

Winter Institute in Data Science and Big Data

Bayesian Inference

JEFF GILL
Distinguished Professor
Departments of Government, and Mathematics & Statistics
Center for Data Science
American University

So What's All This *%\$#@*\$% Bayesian Stuff Anyway?

- ▶ Overt and clear model assumptions.
- ▶ A rigorous way to make *probability* statements about the real quantities of theoretical interest.
- ▶ An ability to update these statements (i.e. learn) as new information is received.
- ▶ Systematic incorporation of *qualitative* knowledge on the subject.
- ▶ Recognition that population quantities are changing over time rather than fixed immemorial.
- ▶ Straightforward assessment of both model quality and sensitivity to assumptions.
- ▶ Freedom from the flawed NHST paradigm.

Typology of Statistics

- ▶ Frequentists: From the Neyman/Pearson/Wald setup. An orthodox view that sampling is infinite and decision rules can be sharp. Estimated quantities usually produced with closed-form statements.
- ▶ Bayesians: From Bayes/Laplace/de Finetti tradition. Unknown quantities are treated probabilistically and the state of the world can always be updated.
- ▶ Likelihoodists: From Fisher. Single sample inference based on finding the parameter value, $\hat{\theta}$, that maximizes the joint distribution of the observed data ($L(\theta|\mathbf{x}) = \prod_{i=1}^n f(x_i|\theta)$), with properties laid-out in Birnbaum (1962). Bayesians that don't know that they are.
- ▶ *So let's look at some critical differences between Frequentists and Bayesians... .*

Critical Differences Between Bayesians and Non-Bayesians, What is Fixed?

Frequentist:

- ▶ Data are an IID random sample from a continuous stream.
- ▶ Parameters are fixed by nature.

Bayesian:

- ▶ Data are observed and therefore fixed by the sample generated.
- ▶ Parameters are unknown and described distributionally.

Critical Differences Between Bayesians and Non-Bayesians, Interpretation of Probability

Frequentist:

- ▶ Probability is observed from the long-run proportion of times that some event occurs in a replicated experiment.
- ▶ Probabilistic quantity of interest is $p(\text{data}|\mathcal{H}_0)$.

Bayesian:

- ▶ Probability is the researcher/observer “degree of belief” before or after the data are observed.
- ▶ Probabilistic quantity of interest is $p(\theta|\text{data})$.

Critical Differences Between Bayesians and Non-Bayesians, General Inference

Frequentist:

- ▶ Point estimates and standard errors or 95% *confidence* intervals.
- ▶ Deduction from $p(\text{data}|\mathcal{H}_0)$, by setting α in advance.
- ▶ Accept \mathcal{H}_1 if $p(\text{data}|\mathcal{H}_0) < \alpha$
- ▶ Accept \mathcal{H}_0 if $p(\text{data}|\mathcal{H}_0) \geq \alpha$

Bayesian:

- ▶ Induction from $p(\theta|\text{data})$, starting with $p(\theta)$.
- ▶ Broad descriptions of the posterior distribution such as means and quantiles.
- ▶ Highest posterior density intervals indicating region of highest posterior probability, regardless of contiguity.

Critical Differences Between Bayesians and Non-Bayesians, Post-hoc Quality Checks

Frequentist:

- ▶ Calculation of Type I and Type II errors, even if there is no setting α in advance.
- ▶ *Sometimes*: effect size and/or power.
- ▶ *Usually*: fixation with small differences in p -values despite large measurement error in the social sciences relative to other scientific disciplines.

Bayesian:

- ▶ Posterior predictive checks from integrating over posterior.
- ▶ Sensitivity checks to forms of the prior, and other assumptions.
- ▶ Bayes factors for model comparison, BIC, DIC.

Reasons *Not* to Use Bayesian Inference in the Social Sciences:

- ▶ The population parameters of interest truly fixed and unchanging under all realistic circumstances.
- ▶ We do *not* have any information prior to the model specification.
- ▶ It is necessary to provide statistical results as if data were from a *controlled experiment*.
- ▶ We care more about “significance” than effect size.
- ▶ Computers are slow and relatively unavailable.
- ▶ We want very automated, “cookbook” type procedures.



CONFORMITY

WHEN PEOPLE ARE FREE TO DO AS THEY PLEASE,
THEY USUALLY IMITATE EACH OTHER.

www.despair.com

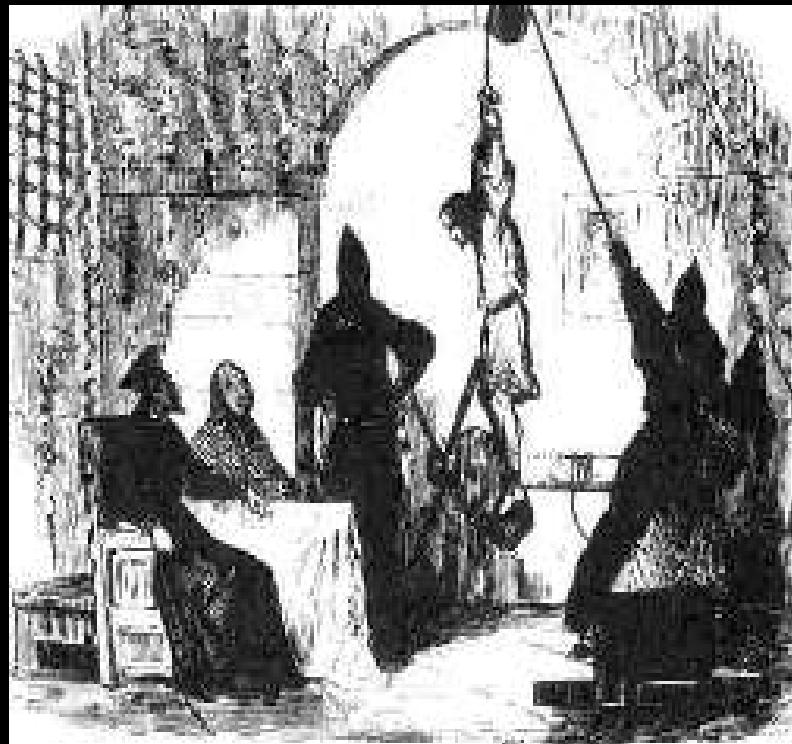
Reasons to Use Bayesian Inference in the Social Sciences:

- ▶ We want to be very careful about stipulating assumptions and are willing to defend them.
- ▶ We view the world probabilistically, rather than as a set of fixed phenomena that are either known or unknown.
- ▶ Every statistical model ever created in the history of the human race is subjective; we are willing to admit it.
- ▶ Prior information abounds in the social sciences and it is important and helpful to use it.

