A large, dark-colored elephant stands behind a small, bright orange ball on a light-colored surface. The elephant's trunk is visible, and it appears to be looking towards the camera. The background is slightly blurred.

S&DS 265 / 565  
Introductory Machine Learning

# Bayes and Topic Models

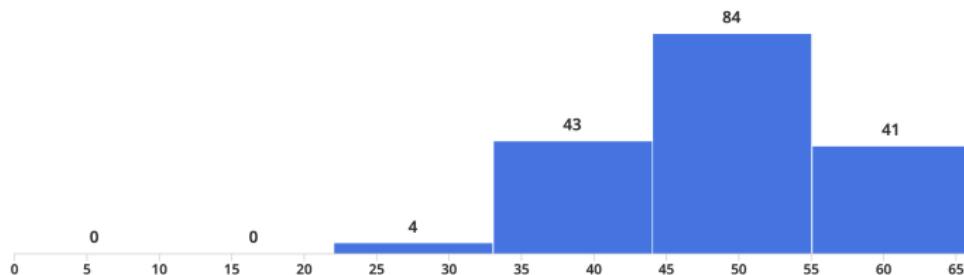
October 24

Yale

# Quick notes

- Assignment 3 due tonight
- Two more to go! Assn4 is posted (language models, embeddings)
- Midterm scores and mid-semester totals posted
- Announcement this morning about mid-semester letter grades
- Quiz 4 next week

# Midterm



Minimum

**24.5**

Median

**50.0**

Maximum

**63.0**

Mean

**48.89**

Std Dev ?

**7.63**

# Reminder

- If you need to email me about any issues with your work for the course, please email `sds265@yale.edu`
- Otherwise, your email may be missed

# For Today

- Mixtures
- Introduction to Bayesian analysis
- Overview of topic models

# Where are we going?

- So far: Mostly “surface representations” & explicit features
  - ▶ Listing information for CA housing
  - ▶ Supernovæ classification — measured light

# Where are we going?

- So far: Mostly “surface representations” & explicit features
  - ▶ Listing information for CA housing
  - ▶ Supernovæ classification — measured light
- Next: “Hidden” representations and latent variables

# George Washington

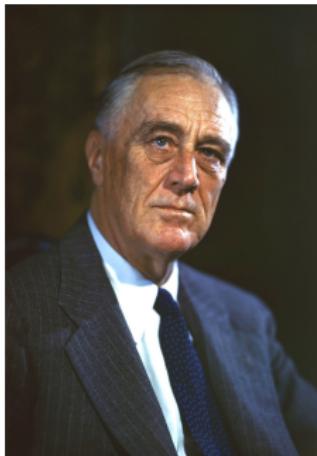
1789



*Among the many interesting objects which will engage your attention that of providing for the common defense will merit particular regard. To be prepared for war is one of the most effectual means of preserving peace.*

# Franklin Roosevelt

1941



*Such a peace would bring no security for us or for our neighbors. As a nation, we may take pride in the fact that we are softhearted; but we cannot afford to be soft-headed.*

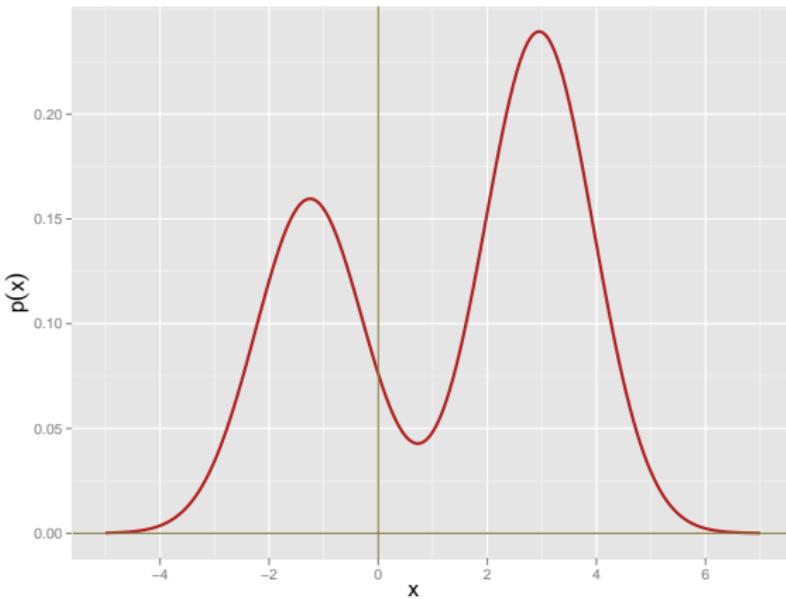
# The elephants in the room



# Mixtures

- Key technique: Mixture models
- Mixtures have latent variables
- Flexible tool
- Simple and difficult at the same time

# Gaussian Mixture



Mixture of two normals, one with mean  $-\frac{5}{2}$ , another with mean 3.

