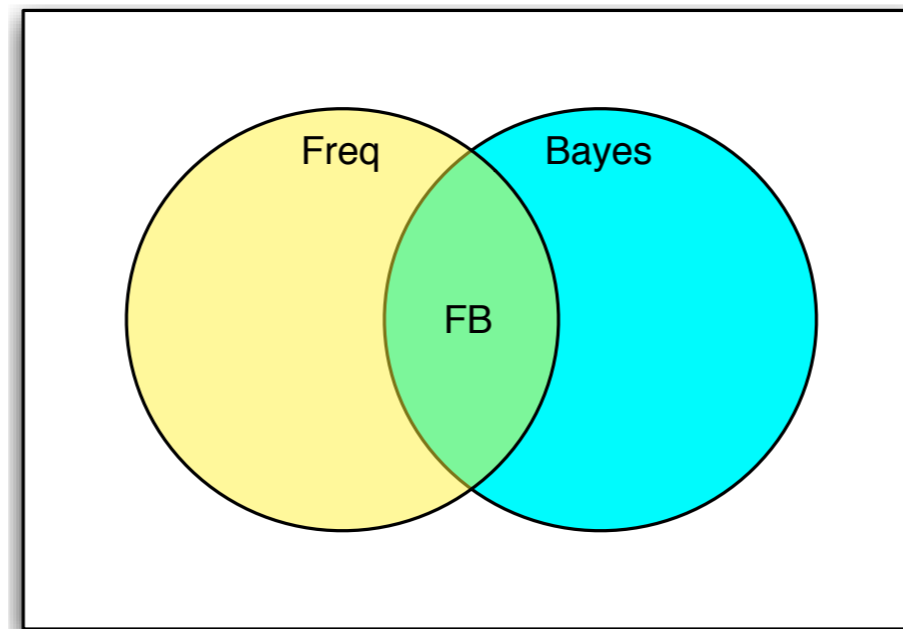


CS 109/Stat 121/AC209/E-109

Data Science Bayesian Methods

Hanspeter Pfister, Joe Blitzstein, Verena Kaynig



This Week

- Form project team if you haven't already. Try your best to have your team formed by next Tuesday, but in any case it is required to fill in the Google form by Tuesday Nov 3, 11:59 pm: <http://goo.gl/forms/CzVRluCZk6>
- HW4 is due Thursday Nov 5, 11:59 pm.

