

Shrinkage for county infection rates



**Yale**

# Goings on

- Assignment 3 due tonight (midnight)
- Assignment 4 is out; due next Thursday
- Quiz 4 this Thursday Nov. 3, on Canvas
  - ▶ PCA, language models, embeddings, Bayes

# For Today

- Embeddings (redux)
- Mixtures
- Bayes

# Core idea of embeddings

- Form a language model but represent words by vectors
- Use PMI-like scores to fit the vectors
- Optimize the embedding vectors using SGD or PCA
- “Two words are similar if they appear with similar words”
- Can be applied whenever have cooccurrence data (Spotify)

# Method 1: Word2Vec

Language model is

$$p_{\text{near}}(w_2 | w_1) = \frac{\exp(\phi(w_2)^T \phi(w_1))}{\sum_w \exp(\phi(w)^T \phi(w_1))}.$$

- Trained on to predict words  $w_2$  that appear near  $w_1$
- SGD over embedding vectors  $\phi(w) \in \mathbb{R}^{100}$

## Method 2: PCA

A different, related approach is to use PCA of PMI:

- Form  $V \times V$  matrix of pointwise mutual information values

$$\text{pmi}(w_1, w_2) = \log \left( \frac{p_{\text{near}}(w_2 | w_1)}{p(w_2)} \right)$$

- Compute top  $k$  eigenvectors  $\phi_1, \dots, \phi_k$
- For each word  $w$ , define embedding as

$$\phi(w) \equiv (\phi_{1w}, \phi_{2w}, \dots, \phi_{kw})^T$$

## Method 3: GloVe

Objective is to minimize squared error

$$\mathcal{O}(\phi) = \sum_{w_1, w_2} f(c_{w_1, w_2}) \left( \phi(w_1)^T \phi(w_2) - \text{pmi}(w_1, w_2) \right)^2$$

where  $c_{w, w'}$  are cooccurrence counts in a window (PMI)

- The parameters are the embedding vectors  $\{\phi(w)\}$
- Unlike least-squares regression, this is a nonconvex minimization problem (wiggly, not bowl-shaped)
- Trained using stochastic gradient descent

---

$f$  is a “fudge factor”. Pennington et al., “GloVe: Global vectors for word representation,” (2015)

# SGD for Word2Vec

- Sample a word pair  $(w_1, w_2)$
- Loss is  $\mathcal{L} = -\log p_{\text{near}}(w_2 | w_1)$
- Gradient update for  $\phi(w_2)$  is

$$\begin{aligned}\phi(w_2) &\leftarrow \phi(w_2) - \eta \frac{\partial \mathcal{L}}{\partial \phi(w_2)} \\ &= \phi(w_2) + \eta (1 - p_{\text{near}}(w_2 | w_1)) \phi(w_1) \\ &= \phi(w_2) + \eta \cdot \text{weight} \cdot \phi(w_1)\end{aligned}$$

Vector  $\phi(w_2)$  is updated to be more similar to  $\phi(w_1)$

# SGD for GloVe

- Sample a word pair ( $w_1, w_2$ )
- Loss is  $\mathcal{L} = \frac{1}{2} (\phi^T(w_1)\phi(w_2) - \text{pmi}(w_1, w_2))^2$
- Gradient update for  $\phi(w_2)$  is

$$\begin{aligned}\phi(w_2) &\leftarrow \phi(w_2) - \eta \frac{\partial \mathcal{L}}{\partial \phi(w_2)} \\ &= \phi(w_2) + \eta \left( \text{pmi}(w_1, w_2) - \phi(w_2)^T \phi(w_1) \right) \phi(w_1) \\ &= \phi(w_2) + \eta \cdot \text{weight} \cdot \phi(w_1)\end{aligned}$$

*Vector  $\phi(w_2)$  is updated to be more similar to  $\phi(w_1)$*

# Where are we going?

- So far: Mostly “surface representations” & explicit features
  - ▶ Listing information for CA housing
  - ▶ Political blogs — word counts

# Where are we going?

- So far: Mostly “surface representations” & explicit features
  - ▶ Listing information for CA housing
  - ▶ Political blogs — word counts
- Next: “Hidden” representations and latent variables