# Mixtures

- For each data point, we have a hidden (latent) variable: Which Gaussian generated the data point?
- $Z = 1$  it was generated from one model;  $Z = 0$  it was generated from the other model
- We don't *observe*  $Z$ , it could be either value

# Mixtures

- *Mixture of  $f$  and  $g$ :*

$$p(x) = \theta f(x) + (1 - \theta)g(x)$$

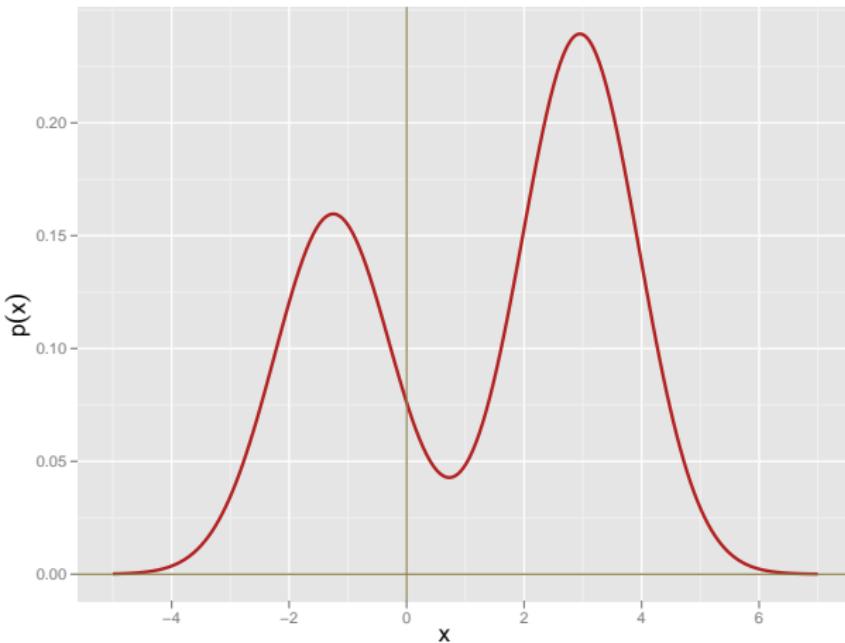
Simplest, most common kind of latent variable model

- *Hidden variable representation:* Define  $Z \sim \text{Bernoulli}(\theta)$  and

$$p(x) = \sum_{z=0,1} p(x | z) p(z)$$

with  $p(x | 1) = f(x)$ ,  $p(x | 0) = g(x)$ ,  $p(z) = \theta^z(1 - \theta)^{(1-z)}$ .

# Gaussian Mixture: All the Key Concepts



# Bayesian Inference

The parameter  $\theta$  of a model is viewed as a random variable.  
Inference usually carried out as follows:

- Choose a *generative model*  $p(x | \theta)$  for the data.
- Choose a *prior distribution*  $\pi(\theta)$  that expresses beliefs about the parameter before seeing any data.
- After observing data  $\mathcal{D}_n = \{x_1, \dots, x_n\}$ , update beliefs and calculate the *posterior distribution*  $p(\theta | \mathcal{D}_n)$ .

# Bayes' Theorem

A simple consequence of conditional probability:

$$\begin{aligned}\mathbb{P}(A | B) &= \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \\ &= \frac{\mathbb{P}(B | A) \mathbb{P}(A)}{\mathbb{P}(B)}\end{aligned}$$

# Bayes' Theorem

The posterior distribution can be written as

$$p(\theta | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \theta) \pi(\theta)}{p(x_1, \dots, x_n)} = \frac{\mathcal{L}_n(\theta) \pi(\theta)}{c_n} \propto \mathcal{L}_n(\theta) \pi(\theta)$$

where  $\mathcal{L}_n(\theta)$  is the *likelihood function* and

$$c_n = p(x_1, \dots, x_n) = \int \mathcal{L}_n(\theta) \pi(\theta) d\theta$$

is the normalizing constant, which is also called *evidence*.

# Important Example

Take model  $X \sim \text{Bernoulli}(\theta)$ .

This is a “coin flip”:  $X = 1$  means “heads” and  $X = 0$  means “tails.”

Natural prior is  $\text{Beta}(\alpha, \beta)$  distribution

$$\pi_{\alpha, \beta}(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

# Important Example

Take model  $X \sim \text{Bernoulli}(\theta)$ .

Natural prior is  $\text{Beta}(\alpha, \beta)$  distribution

$$\pi_{\alpha, \beta}(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

The scaling constant is scary looking:

$$\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}$$

where  $\Gamma(\cdot)$  is the “Gamma function”

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$$

# Important Example

$X \sim \text{Bernoulli}(\theta)$  with data  $\mathcal{D}_n = \{x_1, \dots, x_n\}$ . Prior Beta( $\alpha, \beta$ ) distribution

$$\pi_{\alpha, \beta}(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

Let  $s = \sum_{i=1}^n x_i$  be the number of “heads”

Posterior distribution  $\theta | \mathcal{D}_n$  is another beta distribution!