The Theory That Would Not Die

Copyrighted Material

the theory



that would



not die



how bayes' rule cracked



the enigma code,

hunted down russian

submarines & emerged

triumphant from two

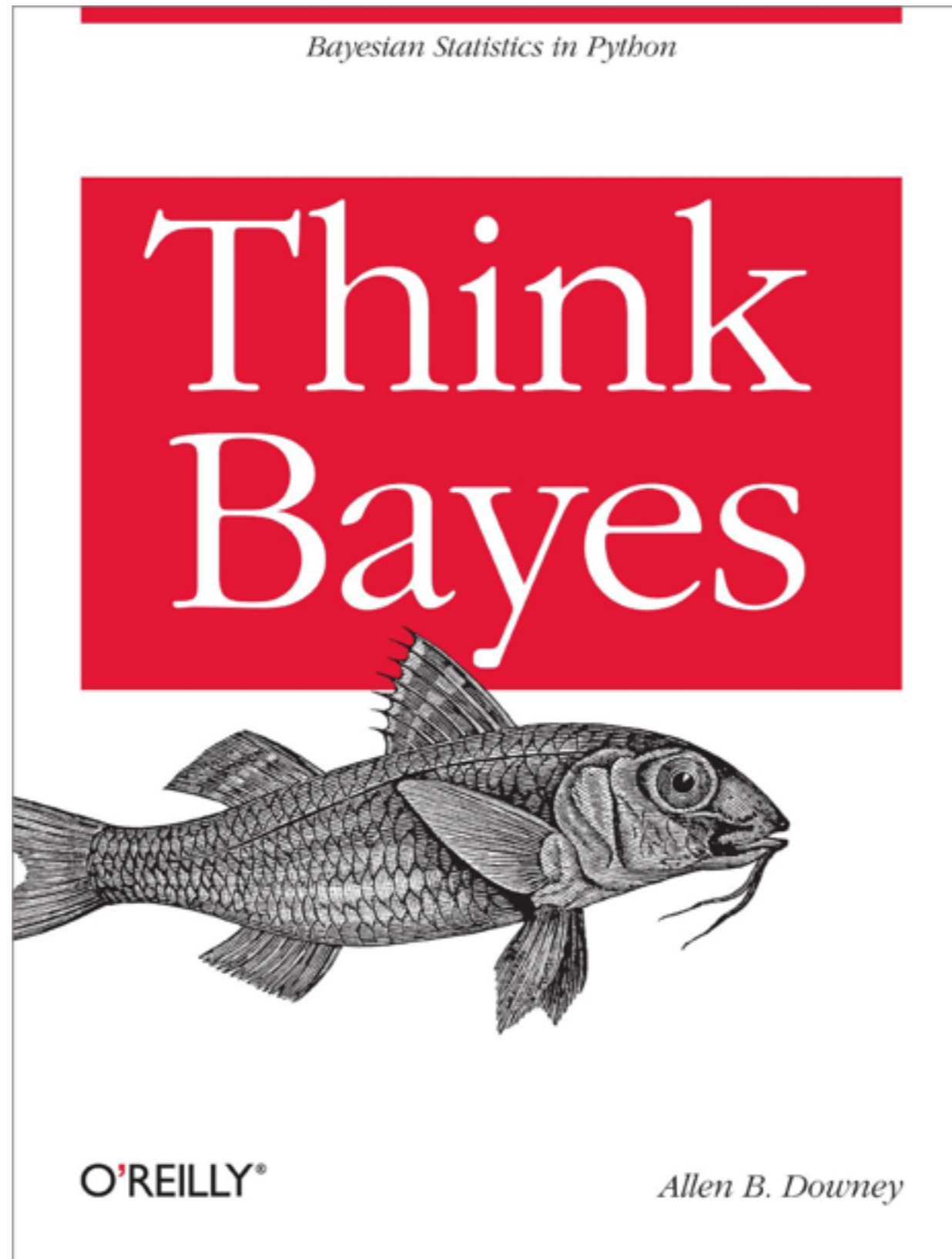


centuries of controversy

sharon bertsch mcgrayne

Copyrighted Material

Think Bayes



<http://greenteapress.com/thinkbayes/>

Probabilistic Programming and Bayesian Methods for Hackers



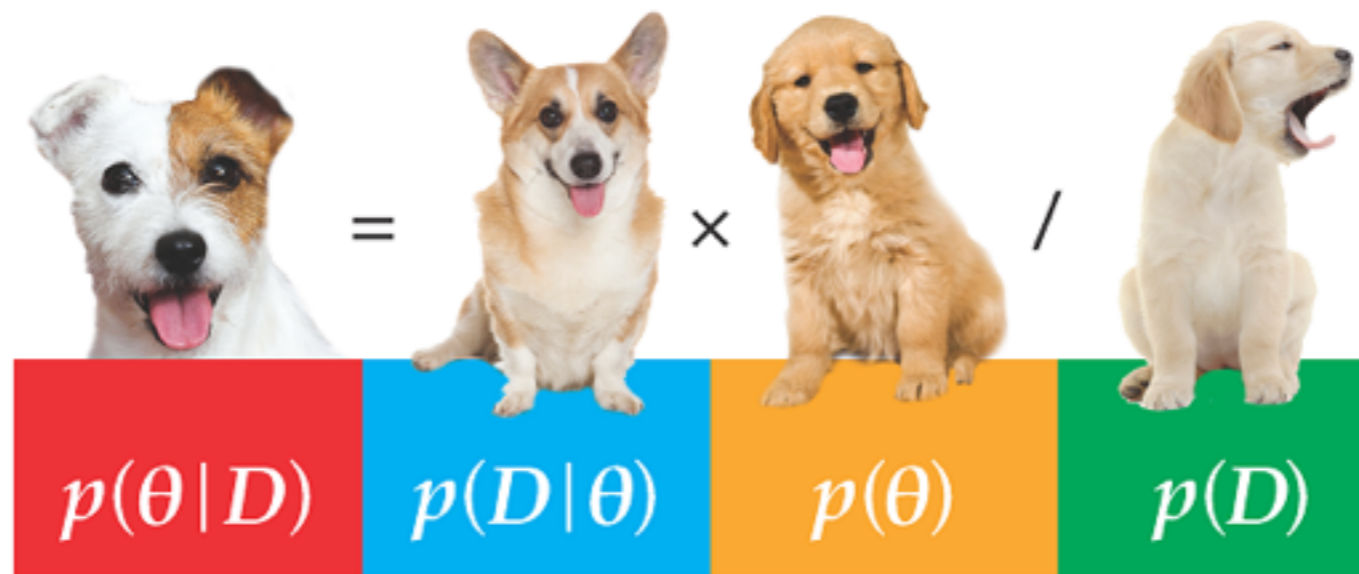
<http://nbviewer.ipython.org/urls/raw.githubusercontent.com/CamDavidsonPilon/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers/master/Prologue/Prologue.ipynb>

