

PREDICTIVE INFERENCE TOOLS FOR RESEARCHERS

by

Voyze G. Harris III

Copyright © Voyze G. Harris III 2021

A Thesis Submitted to the Faculty of the

STATISTICS AND DATA SCIENCE
GRADUATE INTERDISCIPLINARY PROGRAM

In Partial Fulfillment of the Requirements
For the Degree of

MASTER OF SCIENCE

In the Graduate College

THE UNIVERSITY OF ARIZONA

2021

THE UNIVERSITY OF ARIZONA
GRADUATE COLLEGE

As members of the Master's Committee, we certify that we have read the thesis prepared by Voyze Gabriel Harris III, titled *[Enter Thesis Title]* and recommend that it be accepted as fulfilling the dissertation requirement for the Master's Degree.

Dr. Dean Billheimer

Date: _____

Dr. Edward Bedrick

Date: _____

Dr. Walter Piegorsch

Date: _____

Final approval and acceptance of this thesis is contingent upon the candidate's submission of the final copies of the thesis to the Graduate College.

I hereby certify that I have read this thesis prepared under my direction and recommend that it be accepted as fulfilling the Master's requirement.

Dr. Dean Billheimer
Master's Thesis Committee Chair
Biostatistics

Date: _____



ARIZONA

Contents

1	Thesis Abstract	5
2	Introduction: Predictive Inference	5
2.1	Why Predictive Inference?	5
2.2	The Bayesian Parametric Prediction Format	7
3	Chapter 1: Predictive Problems with Conjugate Priors	9
3.1	Prediction of Future Successes: Beta-Binomial (Geisser p. 73)	9
3.1.1	Derivation	9
3.1.2	R Implementation	11
3.1.3	Example	11
3.2	Survival Time: Exponential-Gamma (Geisser p. 74)	11
3.2.1	Derivation	11
3.2.2	R Implementation	13
3.2.3	Example	13
3.3	Poisson-Gamma Model (Hoff p. 43ff)	15
3.3.1	Derivation	15
3.3.2	R Implementation	19
3.3.3	Example	19
3.4	Normal Observation with Normal-Inverse Gamma Prior	21
3.4.1	One sample	21
3.4.1.1	Derivation	21
3.4.1.2	R Implementation	22
3.4.1.3	Example	23
3.4.2	Two samples	26
3.4.2.1	Derivation	26
3.4.2.2	R Implementation	27
3.4.2.3	Example	27
3.4.3	k samples: Comparing multiple groups	28
3.4.3.1	Derivation	29
3.4.3.2	R Implementation	29
3.4.3.3	Example	30
3.4.3.4	Ranking Treatments	32
4	Chapter 2: Normal Regression with Zellner's g-prior	33
4.1	Least Squares Estimation with Example (Hoff p. 149ff.)	33
4.2	Bayesian Estimation for a Regression Model (Hoff p. 154ff)	37
4.2.1	Derivation	37
4.2.1.1	A semiconjugate prior distribution	37
4.2.1.2	Default and weakly informative prior distributions	38
4.2.2	R Implementation	39
4.2.3	Example	39
5	Conclusion	52

1 Thesis Abstract

An obstacle to widespread employment of Bayesian predictive inference in scientific research is the lack of suitable computing tools. In this thesis I document several established useful models, and provide an applicable set of tools for statisticians. For each of the included models, some basic notes on mathematical derivation are presented, and predictive inference is illustrated with examples. Note that throughout this thesis the terms “predictive inference,” “Bayesian inference,” and “Bayesian prediction” are used interchangeably. For the details of the models and some of the examples we relied primarily on Seymour Geisser’s Predictive Inference: An Introduction (1993) and Peter D. Hoff’s A First Course in Bayesian Predictive Methods (2009).

An R package has been developed, the main purpose of which is to provide the researcher with a means of producing random samples from predictive distributions. For all the models, the package includes random sample generators. For those models with analytical solutions, density and distribution functions are also provided. The standard R naming convention has been adopted: density functions are prefixed with the letter “d,” distribution functions with the letter “p,” and random generation functions with the letter “r.” Also included in all function names is the abbreviation “pred” (for predictive) and an initialism or abbreviation identifying the model itself. For example, the density function for the Beta-Binomial model is named “dpredBB().” The R code for each function is included in Appendix [X-insert link to appendix here](#).

2 Introduction: Predictive Inference

2.1 Why Predictive Inference?

The main purpose of statistics is to predict future events based on observed data. Prediction about meaningful quantities that are relevant to the object of study facilitates scientific progress in multiple ways. Advantages include enhancing scientific reproducibility, enabling corroboration or refutation of current hypotheses through future experimentation, informing decision-making by summarizing quantities of direct interest to the researcher, and shifting the focus of statistical analysis from estimation of hypothetical parameters to statements about concrete observables.

It is not the intent of this thesis to suggest that parametric inference should be abandoned in statistical analyses. Conventional statistical inference techniques are useful for summarizing information about large quantities of data in a handful of usable values, and leveraging such summaries to determine whether a particular problem merits continued attention. Indeed, the scientific discipline of statistics developed along frequentist lines, and the evolution of Bayesian methods has occurred atop that foundation.

Prediction is a means of discriminating between scientific hypotheses. Generally, a model may be judged by the quality of its predictions. Given competing models, the better predictor will be given more weight, and a useful model increases in utility as its predictive capability improves.

To illustrate the potential difference between results from Bayesian prediction and using plug-in estimators, consider the game Pass the Pigs[™], a push-your-luck dice game in which the “dice” are actually rubber pig figures. Two pig dice are thrown, and points