Specifically, with

$$\tilde{\alpha} = \alpha + \text{number of heads} = \alpha + s$$

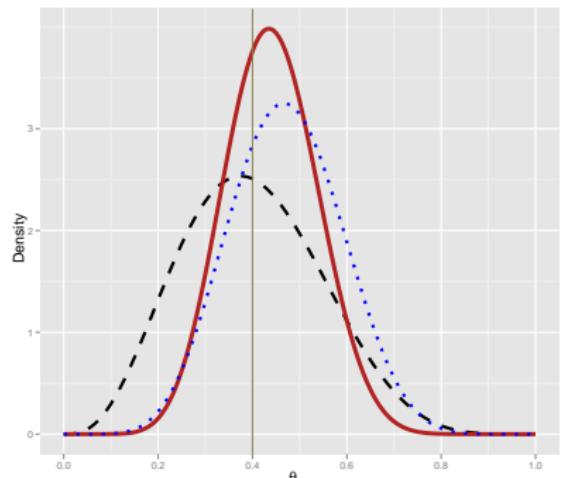
$$\tilde{\beta} = \beta + \text{number of tails} = \beta + n - s$$

---

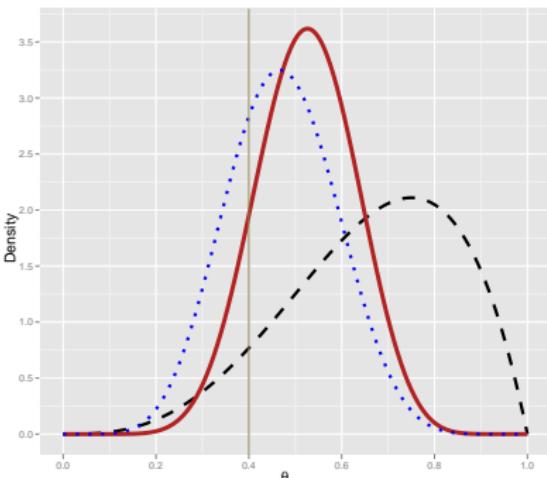
Showing this just uses the simple fact that  $\theta^{\alpha-1} \theta^x = \theta^{x+\alpha-1}$

# Example

$n = 15$  points sampled as  $X \sim \text{Bernoulli}(\theta = 0.4)$ , with  $s = 7$  heads.



prior A



prior B

Prior distribution (black-dashed), likelihood function (blue-dotted), posterior distribution (red-solid).

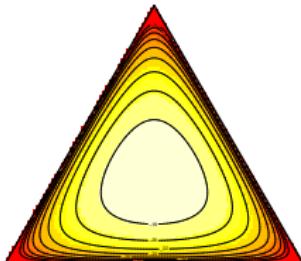
# Dirichlet

Multinomial model with Dirichlet prior is generalization of the Bernoulli/Beta model.

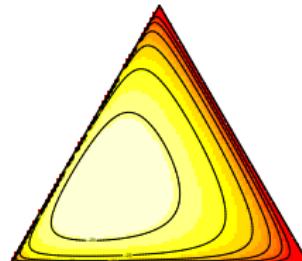
$$\text{Dirichlet}_{\alpha}(\theta) \propto \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \cdots \theta_K^{\alpha_K-1}$$

where  $\alpha = (\alpha_1, \dots, \alpha_K) \in \mathbb{R}_+^K$  is a non-negative vector.

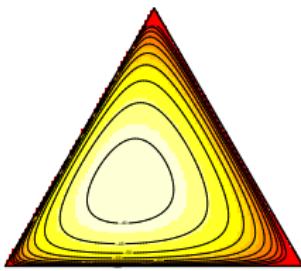
# Example



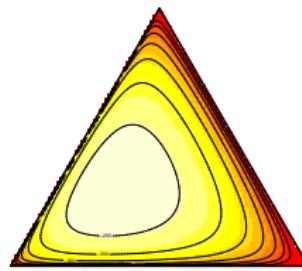
prior with Dirichlet(6,6,6)



likelihood function with  $n = 20$



posterior distribution with  $n = 20$



posterior distribution with  $n = 200$

# Example: Election Forecasting

<https://projects.economist.com/us-2020-forecast/president/how-this-works>

The Economist

Today Weekly edition ≡ Menu



## Forecasting the US elections

*The Economist* is analysing polling, economic and demographic data to predict America's elections in 2020

→ [Read more of our election coverage](#)

---

**President**   [Senate](#)   [House](#)

[National forecast](#)  
[How this works](#)

**COMPETITIVE STATES**

- [Arizona](#)
- [Florida](#)
- [Georgia](#)
- [Iowa](#)
- [Michigan](#)
- [Nevada](#)
- [New Hampshire](#)
- [North Carolina](#)
- [Ohio](#)
- [Pennsylvania](#)
- [Texas](#)
- [Wisconsin](#)

**ALL STATES**

- [Alabama](#)

---

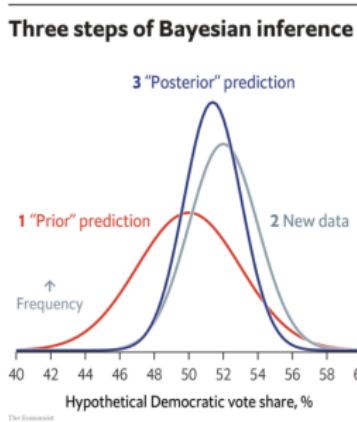
### How The Economist presidential forecast works

**T**HIS YEAR, *The Economist* is publishing its first-ever statistical forecast of an American presidential election. Developed with the assistance of Andrew Gelman and Merlin Heidemanns, political scientists at Columbia University, our model calculates Joe Biden's and Donald Trump's probabilities of winning each individual state and the election overall. Its projections will be updated every day at <https://projects.economist.com/us-2020-forecast/president>.

