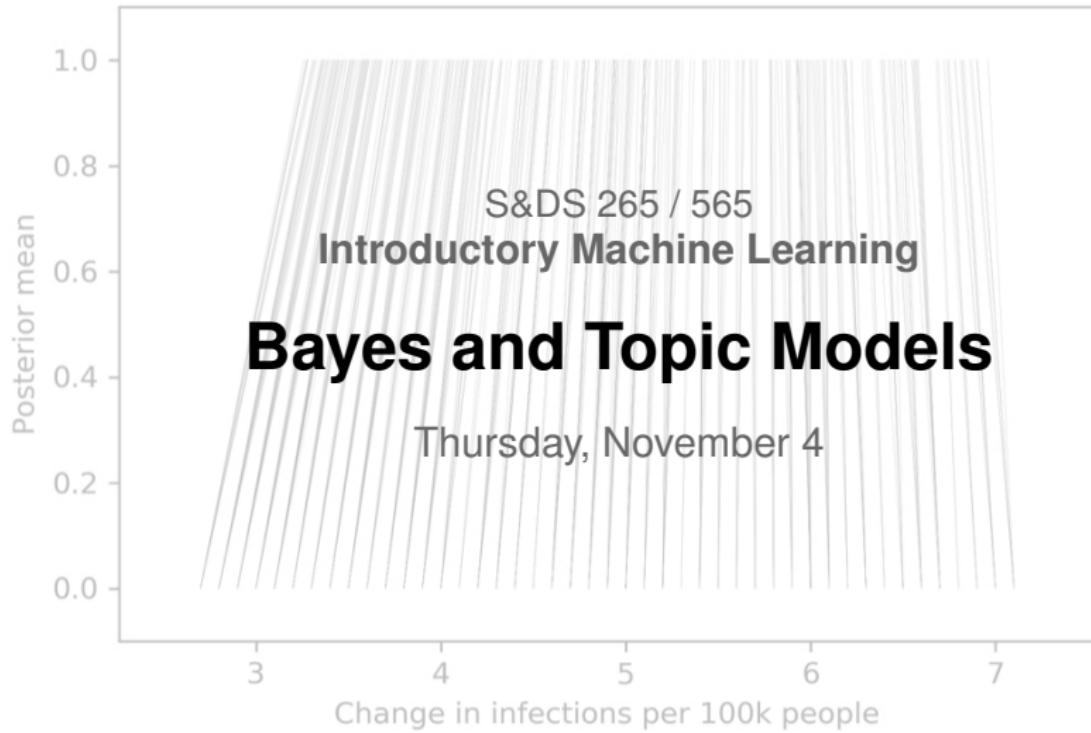


Shrinkage for county infection rates



**Yale**

# Quick notes

- Assignment 5 is out; due week from today (Nov 11)
- Two more to go!
- Midterm solutions posted

# For Today

- Bayes (continued)
- Overview of topic models

# 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

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

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 a little scary:

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

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

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

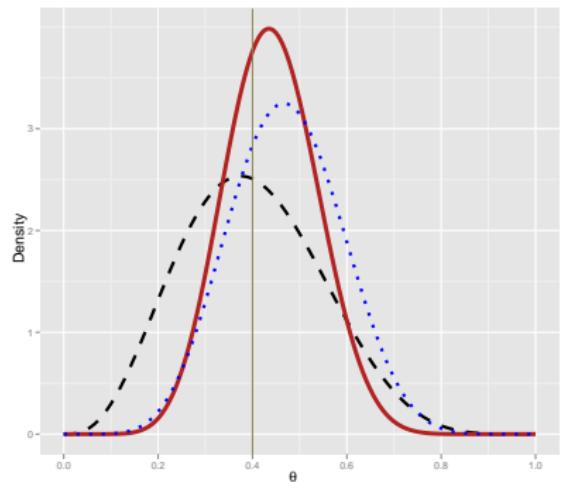
$$\tilde{\beta} = \beta + \text{number of tails} = \beta + n - s$$

---

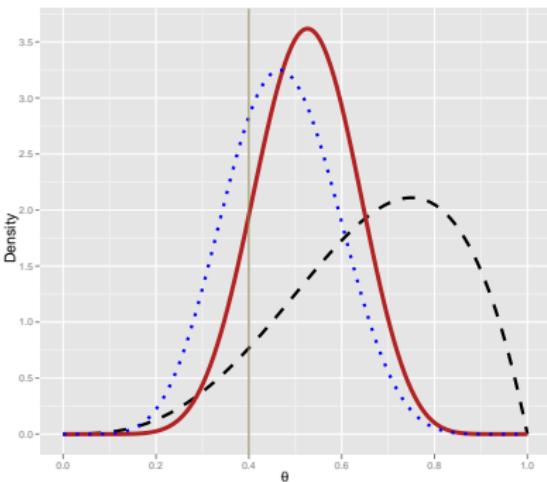
Showing this just uses the simple fact that  $\theta^{\alpha-1} \theta^x = \theta^{x+\alpha-1}$

# 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).

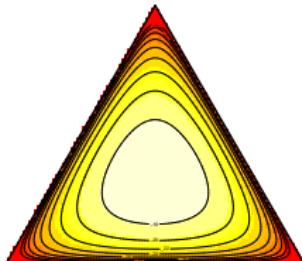
# Dirichlet

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

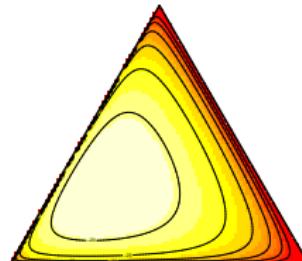
$$\text{Dirichlet}_{\alpha}(\theta) \propto \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \cdots \theta_K^{\alpha_K-1}$$