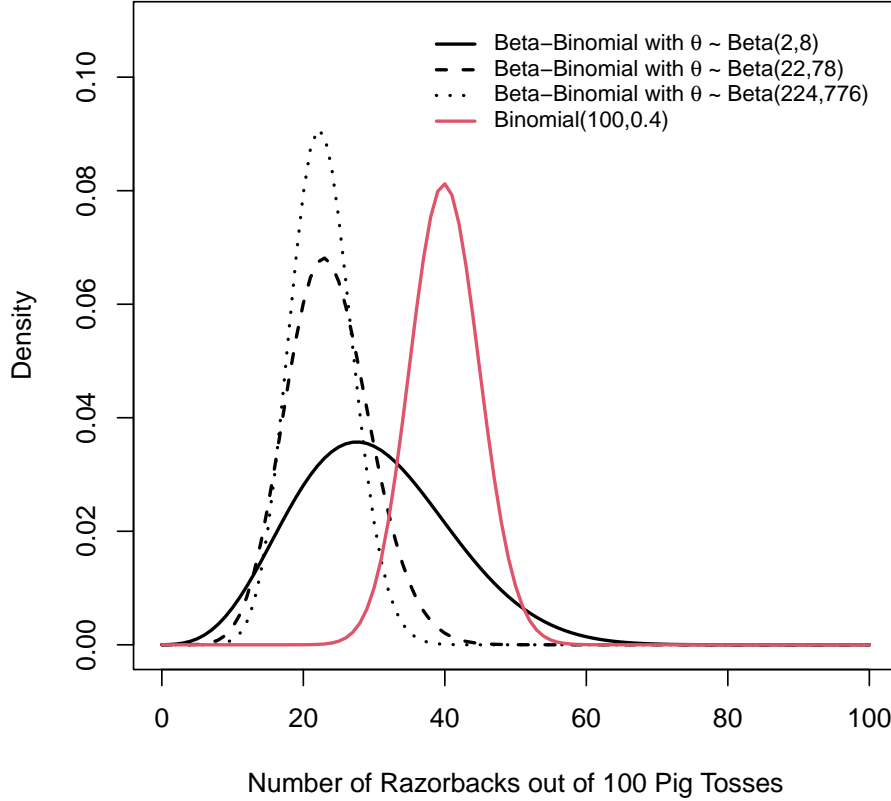
are scored according to the combination of positions in which they come to rest. Details about the game can be found on Wikipedia here: https://en.wikipedia.org/wiki/Pass_the_Pigs

For the purpose of this example, consider the probability of a single pig landing in the “Razorback” position, which occurs when the pig is lying on its back with its legs extended upward. The irregular shape of the pig makes it difficult to assign probabilities to results other than by means of experimentation. Such an experiment was conducted at Duquesne University, and an article describing the experiment as well as Bayesian predictive inference performed on the results appeared in the *Journal of Statistics Education* Volume 14, Number 3, in 2006. The article can be accessed here: <http://jse.amstat.org/v14n3/datasets.kern.html>. Of the 11,954 recorded results for individual pigs, approximately 22.4% were Razorbacks.

Suppose $t = 4$ Razorbacks have been observed out of $N = 10$ tosses of a single pig die, suggesting a straightforward binomial distribution with $\theta = \text{Pr}(\text{Razorback}) = t/N = 0.4$. Taking the Duquesne experiment into consideration, we’ll perform Bayesian prediction using three Beta prior distributions for θ : $\theta \sim \text{Beta}(2, 8)$, $\theta \sim \text{Beta}(22, 78)$, and $\theta \sim \text{Beta}(224, 776)$, and compare these results to predictions obtained from the plug-in estimator $\theta = 0.4$. Any number of prior distributions on θ would satisfy the condition that $E(\theta) \approx 0.224$, suggested by the prior information. The specific choice of a Beta prior is made largely for computational convenience.

In this example the question asked by the researcher is, “For $M = 100$ future observations, how many Razorbacks are predicted?” The density curves in the plot below show the influence of the details of the choice of prior on the location and variance of the predictive distribution. Essentially, each pair of shape parameters (α, β) in the Beta prior reflects the prior knowledge about the results of the Duquesne experiment, with the “weight” given to that knowledge increasing with the shape parameters by orders of magnitude. The choice of shape parameters might be influenced by such things as pig throwing method (perhaps the researcher is throwing them by hand rather than by the carefully controlled method used in the Duquesne experiment, e.g.), or by a need to account for pig-to-pig variation, or anything else the researcher believes introduces a deviation from the events upon which the prior information is based. Notice in the graph and table below that the choice of prior parameters has a significant effect on the variance of the predictive distribution.

iss the Pigs™: Beta–Binomial Prediction vs. Binomial with Plug–in Es



(α, β)	$E(\theta)$	mean(Razorbacks Predicted)	SD(Razorbacks Predicted)
(2,8)	0.2	30.16	10.91
(22,78)	0.22	23.46	5.74
(224,776)	0.224	22.48	4.25

2.2 The Bayesian Parametric Prediction Format

We want to predict future outcomes based on current knowledge. The question is, for observed values $Y_1 = y_1, \dots, Y_n = y_n$, what is likely to be the value of the next observation, $\tilde{Y} = \tilde{y}$? We want to compute $Pr(\tilde{Y}|\theta, y_1, \dots, y_n)$, where θ is a population parameter with a distribution $\pi(\theta)$ based on some prior knowledge or beliefs. Here we are careful to satisfy ourselves that Y_1, \dots, Y_n are *exchangeable*. Exchangeability means that $p(y_1, \dots, y_n) = p(y_{a_1}, \dots, y_{a_n})$ for all permutations $\{a_1, \dots, a_n\}$ of $\{1, \dots, n\}$. In other words the order of the y_i s does not convey any information about their joint density. From de Finetti's representation theorem, then, Y_1, \dots, Y_n are conditionally independent and identically distributed (i.i.d.) given θ (Hoff p. 27), and we can write the joint density

$$p(y_1, \dots, y_n) = \int \left\{ \prod_{i=1}^n p(y_i|\theta) \right\} \pi(\theta) d\theta$$

Get to: We can predict \tilde{Y} using

$$p(\tilde{Y} = \tilde{y} | Y_1 = y_1, \dots, Y_n = y_n) = \int p(\tilde{y} | \theta) p(\theta | y_1, \dots, y_n) d\theta$$

3 Chapter 1: Predictive Problems with Conjugate Priors

[Problems with closed-form solutions. These problems will be what the R package is designed for. Use problems from Geisser, Casella & Berger (Bayesian chapter), other sources. Regression problem—predictive distributions of models that include and exclude some predictor]

write up background for each model: What are they useful for? (2-4 sentence paragraph for each)
need to create “good” examples for Geisser’s models

3.1 Prediction of Future Successes: Beta-Binomial (Geisser p. 73)

3.1.1 Derivation

Let X_i be independent binary variables with $\Pr(X_i = 1) = \theta$, and let $T = \sum X_i$. Then T has probability

$$\binom{N}{t} \theta^t (1 - \theta)^{N-t}.$$

Assume $\theta \sim \text{Beta}(\alpha, \beta)$, so

$$p(\theta) = \frac{\Gamma(\alpha + \beta) \theta^{\alpha-1} (1 - \theta)^{\beta-1}}{\Gamma(\alpha) \Gamma(\beta)}.$$

