

TTIC 31210:
Advanced Natural Language Processing

Kevin Gimpel
Spring 2019

Lecture 13:
Introduction to Bayesian NLP

Roadmap

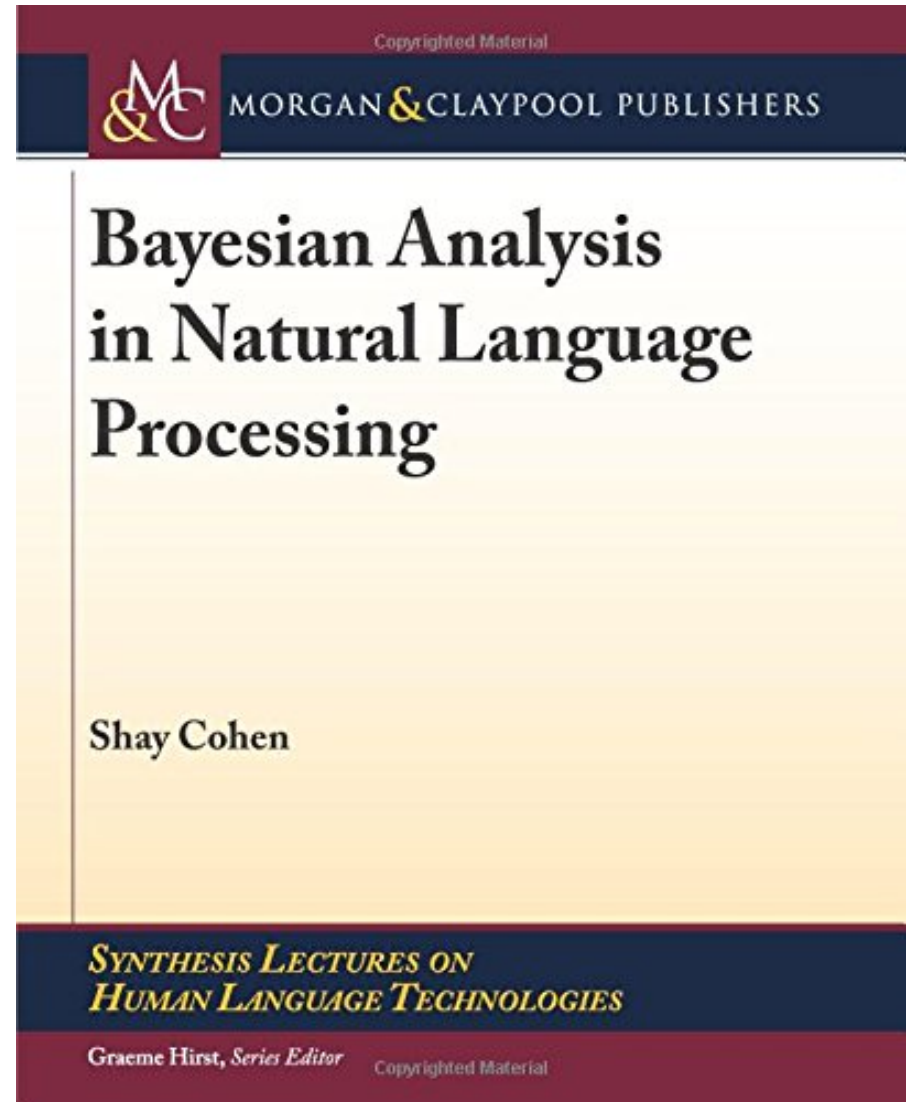
- intro (1 lecture)
- deep learning for NLP (5 lectures)
- structured prediction (4.5 lectures)
- generative models, latent variables, unsupervised learning, variational autoencoders (1.5 lectures)
- **Bayesian methods in NLP (2 lectures)**
- Bayesian nonparametrics in NLP (2 lectures)
- review & other topics (1 lecture)

Assignments

- any questions about Assignment 3?