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

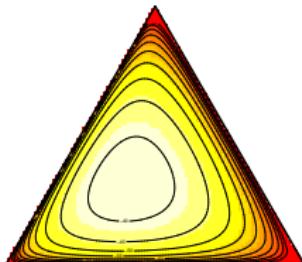
# Example



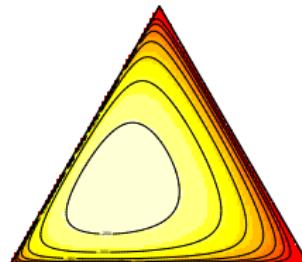
prior with Dirichlet(6,6,6)



likelihood function with  $n = 20$



posterior distribution with  $n = 20$



posterior distribution with  $n = 200$

# Notebook

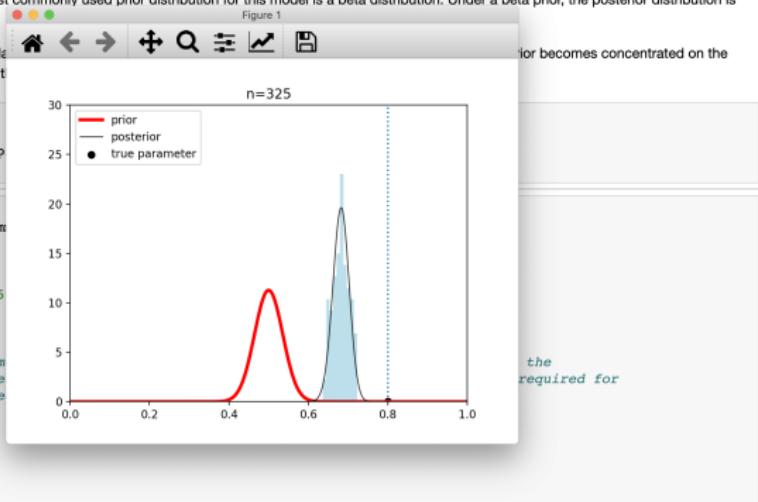
## Demo code for Bayesian analysis

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."

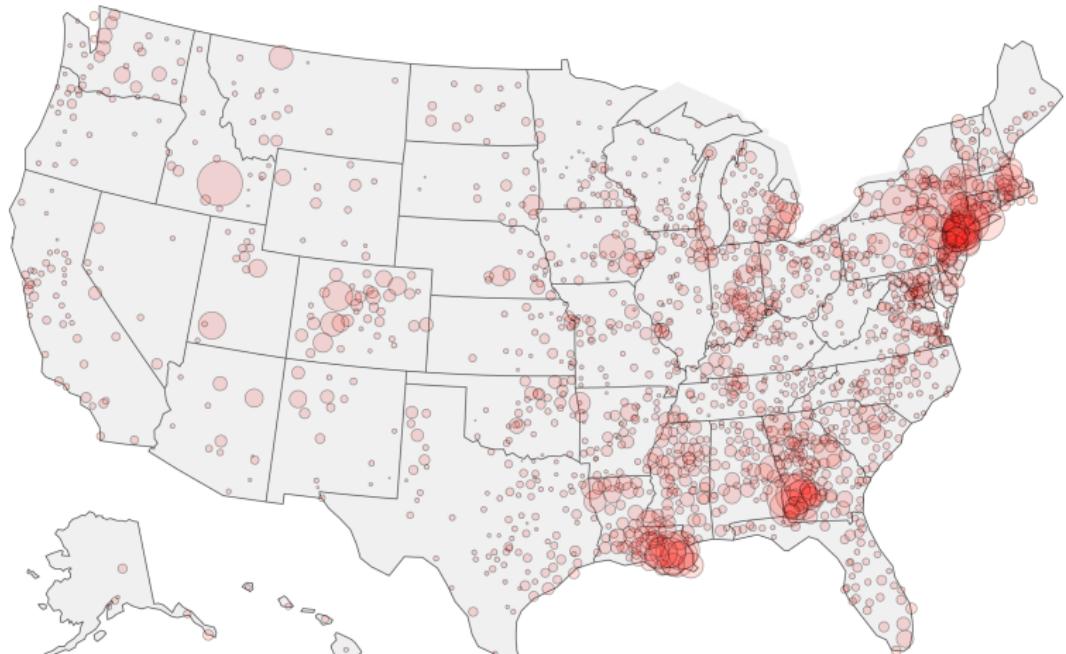
First, we illustrate the situation where the parameter  $\theta$  that we are modeling is a Bernoulli parameter. This can be thought of as the probability that flipping a certain coin comes up heads. The most commonly used prior distribution for this model is a beta distribution. Under a beta prior,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

This is illustrated in the following simulation. We start with a uniform prior for the true parameter. But as the variance of the observed data increases, the posterior distribution becomes concentrated on the true parameter.

```
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 around  
# the true parameter  
  
scale = 100  
a0 = scale*1  
b0 = scale*1  
  
sample_size = 100
```



# Covid Cases per 100,000 people (in April 2020)



per 100,000 people

Data from The New York Times  
[github.com/nytimes/covid-19-data](https://github.com/nytimes/covid-19-data)  
Tuesday April 14, 2020

# Hierarchical models

We are interested in modeling the infection rates across counties

In a Bayesian hierarchical model, we tie together our inferences across regions

In each county  $c$  we test  $n_c$  people (at random) and observe how many people  $Y_c$  have Covid-19.

$Y_c \sim \text{Binomial}(n_c, \theta_c)$ . This is equivalent to  $n_c$  coin flips (Bernoulli) each with probability of heads  $\theta_c$ .

# A simple hierarchical model

For each county  $c$ ,

$$\theta_c \sim F$$

$$Y_c | \theta_c \sim \text{Binomial}(n_c, \theta_c)$$

This ties together the parameters and makes better use of the data

# A simple hierarchical model

A simple transformation makes the use of Gaussians more appropriate.