Then

$$p(\theta | X^{(N)}) = \frac{\Gamma(N + \alpha + \beta) \theta^{t+\alpha-1} (1 - \theta)^{N-t+\beta-1}}{\Gamma(t + \alpha) \Gamma(N - t + \beta)}$$

So for $R = \sum_{i=1}^M X_{N+i}$ we have Beta-Binomial predictive distribution

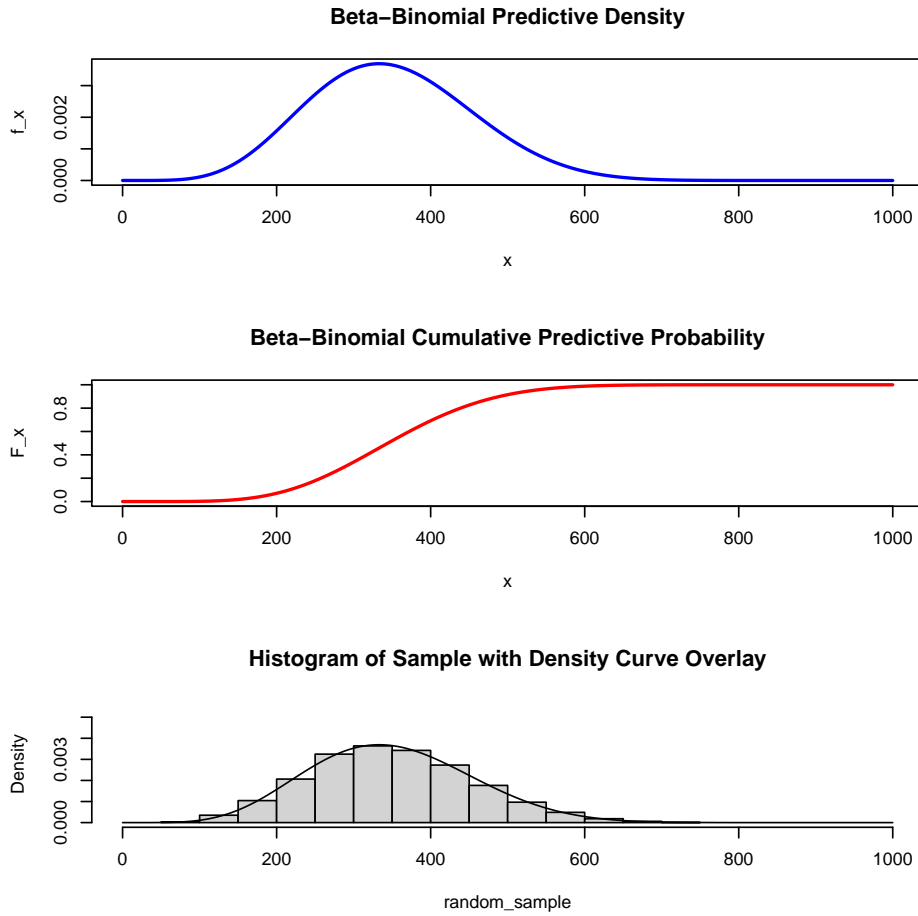
$$\begin{aligned}
\Pr[R = r|t] &= \int \binom{M}{r} \theta^r (1 - \theta)^{M-r} p(\theta|X^{(N)}) d\theta \\
&= \binom{M}{r} \int \theta^r (1 - \theta)^{M-r} \frac{\Gamma(N + \alpha + \beta)}{\Gamma(t + \alpha)\Gamma(N - t + \beta)} \theta^{t+\alpha-1} (1 - \theta)^{N-t+\beta-1} d\theta \\
&= \frac{M!}{r!(M-r)!} \frac{\Gamma(N + \alpha + \beta)}{\Gamma(t + \alpha)\Gamma(N - t + \beta)} \int \theta^{r+t+\alpha-1} (1 - \theta)^{M-r+N-t+\beta-1} d\theta \\
&= \frac{\Gamma(M+1)\Gamma(N + \alpha + \beta)\Gamma(r+t+\alpha)\Gamma(M-r+N-t+\beta)}{\Gamma(r+1)\Gamma(M-r+1)\Gamma(t+\alpha)\Gamma(N-t+\beta)\Gamma(M+N+\alpha+\beta)}
\end{aligned}$$

3.1.2 R Implementation

This result has been used to create “standard” R functions `dpredBB()`, `ppredBB()`, and `rpredBB()` for the Beta-Binomial distribution for density, cumulative probability, and random sampling, respectively (see appendix). These functions are exercised in the following example.

3.1.3 Example

Suppose $t = 5$ successes have been observed out of $N = 10$ binary events, $\alpha = 2$ and $\beta = 8$. For $M = 1000$ future observations, how many successes are predicted? The figures below show the predictive distribution from `dpredBB()`, the cumulative distribution from `ppredBB()`, and a histogram of random draws from `rpredBB()`.



3.2 Survival Time: Exponential-Gamma (Geisser p. 74)

3.2.1 Derivation

Suppose $X^{(N)} = (X^{(d)}, X^{(N-d)})$ where $X^{(d)}$ represents copies fully observed from an exponential survival time density

$$f(x|\theta) = \theta e^{-\theta x}$$

and $X^{(N-d)}$ represents copies censored at x_{d+1}, \dots, x_N , respectively. Hence

$$L(\theta) \propto \theta^d e^{-\theta N\bar{x}}$$

when $N\bar{x} = \sum_{i=1}^N x_i$, as shown below.

The usual exponential likelihood is used for the fully observed copies, whereas for the censored copies we need $\Pr(x > \theta) = 1 - \Pr(x \leq \theta) = 1 - F(x|\theta) = 1 - (1 - e^{-\theta x}) = e^{-\theta x}$. Thus the overall likelihood is

$$L(\theta|x) = \prod_{i=1}^d \theta e^{-\theta x_i} \prod_{i=d+1}^N e^{-\theta x_i} = \theta^d e^{-\theta N\bar{x}}$$

