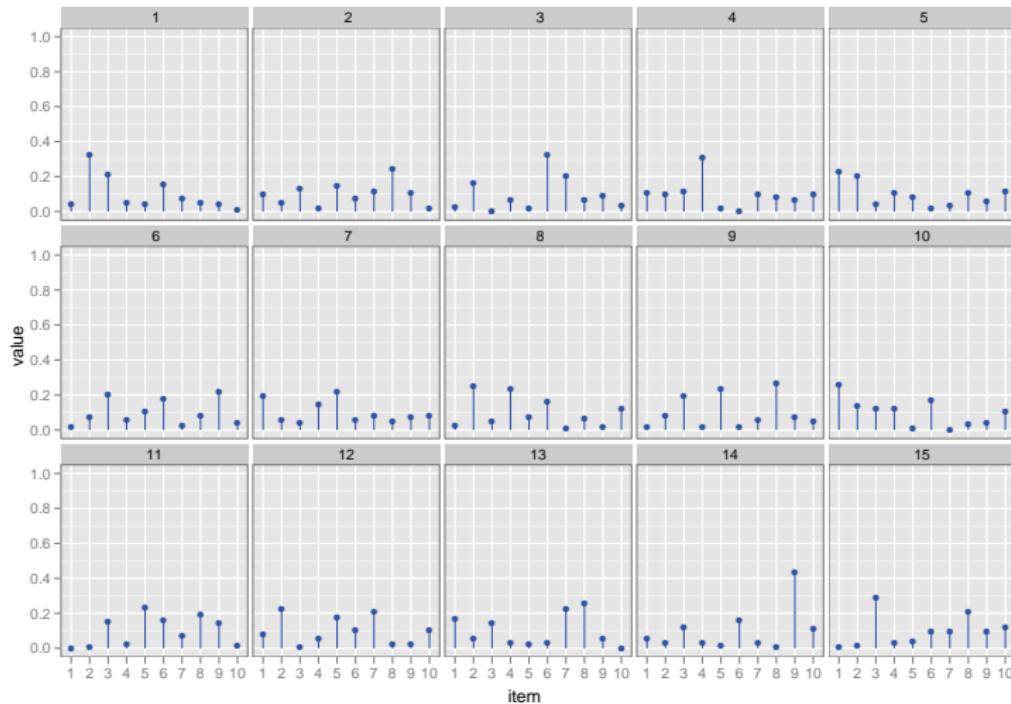
$$\alpha = 10$$



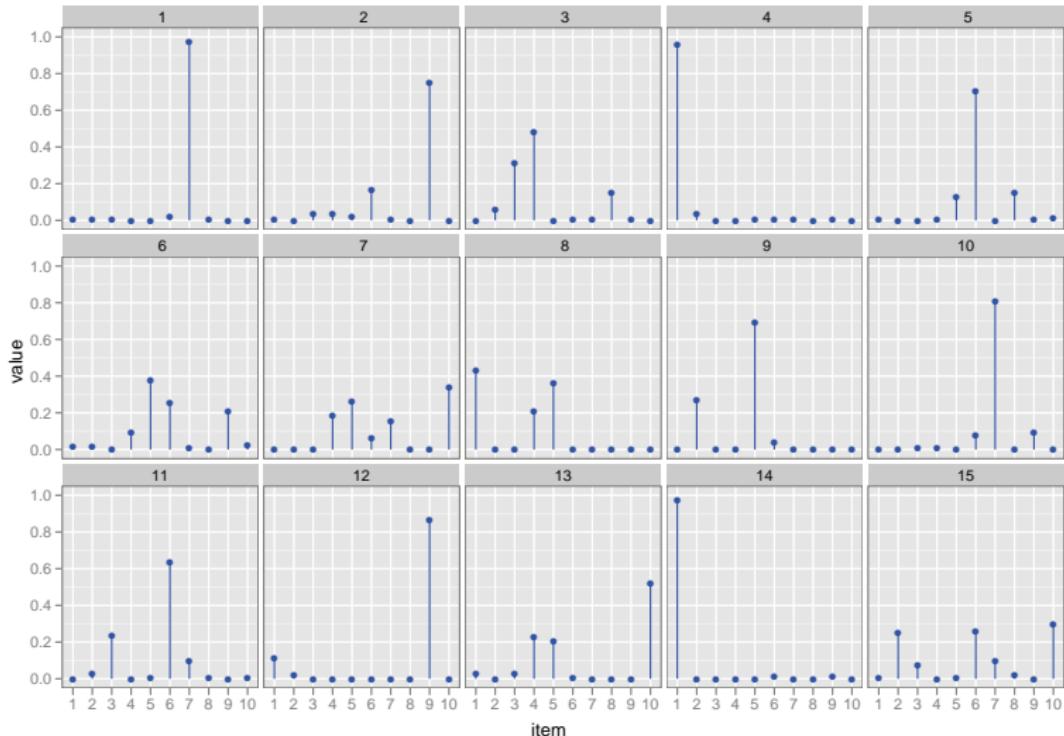
$$\alpha = 100$$



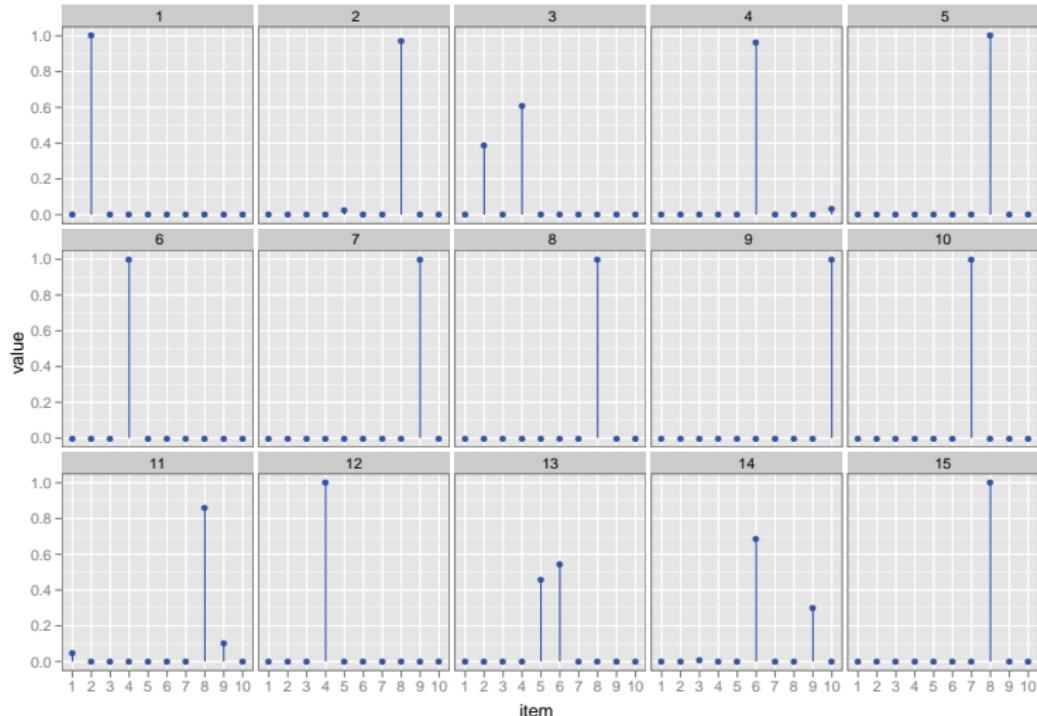
$$\alpha = 1$$



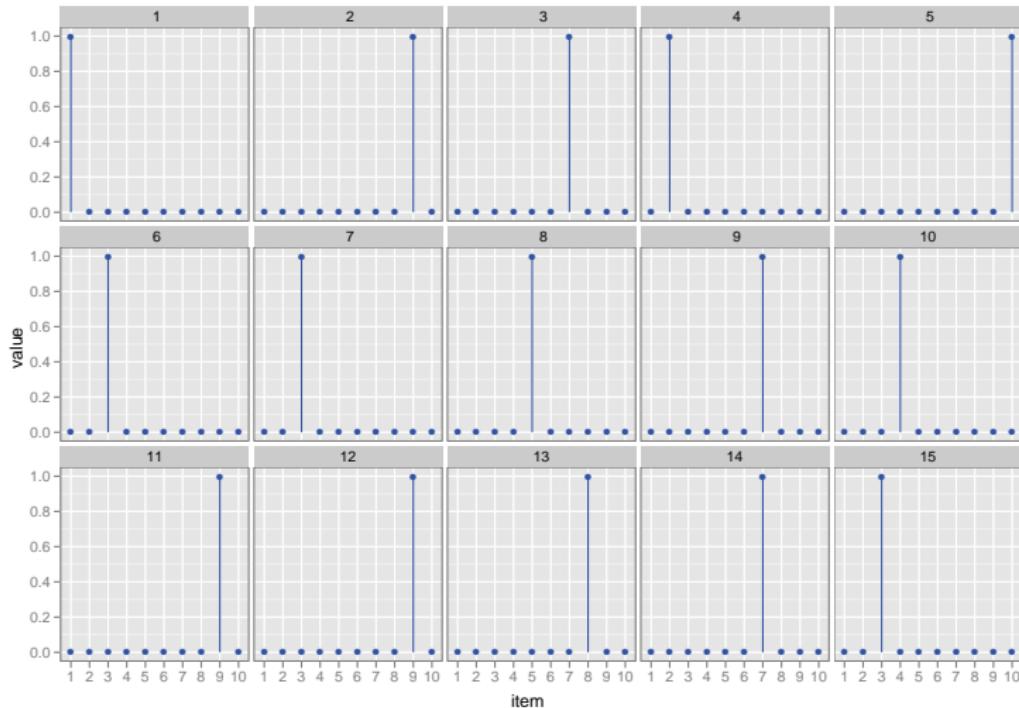
$$\alpha = 0.1$$



$$\alpha = 0.01$$



$$\alpha = 0.001$$



Inference with Gibbs Sampling

- An iterative process.
- Start with random topic assignments for each word.
- In each iteration, *for each word* in the data:
 - Assume you know (from the prev. iteration) the topics of all other words. (pretend they are correct)
 - Determine the probabilities of each topic-assignment given the rest of the data.
 - Choose the most probable assignment.