Transform by  $\psi_c \equiv \log\left(\frac{\theta_c}{1-\theta_c}\right)$  to use Gaussian approximation:

$$\psi_c \sim N(\mu, \tau^2)$$

$$Z_c | \psi_c \sim N(\psi_c, \sigma_c^2)$$

- Posterior distribution may not be written down explicitly
- But we can sample from it
- This approximates the posterior as a mixture of Gaussians

# A run from April 2020

File Edit View Insert Cell Kernel Widgets Help Trusted

In [1]: `import covid19 as cvd  
import covid19_predict as cvd_predict`

covid19: Most recent NY Times data: Tuesday April 14, 2020  
covid19\_predict: initializing for simulation  
covid19\_predict: running Gibbs sampler: n=3193, B=10000  
covid19\_predict: making predictions

In [2]: `cvd_predict.df[cvd_predict.df['county']=='New Haven']`

Out[2]:

	date	county	state	cases	deaths	population	cases_per_100k	delta	delta_bar	delta_95	delta_05	cases_predicted	cases_95_credible
83	2020-04-14	New Haven	Connecticut	3543	151	854757	414.5	21.6	20.36	28.79	14.42	3716	3789

In [3]: `_ = cvd_predict.plot_predictions_for_addr('New Haven, CT', show=True)`

New Haven County, Connecticut

cases  
predicted

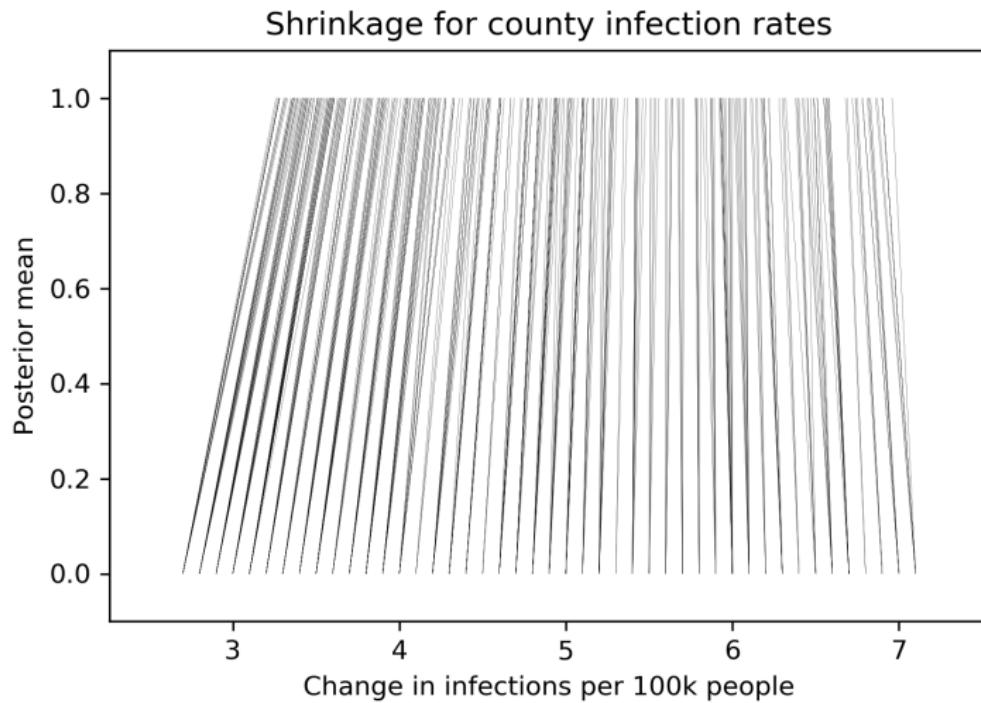
March 25 April 4 April 14

# Shrinkage

Tying the parameters together results in the small probabilities being increased and the larger probabilities being decreased.

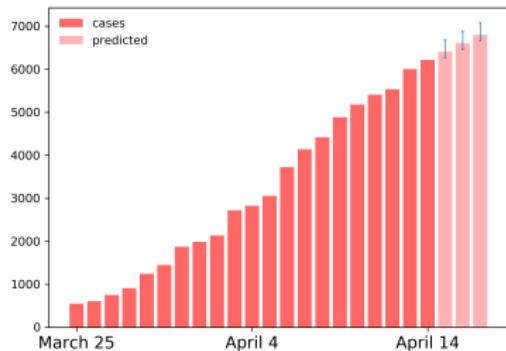
This is called “shrinkage.” The posterior estimates are closer together than the raw frequencies.

# Shrinkage

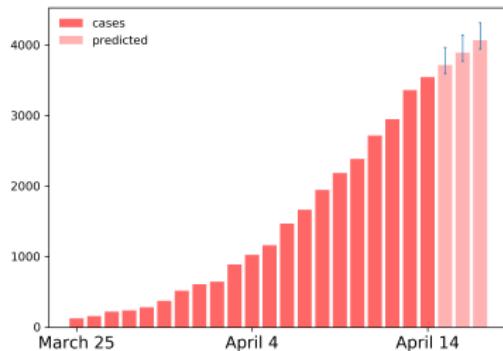


# Examples

Fairfield County, Connecticut



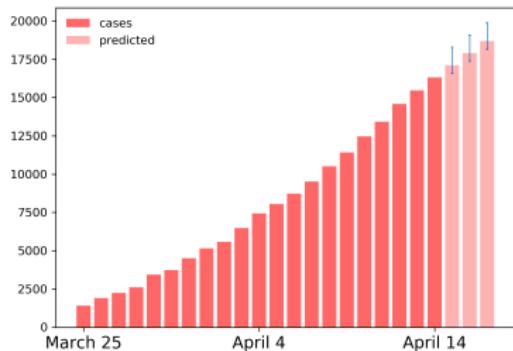
New Haven County, Connecticut



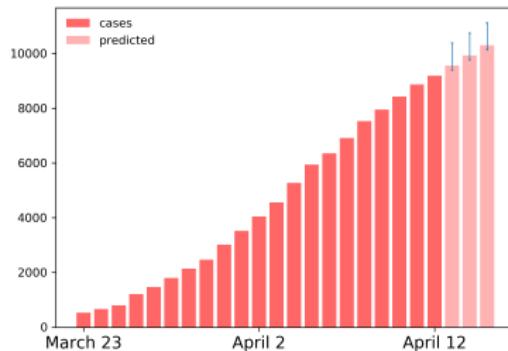
The "credible intervals" trap 95% of the sampled parameters under the simulation.