Assuming a $\text{Gamma}(\delta, \gamma)$ prior for θ ,

$$p(\theta) = \frac{\gamma^\delta \theta^{\delta-1} e^{-\gamma\theta}}{\Gamma(\delta)}$$

we obtain the posterior

$$\begin{aligned} p(\theta|X^{(N)}) &= \frac{p(x^{(N)}|\theta) p(\theta)}{\int p(X^{(N)}|\theta) p(\theta) d\theta} \\ &= \frac{\theta^d e^{-\theta N\bar{x}} \cdot \frac{\gamma^\delta \theta^{\delta-1} e^{-\gamma\theta}}{\Gamma(\delta)}}{\int \left(\theta^d e^{-\theta N\bar{x}} \cdot \frac{\gamma^\delta \theta^{\delta-1} e^{-\gamma\theta}}{\Gamma(\delta)} \right) d\theta} \\ &= \frac{\cancel{\frac{\gamma^\delta}{\Gamma(\delta)}} (\theta^{d+\delta-1} e^{-\theta(\gamma+N\bar{x})})}{\cancel{\frac{\gamma^\delta}{\Gamma(\delta)}} \int (\theta^{d+\delta-1} e^{-\theta(\gamma+N\bar{x})}) d\theta} \\ &= \frac{\frac{(\gamma+N\bar{x})^{d+\delta}}{\Gamma(d+\delta)} (\theta^{d+\delta-1} e^{-\theta(\gamma+N\bar{x})})}{\cancel{\frac{(\gamma+N\bar{x})^{d+\delta}}{\Gamma(d+\delta)}} \int \cancel{(\theta^{d+\delta-1} e^{-\theta(\gamma+N\bar{x})})} d\theta} \\ &= \frac{(\gamma+N\bar{x})^{d+\delta} \theta^{d+\delta-1} e^{-\theta(\gamma+N\bar{x})}}{\Gamma(d+\delta)} \end{aligned}$$

with the $\text{Gamma}(d+\delta, \gamma+N\bar{x})$ density in the next to last step integrating to 1.

Thus the survival time predictive probability is

$$\begin{aligned}
P(X = x|\theta, X^{(N)}) &= \int p(\theta|X^{(N)}) p(x|\theta) d\theta \\
&= \int \frac{(\gamma + N\bar{x})^{d+\delta} \theta^{d+\delta-1} e^{-\theta(\gamma+N\bar{x})}}{\Gamma(d+\delta)} \cdot \theta e^{-\theta x} d\theta \\
&= (d+\delta)(\gamma + N\bar{x})^{d+\delta} \int \frac{\theta^{(d+\delta+1)-1} e^{-\theta(\gamma+N\bar{x}+x)}}{(d+\delta)\Gamma(d+\delta)} d\theta \\
&= \frac{(d+\delta)(\gamma + N\bar{x})^{d+\delta}}{(\gamma + N\bar{x} + x)^{d+\delta+1}} \int \frac{(\gamma + N\bar{x} + x)^{d+\delta+1} \theta^{(d+\delta+1)-1} e^{-\theta(\gamma+N\bar{x}+x)}}{\Gamma(d+\delta+1)} d\theta \\
&= \frac{(d+\delta)(\gamma + N\bar{x})^{d+\delta}}{(\gamma + N\bar{x} + x)^{d+\delta+1}}
\end{aligned}$$

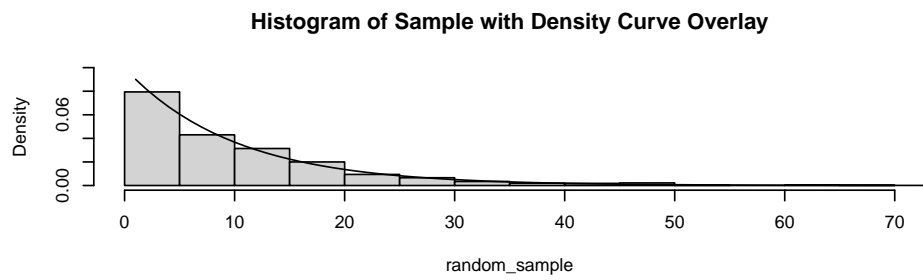
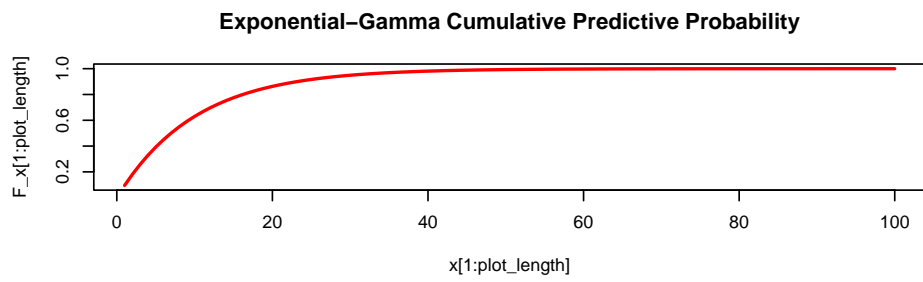
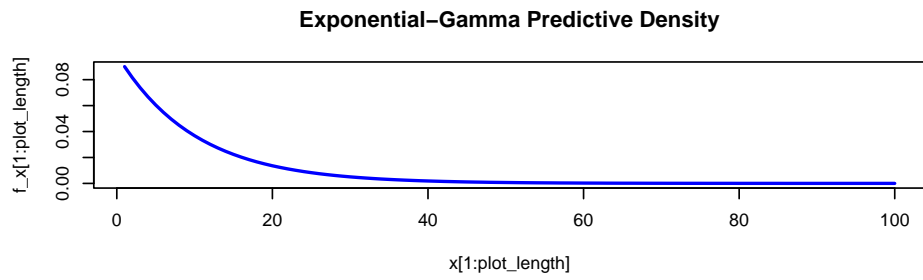
(simplifying by constructing a $\text{Gamma}(d + \delta + 1, \gamma + N\bar{x} + x)$ density in the final integrand.)

3.2.2 R Implementation

This result has been used to create standard format R functions `dpredEG()`, `ppredEG()`, and `rpredEG()` for the Gamma-Exponential distribution for density, cumulative probability, and random sampling, respectively (see appendix). These functions are exercised in the following example.

3.2.3 Example

Suppose $d = 800$ out of $N = 1000$ copies have been observed, and the remaining 200 censored. Say $\delta = 20$, $\gamma = 5$, and we are interested in the number of survivors out of $M = 1000$ future observations. The figures below illustrate the predictive probability using `dpredEG()` and `rpredEG()`, along with a histogram of a random sample taken using `rpredEG()`.



3.3 Poisson-Gamma Model (Hoff p. 43ff)

3.3.1 Derivation

[using Hoff's notation and variable names below. Should I convert this to Geisser's $x^{(N)}, x_{(M)}$ convention for uniformity throughout my thesis?]