Additional Reading

- for this segment of the course, the optional text is Cohen (2016, 2019)
- there is a copy of the second edition (2019) in the TTIC library
- readings will be drawn from this book for the next few lectures



Motivation

- in most neural NLP, we assume parameters and architectures are fixed
- consider a one-hidden-layer MLP:

$$p(Y = y \mid \mathbf{x}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

Motivation

- in most neural NLP, we assume parameters and architectures are fixed
- consider a one-hidden-layer MLP:

$$p(Y = y \mid \mathbf{x}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

- now let's be more explicit about what we're conditioning on:

$$p(Y = y \mid \mathbf{x}, \Theta = \{\mathbf{w}, \mathbf{W}\}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

Motivation

- in most neural NLP, we assume parameters and architectures are fixed
- consider a one-hidden-layer MLP:

$$p(Y = y \mid \mathbf{x}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

- now let's be more explicit about what we're conditioning on:

$$p(Y = y \mid \mathbf{x}, \Theta = \{\mathbf{w}, \mathbf{W}\}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

note: the notation above suggests that we can think of parameters as random variables; this is not uncontroversial.

Motivation

$$p(Y = y \mid \mathbf{x}, \Theta = \{\mathbf{w}, \mathbf{W}\}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

- how do we get back to $p(Y = y \mid \mathbf{x})$?

Motivation

$$p(Y = y \mid \mathbf{x}, \Theta = \{\mathbf{w}, \mathbf{W}\}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

- how do we get back to $p(Y = y \mid \mathbf{x})$?
- marginalize over new random variables:

$$p(Y = y \mid \mathbf{x}) = \int_{\Theta} p(Y = y, \Theta = \{\mathbf{w}, \mathbf{W}\} \mid \mathbf{x}) d\Theta$$

Motivation

$$p(Y = y \mid \mathbf{x}, \Theta = \{\mathbf{w}, \mathbf{W}\}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

- how do we get back to $p(Y = y \mid \mathbf{x})$?
- marginalize over new random variables:

$$p(Y = y \mid \mathbf{x}) = \int_{\Theta} p(Y = y, \Theta = \{\mathbf{w}, \mathbf{W}\} \mid \mathbf{x}) d\Theta$$

$$p(Y = y, \Theta = \{\mathbf{w}, \mathbf{W}\} \mid \mathbf{x}) = p(Y = y \mid \Theta = \{\mathbf{w}, \mathbf{W}\}, \mathbf{x}) p(\Theta = \{\mathbf{w}, \mathbf{W}\} \mid \mathbf{x})$$

Motivation

$$p(Y = y \mid \mathbf{x}, \Theta = \{\mathbf{w}, \mathbf{W}\}) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

- how do we get back to $p(Y = y \mid \mathbf{x})$?
- marginalize over new random variables:

$$p(Y = y \mid \mathbf{x}) = \int_{\Theta} p(Y = y, \Theta = \{\mathbf{w}, \mathbf{W}\} \mid \mathbf{x}) d\Theta$$

- intuitively: don't commit to a single set of parameter values; use them all (with a suitable prior distribution)

Going Further...

- marginalize over architectures & parameters:

$$p(Y = y \mid \mathbf{x}, \Lambda = \text{MLP}(\mathbf{w}, \mathbf{W})) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

$$p(Y = y \mid \mathbf{x}) = \int_{\Lambda} p(Y = y, \Lambda = \text{MLP}(\mathbf{w}, \mathbf{W}) \mid \mathbf{x}) d\Lambda$$

Going Further...

- marginalize over architectures & parameters:

$$p(Y = y \mid \mathbf{x}, \Lambda = \text{MLP}(\mathbf{w}, \mathbf{W})) = \frac{\exp\{\mathbf{w}_y^\top \tanh(\mathbf{W}g(\mathbf{x}))\}}{Z}$$

$$p(Y = y \mid \mathbf{x}) = \int_{\Lambda} p(Y = y, \Lambda = \text{MLP}(\mathbf{w}, \mathbf{W}) \mid \mathbf{x}) d\Lambda$$

- the Bayesian framework gives us a vocabulary to discuss this kind of thing and methods for approximating these computations

Why “Bayesian”?

Likelihood

Probability of collecting this data when our hypothesis is true

Bill Howe, UW

Prior

The probability of the hypothesis being true before collecting data

$$P(H|D) = \frac{P(D|H) P(H)}{P(D)}$$

Posterior

The probability of our hypothesis being true given the data collected

Marginal

What is the probability of collecting this data under all possible hypotheses?