In another first, we are [publishing the source code](#) for what we believe to be the most innovative section of the model. All readers are welcome to download it, explore how it works, tweak its parameters and run it

# Example: Election Forecasting

<https://projects.economist.com/us-2020-forecast/president/how-this-works>



## Back to Bayes-ics

Readers acquainted with the workings of similar forecasting models may be surprised that the phrase "state polls" has not yet entered the equation. This exclusion is by design. Our model follows a logical structure first developed by Thomas Bayes, an 18th-century reverend whose ideas have shaped a large and growing family of statistical techniques. His approach works in two stages. First, before conducting a study, researchers

explicitly state what they believe to be true, and how confident they are in that belief. This is called a "prior". Next, after acquiring data, they update this prior to reflect the new information—gaining more confidence if it confirms the prior, and generally becoming more uncertain if it refutes the prior (though not if the new numbers are so definitive that leave little room for doubt). In this framework, the expected distribution of potential vote shares in each state derived above is the prior, and state polls that trickle in during the course of the campaign are the new data. The result—a "posterior", in Bayesian lingo—is our forecast.

# Let's go to the notebook!

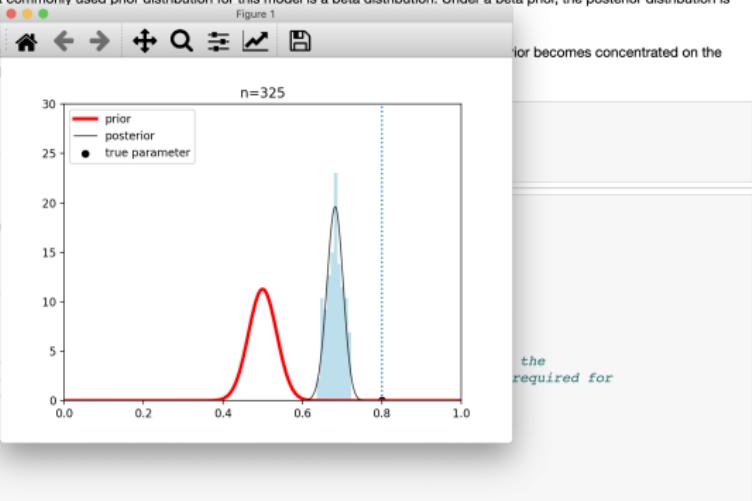
## Demo code for Bayesian analysis

In this notebook we illustrate some of the basic models and priors for Bayesian inference. These concepts will be important for our discussions about "topic models."

First, we illustrate the situation where the parameter  $\theta$  that we are modeling is a Bernoulli parameter. This can be thought of as the probability that flipping a certain coin comes up heads. The most commonly used prior distribution for this model is a beta distribution. Under a beta prior, the posterior distribution is again a beta distribution.

This is illustrated in the following simulation. We start with a true parameter. But as the variance of the

```
In [5]: import os, gzip  
import numpy as np  
import matplotlib.pyplot as plt  
  
In [*]: %matplotlib qt  
from scipy.special import gammainc  
from scipy import random  
from scipy.stats import beta  
  
theta = np.linspace(0,1,num=5  
fig = plt.figure(1)  
plt.ion()  
  
# The following are the parameters  
# variance of the prior decreases  
# the posterior to be centered  
  
scale = 100  
a0 = scale*1  
b0 = scale*1  
  
sample_size = 100
```



# Reading

- See notes on Bayesian inference on iML site  
(bayes-notes.pdf)
- Not necessary to understand everything—but you should be able to do some of the basic “coin flipping” calculations
- We’ll build on this when discussing topic models

# Summary

- Mixtures are latent variable models
- The mixing weight encodes a hidden variable
- Computing with mixtures uses basic probabilistic reasoning
- Bayesian inference is popular in ML
- In a Bayesian approach, the parameters are random, and the data are fixed.
- Three ingredients: Prior, likelihood, posterior

## Next up: Topic models

- High level intro to topic models
- Use of latent variables, mixtures
- More details and examples next week

Readings:

<http://www.cs.columbia.edu/~blei/topicmodeling.html>

A survey paper describing many of these ideas in more detail is here:

<https://www.cs.columbia.edu/~blei/papers/Blei2012.pdf>

# Topic modeling



Topic modeling provides methods for automatically organizing, understanding, searching, and summarizing large electronic archives.

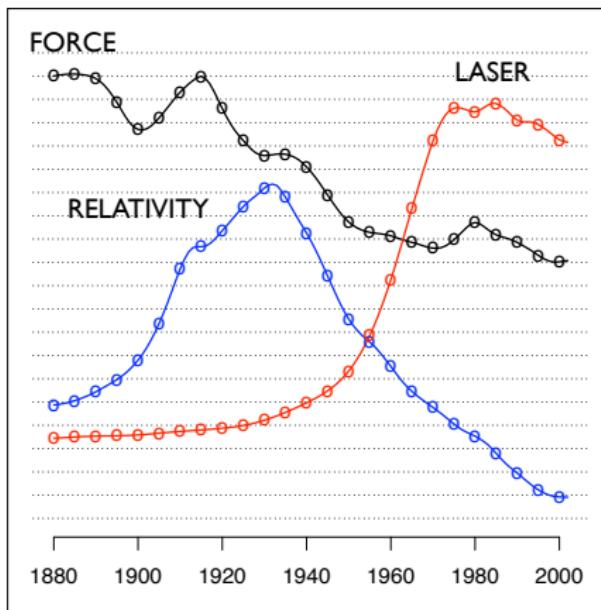
- ① Discover the hidden themes that pervade the collection.
- ② Annotate the documents according to those themes.
- ③ Use annotations to organize, summarize, and search the texts.