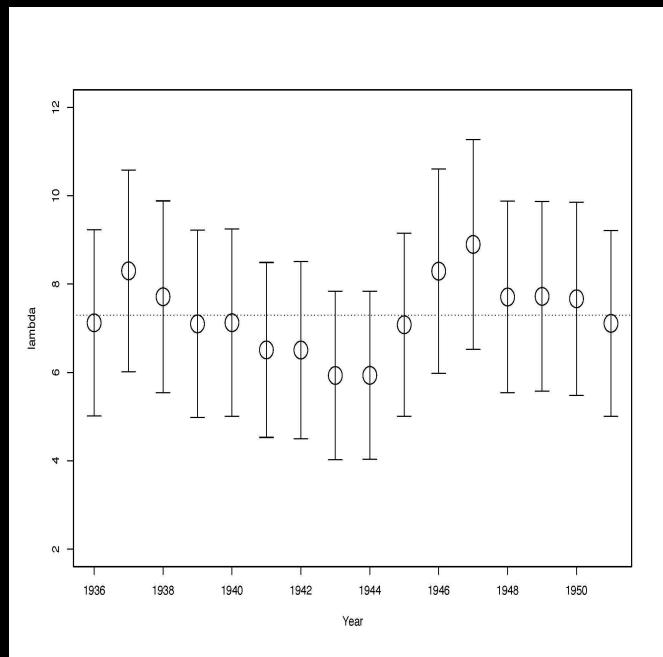
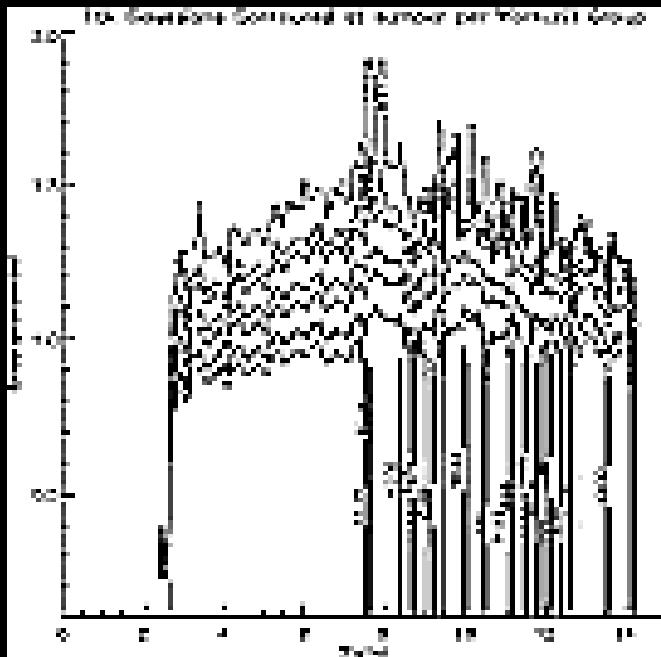


Some Problems with Traditional Statistical Thinking in the Social Sciences

- ▶ Small-n inference.
- ▶ Significance through sample size.
- ▶ Confidence.
- ▶ Contrived ignorance and buried assumptions.
- ▶ Null Hypothesis Testing/Star-gazing.



Large and Small Sample Inference

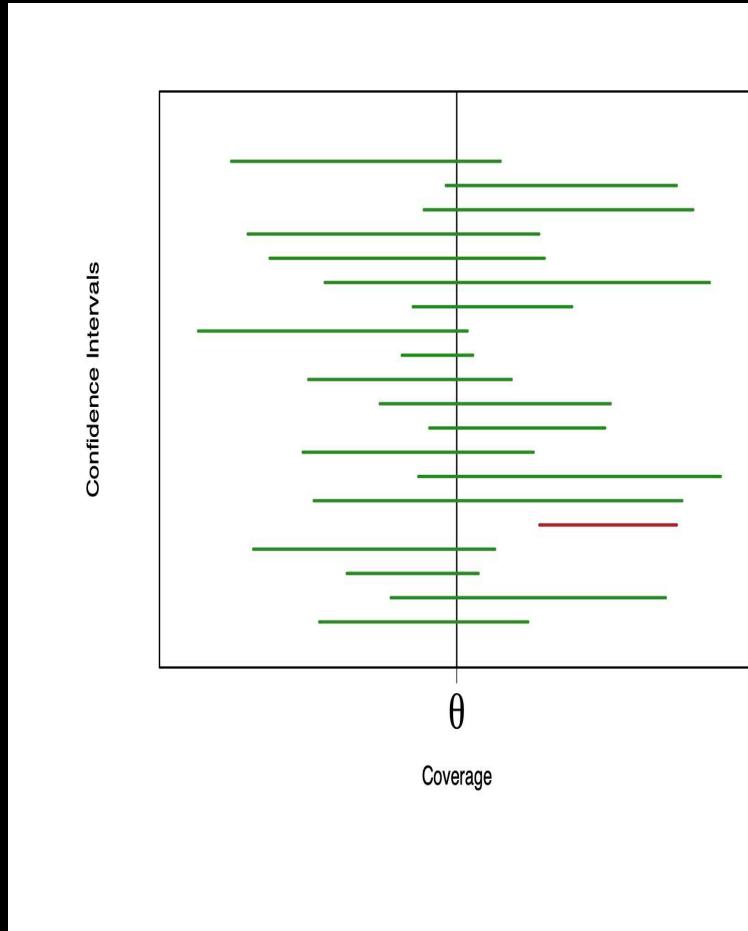


<http://setiathome.ssl.berkeley.edu/>

Marriage Rates per 1000 in Italy 1936 to 1951.

Confidence

- ▶ Which of these is the correct interpretation of a $(1 - \alpha)$ confidence interval?
 - ▷ An interval that has a $1 - \alpha\%$ chance of containing the true value of the parameter.
 - ▷ An interval that over $1 - \alpha\%$ of replications contains the true value of the parameter, *on average*.
- ▶ What interpretation do people really *want*.



Contrived Ignorance, Buried Assumptions

- ▶ Models with uniform priors.
- ▶ Normality.
- ▶ Correlation coefficient.
- ▶ Only two models tested.
- ▶ No such thing as specification searches.