Doing Bayesian Data Analysis

Second Edition

Doing Bayesian Data Analysis

A Tutorial with R, JAGS, and Stan



John K. Kruschke



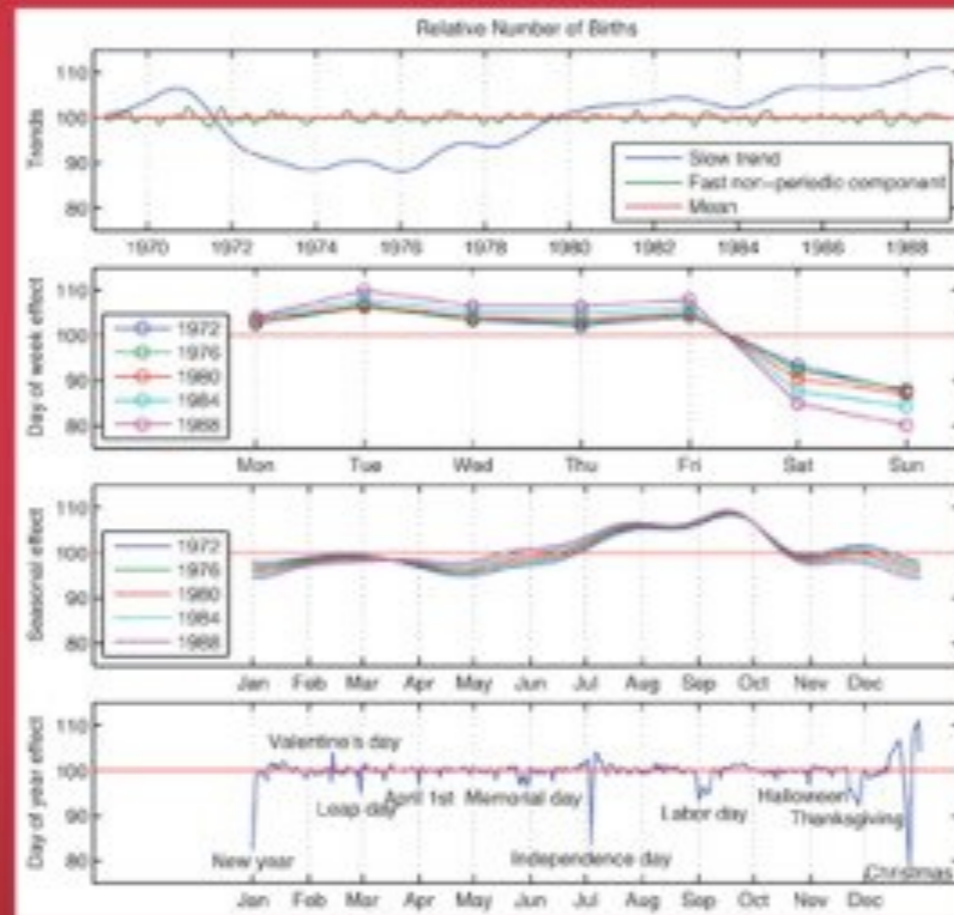
<https://sites.google.com/site/doingbayesiandataanalysis/>