# Discover topics from a corpus

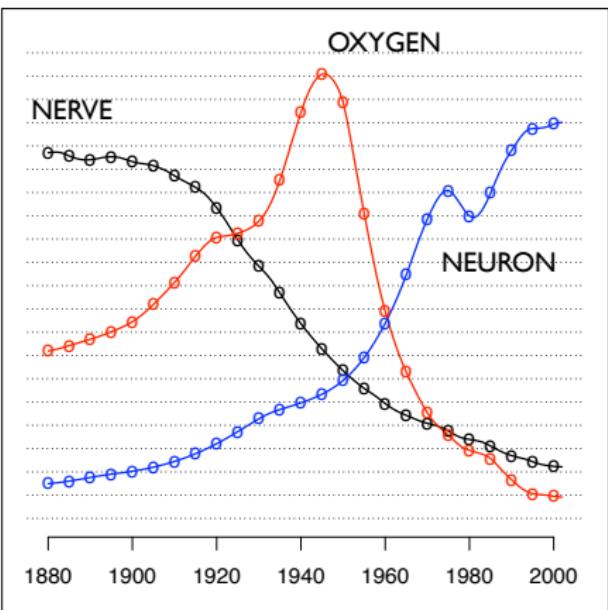
human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living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

# Model the evolution of topics over time

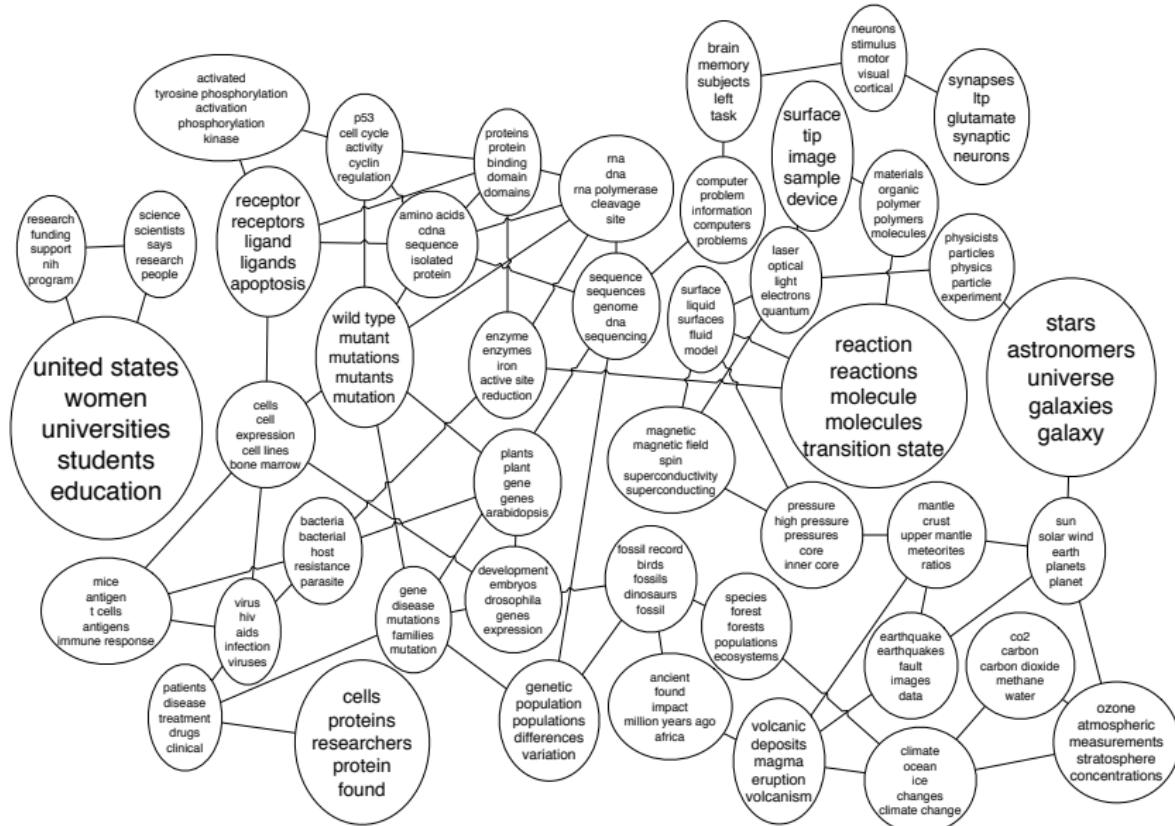
"Theoretical Physics"