# Examples

Cook County, Illinois



Los Angeles County, California



The “credible intervals” trap 95% of the sampled parameters under the simulation.

# Another 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

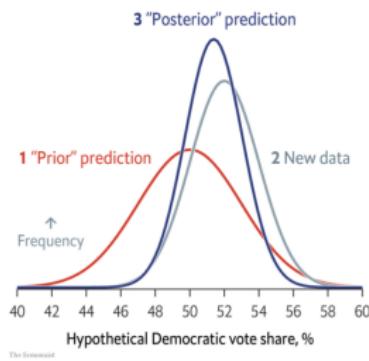
**T**HIS YEAR, *The Economist* is publishing its first-ever statistical forecast of an American presidential election. Developed with the assistance of Andrew Gelman and Merlin Heidemanns, political scientists at Columbia University, our model calculates Joe Biden's and Donald Trump's probabilities of 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

# Another Example: Election Forecasting

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

## Three steps of Bayesian inference



The Economist

## 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.

# Reading

- See notes on Bayesian inference on iML site  
(bayes-notes.pdf)
- These are more advanced, from S&DS 365.
- 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: Bayes

- 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.
- Bayesian hierarchical models tie together data using latent variables

## Next up: Topic models

- High level intro to topic models
- Use of latent variables, mixtures
- Work through a Jupyter notebook
- More details and examples on Thursday

# Intro to Topic Modeling

Readings:

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

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

[https://cacm.acm.org/magazines/2012/4/  
147361-probabilistic-topic-models/fulltext](https://cacm.acm.org/magazines/2012/4/147361-probabilistic-topic-models/fulltext)

# Topic modeling



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

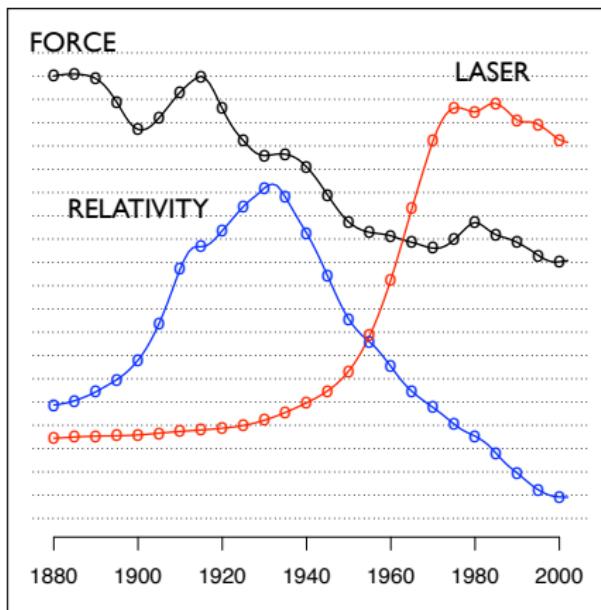
- ① Discover the hidden themes that pervade the collection.
- ② Annotate the documents according to those themes.
- ③ Use annotations to organize, summarize, and search the texts.

# Discover topics from a corpus

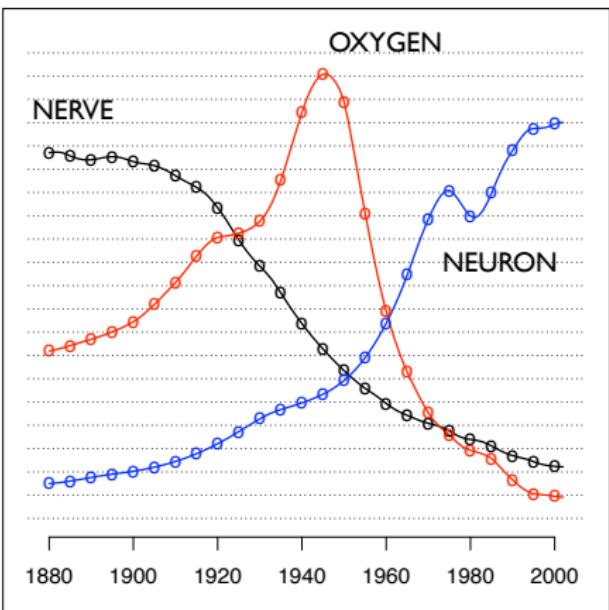
human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	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