The pseudo-Frequentist NHST is wrong

- A few authors have noted this (**just a sample**): *Barnett 1973, Berger, Boukai, and Wang 1997, Berger Thomas Sellke 1987, Berkhardt and Schoenfeld 2003, Bernardo 1984, Brandstätter 1999, Carver 1978, 1993, Dar, Serlin and Omar 1994, Cohen 1988, 1994, 1992, 1977, 1962, Denis 2005, Falk and Greenbaum 1995, Gelman, Carlin, Stern, and Rubin 1995, Gigerenzer 1987, 1993, 1998, Gigerenzer and Murray 1987, Gill 1999, 2005, Gliner, Leech and Morgan 2002, Grayson 1998, Greenwald 1975, Greenwald, Gonzalez, Harris and Guthrie 1996, Hager 2000, Howson and Urbach 1993, Hunter 1997, Hunter and Schmidt 1990, Jeffreys 1961, Kirk 1996, Krueger 1999, 2001, Lindsay 1995, Loftus 1991, 1993a, 1993b, 1994, 1996, Loftus and Bamber 1990, Macdonald 1997, Meehl 1967, 1978, 1990, 1978, Nickerson 2000, Oakes 1986, Pollard 1993, Pollard and Richardson 1987, Robinson and Levin 1997, Rosnow and Rosenthal 1989, Rozeboom 1960, 1997, Schmidt 1996, Schmidt and Hunter 1977, Sedlmeier and Gigerenzer 1989, Thompson 2002, Wilkinson 1999.*
 - 1. Artificial Model Selection Criteria
 - 2. The Arbitrariness of Alpha
 - 3. Replication Fallacy
 - 4. Asymmetry and Accepting the Null Hypothesis
 - 5. Probabilistic Modus Tollens
 - 6. Inverse Probability Problem
- Why?

Regular Modus Tollens

If A then B	If H_0 is true then the data will follow an expected pattern
Not B observed	The data do not follow the expected pattern
Therefore not A	Therefore H_0 is false.

Probabilistic Modus Tollens

If A then B is highly likely	If H_0 is true then the data are highly likely to follow an expected pattern
Not B observed	The data do not follow the expected pattern
Therefore A is highly unlikely	Therefore H_0 is highly unlikely.

Probabilistic Modus Tollens Example

If A then B is
highly likely

Not B observed

Therefore A is highly unlikely

If a person is an
American, then it is
highly unlikely she is
a member of Congress.

The person is a member
of Congress

Therefore it is highly
unlikely she is
an American.

Misconceptions about Inverse Probability

- ▶ The inferential mechanism of the null hypothesis significance test is based on conditional probability.
- ▶ The test looks at: $p(\text{data}|H_0)$, “how likely is it to observe these data, given that the null hypothesis of no effect is *true*.”
- ▶ It is commonly (mis)interpreted as: $p(H_0|\text{data})$, “how probable is the null hypothesis, given these observed data.”
- ▶ These (the right and the wrong) statements are fundamentally different quantities and can only be related with Bayes’ Law:

$$p(H_0|\text{data}) = \frac{p(H_0)}{p(\text{data})} p(\text{data}|H_0).$$

- ▶ The problem comes from an unholy blending of Fisher and Neyman/Pearson.

Misconceptions about Inverse Probability

- The order of conditionality can be really important.
- suspected probability of AIDS in risk group: $P(A) = 0.02$
probability of correct positive classification: $P(C|A) = 0.95$
probability of correct negative classification: $P(C^c|A^c) = 0.97$

Misconceptions about Inverse Probability

- ▶ The order of conditionality can be really important.
- ▶ suspected probability of AIDS in risk group: $P(A) = 0.02$
probability of correct positive classification: $P(C|A) = 0.95$
probability of correct negative classification: $P(C^c|A^c) = 0.97$
- ▶ Suppose we want $P(A|C)$, from:
$$P(A|C) = \frac{P(A)}{P(C)} P(C|A)$$

Misconceptions about Inverse Probability

- The order of conditionality can be really important.

- suspected probability of AIDS in risk group: $P(A) = 0.02$

probability of correct positive classification: $P(C|A) = 0.95$

probability of correct negative classification: $P(C^c|A^c) = 0.97$

- Suppose we want $P(A|C)$, from:

$$P(A|C) = \frac{P(A)}{P(C)} P(C|A)$$

$$\begin{aligned} P(C) &= P(C \cap A) + P(C \cap A^c) \\ &= P(C|A)P(A) + P(C|A^c)P(A^c) \\ &= P(C|A)P(A) + [1 - P(C^c|A^c)]P(A^c) \\ &= (0.95)(0.02) + (1 - 0.97)(0.98) \cong 0.05 \end{aligned}$$

- Getting the unconditional:

Misconceptions about Inverse Probability

- ▶ The order of conditionality can be really important.
- ▶ suspected probability of AIDS in risk group: $P(A) = 0.02$

probability of correct positive classification: $P(C|A) = 0.95$

probability of correct negative classification: $P(C^c|A^c) = 0.97$

- ▶ Suppose we want $P(A|C)$, from:

$$P(A|C) = \frac{P(A)}{P(C)} P(C|A)$$

$$\begin{aligned} P(C) &= P(C \cap A) + P(C \cap A^c) \\ &= P(C|A)P(A) + P(C|A^c)P(A^c) \\ &= P(C|A)P(A) + [1 - P(C^c|A^c)]P(A^c) \\ &= (0.95)(0.02) + (1 - 0.97)(0.98) \cong 0.05 \end{aligned}$$

- ▶ So now we can calculate:

$$P(A|C) = \frac{P(A)}{P(C)} P(C|A) = \frac{0.02}{0.05}(0.95) = 0.38$$

The History of Bayesian Statistics—Milestones



- ▶ Reverend Thomas Bayes (1702-1761).
- ▶ Pierre Simon Laplace.
- ▶ Pearson (Karl), Fisher, Neyman and Pearson (Egon), Wald.
- ▶ Jeffreys, de Finetti, Good, Savage, Lindley, Zellner.
- ▶ A world divided.
- ▶ The revolution: Gelfand and Smith (1990).
- ▶ Today...

Two Primary Principles of Bayesian Inference

Principle I.

Principle II.

Two Primary Principles of Bayesian Inference

Principle I.

Explicit and direct use of probability for describing uncertainty:

- ▷ probability models (likelihood fn.) for data given parameters,
- ▷ probability distributions (PDF,PMF) for parameters.

Principle II.

Two Primary Principles of Bayesian Inference

Principle I.

Explicit and direct use of probability for describing uncertainty:

- ▷ probability models (likelihood fn.) for data given parameters,
- ▷ probability distributions (PDF,PMF) for parameters.

Principle II.

Inference for unknown values conditioned on observed data:

- ▷ use of inverse probability,
- ▷ Bayes theorem,
- ▷ description of full posterior.

The Three General Steps

Step I.

Step II.

Step III.

The Three General Steps

Step I.