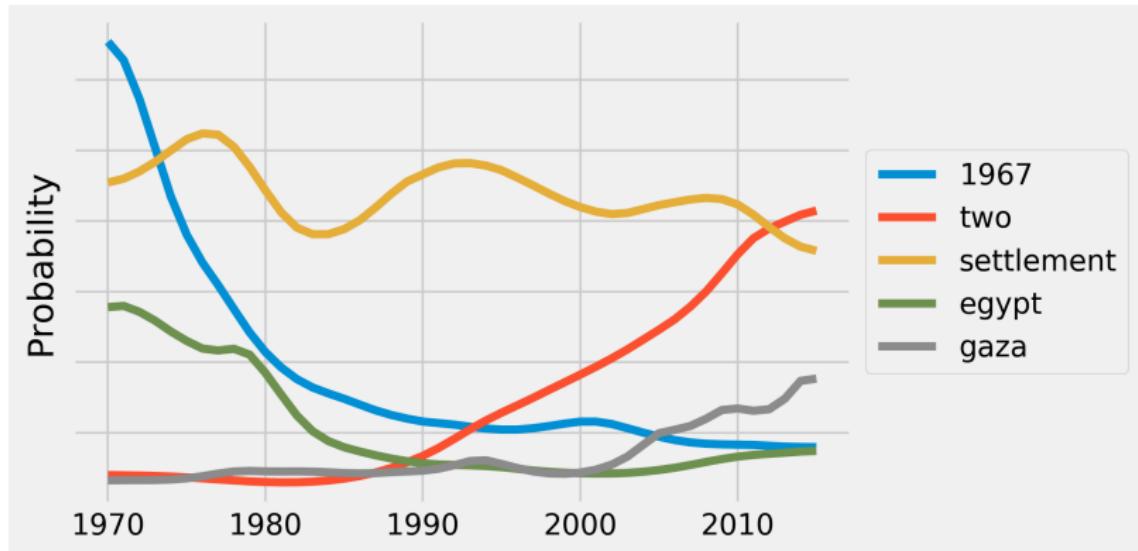
"Neuroscience"