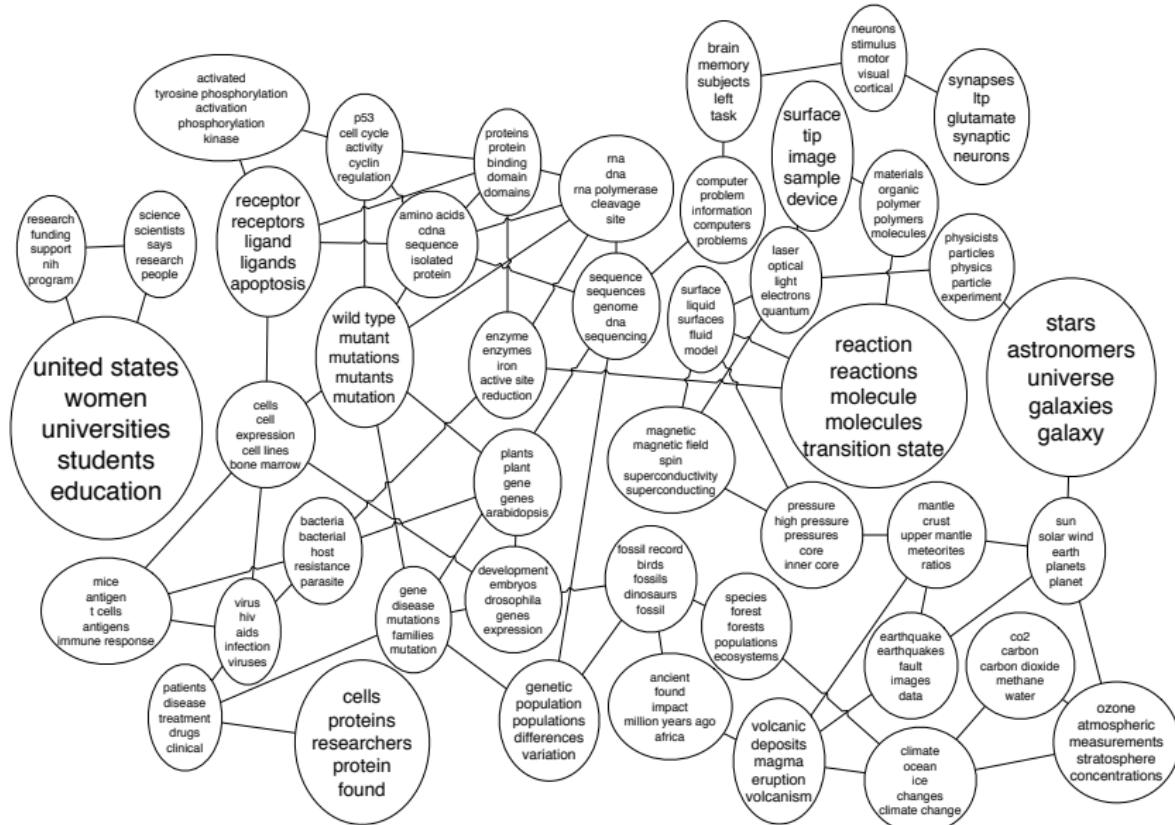
"Theoretical Physics"



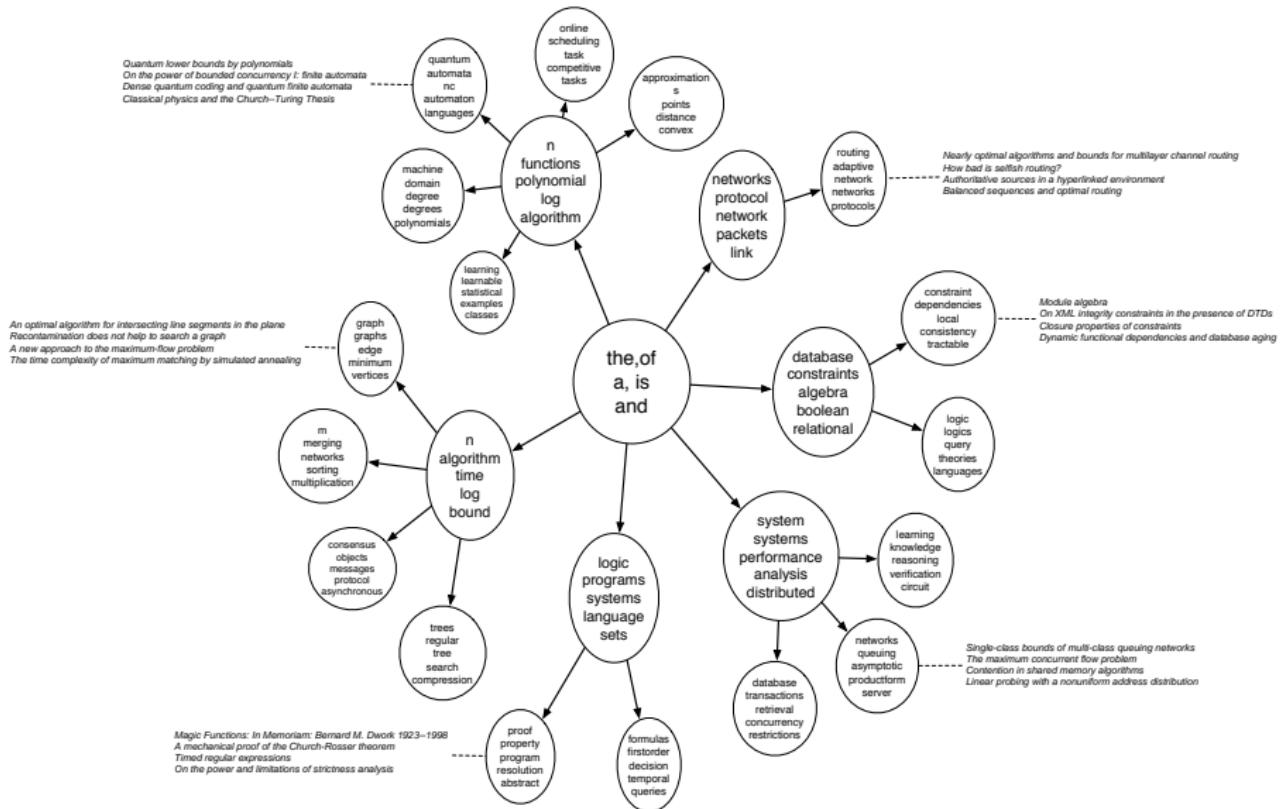
"Neuroscience"