Suppose $Y_1, \dots, Y_n | \theta \stackrel{i.i.d.}{\sim} \text{Poisson}(\theta)$ with Gamma prior $\theta \sim \text{Gamma}(\alpha, \beta)$. That is,

$$\begin{aligned} P(Y_1 = y_1, \dots, Y_n = y_n | \theta) &= \prod_{i=1}^n p(y_i | \theta) \\ &= \prod_{i=1}^n \frac{1}{y_i!} \theta^{y_i} e^{-\theta} \\ &= \left(\prod_{i=1}^n \frac{1}{y_i!} \right) \theta^{\sum y_i} e^{-n\theta} \\ &= c(y_1, \dots, y_n) \theta^{\sum y_i} e^{-n\theta} \end{aligned}$$

and

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \theta, \alpha, \beta > 0.$$

Then we have posterior distribution

$$\begin{aligned} p(\theta | y_1, \dots, y_n) &= \frac{p(y_1, \dots, y_n | \theta) p(\theta)}{\int_{\theta} p(y_1, \dots, y_n | \theta) p(\theta)} \\ &= \frac{p(y_1, \dots, y_n | \theta) p(\theta)}{p(y_1, \dots, y_n)} \\ &= \frac{1}{p(y_1, \dots, y_n)} \theta^{\sum y_i} e^{-n\theta} \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} \\ &= C(y_1, \dots, y_n, \alpha, \beta) \theta^{\alpha + \sum y_i - 1} e^{-(\beta + n)\theta} \\ &\sim \text{Gamma}\left(\alpha + \sum y_i, \beta + n\right). \end{aligned}$$

Here

$$\begin{aligned}
C(y_1, \dots, y_n, \alpha, \beta) &= \frac{1}{p(y_1, \dots, y_n)} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \\
&= \frac{1}{\int_\theta p(y_1, \dots, y_n | \theta) p(\theta)} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} \\
&= \frac{1}{\int_\theta \left(\prod \frac{1}{y_i!} \right) \theta^{\sum y_i} e^{-n\theta} \left(\frac{\beta^\alpha}{\Gamma(\alpha)} \right) \theta^{\alpha-1} e^{-\beta\theta} \cancel{\left(\frac{\beta^\alpha}{\Gamma(\alpha)} \right)}} \cdot \left(\frac{\beta^\alpha}{\Gamma(\alpha)} \right) \\
&= \frac{1}{\left(\prod \frac{1}{y_i!} \right) \frac{\Gamma(\alpha + \sum y_i)}{(\beta + n)^{\alpha + \sum y_i}} \int_\theta \frac{(\beta + n)^{\alpha + \sum y_i}}{\Gamma(\alpha + \sum y_i)} \theta^{\sum y_i + \alpha - 1} e^{-(\beta + n)\theta}} \\
&= \frac{\prod_{i=1}^n y_i! (\beta + n)^{\alpha + \sum y_i}}{\Gamma(\alpha + \sum y_i)}
\end{aligned}$$

Call this constant C_n (for n observations).

Note that an additional observation $y_{n+1} = \tilde{y}$ the constant becomes

$$C_{n+1} = \frac{\prod_{i=1}^{n+1} y_i! (\beta + n + 1)^{\alpha + \sum_{i=1}^{n+1} y_i}}{\Gamma(\alpha + \sum_{i=1}^{n+1} y_i)}.$$

Also note that the marginal joint distribution of k observations is

$$p(\tilde{y} | y_1, \dots, y_k) = \frac{1}{C_k} \frac{\beta^\alpha}{\Gamma(\alpha)}.$$

For future observation \tilde{y} , then, we compute predictive distribution

$$\begin{aligned}
p(\tilde{y}|y_1, \dots, y_n) &= \frac{p(y_1, \dots, y_n, \tilde{y})}{p(y_1, \dots, y_n)} = \frac{p(y_1, \dots, y_{n+1})}{p(y_1, \dots, y_n)} = \frac{\frac{1}{C_{n+1}} \frac{\beta^\alpha}{\Gamma(\alpha)}}{\frac{1}{C_n} \frac{\beta^\alpha}{\Gamma(\alpha)}} = \frac{C_n}{C_{n+1}} \\
&= \frac{\prod_{i=1}^n y_i! (\beta + n)^{\alpha + \sum_{i=1}^n y_i}}{\Gamma(\alpha + \sum_{i=1}^n y_i)} \\
&= \frac{\prod_{i=1}^{n+1} y_i! (\beta + n + 1)^{\alpha + \sum_{i=1}^{n+1} y_i}}{\Gamma(\alpha + \sum_{i=1}^{n+1} y_i)} \\
&= \frac{\Gamma(\alpha + \sum_{i=1}^{n+1} y_i) (\beta + n)^{\alpha + \sum_{i=1}^n y_i}}{(y_{n+1}!) \Gamma(\alpha + \sum_{i=1}^n y_i) (\beta + n + 1)^{\alpha + \sum_{i=1}^{n+1} y_i}} \\
&= \frac{\Gamma(\alpha + \sum_{i=1}^n y_i + \tilde{y}) (\beta + n)^{\alpha + \sum_{i=1}^n y_i}}{(\tilde{y}!) \Gamma(\alpha + \sum_{i=1}^n y_i) (\beta + n + 1)^{\alpha + \sum_{i=1}^n y_i + \tilde{y}}} \\
&= \frac{\Gamma(\alpha + \sum y_i + \tilde{y})}{\Gamma(\tilde{y} + 1) \Gamma(\alpha + \sum y_i)} \cdot \left(\frac{\beta + n}{\beta + n + 1} \right)^{\alpha + \sum y_i} \cdot \left(\frac{1}{\beta + n + 1} \right)^{\tilde{y}}
\end{aligned}$$

This is a negative binomial distribution: $\tilde{y} \sim NB(\alpha + \sum y_i, \beta + n)$, for which

$$\begin{aligned}
E[\tilde{Y}|y_1, \dots, y_n] &= \frac{a + \sum y_i}{b + n} = E[\theta|y_1, \dots, y_n]; \\
\text{Var}[\tilde{Y}|y_1, \dots, y_n] &= \frac{a + \sum y_i}{b + n} \frac{b + n + 1}{b + n} \\
&= \text{Var}[\theta|y_1, \dots, y_n] \times (b + n + 1) \\
&= E[\theta|y_1, \dots, y_n] \times \frac{b + n + 1}{b + n}
\end{aligned}$$