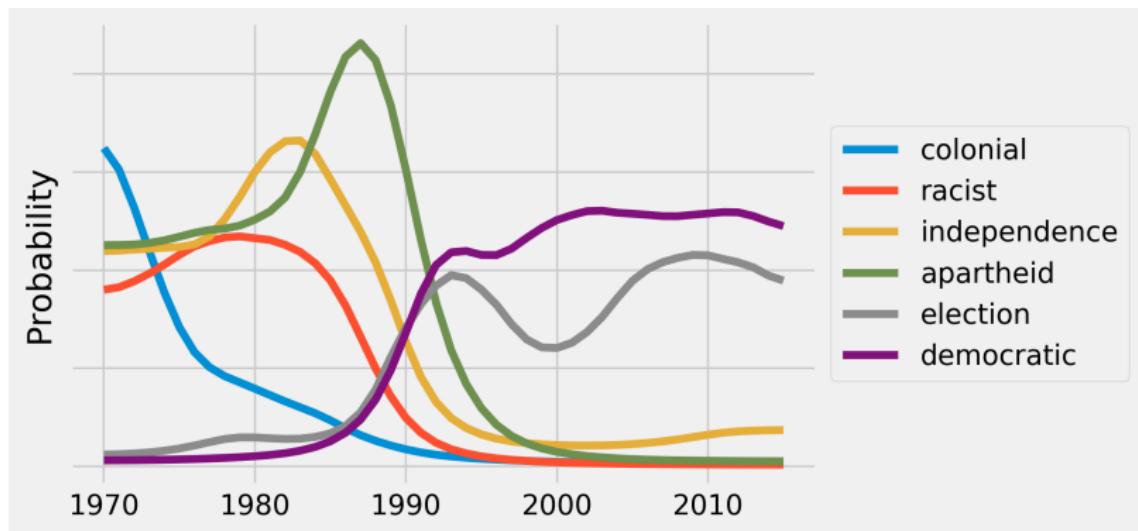
# Model connections between topics



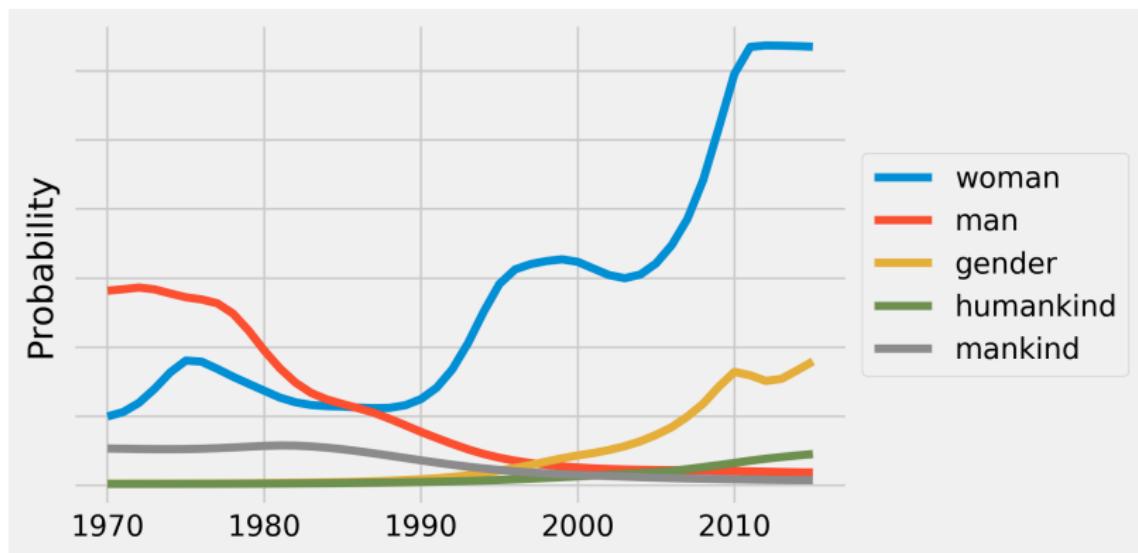
# Proceedings of the United Nations



# Proceedings of the United Nations



# Proceedings of the United Nations



# Topics in scientific texts

<b>Quantum physics</b>	spin energy field electron magnetic state states hamiltonian
<b>Particle physics</b>	higgs neutrino coupling decay scale masses mixing quark
<b>Astrophysics</b>	mass gas star stellar galaxies disk halo radius luminosity
<b>Relativity</b>	black metric hole schwarzschild gravity holes einstein
<b>Number theory</b>	prime integer numbers conjecture integers degree modulo
<b>Graph theory</b>	graph vertex vertices edges node edge number set tree
<b>Linear algebra</b>	matrix matrices vector basis vectors diagonal rank linear
<b>Optimization</b>	problem optimization algorithm function solution gradient
<b>Probability</b>	random probability distribution process measure time
<b>Machine learning</b>	layer word image feature sentence model cnn lstm training

In all of these examples, the topics are automatically inferred (unsupervised) through Bayesian inference

# Topic modeling for equations

Topic	Generated Equations
Quantum physics	<ul style="list-style-type: none"><li><math>E = \hbar \frac{\partial^2 S}{\partial t^2} \left( \frac{\partial \varphi}{\partial c} \right) - \frac{k}{\hbar^2} \frac{\partial B}{\partial t} (t + \partial_t \delta).</math></li><li><math>\Psi_{\text{pr}} = \sum_{\mathbf{l}} (\psi_{\mathbf{r}+\uparrow} - \psi_{\mathbf{r}\downarrow}^\dagger) + \sum_{\mathbf{r}'} (\psi_{\mathbf{r}',\uparrow}^\dagger - \psi_{\mathbf{r}'\downarrow} \sigma^\dagger).</math></li></ul>
Particle physics	<ul style="list-style-type: none"><li><math>\mathcal{H} = \frac{1}{4}(\partial_\mu \phi)^2 + 2m\phi_\nu(\phi) + \frac{1}{2}m^2(\phi)(1-\phi^2)^2.</math></li><li><math>m_{\text{eff}}(M) = 1.4 \cdot 10^{-13} \text{ GeV}.</math></li></ul>
Relativity	<ul style="list-style-type: none"><li><math>\mathcal{M} = \frac{1}{2}g^{\mu\nu}(f_{\mu\nu,\mu} - g_{\mu\nu,\nu} + g_{\nu\nu,b}f_{\mu,\nu}) + \frac{1}{2}g^{\mu\nu}.</math></li><li><math>T_{\mu\nu} = \int_0^\infty ds_{\mu\nu} ds^2 + a_\mu^2 dr^2 + r^2 d\Omega^2.</math></li></ul>

# Uber Topics

From a former student:

*Hi Prof. Lafferty,*

*I took your course last spring. I just wanted to let you know that I'm currently topic modeling at my current job at Uber!*

*We're using topic models at Uber to discover topics in rider feedback – when riders write comments about their driver after the trip. We're trying to find topics such as 'unprofessional driver', 'driver no-show', 'sexual harassment', etc.*

# Probabilistic modeling

- ① Data are assumed to be observed from a generative probabilistic process that includes hidden variables.
  - *In text, the hidden variables are the thematic structure.*
- ② Infer the hidden structure using posterior inference
  - *What are the topics that describe this collection?*
- ③ Situate new data into the estimated model.
  - *How does a new document fit into the topic structure?*

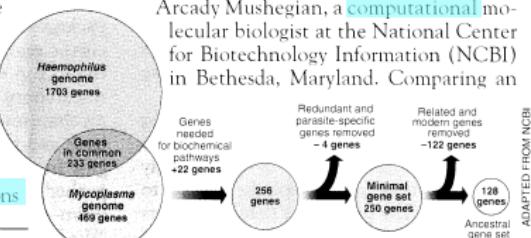
# Latent Dirichlet allocation (LDA)