Specify a probability model for unknown parameter values that includes some prior knowledge about the parameters if available.

Step II.

Step III.

The Three General Steps

Step I.

Specify a probability model for unknown parameter values that includes some prior knowledge about the parameters if available.

Step II.

Update knowledge about the unknown parameters by conditioning this probability model on observed data.

Step III.

The Three General Steps

Step I.

Specify a probability model for unknown parameter values that includes some prior knowledge about the parameters if available.

Step II.

Update knowledge about the unknown parameters by conditioning this probability model on observed data.

Step III.

Evaluate the fit of the model to the data and the sensitivity of the conclusions to the assumptions.

Simple Mechanics

$$\pi(\theta|\mathbf{x}) = \frac{p(\theta)L(\theta|\mathbf{x})}{\int_{\Theta} p(\theta)L(\theta|\mathbf{x})d\theta}$$
$$\propto p(\theta)L(\theta|\mathbf{x})$$

Posterior Probability \propto Prior Probability \times Likelihood Function

Views On Priors Determine Types of Bayesians

- ▶ **Empirical Bayes**: prior distributions are produced from other parts of the data, or possibly from the same data. Results are reported like Frequentists.
- ▶ **Proper Bayes**:
- ▶ **Reference Bayes**:
- ▶ **Decision-Theoretic Bayes**:
- ▶ **Bayesians of Convenience**:

Views On Priors Determine Types of Bayesians

- ▶ **Empirical Bayes**: prior distributions are produced from other parts of the data, or possibly from the same data. Results are reported like Frequentists.
- ▶ **Proper Bayes**: prior distributions come from previously compiled evidence, such earlier studies or published work, researcher intuition, or substantive experts. Results are reported without utility or loss functions.
- ▶ **Reference Bayes**:
- ▶ **Decision-Theoretic Bayes**:
- ▶ **Bayesians of Convenience**:

Views On Priors Determine Types of Bayesians

- ▶ **Empirical Bayes**: prior distributions are produced from other parts of the data, or possibly from the same data. Results are reported like Frequentists.
- ▶ **Proper Bayes**: prior distributions come from previously compiled evidence, such earlier studies or published work, researcher intuition, or substantive experts. Results are reported without utility or loss functions.
- ▶ **Reference Bayes**: prior distributions are created to influence the posterior as little as mathematically possible (“objective”). Results are reported without utility or loss functions.
- ▶ **Decision-Theoretic Bayes**:
- ▶ **Bayesians of Convenience**:

Views On Priors Determine Types of Bayesians

- ▶ **Empirical Bayes**: prior distributions are produced from other parts of the data, or possibly from the same data. Results are reported like Frequentists.
- ▶ **Proper Bayes**: prior distributions come from previously compiled evidence, such earlier studies or published work, researcher intuition, or substantive experts. Results are reported without utility or loss functions.
- ▶ **Reference Bayes**: prior distributions are created to influence the posterior as little as mathematically possible (“objective”). Results are reported without utility or loss functions.
- ▶ **Decision-Theoretic Bayes**: prior distributions are from either of the last two sources. Results are presented in a full decision-theoretic framework where utility functions determine decision losses, which are minimized according to different probabilistic criteria.
- ▶ **Bayesians of Convenience**:

Views On Priors Determine Types of Bayesians

- ▶ **Empirical Bayes**: prior distributions are produced from other parts of the data, or possibly from the same data. Results are reported like Frequentists.
- ▶ **Proper Bayes**: prior distributions come from previously compiled evidence, such earlier studies or published work, researcher intuition, or substantive experts. Results are reported without utility or loss functions.
- ▶ **Reference Bayes**: prior distributions are created to influence the posterior as little as mathematically possible (“objective”). Results are reported without utility or loss functions.
- ▶ **Decision-Theoretic Bayes**: prior distributions are from either of the last two sources. Results are presented in a full decision-theoretic framework where utility functions determine decision losses, which are minimized according to different probabilistic criteria.
- ▶ **Bayesians of Convenience**: conjugate diffuse priors or uniform priors on all parameters. Results are reported without utility or loss functions.

Important Observations On Priors

- ▶ Not considering prior distributions can lead to the wrong conclusion: consider alternative probabilities that a randomly chosen US citizen is a member of Congress: $1/300M$, $535/300M$, and 0.5 .
- ▶ Stipulating priors is an overt public statement made to a naturally skeptical scientific audience, so priors cannot be covertly adjusted to “cook” the conclusions.
- ▶ When there are multiple theories or empirical observations, then the associated prior distributions can be used to test the efficacy of the different subsequent posterior distributions.
- ▶ Prior to posterior inference is how scientific knowledge accumulates.
- ▶ With very large data stated priors are typically not important.
- ▶ Producing even multiple priors is a straightforward and easy process.

Spreading the Bayesian Gospel?



- ▶ Preaching: $\pi(\theta|\mathbf{x}) \propto p(\theta)L(\theta|\mathbf{x})$
- ▶ Practicing: $p(\theta) = 1, -\infty < \theta < \infty; p(\theta) \sim \mathcal{N}(0, \sigma^2), \sigma^2 \gg 0;$ or $p(\theta) = -E_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) \right]$

Example: the Beta-Binomial

- X_1, X_2, \dots, X_n iid Bernoulli, $p \sim \text{beta}(A, B)$ prior.
- Standard trick: $Y = \sum_{i=1}^n X_i \sim \text{binomial}(n, p)$.

Example: the Beta-Binomial

- X_1, X_2, \dots, X_n iid Bernoulli, $p \sim \text{beta}(A, B)$ prior.
- Standard trick: $Y = \sum_{i=1}^n X_i \sim \text{binomial}(n, p)$.

$$\begin{aligned}
 f(y, p) &= f(y|p)f(p) \\
 &= \left[\binom{n}{y} p^y (1-p)^{n-y} \right] \times \left[\frac{\Gamma(A+B)}{\Gamma(A)\Gamma(B)} p^{A-1} (1-p)^{B-1} \right] \\
 &= \frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)} p^{y+A-1} (1-p)^{n-y+B-1}
 \end{aligned}$$

Example: the Beta-Binomial

- X_1, X_2, \dots, X_n iid Bernoulli, $p \sim \text{beta}(A, B)$ prior.
- Standard trick: $Y = \sum_{i=1}^n X_i \sim \text{binomial}(n, p)$.