# Model connections between topics



# Find hierarchies of topics



# Annotate images



SKY WATER TREE  
MOUNTAIN PEOPLE



SCOTLAND WATER  
FLOWER HILLS TREE



SKY WATER BUILDING  
PEOPLE WATER



FISH WATER OCEAN  
TREE CORAL

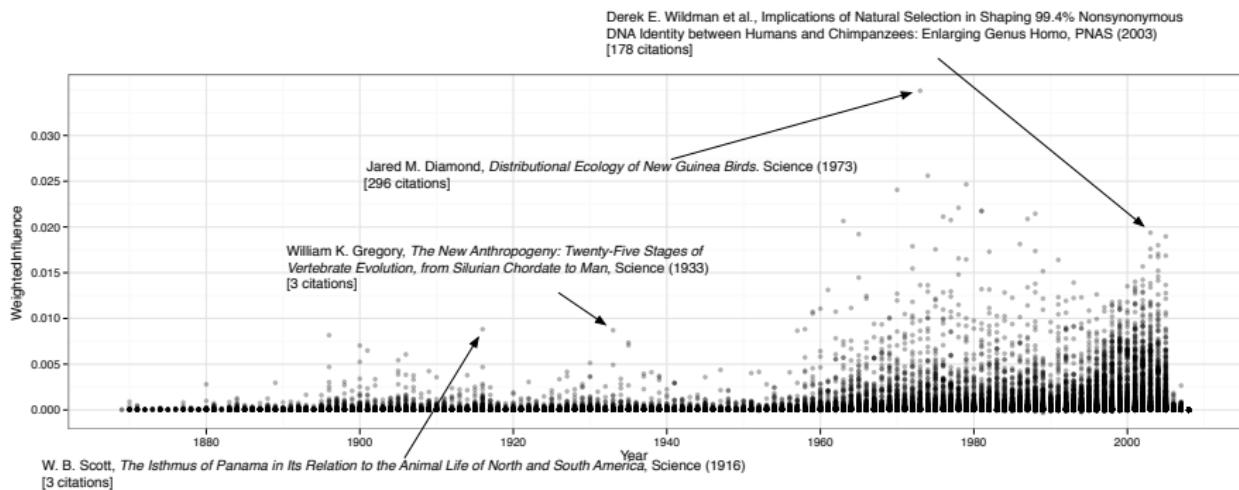


PEOPLE MARKET PATTERN  
TEXTILE DISPLAY



BIRDS NEST TREE  
BRANCH LEAVES

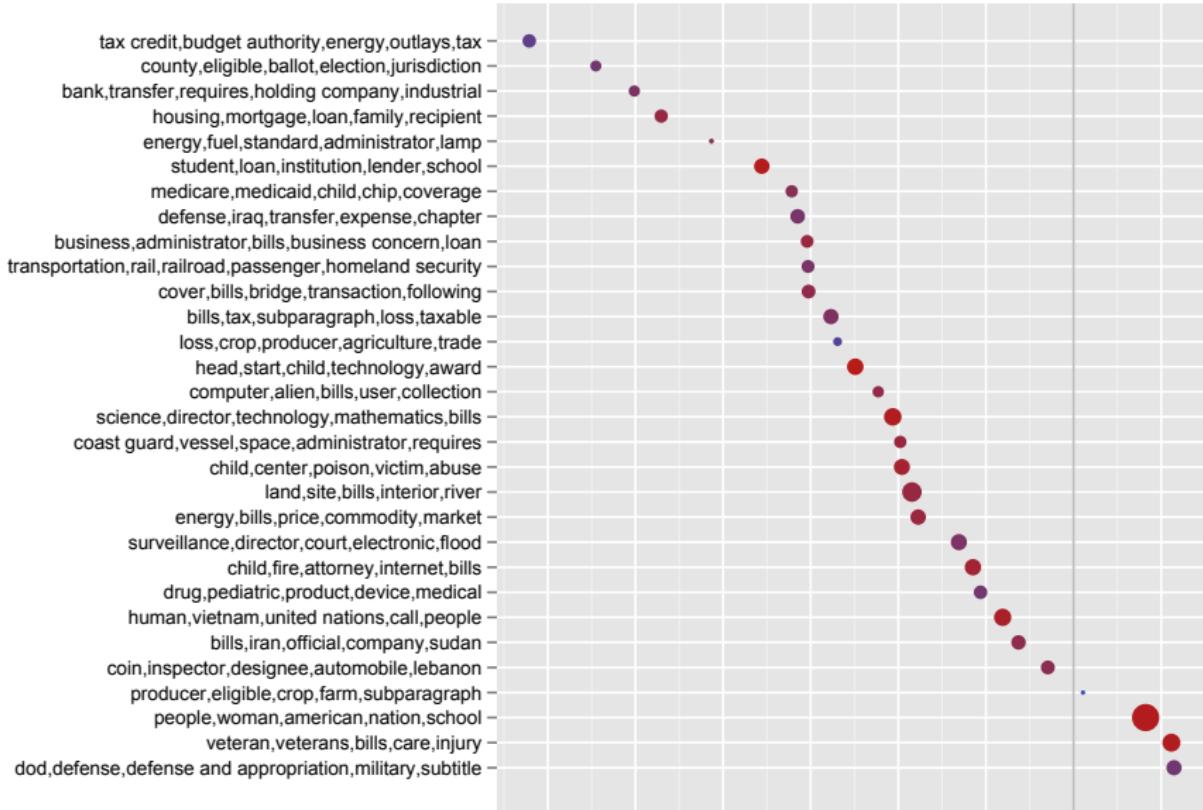
# Discover influential articles



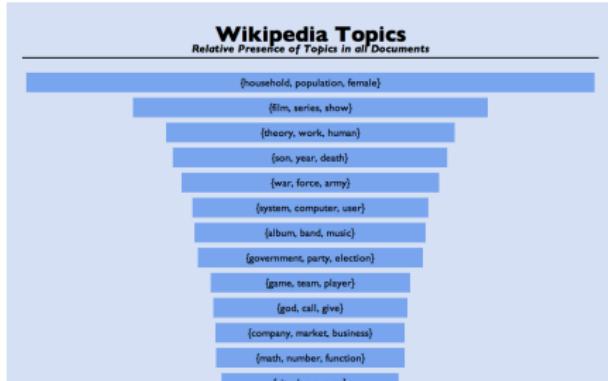
# Predict links between articles

<p><i>Markov chain Monte Carlo convergence diagnostics: A comparative review</i></p> <p><b>Minorization conditions and convergence rates for Markov chain Monte Carlo</b></p> <ul style="list-style-type: none"><li>Rates of convergence of the Hastings and Metropolis algorithms</li></ul> <p><b>Possible biases induced by MCMC convergence diagnostics</b></p> <ul style="list-style-type: none"><li>Bounding convergence time of the Gibbs sampler in Bayesian image restoration</li><li>Self regenerative Markov chain Monte Carlo</li></ul> <p>Auxiliary variable methods for Markov chain Monte Carlo with applications</p> <p><b>Rate of Convergence of the Gibbs Sampler by Gaussian Approximation</b></p> <ul style="list-style-type: none"><li>Diagnosing convergence of Markov chain Monte Carlo algorithms</li></ul>	<p>RTM (<math>\psi_e</math>)</p>
<p>Exact Bound for the Convergence of Metropolis Chains</p> <p>Self regenerative Markov chain Monte Carlo</p> <p><b>Minorization conditions and convergence rates for Markov chain Monte Carlo</b></p> <ul style="list-style-type: none"><li>Gibbs-markov models</li></ul> <p>Auxiliary variable methods for Markov chain Monte Carlo with applications</p> <p>Markov Chain Monte Carlo Model Determination for Hierarchical and Graphical Models</p> <ul style="list-style-type: none"><li>Mediating instrumental variables</li><li>A qualitative framework for probabilistic inference</li><li>Adaptation for Self Regenerative MCMC</li></ul>	<p>LDA + Regression</p>

# Characterize political decisions



# Organize and browse large corpora



## {film, series, show}

words	related documents	related topics
film	The X-Files	{son, year, death}
series	Orson Welles	{work, book, publish}
show	Stanley Kubrick	{album, band, music}