Bayesian Data Analysis

Texts in Statistical Science

Bayesian Data Analysis

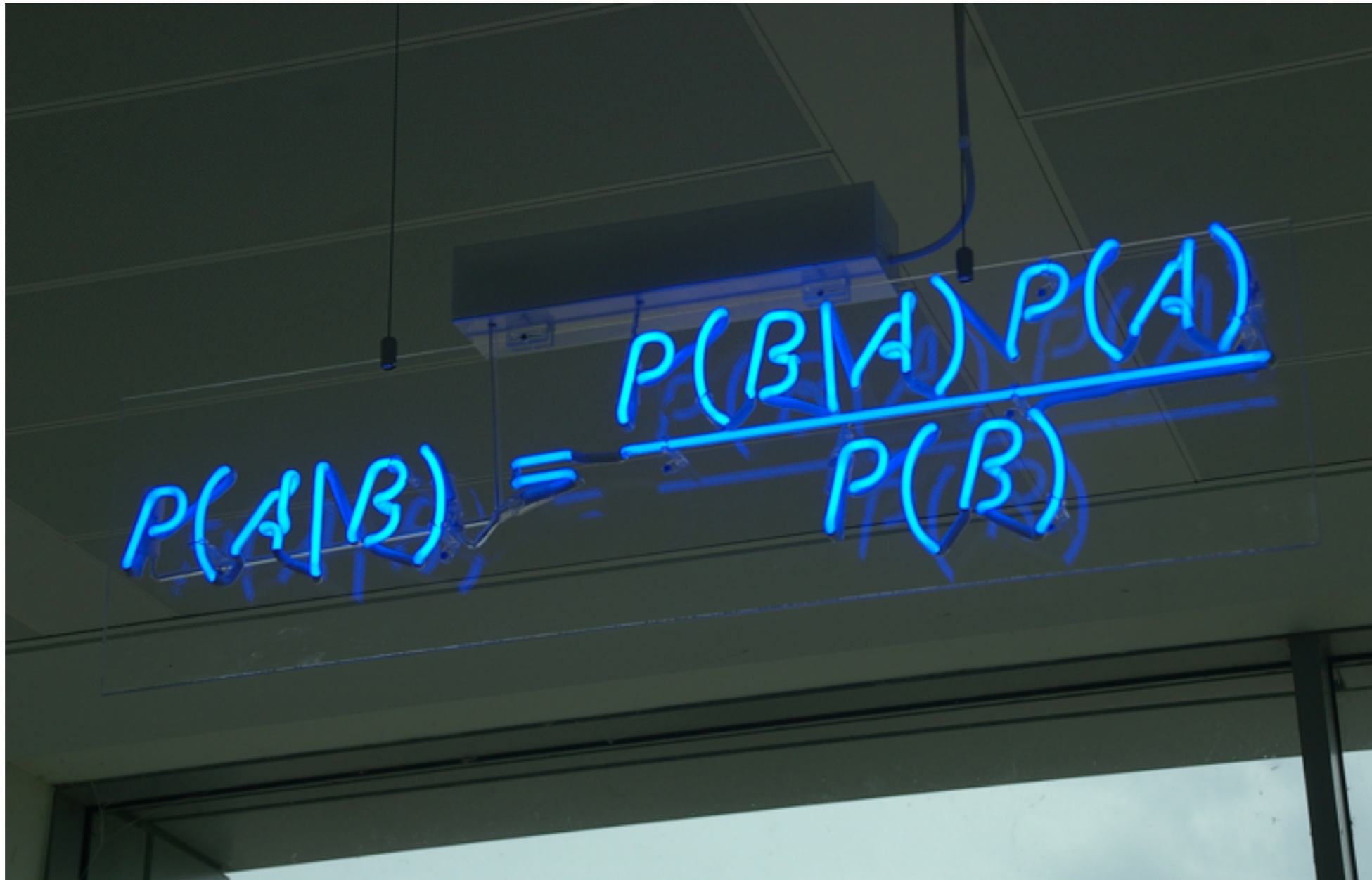
Third Edition



Andrew Gelman, John B. Carlin, Hal S. Stern,
David B. Dunson, Aki Vehtari, and Donald B. Rubin

 CRC Press
Taylor & Francis Group
A CHAPMAN & HALL BOOK

Bayes' rule

A photograph of a whiteboard with the Bayes' rule formula written in blue marker. The formula is $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$. The whiteboard is mounted on a wall, and the lighting is somewhat dim, with the blue marker providing the primary source of color in the image.
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Bayes' rule

$$\overset{\text{posterior probability for A}}{P(A|B)} = \frac{\overset{\text{prior probability for A}}{P(B|A)P(A)}}{P(B)}$$

Bayes' rule, likelihood version

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}$$

Treating the data y as fixed,

$$p(\theta|y) \propto L(\theta)p(\theta)$$

Bayes' rule says the posterior density is proportional to the likelihood function times the prior density.

Discriminative vs. Generative Classifiers

What to model and what not to model?

discriminative: just model $p(y|x)$

generative: give a full probability model

$$p(x,y) = p(x)p(y|x) = p(y)p(x|y)$$

Generative Models

$$P(Y = 1|X = x) = \frac{f(x|Y = 1)P(Y = 1)}{f(x|Y = 1)P(Y = 1) + f(x|Y = 0)P(Y = 0)}$$

(by Bayes' Rule)

Then can model the densities $f(x|Y=1)$, $f(x|Y=0)$.

Naive Bayes Spam Filter

Consider 10 words that occur frequently in spam, and let W_j be the event that the j th word appears in the email.

A certain email uses the 1st and 10th words but not the rest. What's the probability that it is spam?

$$P(\text{spam} | W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10}) = \frac{P(W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10} | \text{spam}) P(\text{spam})}{P(W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10})}$$

Expand denominator with law of total probability

$$P(W) = P(W | \text{spam}) P(\text{spam}) + P(W | \text{not spam}) P(\text{not spam})$$

Naive Bayes Spam Filter

$$P(\text{spam}|W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10}) = \frac{P(W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10}|\text{spam})P(\text{spam})}{P(W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10})}$$

Naive Bayes assumption: *conditional independence* given spam, and also *conditional independence* given not spam.

$$P(W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10}|\text{spam}) = P(W_1|\text{spam})P(W_2^c|\text{spam}) \dots P(W_{10}|\text{spam})$$

$$P(W_1, W_2^c, W_3^c, \dots, W_9^c, W_{10}|\text{not spam}) = P(W_1|\text{not spam})P(W_2^c|\text{not spam}) \dots P(W_{10}|\text{not spam})$$

Huge assumption but huge simplification in statistical and computational complexity.

Naive Bayes

Naive conditional independence assumption:

$$f_j(x_1, \dots, x_d) = f_{j1}(x_1) f_{j2}(x_2) \dots f_{jd}(x_d)$$

Often unrealistic, but still may be *useful* esp. since it leads to a drastic reduction in the number of parameters to estimate.