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

$$= \left[\binom{n}{y} p^y (1-p)^{n-y} \right] \times \left[\frac{\Gamma(A+B)}{\Gamma(A)\Gamma(B)} p^{A-1} (1-p)^{B-1} \right]$$

- Joint Distribution:

$$= \frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)} p^{y+A-1} (1-p)^{n-y+B-1}$$

$$f(y) = \int_0^1 \frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)} p^{y+A-1} (1-p)^{n-y+B-1} dp$$

- Marginal Distribution for y :

$$= \frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)} \frac{\Gamma(y+A)\Gamma(n-y+B)}{\Gamma(n+A+B)}$$

Example: the Beta-Binomial, Cont.

- Posterior Distribution for p :

$$\begin{aligned} f(p|y) &= \frac{f(y, p)}{f(y)} = \frac{\frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)}p^{y+A-1}(1-p)^{n-y+B-1}}{\frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)}\frac{\Gamma(y+A)\Gamma(n-y+B)}{\Gamma(n+A+B)}} \\ &= \frac{\Gamma(n + A + B)}{\Gamma(y + A)\Gamma(n - y + B)}p^{(y+A)-1}(1-p)^{(n-y+B)-1} \end{aligned}$$
$$p|y \sim \text{beta}(y + A, n - y + B)$$

Example: the Beta-Binomial, Cont.

- Posterior Distribution for p :

$$\begin{aligned}
 f(p|y) &= \frac{f(y, p)}{f(y)} = \frac{\frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)}p^{y+A-1}(1-p)^{n-y+B-1}}{\frac{\Gamma(n+1)\Gamma(A+B)}{\Gamma(y+1)\Gamma(n-y+1)\Gamma(A)\Gamma(B)}\frac{\Gamma(y+A)\Gamma(n-y+B)}{\Gamma(n+A+B)}} \\
 &= \frac{\Gamma(n+A+B)}{\Gamma(y+A)\Gamma(n-y+B)}p^{(y+A)-1}(1-p)^{(n-y+B)-1} \\
 p|y &\sim \text{beta}(y+A, n-y+B)
 \end{aligned}$$

- An implication:

$$\bar{p} = \frac{(y+A)}{(y+A)+(n-y+B)} = \left[\frac{n}{A+B+n} \right] \left(\frac{y}{n} \right) + \left[\frac{A+B}{A+B+n} \right] \left(\frac{A}{A+B} \right)$$

Example: the Beta-Binomial, Cont.

- The Data (Romney 1999):

Response:	1	1	1	1	0	1	1	0	1	0	1	1
	1	0	1	1	1	1	1	1	0	0	0	1

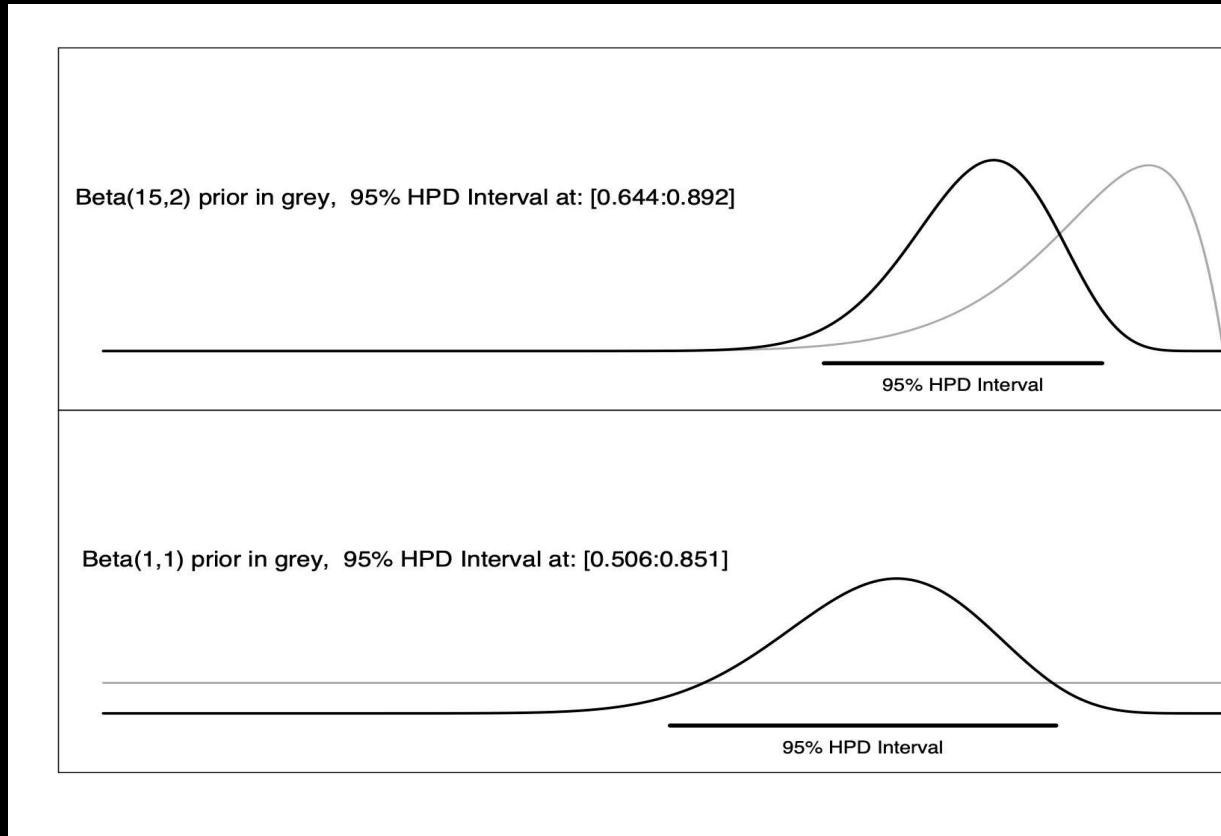
- Two Priors: $\text{BE}(p|15, 2)$, $\text{BE}(p|1, 1)$