character	B movie	{woman, child, man}
play	Mystery Science Theater 3000	{law, state, case}
make	Monty Python	{black, white, people}
episode	Doctor Who	{theory, work, human}
movie	Sam Peckinpah	{@card@, make, design}
good	Married... with Children	{war, force, army}
release	History of film	{god, call, give}
feature	The A-Team	{game, team, player}
television	Pulp Fiction (film)	{day, year, event}
star	Mad (magazine)	{company, market, business}

## Stanley Kubrick

A pie chart illustrating the relative presence of various topics associated with Stanley Kubrick. The largest segment is 'film, series, show'.

Topic	Relative Presence (approx.)
film, series, show	0.45
theory, work, human	0.15
{son, year, death}	0.10
{black, white, people}	0.08
{god, call, give}	0.05
{math, energy, light}	0.05

### related topics

- {film, series, show}
- {theory, work, human}
- {son, year, death}
- {black, white, people}
- {god, call, give}
- {math, energy, light}

### related documents

- Orson Welles
- B movie
- Mystery Science Theater 3000
- Monty Python
- Doctor Who
- Sam Peckinpah
- The A-Team
- Pulp Fiction (film)
- Buffy the Vampire Slayer (TV series)
- The X-Files
- Sunset Boulevard (film)
- Jack Benny

## {theory, work, human}

words	related documents	related topics
theory	Meme	{work, book, publish}
work	Intelligent design	{law, state, case}
human	Immanuel Kant	{son, year, death}
ideas	Philosophy of mathematics	{woman, child, man}
term	History of science	{god, call, give}
study	Free will	{black, white, people}
view	Truth	{film, series, show}
science	Psychoanalysis	{war, force, army}
concept	Charles Peirce	{language, word, form}
form	Existentialism	{@card@, make, design}
world	Deconstruction	{church, century, christian}
argue	Social sciences	{rate, high, increase}
social	Idealism	{company, market, business}

# 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

# 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:

*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. LDA has worked really well with this...*

# **Introduction to Topic Modeling**

# 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?*

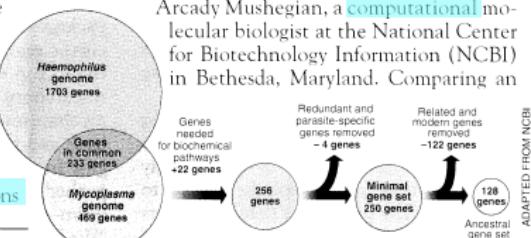
# Latent Dirichlet allocation (LDA)