# George Washington

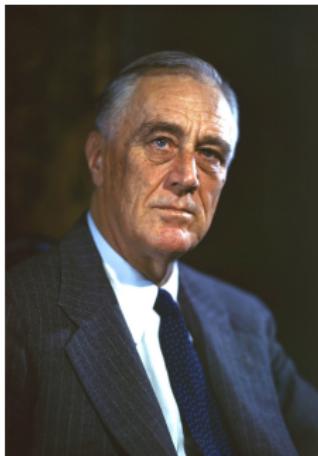
1789



*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.*

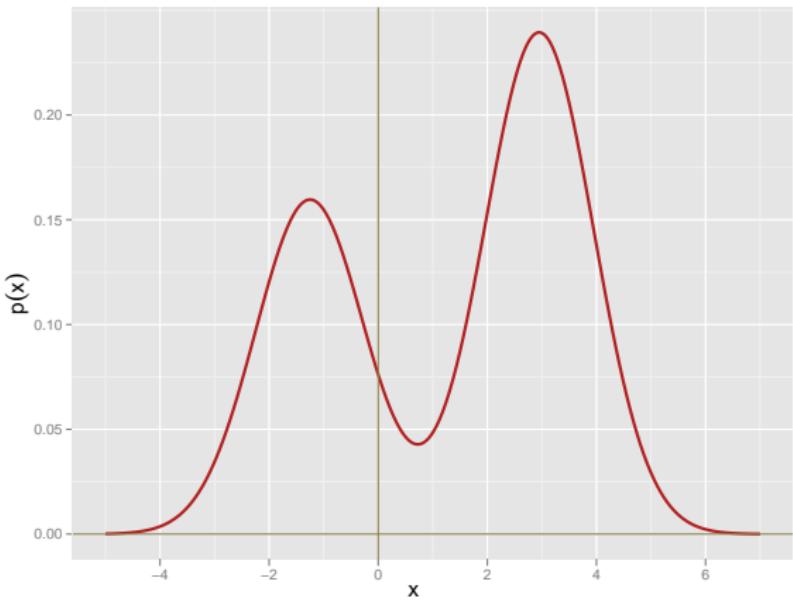
# The elephants in the room



# 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

- *Mixture of f and g:*

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

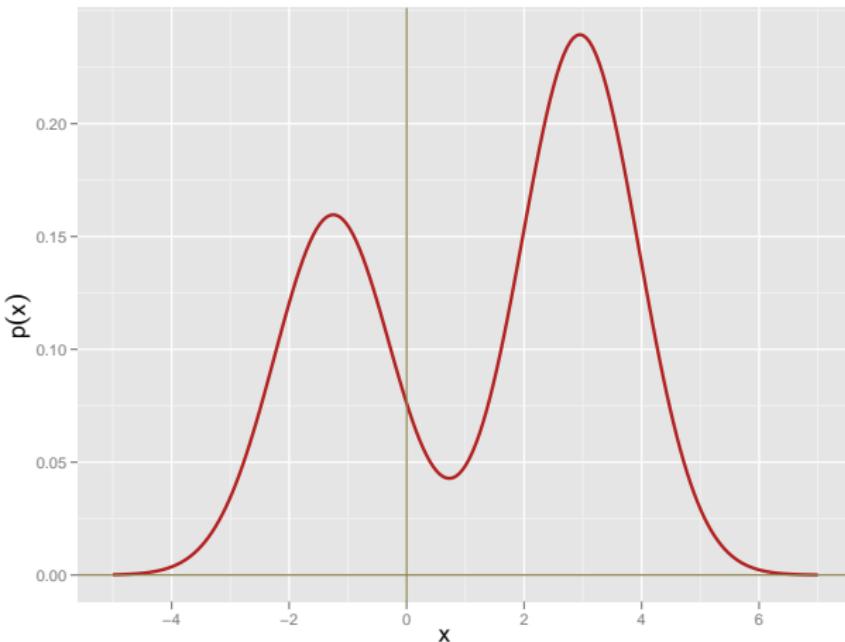
Simplest, most common kind of latent variable model