- Resulting Posteriors:

$$\text{BE}(\sum x_i + 15, n - \sum x_i + 2) = \text{BE}(32, 9),$$

and $\text{BE}(\sum x_i + 1, n - \sum x_i + 1) = \text{BE}(18, 8)$

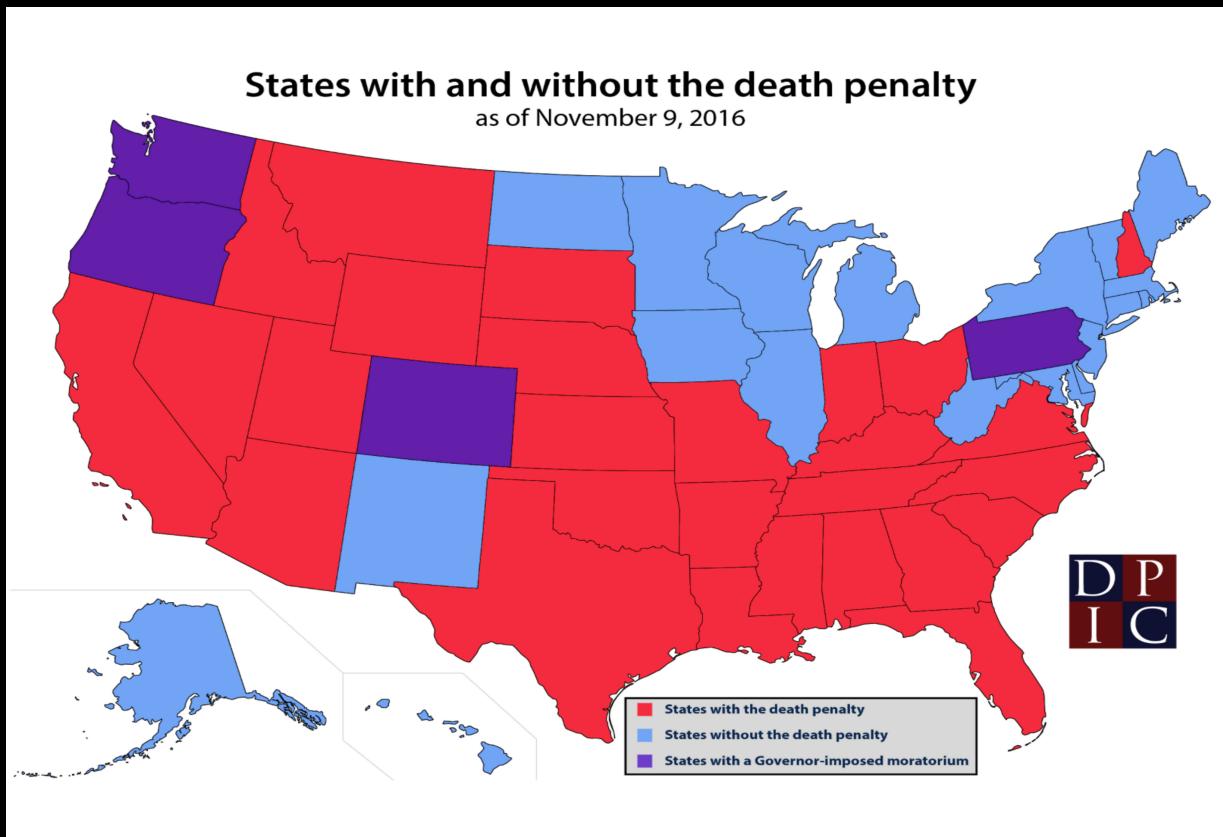
Example: the Beta-Binomial, Cont.



Point Estimates Versus Intervals

- ▶ Bayesians often describe posterior results with point estimates (posterior means, etc.) but generally prefer more full descriptions of the posterior distribution like HPD intervals and quantiles.
- ▶ Also, no point estimator can be absolutely correct and no point estimator can describe the full posterior.
- ▶ So the vast literature on what point estimator is better is not a concern to Bayesians and Bayesians give more information anyways.

Bayesian Tobit Model for Death Penalty Support



Bayesian Tobit Model for Death Penalty Support

- ▶ Use the Tobit model (Tobin 1958) to look at social and political influences on U.S. state decisions to impose the death penalty since the Supreme Court ruled the practice constitutional in *Furman v. Georgia* 1972.
- ▶ Does the ideological, racial and religious makeup, political culture, and urbanization are causal effects for state-level death sentences from 1993 to 1995.
- ▶ The Tobit model is necessary to account for censoring here because 15 states did not have capital punishment provisions on the books in the studied period.



Bayesian Tobit Model for Death Penalty Support

- If \mathbf{z} is a latent outcome variable in this context with the assumptions

$$\mathbf{z} = \mathbf{x}\boldsymbol{\beta} + \mathbf{E}$$

and

$$z_i \sim \mathcal{N}(\mathbf{x}\boldsymbol{\beta}, \sigma^2),$$

then the observed outcome variable is produced according to:

$$y_i = z_i \quad \text{if} \quad z_i > 0,$$

and

$$y_i = 0 \quad \text{if} \quad z_i \leq 0.$$

- The likelihood function is then:

$$L(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) = \prod_{y_i=0} \left[1 - \Phi\left(\frac{x_i\boldsymbol{\beta}}{\sigma}\right) \right] \prod_{y_i>0} (\sigma^{-1}) \exp\left[-\frac{1}{2\sigma^2}(y_i - x_i\boldsymbol{\beta})^2\right].$$

Bayesian Tobit Model for Death Penalty Support

- A flexible parameterization for the priors is given by

$$\boldsymbol{\beta} | \sigma^2 \sim \mathcal{N}(\boldsymbol{\beta}_0, \mathbf{I}\sigma^2 B_0^{-1}) \quad \sigma^2 \sim \mathcal{IG}\left(\frac{\gamma_0}{2}, \frac{\gamma_1}{2}\right)$$

with vector hyperparameter $\boldsymbol{\beta}_0$, scalar hyperparameters B_0 , $\gamma_0 > 2$, $\gamma_1 > 0$, and appropriately sized identity matrix \mathbf{I} .

- Substantial prior flexibility can be achieved with varied levels of these parameters, although values far from those implied by the data will make the Gibbs sampler algorithm run very slowly.

Bayesian Tobit Model for Death Penalty Support

- The resulting joint posterior, $\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X})$, is now analytically *intractable*, even with this basic model:

$$\begin{aligned} \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &= \prod_{y_i=0} \left[1 - \Phi\left(\frac{x_i \boldsymbol{\beta}}{\sigma}\right) \right] \prod_{y_i>0} (\sigma^{-1}) \exp\left[-\frac{1}{2\sigma^2}(y_i - x_i \boldsymbol{\beta})^2\right] \\ &\quad \times (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left[-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)\right] \frac{\left(\frac{\gamma_1}{2}\right)^{\frac{\gamma_0}{2}}}{\Gamma\left(\frac{\gamma_0}{2}\right)} (\sigma^2)^{-(\frac{\gamma_0}{2}+1)} \exp\left[-\left(\frac{\gamma_1}{2}\right)/\sigma^2\right] \end{aligned}$$

- To produce a regression table we now need to solve seven six-dimensional integrals (6 parameters in $\boldsymbol{\beta}$ plus σ^2) to get the marginal posteriors, then seven times two more one-dimensional integrals to get the first two moments for each parameter.
- Better solution: Gibbs sampling (MCMC) which cycles through iterative draws of the full conditional distributions for each model parameter.