[Showing here that it is indeed a NB distribution]

$$\theta \sim NB(\alpha, \beta) \Rightarrow p(\theta) = \binom{\theta + \alpha - 1}{\alpha - 1} \left(\frac{\beta}{\beta + 1} \right)^\alpha \left(\frac{1}{\beta + 1} \right)^\theta$$

so

$$\begin{aligned}
\tilde{y} \sim NB\left(\alpha + \sum y_i, \beta + n\right) &\Rightarrow p(\tilde{y}) = \binom{\tilde{y} + \alpha + \sum y_i - 1}{\alpha + \sum y_i - 1} \left(\frac{\beta + n}{\beta + n + 1}\right)^{\alpha + \sum y_i} \left(\frac{1}{\beta + n + 1}\right)^{\tilde{y}} \\
&= \frac{(\alpha + \sum y_i + \tilde{y} - 1)!}{(\alpha + \sum y_i - 1)! (\tilde{y})!} \left(\frac{\beta + n}{\beta + n + 1}\right)^{\alpha + \sum y_i} \left(\frac{1}{\beta + n + 1}\right)^{\tilde{y}} \\
&= \frac{\Gamma(\alpha + \sum y_i + \tilde{y})}{\Gamma(\alpha + \sum y_i) \Gamma(\tilde{y} + 1)} \left(\frac{\beta + n}{\beta + n + 1}\right)^{\alpha + \sum y_i} \left(\frac{1}{\beta + n + 1}\right)^{\tilde{y}}
\end{aligned}$$

[This is the result in Hoff. The straightforward derivation below is off by a constant multiple. Need to figure out what went awry.]

$$\begin{aligned}
p(\tilde{y}|y_1, \dots, y_n) &= \int_0^\infty p(\tilde{y}|\theta, y_1, \dots, y_n) p(\theta|y_1, \dots, y_n) d\theta \\
&= \int p(\tilde{y}|\theta) p(\theta|y_1, \dots, y_n) d\theta \\
&= C \int \left(\frac{1}{\tilde{y}!} \theta^{\tilde{y}} e^{-\theta}\right) \theta^{\alpha + \sum y_i - 1} e^{-(\beta + n)\theta} d\theta \\
&= \frac{C}{\tilde{y}!} \int \theta^{\tilde{y} + \alpha + \sum y_i - 1} e^{-(\beta + n + 1)\theta} d\theta \\
&= \frac{C \Gamma(\tilde{y} + \alpha + \sum y_i)}{\Gamma(\tilde{y} + 1) (\beta + n + 1)^{\tilde{y} + \alpha + \sum y_i}} \int \frac{(\beta + n + 1)^{\tilde{y} + \alpha + \sum y_i}}{\Gamma(\tilde{y} + \alpha + \sum y_i)} \theta^{\tilde{y} + \alpha + \sum y_i - 1} e^{-(\beta + n + 1)\theta} d\theta \\
&= C \cdot \frac{\Gamma(\tilde{y} + \alpha + \sum y_i)}{\Gamma(\tilde{y} + 1) (\beta + n + 1)^{\tilde{y} + \alpha + \sum y_i}} \\
&= \frac{\prod_{i=1}^n y_i! (\beta + n)^{\alpha + \sum y_i}}{\Gamma(\alpha + \sum y_i)} \cdot \frac{\Gamma(\tilde{y} + \alpha + \sum y_i)}{\Gamma(\tilde{y} + 1) (\beta + n + 1)^{\tilde{y} + \alpha + \sum y_i}} \\
&= \prod_{i=1}^n y_i! \cdot \frac{\Gamma(\tilde{y} + \alpha + \sum y_i)}{\Gamma(\tilde{y} + 1) \Gamma(\alpha + \sum y_i)} \cdot \left(\frac{\beta + n}{\beta + n + 1}\right)^{\alpha + \sum y_i} \cdot \left(\frac{1}{\beta + n + 1}\right)^{\tilde{y}}
\end{aligned}$$

Hoff p.47:

- b is interpreted as the number of prior observations
- a is interpreted as the sum of counts from b prior observations

Hoff p. 49 (Birth rate example): $a = 2, b = 1$.

3.3.2 R Implementation

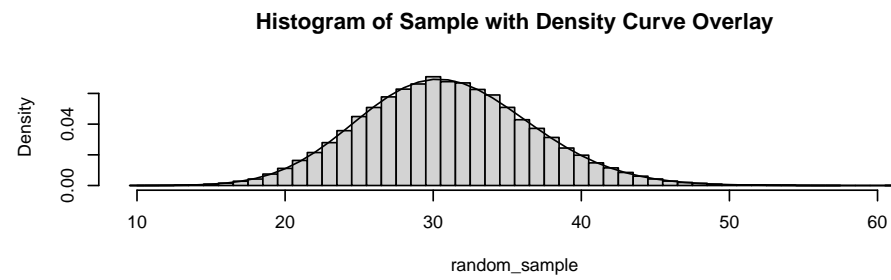
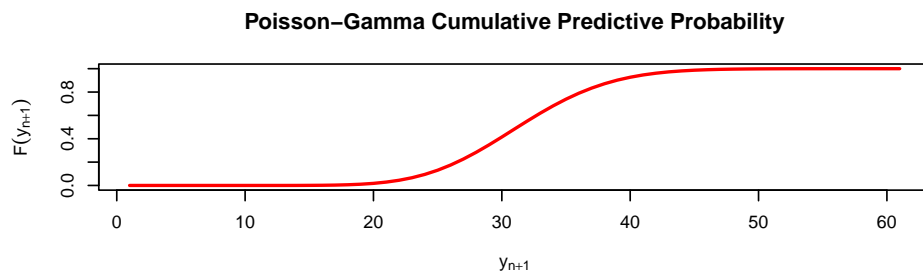
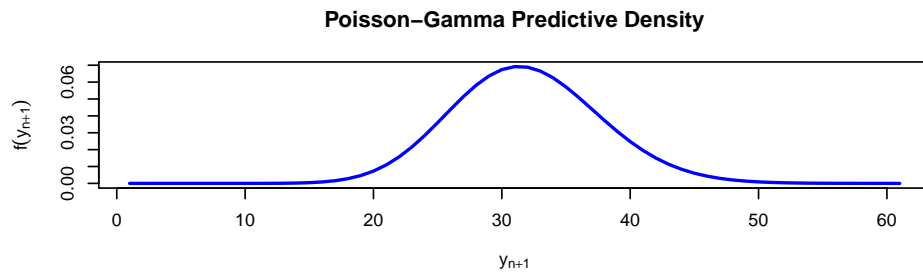
This result has been used to create standard format R functions `dpredPG()`, `ppredPG()`, and `rpredPG()` for the Poisson-Gamma distribution for density, cumulative probability, and random sampling, respectively (see appendix). These functions are exercised in the following example.