- *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

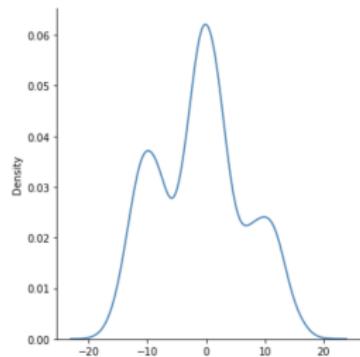


# Let's go to the notebook

```
In [3]: from scipy.stats import multinomial
n = 10000
eta = [.3, .5, .2]
sigma = 3
z = multinomial.rvs(p=eta, n=1, size=n)

x0 = norm.rvs(loc=-10, scale=sigma, size=n)
x1 = norm.rvs(loc=0, scale=sigma, size=n)
x2 = norm.rvs(loc=10, scale=sigma, size=n)

x = z[:,0]*x0 + z[:,1]*x1 + z[:,2]*x2
_ = sns.kdeplot(x, kind='kde')
```



# 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 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

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

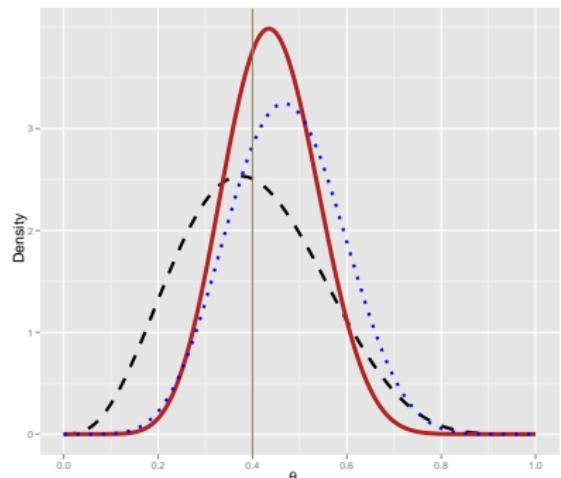
$$\tilde{\beta} = \beta + \text{number of tails} = \beta + n - s$$

---

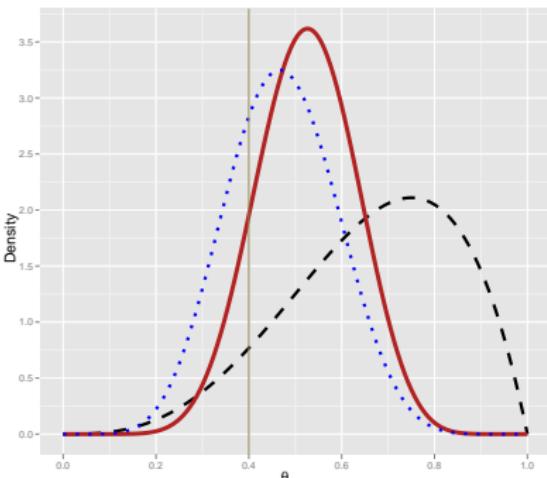
Showing this just uses the simple fact that  $\theta^{\alpha-1} \theta^x = \theta^{x+\alpha-1}$

# 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).

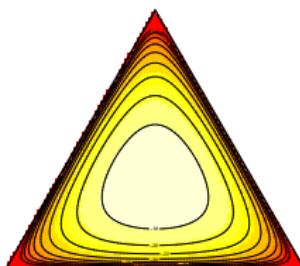
# Dirichlet

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

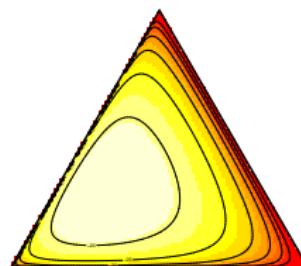
$$\text{Dirichlet}_{\alpha}(\theta) \propto \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \cdots \theta_K^{\alpha_K-1}$$

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

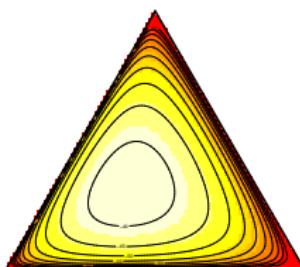
# Example



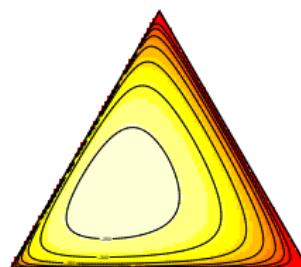
prior with Dirichlet(6,6,6)