- Iterate until convergence.

Gibbs Sampling

- Collection of documents is a set of word indices w_i and document indices d_i , for each word token i .

$$P(z_i = j | z_{-i}, w_i, d_i, \cdot) = \frac{C_{d_i j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i t}^{DT} + T\alpha} \frac{C_{w_i j}^{WT} + \beta}{\sum_{w=1}^W C_{w j}^{WT} + W\beta}$$

- From this conditional distribution, a topic is sampled and stored as the new topic assignment for this word token.

$z_i = j$: topic assignment of token i to topic j

z_{-i} : the topic assignments of all other word tokens

\cdot : all other known information (w_{-i} , d_{-i} , α , β).

C^{WT} : matrix of counts with dimensions $W \times T$

C^{DT} : matrix of counts with dimensions $D \times T$

Posterior Estimates of β and θ

$$\phi_{ij} = \frac{C_{ij}^{WT} + \beta}{\sum_{t=1}^w C_{w;j}^{WT} + W\beta}$$

$$\theta_{dj} = \frac{C_{d;j}^{DT} + \alpha}{\sum_{t=1}^T C_{d;t}^{DT} + T\alpha}$$

- ϕ_{ij} : normalizes word-topic count matrix (probability of word type i for topic j)
- θ_{dj} : normalizes the counts in the document-topic count matrix (probability that document d belongs to topic j).
- technically the posterior mean of the parameters ϕ and θ

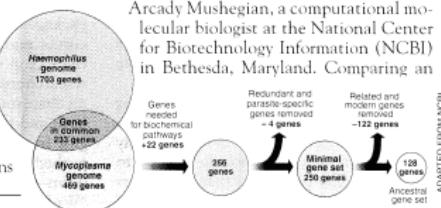
Example Inference

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 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.