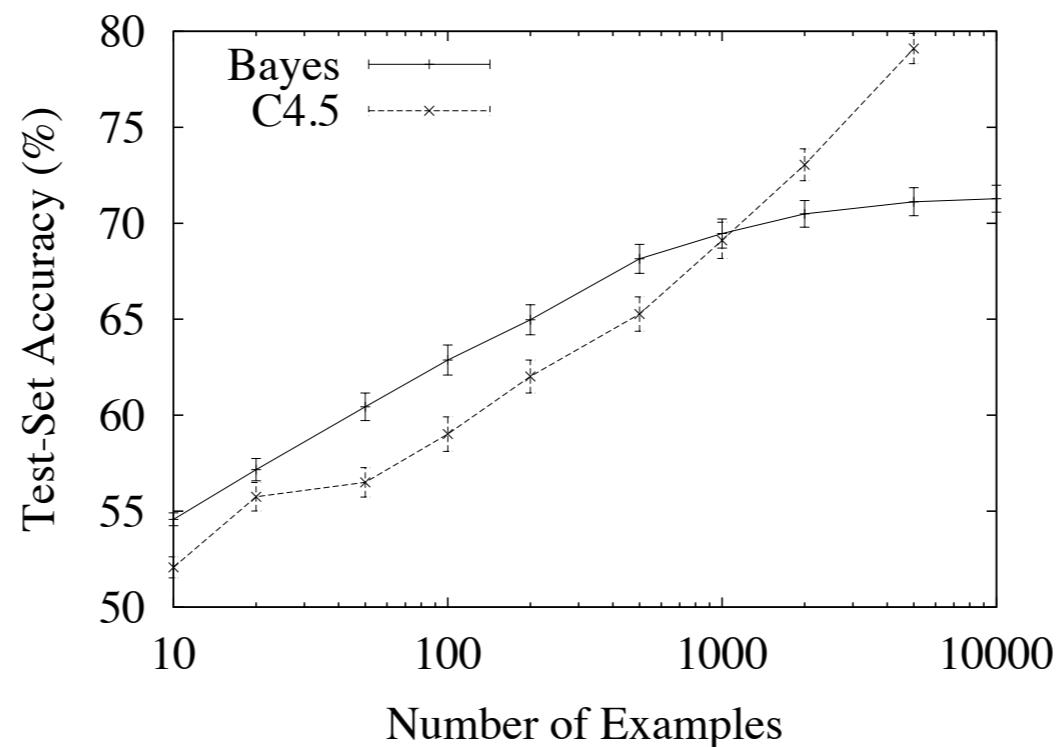


Figure 2: Naive Bayes can outperform a state-of-the-art rule learner (C4.5rules) even when the true classifier is a set of rules.

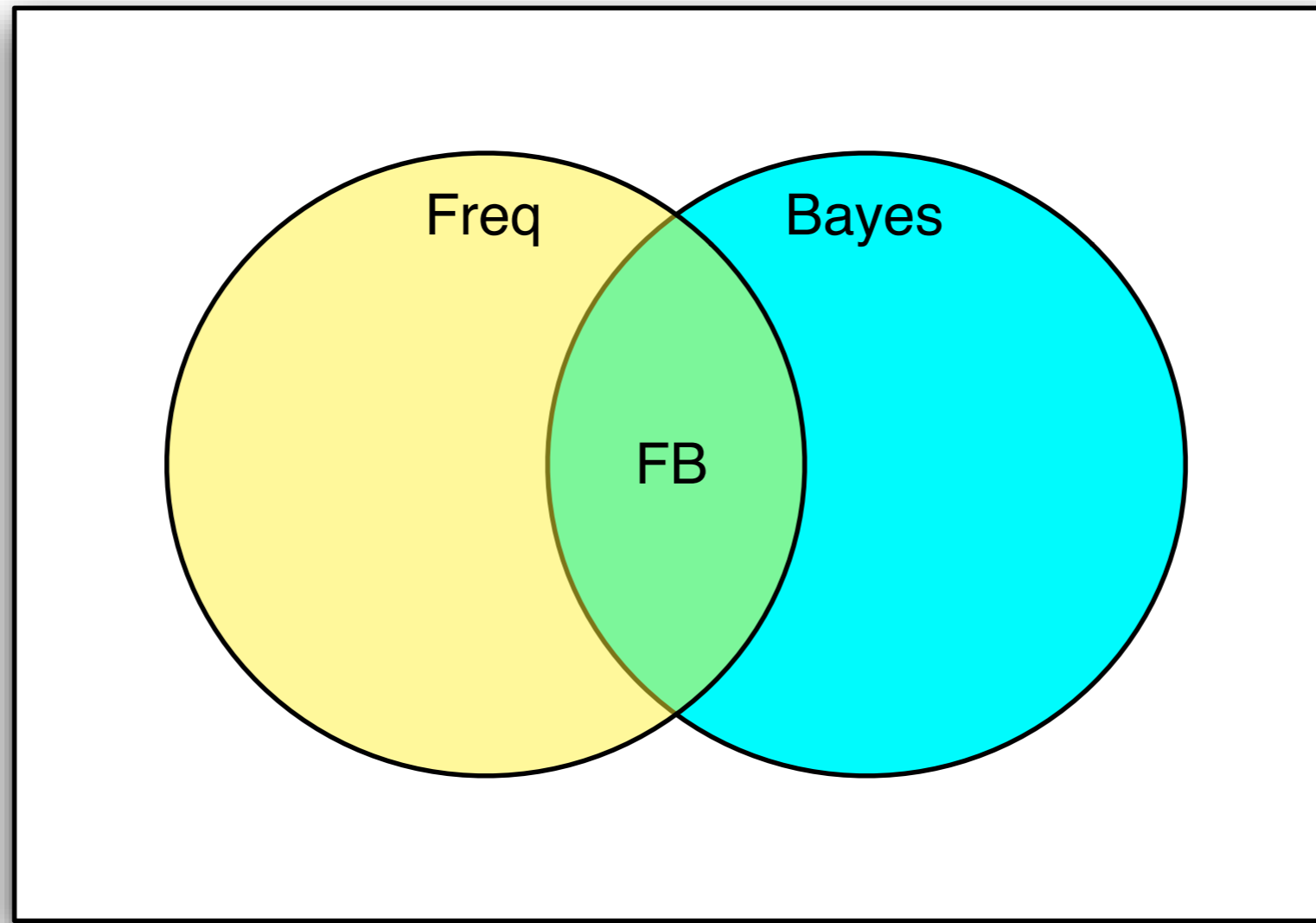
Full Probability Modeling

“The process of Bayesian data analysis can be idealized by dividing it into the following three steps:

1. Setting up a full probability model – a joint probability distribution for all observable and unobservable quantities in a problem...
2. Conditioning on observed data: calculating and interpreting the appropriate posterior distribution – the conditional probability distribution of the unobserved quantities of ultimate interest, given the observed data.
3. Evaluating the fit of the model and the implications of the resulting posterior distribution...”

-- Gelman et al, Bayesian Data Analysis

Bayes-Frequency Reconciliation



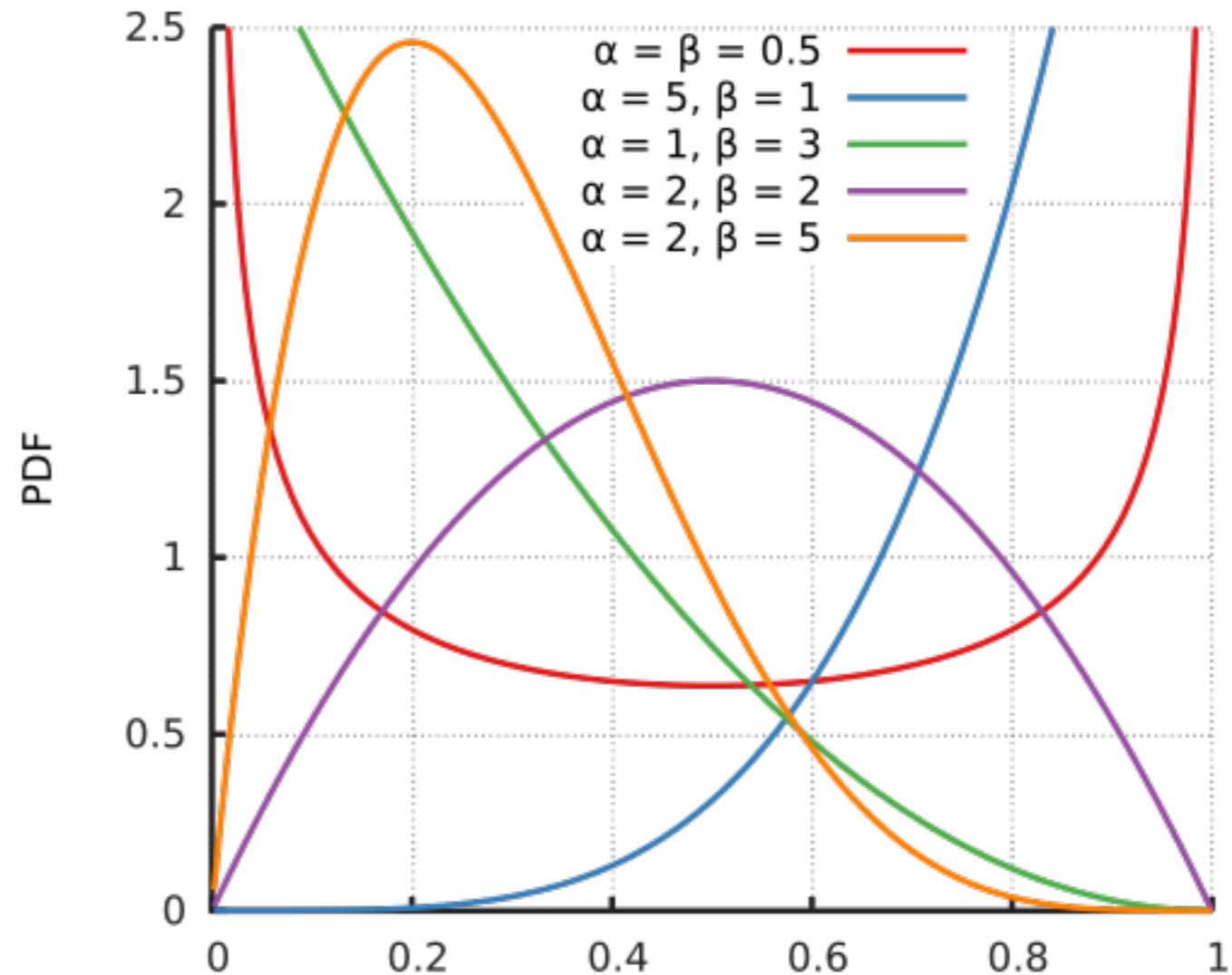
Think like a Bayesian, check like a frequentist.