likelihood function with  $n = 20$



posterior distribution with  $n = 20$



posterior distribution with  $n = 200$

# Let's go to the notebook!

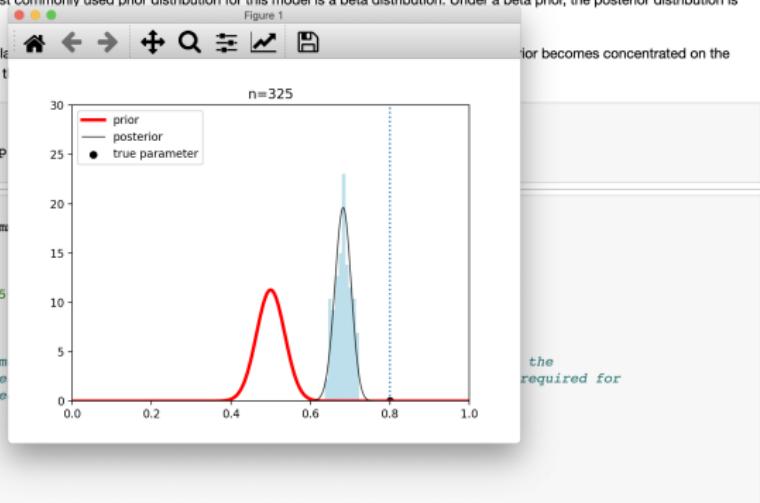
## Demo code for Bayesian analysis

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.

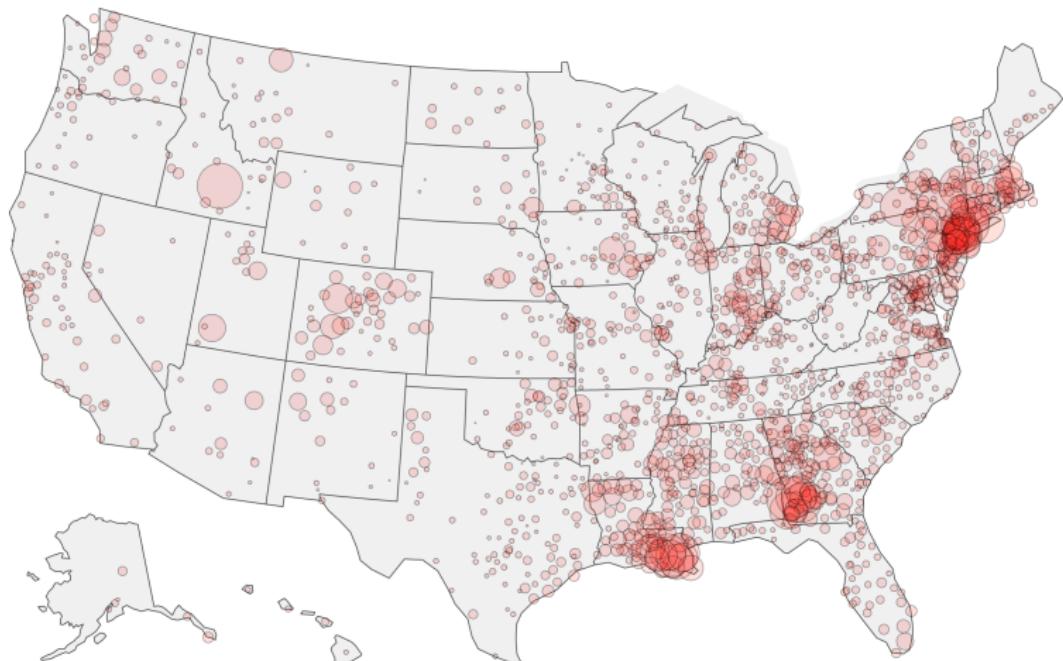
This is illustrated in the following simulation. As the sample size increases, the posterior distribution becomes concentrated on the true parameter. But as the variance of the prior decreases, the posterior distribution becomes more spread out.

```
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 of the prior.  
# variance of the prior decreases as the sample size increases.  
# the posterior to be centered around the true parameter.  
  
scale = 100  
a0 = scale*1  
b0 = scale*1  
  
sample_size = 100
```



The next couple of examples are just covered briefly, as motivation and context. Don't worry about the details!