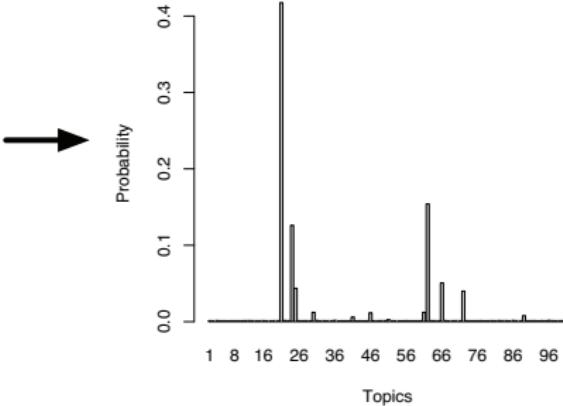
Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of 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.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.



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

Example Inference

Chaotic Beetles

Charles Godfray and Michael Hassell

Ecologists have known since the pioneering work of May in the mid-1970s (1) that the population dynamics of animals and plants can be exceedingly complex. This complexity arises from two sources: The tangled web of interactions that constitute any natural community provide a myriad of different pathways for species to interact, both directly and indirectly. And even in isolated populations the nonlinear feedback processes present in all natural populations can result in complex dynamic behavior. Natural populations can show persistent oscillatory dynamics and chaos, the latter characterized by extreme sensitivity to initial conditions. If such chaotic dynamics were common in nature, then this would have important ramifications for the management and conservation of natural resources. On page 389 of this issue, Costantino *et al.* (2) provide the most

convincing evidence to date of complex dynamics and chaos in a biological population—of the flour beetle, *Tribolium castaneum* (see figure).

It has proven extremely difficult to demonstrate complex dynamics in populations in the field. By its very nature, a chaotically fluctuating population will superficially resemble a stable or cyclic population buffeted by the normal random perturbations experienced by all species. Given a long enough time series, diagnostic tools from nonlinear mathematics can be used to identify the tell-tale signatures of chaos. In phase space, chaotic trajectories come to lie on “strange attractors,” curious geometric objects with fractal structure and hence noninteger dimension. As they

move over the surface of the attractor, sets of adjacent trajectories are pulled apart, then stretched and folded, so that it becomes impossible to predict exact population densities into the future. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured mathematically estimating the Liapunov exponent,

which is positive for chaotic dynamics and nonpositive otherwise. There have been many attempts to estimate attractor dimension and Liapunov exponents from time series data, and some candidate chaotic population have been identified (some insects, rodents, and most convincingly, human childhood diseases), but the statistical difficulties preclude any broad generalization (3).

An alternative approach is to parameterize population models with data from natural populations and then compare their predictions with the dynamics in the field. This technique has been gaining popularity in recent years, helped by statistical advances in parameter estimation. Good ex-

The authors are in the Department of Biology, Imperial College at Silwood Park, Ascot, Berks, SL5 7PZ, UK. E-mail: m.hassell@ic.ac.uk

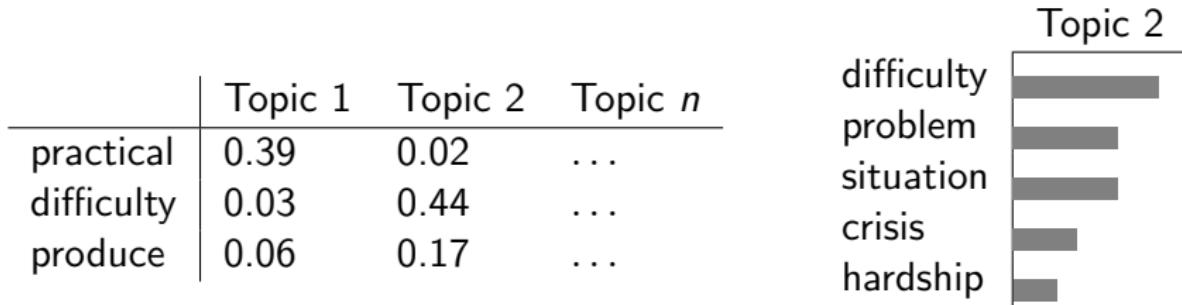


Cannibalism and chaos. The flour beetle, *Tribolium castaneum*, exhibits chaotic population dynamics when the amount of cannibalism is altered in a mathematical model.