## 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



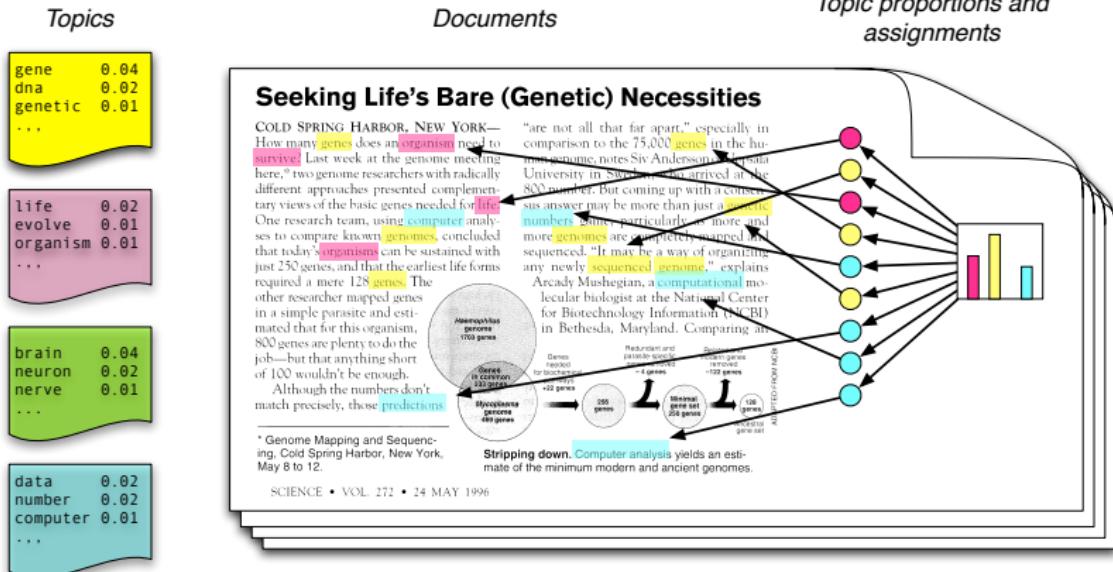
ADAPTED FROM NCBI

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

SCIENCE • VOL. 272 • 24 MAY 1996

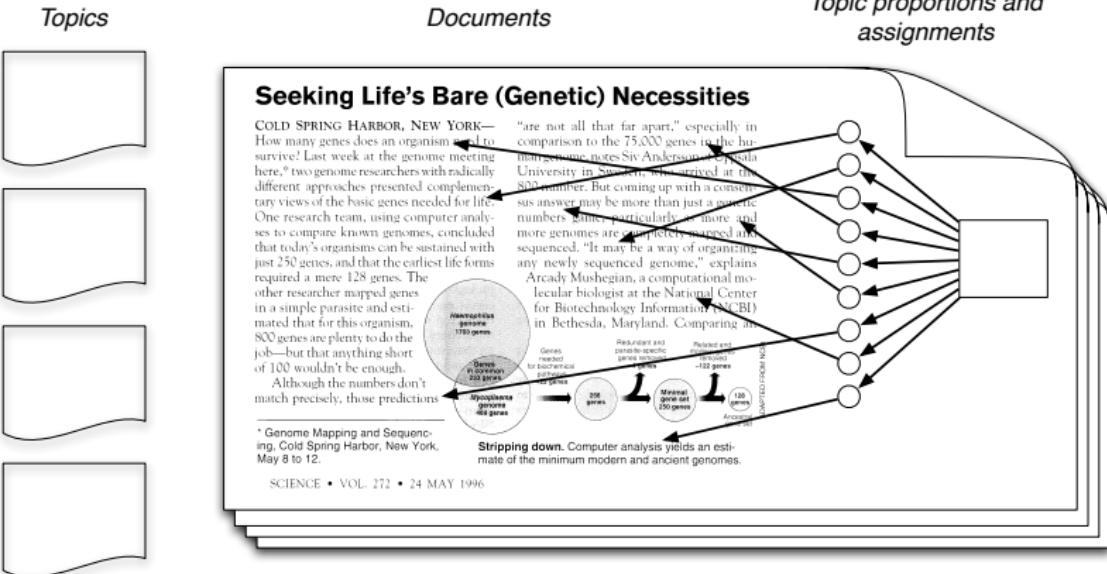
**Simple intuition:** Documents exhibit multiple topics.

# Generative model for LDA



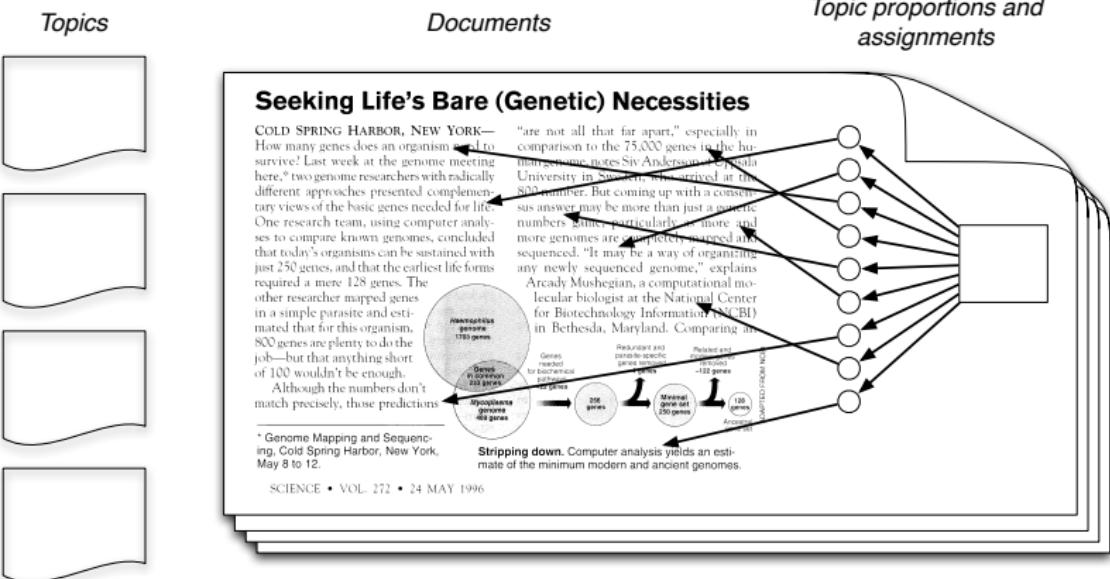
- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

# The posterior distribution



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

# The posterior distribution



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

# Summary: Topic models

- Topic models automatically extract “semantic themes” from large document collections
- Based on latent variables, mixtures, and Bayesian inference
- Can be useful for a wide variety of data