Bayesian NLP

- typically used with unsupervised learning:
 - we have data
 - we hypothesize some latent variables through which the data are generated
 - we define the “generative story” that describes how latent variables are generated, then how data is generated using latent variables
 - **new**: we parameterize the distributions & add the parameters themselves to the generative story

Generative Story Template

- 1: Draw a set of parameters θ from $p(\Theta)$
- 2: Draw a latent structure z from $p(Z | \theta)$
- 3: Draw the observed data x from $p(X | z, \theta)$

the above generative story implies the following factorization of the joint distribution:

$$p(x, z, \theta) = p(\theta)p(z | \theta)p(x | z, \theta)$$

Categorical Distribution

- parameterized by a vector of probabilities, one for drawing each outcome
- i.e., prob. of drawing outcome i for variable X :

$$p(X = x_i | \theta) = \theta_i \quad i \in \{1, \dots, K\}$$

- when we want to draw from this distribution, we will write:

$$x \sim \text{Categorical}(\theta)$$

Categorical vs. Multinomial

- “multinomial” is used frequently to mean categorical in this literature, so we’ll use them interchangeably
- a multinomial is actually more general (permits more than 1 instance of an event)

Vector Form of Categorical Distribution

- form we saw earlier:

$$p(X = x_i | \theta) = \theta_i \quad i \in \{1, \dots, K\}$$

- we can also write the categorical distribution as a (one-hot) vector random variable Y :

$$Y_i = \mathbb{I}[X = i] \quad Y \in \{0, 1\}^K$$

$$p(Y = y | \theta) = ?$$

Vector Form of Categorical Distribution

- form we saw earlier:

$$p(X = x_i | \theta) = \theta_i \quad i \in \{1, \dots, K\}$$

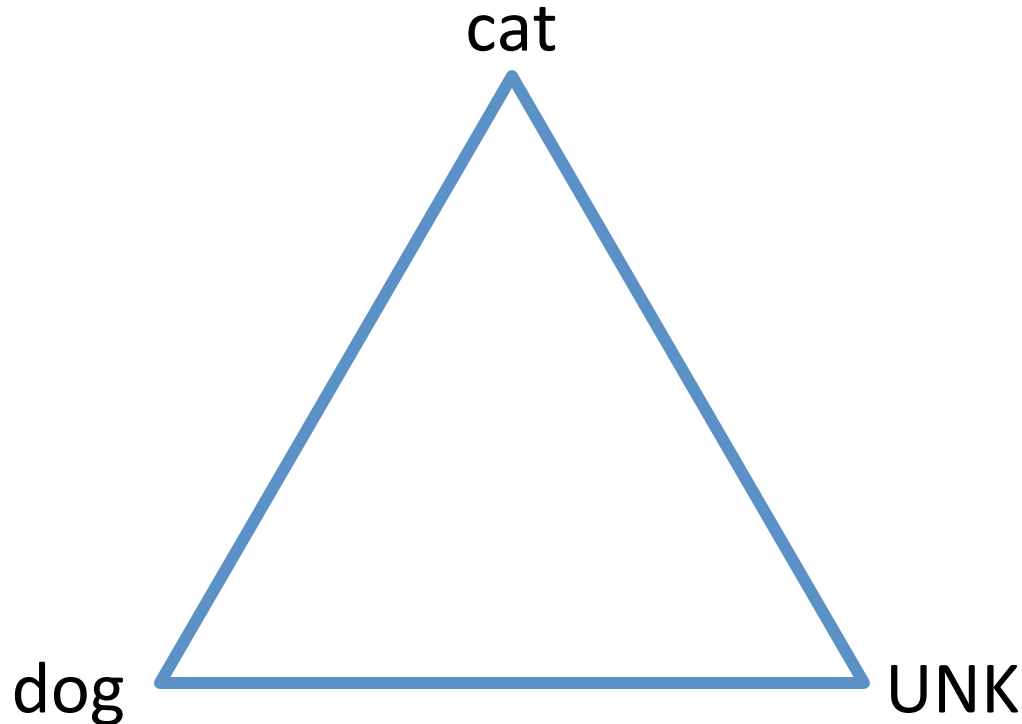
- we can also write the categorical distribution as a (one-hot) vector random variable Y :