Conjugate Priors: Beta-Binomial

$$X|p \sim \text{Bin}(n, p)$$

$$p \sim \text{Beta}(a, b)$$

$$f(p) \propto p^{a-1} (1-p)^{b-1}$$



Conjugate Priors: Beta-Binomial

$$X|p \sim \text{Bin}(n, p)$$

$$p \sim \text{Beta}(a, b)$$

Posterior is then $p|X = x \sim \text{Beta}(a + x, b + n - x)$

Conjugate Priors: Normal-Normal

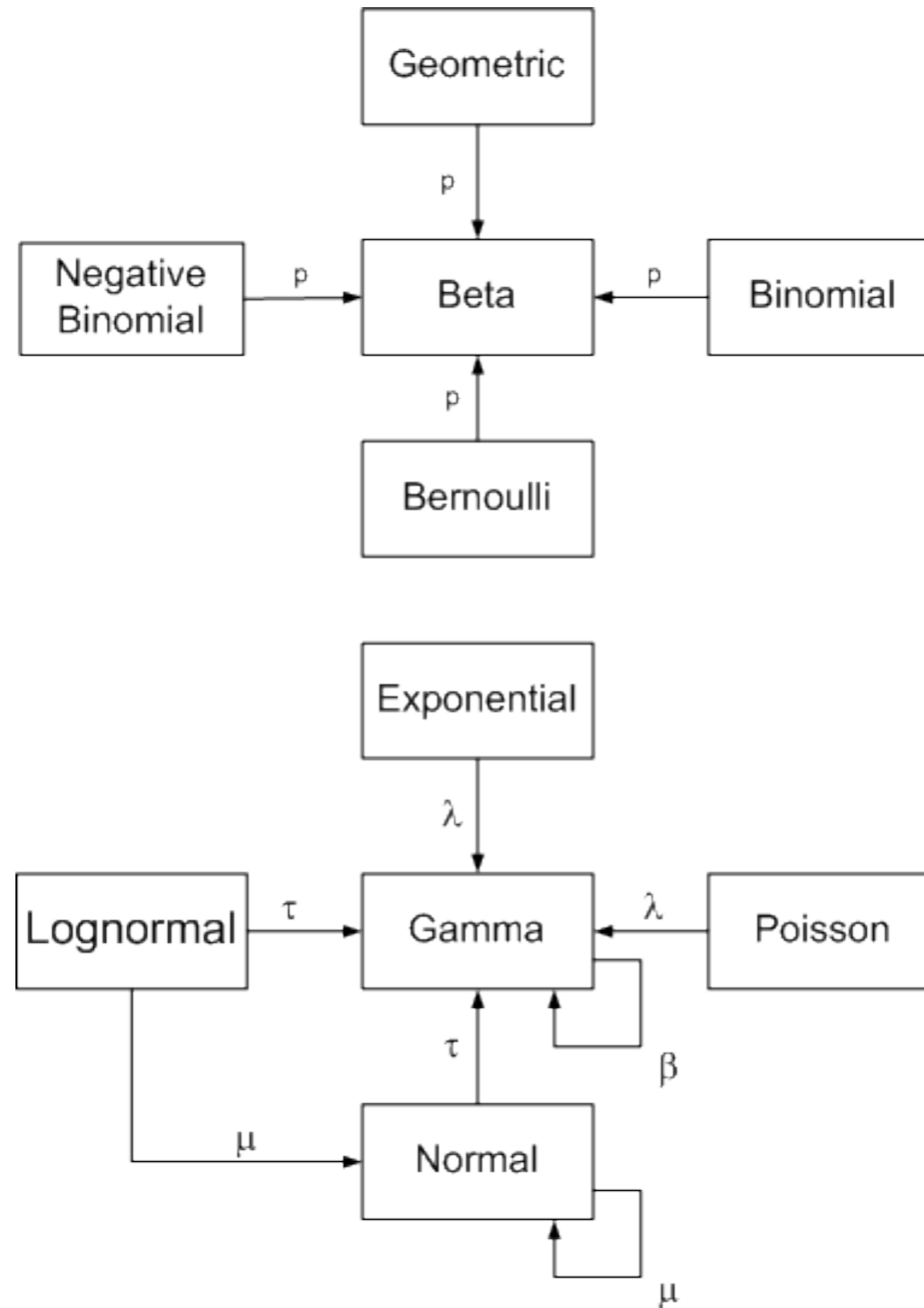
$$y|\mu \sim \mathcal{N}(\mu, \sigma^2)$$

$$\mu \sim \mathcal{N}(\mu_0, \tau^2)$$

$$\text{Then } \mu|y \sim \mathcal{N}\left((1 - B)y + B\mu_0, \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}\right)$$

where $B = \frac{\sigma^2}{\sigma^2 + \tau^2}$ is the *shrinkage factor*

Conjugate Priors



Ranking Reddit Comments:

Example from Probabilistic Programming and Bayesian Methods for Hackers

↑ [-] **Cocky_All_Day** 1957 points 2 days ago (2355|401)
↓ If I ever found myself being attacked by a bear, what advice would you give me?
permalink source report give gold save reply hide child comments

↑ [-] **allenahansen** [S] 4848 points 2 days ago (9861|5007)
↓ If it's a Grizzly Bear, play dead. If you're in California, it's a Black Bear. Fight back with everything you've got because it's trying to kill you. If it's a Polar Bear, you're fucked.
permalink source parent report save give gold reply