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



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

per 100,000 people

# 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

Legend: cases (red), predicted (pink)

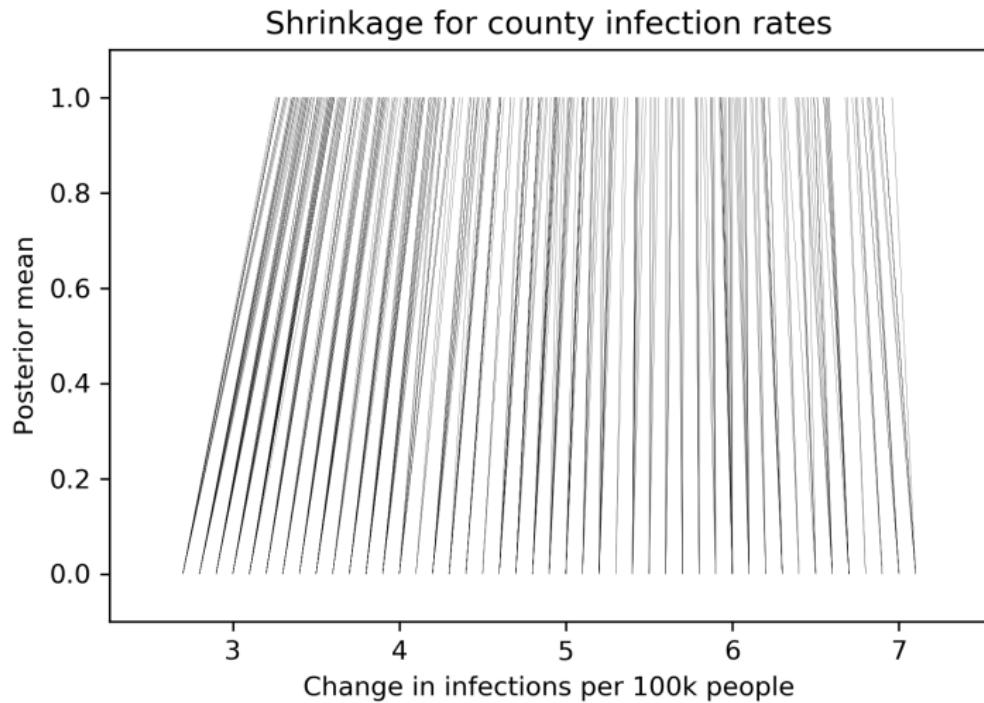
Date	Actual Cases (Red Bar)	Predicted Cases (Pink Bar)
March 25	~100	~100
April 1	~1000	~1000
April 4	~1200	~1200
April 7	~2000	~2000
April 10	~2800	~2800
April 13	~3500	~3500
April 14	~4000	~4000

# Shrinkage

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

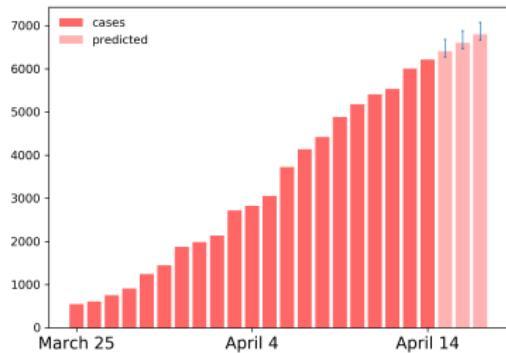
This is called “shrinkage.” The posterior estimates are closer together than the raw frequencies.