$$Y_i = \mathbb{I}[X = i] \quad Y \in \{0, 1\}^K$$

$$p(Y = y | \theta) = \prod_{i=1}^K \theta_i^{y_i}$$

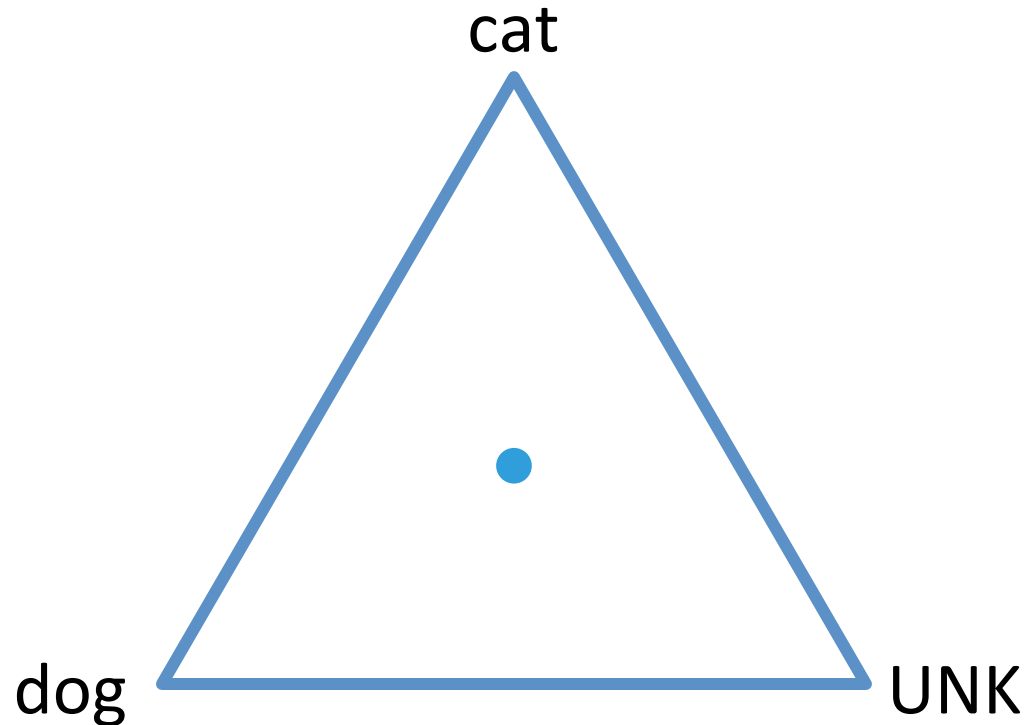
The 2-Simplex

- consider a categorical distribution with 3 outcomes
- e.g., a distribution over words using a vocabulary of size 3:



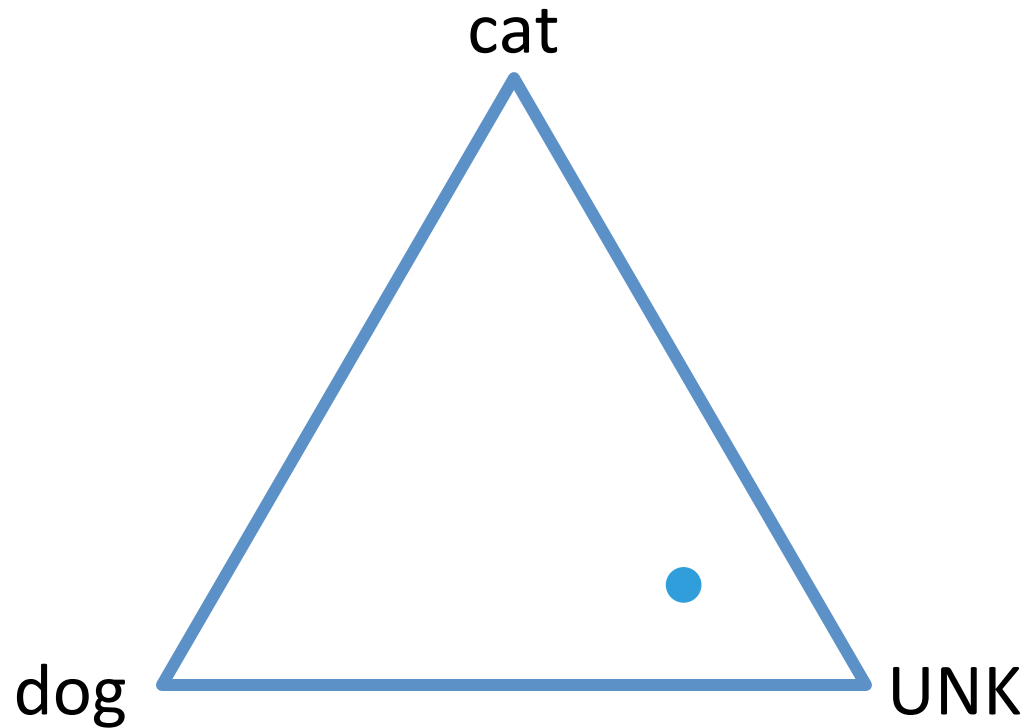
The 2-Simplex

- a point on this simplex represents a categorical distribution over the 3 outcomes
- a uniform distribution:



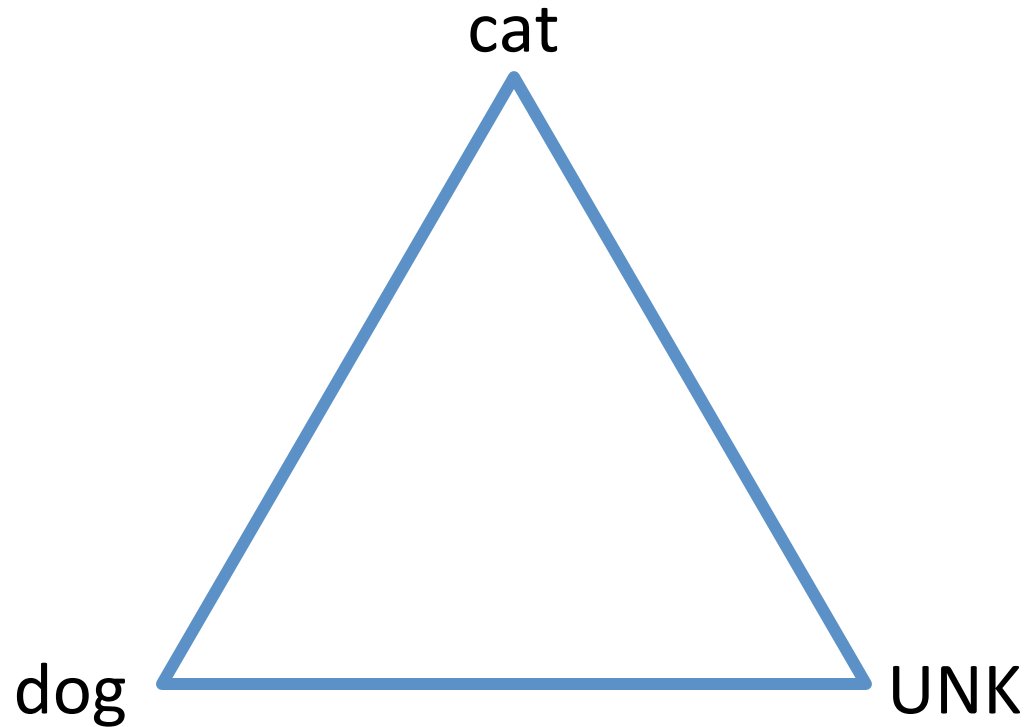
The 2-Simplex

- a distribution that puts most probability on UNK:

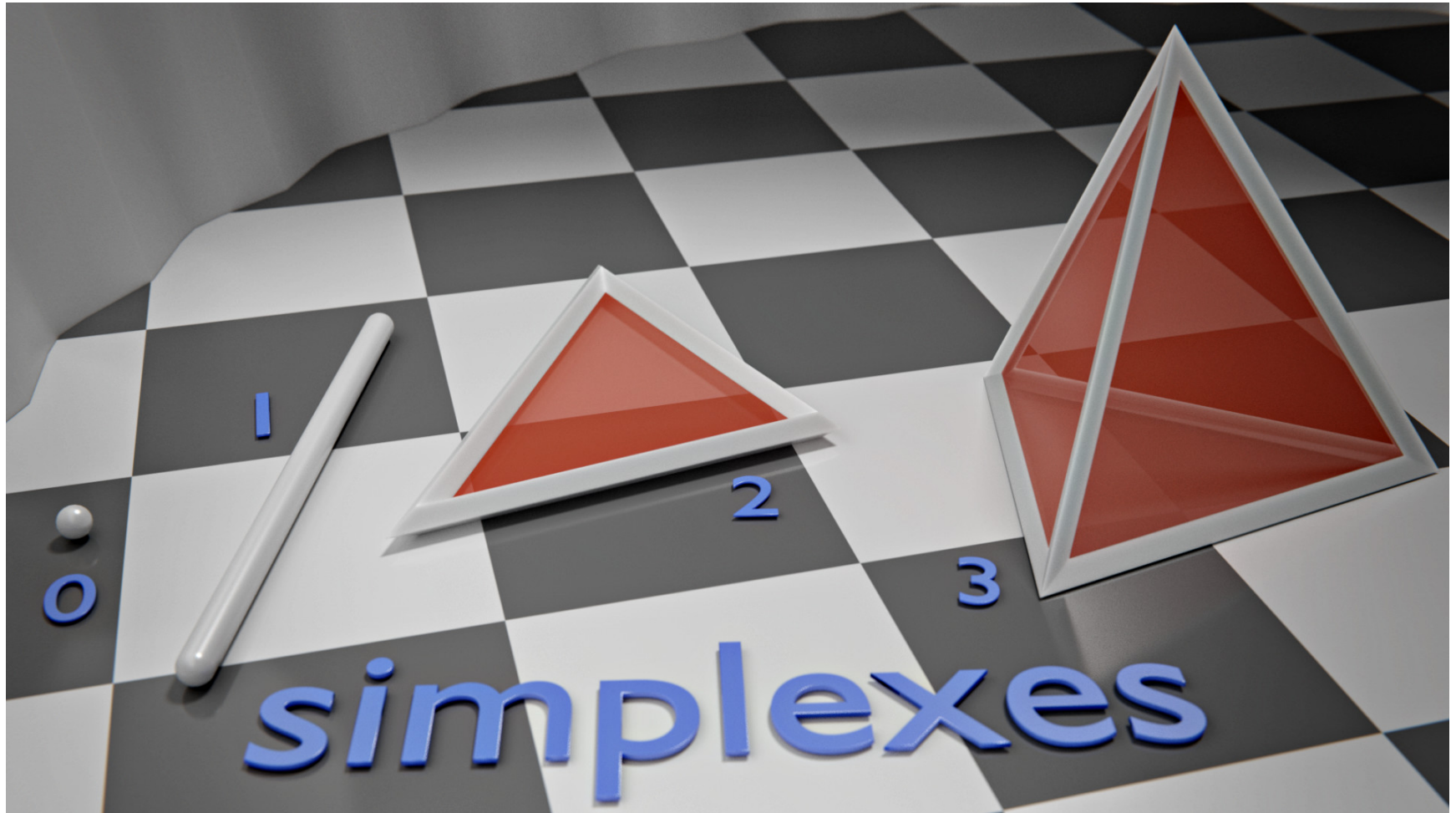


The 2-Simplex

- a categorical distribution with K outcomes has $K-1$ parameters



Categorical Parameters lie in $(K-1)$ -Simplex



Latent Dirichlet Allocation

David M. Blei

*Computer Science Division
University of California
Berkeley, CA 94720, USA*

BLEI@CS.BERKELEY.EDU

Andrew Y. Ng

*Computer Science Department
Stanford University
Stanford, CA 94305, USA*

ANG@CS.STANFORD.EDU

Michael I. Jordan

*Computer Science Division and Department of Statistics
University of California
Berkeley, CA 94720, USA*

JORDAN@CS.BERKELEY.EDU

- generative model for document collections using latent variables that can be interpreted as “topics”
- learns a multinomial distribution over words for each topic

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
...