Developing the random sample function `rpredPG()`: I need to establish the support of the predictive distribution f_x from which to sample. the `uniroot()` function is not working because it keeps feeding non-integer values to `dnbinom()`. Strategy: a modified bisection method as follows:

1. set a desired tolerance ϵ .
2. Find the expected value E_x (closed formula, see above).
3. Step to the right of E_x by whole integers, in the sequence $E_x + \{1, 2, 4, \dots, 2^n\}$, stopping at $U = f_x(E_x + 2^n) < 0$. This is the upper bound for the bisection method.
4. Bisect the interval, rounding to the nearest integer. Call the resulting mid-interval number B .
5. If B is positive, test whether $0 \leq f_x(B) \leq \epsilon$. If so, DONE. If not:
6. Establish new interval, choosing endpoints from E_x , B , and U so that the interval straddles 0, and repeat the steps until the condition in step 5 is reached.

3.3.3 Example

Suppose we have 10 prior observations with counts 27, 79, 21, 100, 8, 4, 37, 15, 3, 97. Let $\alpha = 11$ and $\beta = 3$. For $\tilde{y} = 1 : 100$ possible future occurrences, the figures below show the predictive distribution from `dpredPG()`, the cumulative distribution from `ppredPG()`, and a histogram of random draws from `rpredPG()`.



3.4 Normal Observation with Normal-Inverse Gamma Prior

also created `rpredNormIG2()` and `rpredNormIGk()` for samples comparing 2 groups and k groups, respectively. I need to look into combining them into one function. maybe do away with `dpredNormIG()` and `ppredNormIG()`?

3.4.1 One sample

3.4.1.1 Derivation [Hoff p. 69ff]

Let $\{Y_1, \dots, Y_n | \theta, \sigma^2\} \stackrel{i.i.d.}{\sim} N(\theta, \sigma^2)$. Then the joint sampling density is

$$\begin{aligned} p(y_1, \dots, y_n | \theta, \sigma^2) &= \prod_{i=1}^n p(y_i | \theta, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{y_i - \theta}{\sigma}\right)^2} \\ &= (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2}\sum_{i=1}^n \left(\frac{y_i - \theta}{\sigma}\right)^2}. \end{aligned}$$

Following Hoff (p. 74ff), for joint inference on both θ and σ , assume priors

$$\begin{aligned} \frac{1}{\sigma^2} &\sim \text{gamma}(\nu_0/2, \nu_0\sigma_0^2/2) \\ \theta | \sigma^2 &\sim \text{normal}(\mu_0, \sigma^2/\kappa_0) \end{aligned}$$

where (σ_0^2, ν_0) are the sample variance and sample size of prior observations, and (μ_0, κ_0) are the sample mean and sample size of prior observations.

Note: μ_0 , κ_0 , ν_0 , and σ_0^2 come from prior knowledge. [in the Hoff example (Midge Wing Length), κ_0 and ν_0 are both set to 1 so that “our prior distributions are only weakly centered around these estimates from other populations.”]

From this we derive joint posterior distribution

$$\begin{aligned} \{\theta | y_1, \dots, y_n, \sigma^2\} &\sim \text{normal}(\mu_n, \sigma^2/\kappa_n) \\ \{\sigma^2 | y_1, \dots, y_n\} &\sim \text{inverse-gamma}(\nu_n/2, \sigma_n^2\nu_n/2). \end{aligned}$$

where

$$\kappa_n = \kappa_0 + n$$

$$\mu_n = \frac{\kappa_0 \mu_0 + n \bar{y}}{\kappa_n}$$

$$\nu_n = \nu_0 + n$$

$$\sigma_n^2 = \frac{1}{\nu_n} \left[\nu_0 \sigma_0^2 + (n-1) s^2 + \frac{\kappa_0 n}{\kappa_n} (\bar{y} - \mu_0)^2 \right].$$

Here $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is the sample mean and $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$ is the sample variance.

From the joint posterior distribution we generate marginal samples by means of the Monte Carlo method (Hoff, p. 77):

$$\begin{aligned} \sigma^{2(1)} &\sim \text{inverse-gamma}(\nu_n/2, \sigma_n^2 \nu_n/2), & \theta^{(1)} &\sim \text{normal}(\mu_n, \sigma^{2(1)}/\kappa_n) \\ &\vdots & &\vdots \\ \sigma^{2(S)} &\sim \text{inverse-gamma}(\nu_n/2, \sigma_n^2 \nu_n/2), & \theta^{(S)} &\sim \text{normal}(\mu_n, \sigma^{2(S)}/\kappa_n) \end{aligned}$$

For prediction of future $\tilde{y}|y_1, \dots, y_n, \theta, \sigma^2$, generate $\tilde{y}_i \sim \text{normal}(\theta^{(i)}, \sigma^{2(i)})$.

For prediction without the influence of any previous knowledge (Hoff p. 79), we can employ Jeffreys prior $\tilde{p}(\theta, \sigma^2) = 1/\sigma^2$. This leads to the same conditional distribution for θ but a $\text{gamma}(\frac{n-1}{2}, \frac{1}{2} \sum (y_i - \bar{y})^2)$ distribution for $1/\sigma^2$. This joint posterior distribution can be used to predict future \tilde{y} by first drawing θ, σ^2 and then simulating $\tilde{y} \sim \text{normal}(\theta, \sigma^2)$. Alternatively, the joint posterior can be integrated to show that

$$\frac{\theta - \bar{y}}{s/\sqrt{n}} | y_1, \dots, y_n \sim t_{n-1}.$$

The resulting predictive distribution for \tilde{y} is a t-distribution with location \bar{y} and scale $s\sqrt{1 + 1/n}$ and $n - 1$ degrees of freedom (Gelman et. al. p. 66).

3.4.1.2 R Implementation Standard format R functions `dpredNormIG()`, `ppredNormIG()`, and `rpredNormIG()` have been created for the Normal-Inverse Gamma distribution for density, cumulative probability, and random sampling, respectively (see appendix). These functions all include options for implementation with or without previous knowledge as desired. If Jeffreys prior is used, the functions simply implement R's Student's t-distribution functions `rt()`, `dt()`, and `pt()`, applying the location and scale parameters as described above. For predictions using previous knowledge, the functions work as follows: For the random sampler `rpredNormIG()`, the Monte-Carlo method described above

is directly employed. The predictive density and cumulative predictive density functions (`dpredNormIG()` and `ppredNormID()`, respectively) depend on the random sample. `ppredNormIG()` utilizes the empirical cumulative density function `ecdf()` from R's `stats` package. `dpredNormIG()` utilizes a Kernel Density Estimation (KDE) method and R's built-in `density()` function. The KDE is computed by definition, using a normal kernel:

$$\hat{f}_K(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right),$$

where

X_i is the random sample generated using `rpredNormIG()`

K is `Normal(0,1)`

h is the bandwidth from R's `density()` function (that is, $h = \text{density}(X_i)\$bw$)

These functions are exercised in the following example.