# Shrinkage

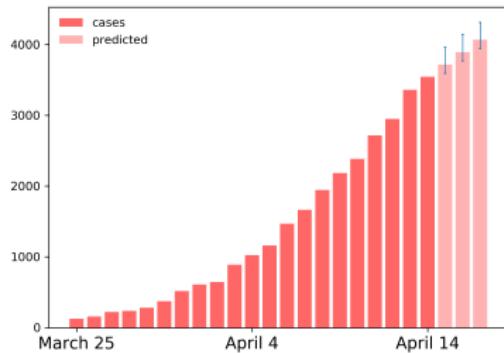


# Examples

Fairfield County, Connecticut



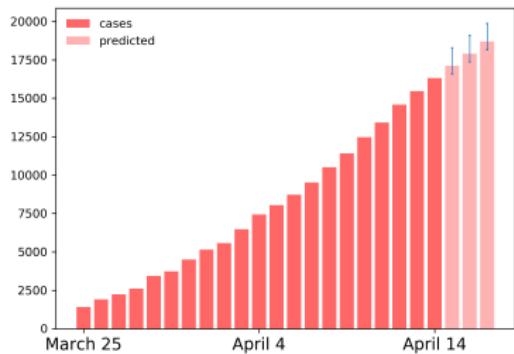
New Haven County, Connecticut



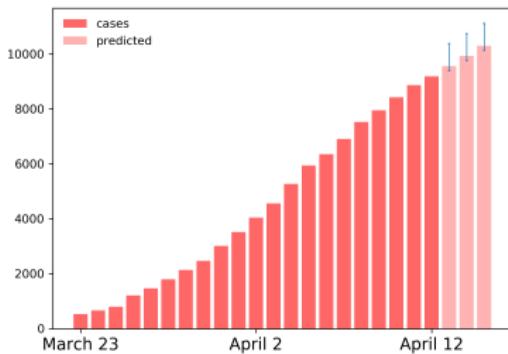
The “credible intervals” trap 95% of the sampled parameters under the simulation.

# 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

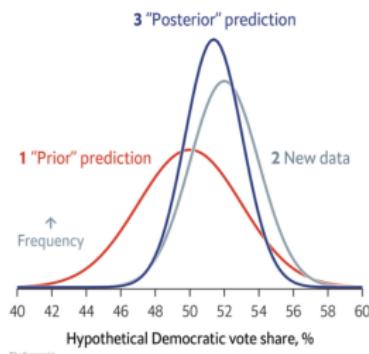
**T**HIS YEAR, *The Economist* is publishing its first-ever statistical forecast of an American presidential election. Developed with the assistance of Andrew Gelman and Merlin Heidemanns, political scientists at Columbia University, our model calculates Joe Biden's and Donald Trump's probabilities of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



## 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

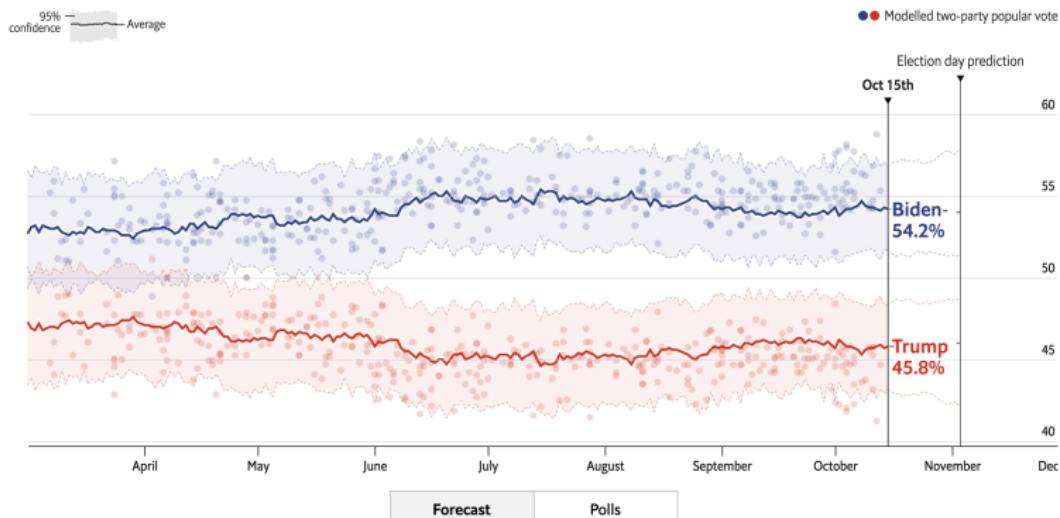
explicitly state what they believe to be true, and how confident they are in that belief. This is called a "prior". Next, after acquiring data, they update this prior to reflect the new information—gaining more confidence if it confirms the prior, and generally becoming more uncertain if it refutes the 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.

# Another Example: Election Forecasting

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

## Modelled popular vote on each day

The model first averages the polls, weighting them by their sample sizes and correcting them for tendencies to overestimate support for one party. It then combines this average with our forecast based on non-polling data, pulling vote shares on each day slightly towards the final election-day projection.



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