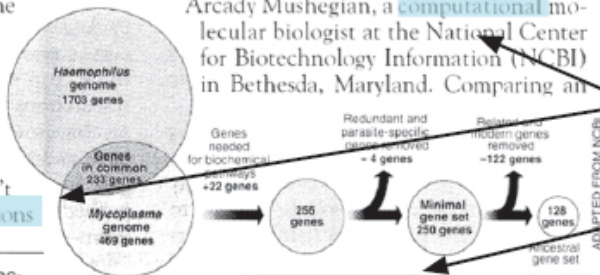
Documents

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,* 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

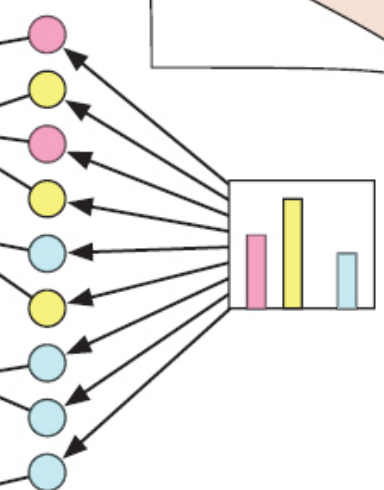


* Genome Mapping and Sequencing. Cold Spring Harbor, New York, May 8 to 12.

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

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Blei (2012): Probabilistic Topic Models

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or fee. Request permission to publish from permissions@acm.org or fax (212) 869-0481.

The Digital Library is published by the Association for Computing Machinery. Copyright © 2012 ACM, Inc.

Latent Dirichlet Allocation

(Blei et al., 2003)

categorical distributions over words for four topics:

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT

Generative Story for Simple LDA

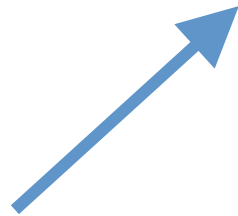
simplified LDA, and only showing generative story for 1 document:

- 1: Draw a multinomial topic distribution θ from some distribution $p(\Theta)$
- 2: For each position i in document:
 - a: Draw a topic $z_i \sim \text{Multinomial}(\theta)$
 - b: Draw a word $w_i \sim \text{Multinomial}(\beta_{z_i})$

Generative Story for Simple LDA

simplified LDA, and only showing generative story for 1 document:

- 1: Draw a multinomial topic distribution θ from some distribution $p(\Theta)$
- 2: For each position i in document:
 - a: Draw a topic $z_i \sim \text{Multinomial}(\theta)$
 - b: Draw a word $w_i \sim \text{Multinomial}(\beta_{z_i})$

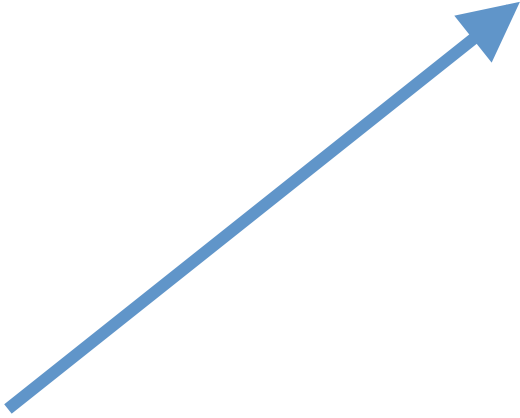


multinomial distribution over words for topic z_i

Generative Story for Simple LDA

simplified LDA, and only showing generative story for 1 document:

- 1: Draw a multinomial topic distribution θ from some distribution $p(\Theta)$
- 2: For each position i in document:
 - a: Draw a topic $z_i \sim \text{Multinomial}(\theta)$
 - b: Draw a word $w_i \sim \text{Multinomial}(\beta_{z_i})$

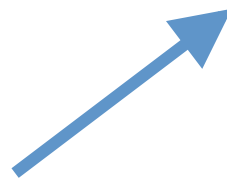


what should we keep in mind
when choosing this distribution?