Example Inference

problem	model	selection	species
problems	rate	male	forest
mathematical	constant	males	ecology
number	distribution	females	fish
new	time	sex	ecological
mathematics	number	species	conservation
university	size	female	diversity
two	values	evolution	population
first	value	populations	natural
numbers	average	population	ecosystems
work	rates	sexual	populations
time	data	behavior	endangered
mathematicians	density	evolutionary	tropical
chaos	measured	genetic	forests
chaotic	models	reproductive	ecosystem

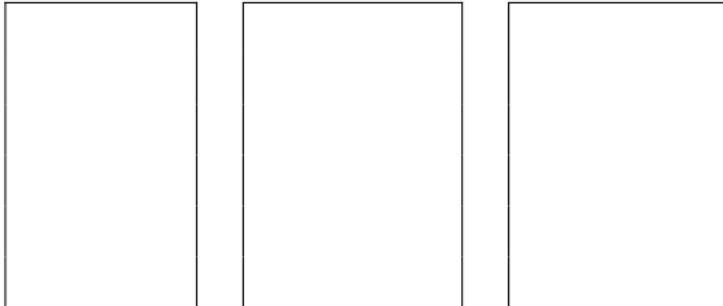
What does this have to do with semantics?



- Topics are the dimensions of the space (500, 1000)
- Vector components: probability of word given topic
- Topics correspond to coarse-grained sense distinctions
- Cosine similarity can be used (probabilistic alternatives)

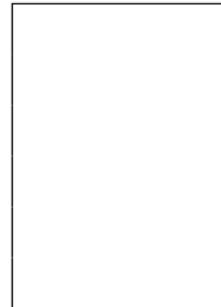
Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, free association, ambiguity, semantic priming, reading time, free recall



Model Evaluation

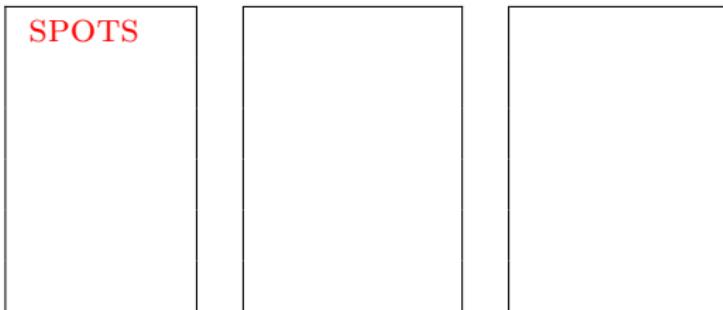
- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.



Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.

SPOTS



Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.

SPOTS

DOG

DIRTY

DIRT

STRIPES

DARK

Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.

SPOTS

DOG

DIRTY

DIRT

STRIPES

DARK

RICE

Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.

SPOTS

DOG

DIRTY

DIRT

STRIPES

DARK

RICE

CHINESE

WEDDING

FOOD

WHITE

CHINA

Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.

SPOTS

DOG

DIRTY

DIRT

STRIPES

DARK

RICE

CHINESE

WEDDING

FOOD

WHITE

CHINA

SEPARATE

Model Evaluation

- Griffiths et al. present a range of evaluations and modeling tasks comparing LDA and LSA.
- Synonym tests from TOEFL, **free association**, ambiguity, semantic priming, reading time, free recall
- Free-association:** produce the first word that comes to mind in response to a cue word.