## 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



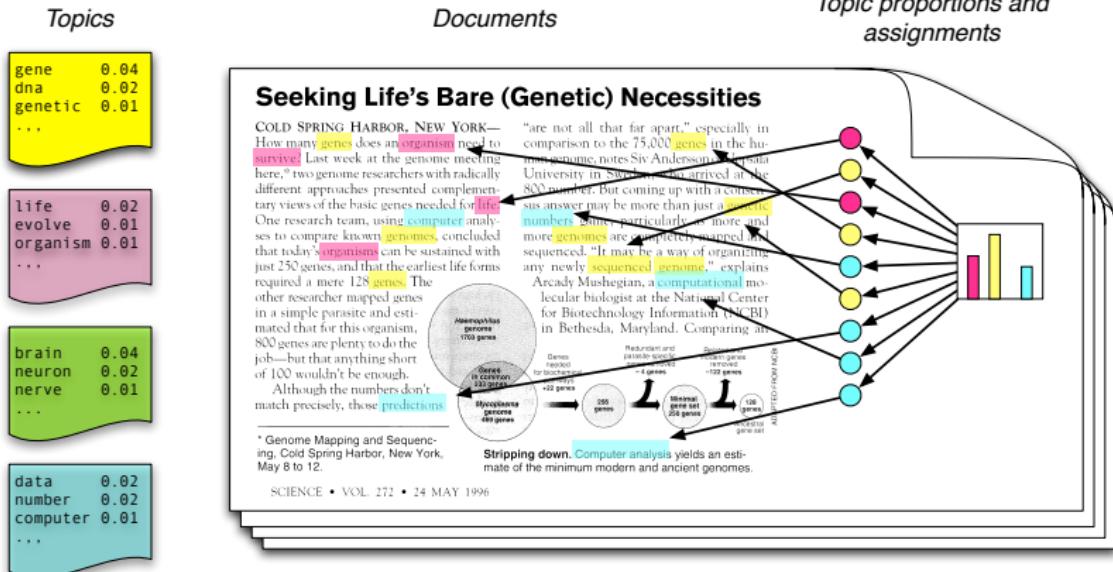
ADAPTED FROM NCBI

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

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

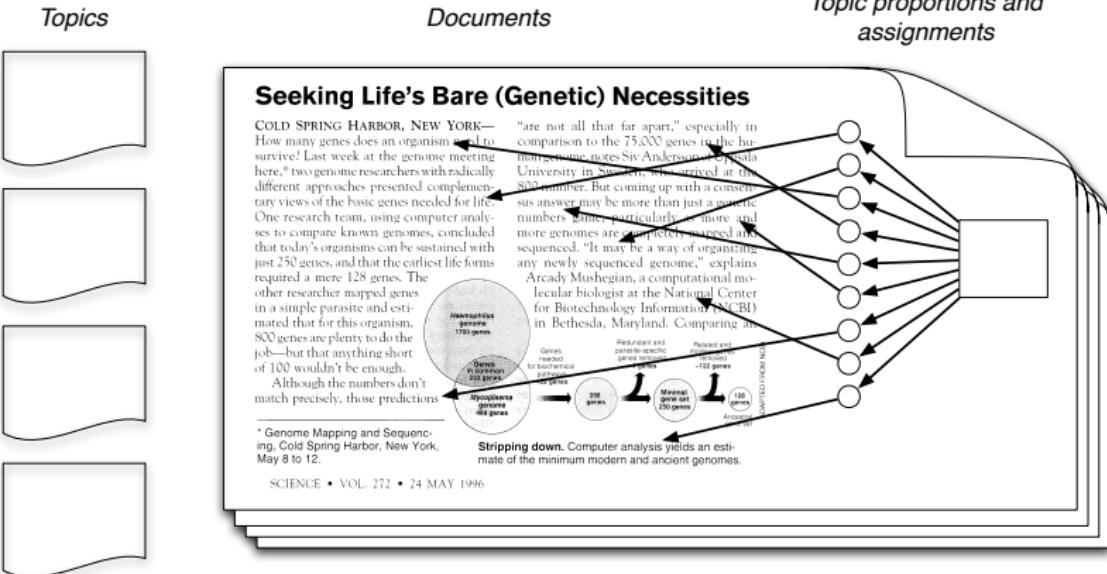
**Simple intuition:** Documents exhibit multiple topics.

# Generative model for LDA



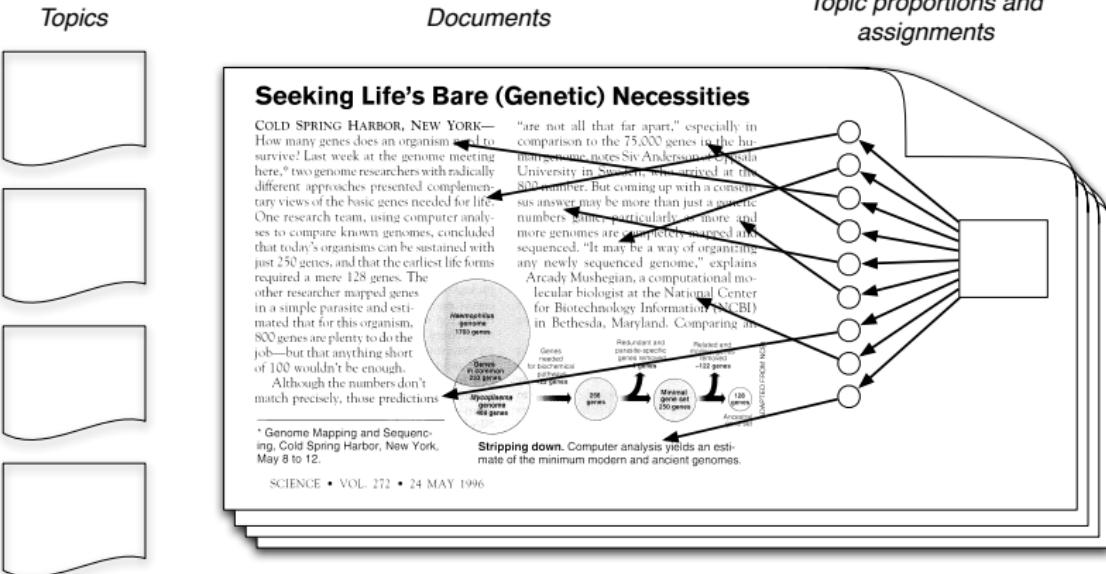
- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

# 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