↑ [-] **ExBoop** 4052 points 2 days ago (6960|2910)
↓ If it's brown, stay down. If it's black, attack. Now confirmed to be true.
permalink source parent report save give gold reply

↑ [-] **IrkenInvaderGir** 4439 points 2 days ago (8395|3948)
↓ If it's white, good night.
permalink source parent report save give gold reply

http://nbviewer.ipython.org/urls/raw.githubusercontent.com/CamDavidsonPilon/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers/master/Chapter4_TheGreatestTheoremNeverTold/LawOfLargeNumbers.ipynb

Ranking Reddit Comments: A Simple Model

number of upvotes $\sim \text{Bin}(n, p)$

conjugate prior: $p \sim \text{Beta}(a, b)$, pdf $\propto p^{a-1} (1-p)^{b-1}$

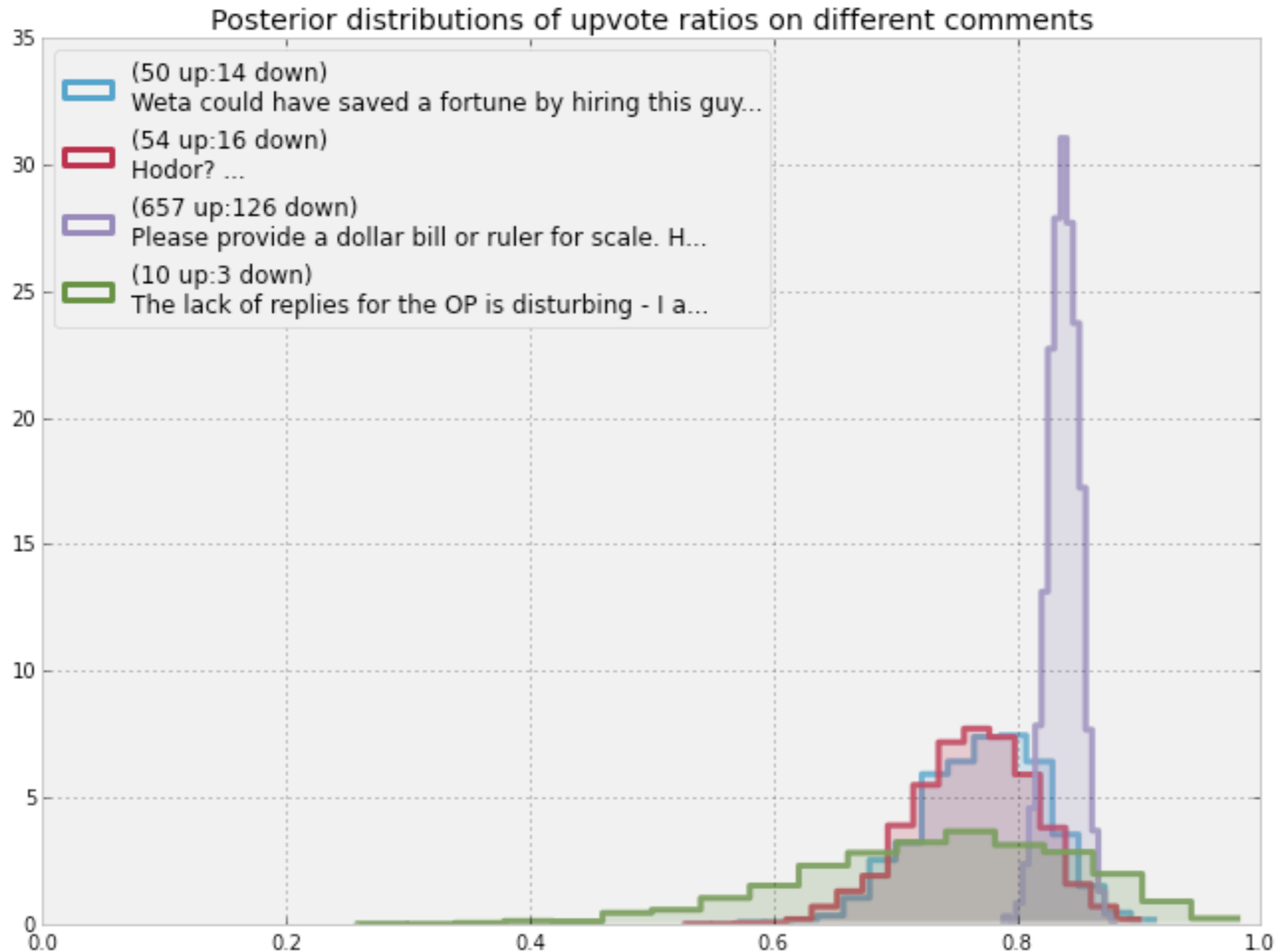
posterior: $p|\text{data} \sim \text{Beta}(a + \#\text{upvotes}, b + \#\text{downvotes})$

Ranking Reddit Comments

Why not just add “pseudocounts” and then use proportion? Why bother with Bayes?

For example, the Agresti-Coull method adds 2 successes and 2 failures.

Posterior Distributions for Reddit Comments



Ranking Reddit Comments by Posterior Quantiles