Dirichlet Distribution

- distribution over vectors with entries that are all positive and sum to 1
- so it's kind of like a “distribution over (categorical) distributions”

$$p(\Theta = \theta \mid \alpha) = \frac{1}{B(\alpha)} \prod_i \theta_i^{\alpha_i - 1}$$



normalization term that depends on α

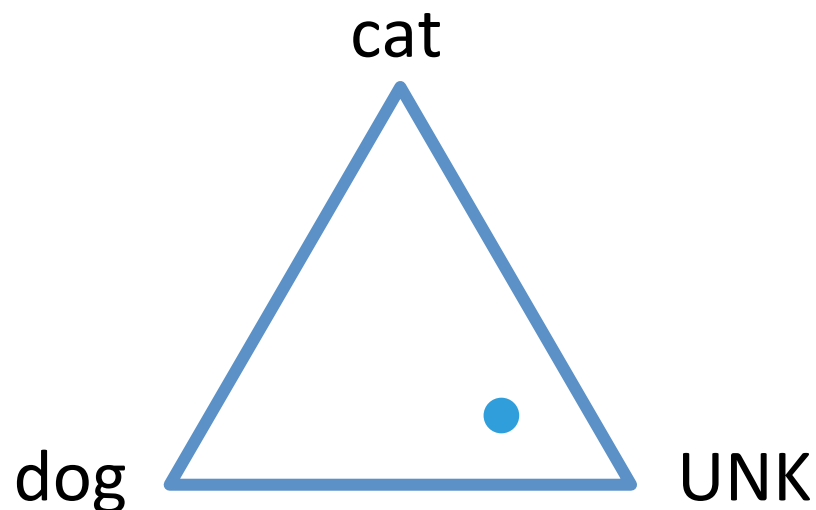
Dirichlet Distribution

- parameterized by a positive vector α

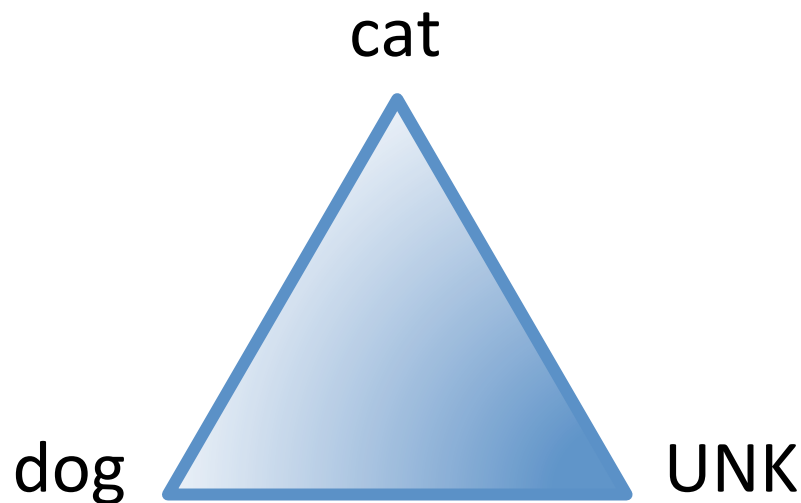
$$p(\Theta = \theta \mid \alpha) = \frac{1}{B(\alpha)} \prod_i \theta_i^{\alpha_i - 1}$$

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

- categorical = point on the simplex



- Dirichlet = distribution over the simplex



[see Jupyter Notebook]

Generative Story for Simple LDA

simplified LDA, and only showing generative story for 1 document:

- 1: Draw a multinomial topic distribution θ from some distribution $p(\Theta)$
- 2: For each position i in document:
 - a: Draw a topic $z_i \sim \text{Multinomial}(\theta)$
 - b: Draw a word $w_i \sim \text{Multinomial}(\beta_{z_i})$

Generative Story for Simple LDA

simplified LDA, and only showing generative story for 1 document:

- 1: Draw a multinomial topic distribution $\theta \sim \text{Dirichlet}(\alpha)$
- 2: For each position i in document:
 - a: Draw a topic $z_i \sim \text{Multinomial}(\theta)$
 - b: Draw a word $w_i \sim \text{Multinomial}(\beta_{z_i})$

Generative Story for LDA

- 1: For each topic, draw a multinomial word distribution $\beta_i \sim \text{Dirichlet}(\eta)$
 - 2: For each document d :
 - a: Draw a multinomial topic distribution $\theta \sim \text{Dirichlet}(\alpha)$
 - b: For each position i in document d :
 - i: Draw a topic $z_i \sim \text{Multinomial}(\theta)$
 - ii: Draw a word $w_i \sim \text{Multinomial}(\beta_{z_i})$
- now we show explicitly the generation of the word multinomials (once for the document collection)
 - where should the hyperparameters (alpha and psi) come from?

Graphical Model for LDA