Bayesian Tobit Model for Death Penalty Support

- The full conditional distributions for Gibbs sampling are given for the β block, σ^2 , and the individual $z_i|y_i = 0$ as:

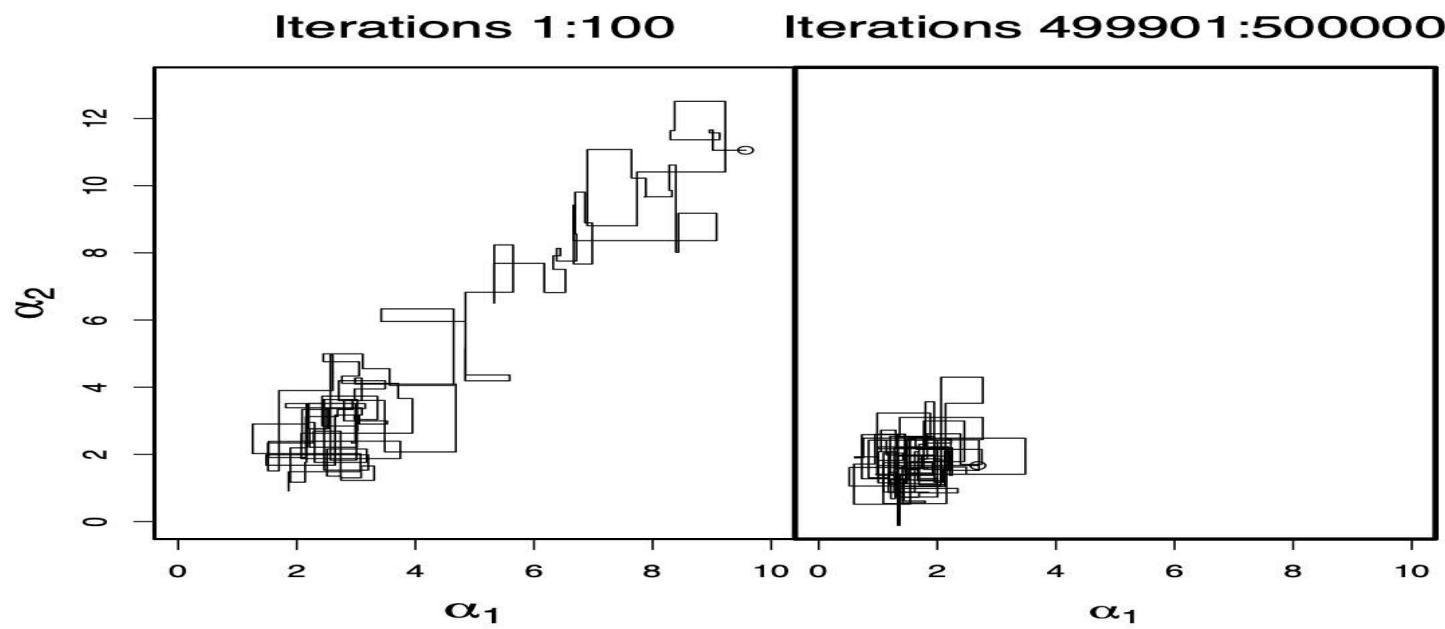
$$\beta|\sigma^2, \mathbf{z}, \mathbf{y}, \mathbf{X} \sim \mathcal{N}\left((B_0 + \mathbf{X}'\mathbf{X})^{-1})(\beta_0 B_0 + \mathbf{X}'\mathbf{z}), (\sigma^{-2}B_0 + \sigma^{-2}\mathbf{X}'\mathbf{X})^{-1}\right)$$

$$\sigma^2|\beta, \mathbf{z}, \mathbf{y}, \mathbf{X} \sim \mathcal{IG}\left(\frac{\gamma_0 + n}{2}, \frac{\gamma_1 + (\mathbf{z} - \mathbf{X}\beta)'(\mathbf{z} - \mathbf{X}\beta)}{2}\right)$$

$$z_i|y_i = 0, \beta, \sigma, \mathbf{X} \sim \mathcal{TN}(\mathbf{X}\beta, \sigma^2)I_{(-\infty, 0)},$$

where $\mathcal{TN}()$ denotes the truncated normal and the indicator function $I_{(-\infty, 0)}$ provides the bounds of truncation.

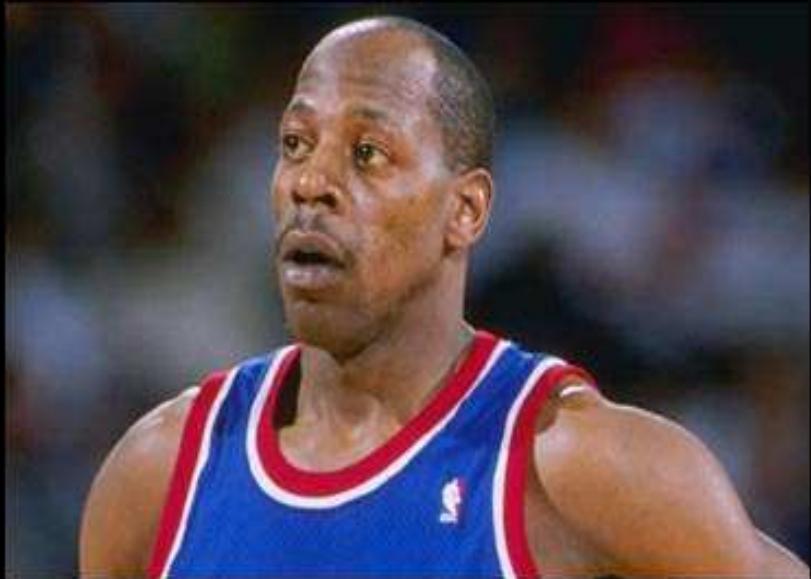
Gibbs Sampling Illustration



Posterior Summary, Tobit Model for Death Penalty Support

	$\bar{\beta}$	σ_{β}
Constant	-6.7600	3.5630
Past Rates	25.5586	8.0697
Political Culture	0.7919	0.1398
Current Opinion	5.9499	1.0805
Ideology	0.2638	1.0961
Murder Rate	0.1800	0.0764

Important Application: Did Vinnie Johnson Have a Hot Hand?



- ▶ On May 5, 1985, the Pistons trailed the Celtics by 87-76 after three periods at home during the playoffs (Eastern Conference Semi-Finals). Then in an amazing scoring display off the bench, Johnson scored 22 of 26 Piston's points in the period to pace an incredible 102-99 win, and tie the series 2-2.
- ▶ His overall shooting rate for 380 games was $\hat{\xi} = 0.43$.
- ▶ Question: is there extra game-to-game extra-binomial variability?