SPOTS

DOG

DIRTY

DIRT

STRIPES

DARK

RICE

CHINESE

WEDDING

FOOD

WHITE

CHINA

SEPARATE

DIVIDE

DIVORCE

PART

SPLIT

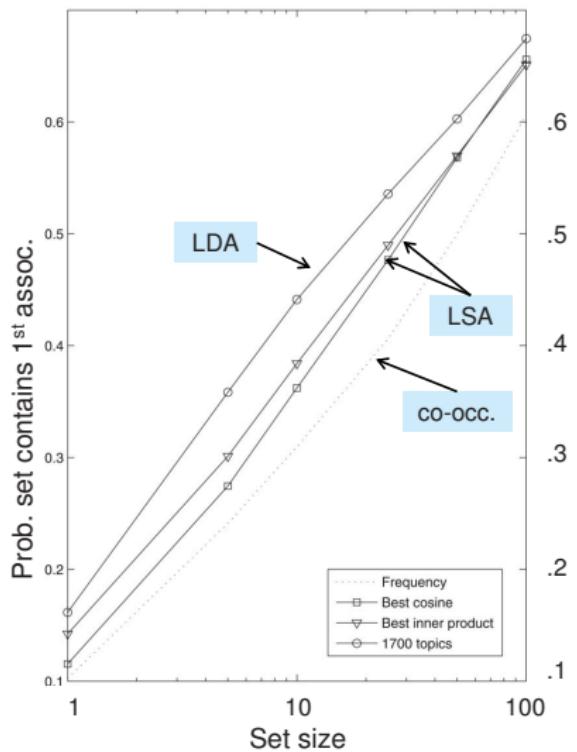
REMOVE

Computing Word Similarity

- Build LSA and LDA representations from corpus of educational materials.
- Compute associates predicted by each model.
- LSA: top associate of w_1 is word with closest cosine (or dot product) similarity.
- LDA: top associate of w_1 is word with highest $P(w_2|w_1)$.

$$P(w_2|w_1) = \sum_z P(w_2|z)P(z|w_1)$$

Results

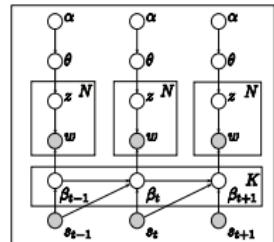
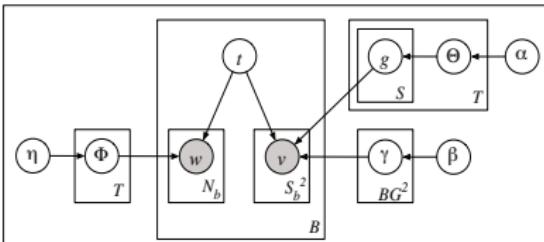
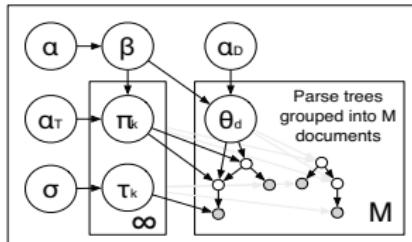
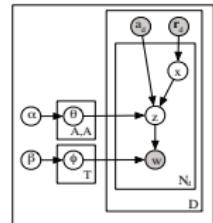
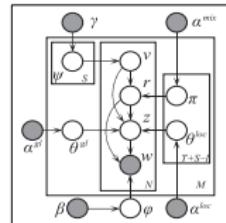
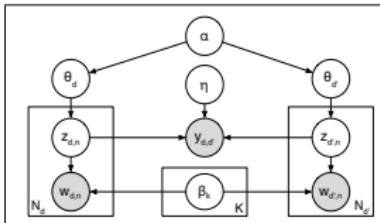
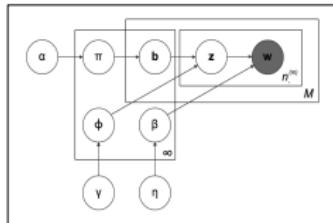


If we take the top n associates returned by the model (set size), what is the probability that this set contains the true first associate?

Discussion

- LSA is useful in practice: IR, unsupervised morphology, essay scoring, text coherence, language modeling, summarization.
- LDA gave rise to many machine learning papers: Hierarchical LDA, nested LDA, author-topic model, topics over time, topics + syntax, etc.
- So far, limited work showing these models useful in NLP.
- LDA for discourse segmentation, extensions for text classification, WSD, tracking social phenomena.
- In comparisons, LSA often performs as well or better.

Discussion



References

Griffiths, Thomas L., Mark Steyvers, and Joshua B. Tenenbaum. 2007. Topics in semantic representation. *Psychological Review* 114(2):211–244.