Vinnie Johnson (“microwave”)

- Null: success probability constant across games

$$H_0 : Y_i \stackrel{iid}{\sim} \text{Bin}(n_i, p), \quad i = 1, \dots, 380.$$

- Alternative: success probability has extra-variability

$$H_1 : Y_i \sim \text{Bin}(n_i, p_i) \text{ independently,}$$

where:

$$p_i \sim \text{Beta}(\xi/\omega, (1 - \xi)/\omega),$$

so from beta-binomial $E[Y_i/n_i] = \xi$ since $\frac{\xi/\omega}{\xi/\omega + (1 - \xi)/\omega} = \xi$, ω is an assigned parameter, and the prior on ξ is $\mathcal{U}(0, 1)$.

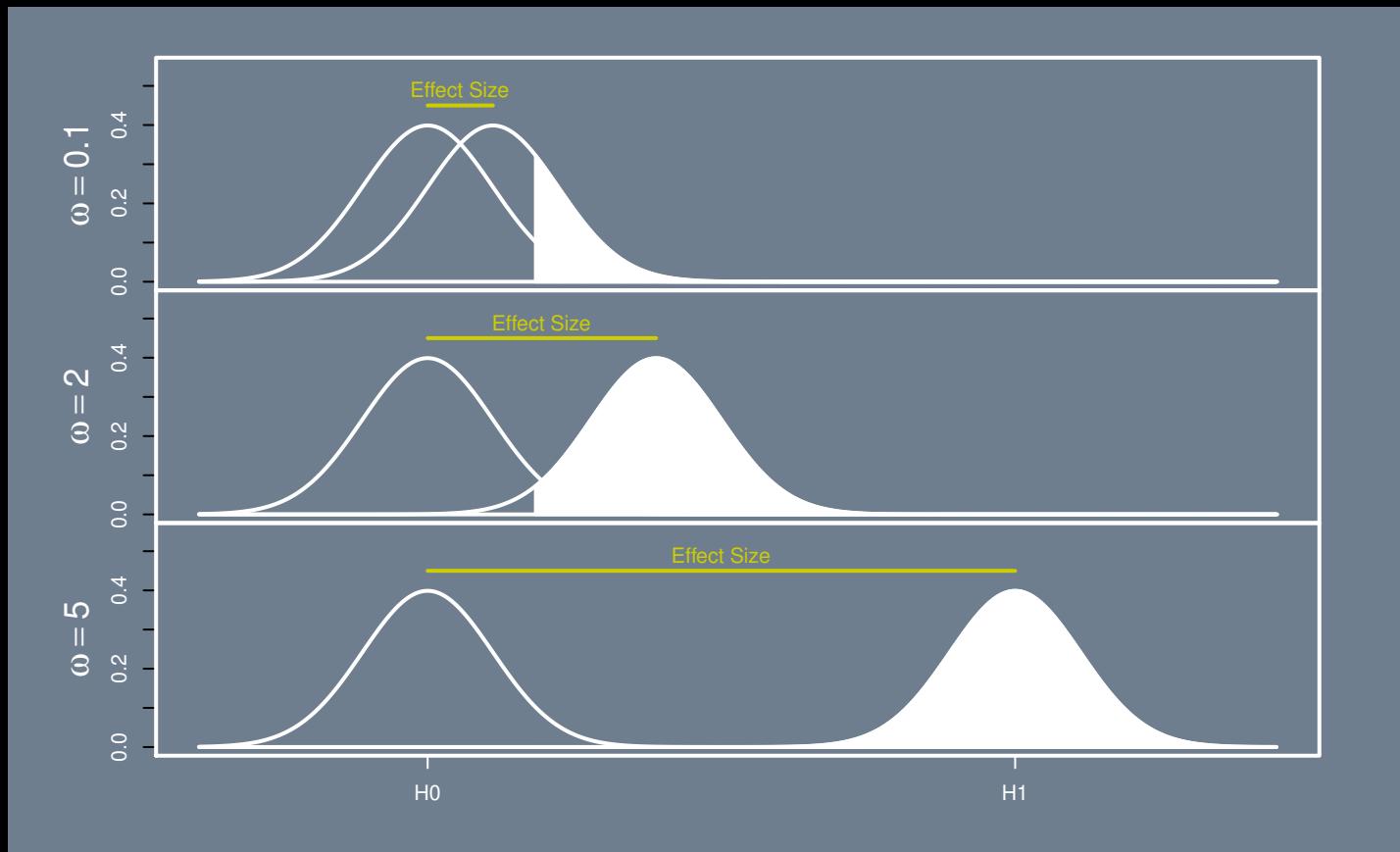
Vinnie Johnson (6'2"; 200 lb.)

- Under H_1 the variance of p_i is large causing “swings” across games inducing the Hot Hand analogy:

$$Var(p_i) = \frac{\left(\frac{\xi}{\omega}\right)\left(\frac{1-\xi}{\omega}\right)}{\left(\frac{\xi}{\omega} + \frac{1-\xi}{\omega}\right)^2 \left(\frac{\xi}{\omega} + \frac{1-\xi}{\omega} + 1\right)} = \frac{\xi(1-\xi)}{1 + \frac{1}{\omega}}$$

- Interpretation of ω parameter: as $\omega \rightarrow 0$, $Var(p_i) \rightarrow 0$, meaning $H_1 \rightarrow H_0$.
- So it is interesting to compare models with different ω values: how far does the data support a difference from the binomial?

Flexible Hypothesis Testing Alternatives, $\alpha = 0.05$, One-Tail



Vinnie Johnson (#15)

- One way to directly test competing models is the Bayes Factor:

$$B_{10}(\mathbf{y}) = \frac{\pi(M_1|\mathbf{y})/p(M_1)}{\pi(M_0|\mathbf{y})/p(M_0)}$$

- Results here:

ω	$Var(p)$	$B_{10}(\mathbf{y})$
0.001	0.007	0.21
0.005	0.035	0.16
0.010	0.049	0.017
0.030	0.085	0.0000003

- So as $\omega \rightarrow 0$, $B_{10}(\mathbf{y}) \rightarrow 1$, and as $\omega \uparrow$, $Var(p) \uparrow$, but $B_{10}(\mathbf{y}) \rightarrow 0$.
- For further exploration the data can be downloaded at
<http://www.stat.washington.edu/raftery/software.html>, and there is a brief discussion in Kass and Raftery (1995).