3.4.1.3 Example *Example (Hoff p. 72ff, using data from Grogan and Wirth (1981)): Midge wing length*

Grogan and Wirth (1981) provide 9 measurements of midge wing length, in millimeters: $y = \{1.64, 1.7, 1.72, 1.74, 1.82, 1.82, 1.82, 1.90, 2.08\}$. Previous studies suggest values $\mu_0 = 1.9$ and $\sigma_0^2 = 0.01$. We choose $\kappa_0 = \nu_0 = 1$ “...so that our prior distributions are only weakly centered around these estimates from other populations” (Hoff p. 76). We compute

$$\bar{y} = 1.804$$

$$\text{var}(y) = 0.0169$$

$$\kappa_n = 1 + 9 = 10$$

$$\mu_n = \frac{1 \cdot 1.9 + 9 \cdot 1.804}{10} = 1.814$$

$$\nu_n = 1 + 9 = 10$$

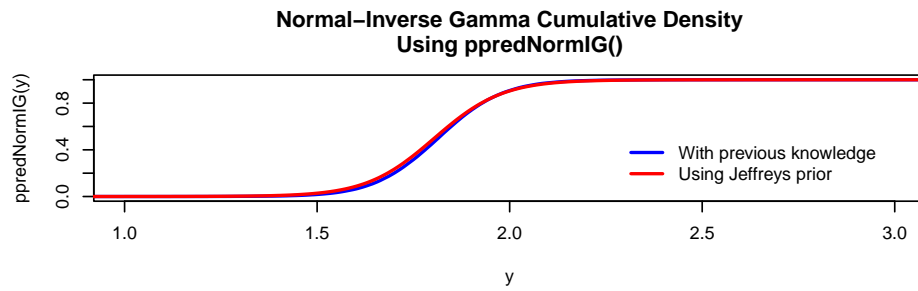
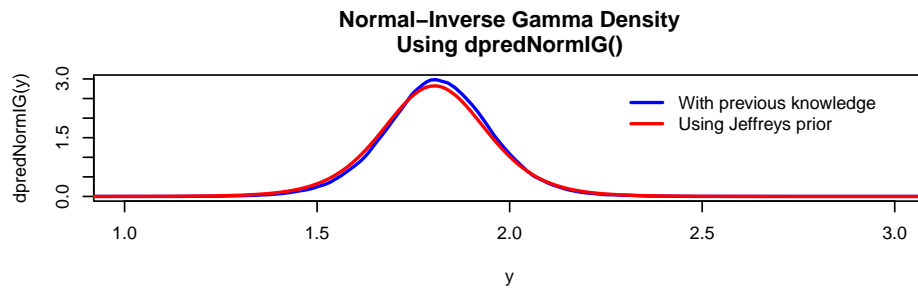
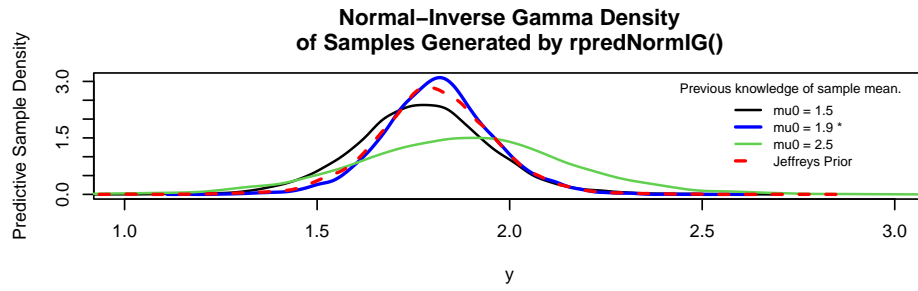
$$\sigma_n^2 = \frac{1}{10} \left[1 \cdot 0.01 + (9 - 1) \cdot 0.0169 + \frac{1 \cdot 9}{10} (1.804 - 1)^2 \right] = 0.0153$$

Thus $\nu_n/2 = 5$ and $\nu_n \sigma_n^2/2 = 0.7662$ and we have posteriors

$$\{\theta|y_1, \dots, y_n, \sigma^2\} \sim \text{normal}(1.814, \sigma^2/10)$$

$$\{\sigma^2|y_1, \dots, y_n\} \sim \text{inverse-gamma}(5, 0.7662)$$

The plot below illustrates the influence of previous knowledge of the population mean, and compares to the predictions resulting from Jeffreys prior.



3.4.2 Two samples

3.4.2.1 Derivation For a Bayesian analysis comparing two groups we use the following sampling model (Hoff p. 127):

$$\begin{aligned} Y_{i,1} &= \mu + \delta + \epsilon_{i,1} \\ Y_{i,2} &= \mu - \delta + \epsilon_{i,2} \\ \{\epsilon_{i,j}\} &\sim \text{i.i.d. normal}(0, \sigma^2). \end{aligned}$$

Letting $\theta_1 = \mu + \delta$ and $\theta_2 = \mu - \delta$ we see that $\delta = (\theta_1 - \theta_2)/2$ is half the population difference in means, and $\mu = (\theta_1 + \theta_2)/2$ is the pooled average. We'll assume conjugate prior distributions

$$\begin{aligned} p(\mu, \delta, \sigma^2) &= p(\mu) \times p(\delta) \times p(\sigma^2) \\ \mu &\sim \text{normal}(\mu_0, \gamma_0^2) \\ \delta &\sim \text{normal}(\delta_0, \tau_0^2) \\ \sigma^2 &\sim \text{inverse-gamma}(\nu_0/2, \nu_0\sigma_0^2/2), \end{aligned}$$

where ν_0 as before is the assumed prior sample size. The full conditional distributions follow:

$$\{\mu | \mathbf{y}_1, \mathbf{y}_2, \delta, \sigma^2\} \sim \text{normal}(\mu_n, \gamma_n^2), \text{ where}$$

$$\mu_n = \gamma_n^2 \times \left[\frac{\mu_0}{\gamma_0^2} + \frac{\sum_{i=1}^{n_1} (y_{i,1} - \delta) + \sum_{i=1}^{n_2} (y_{i,2} + \delta)}{\sigma^2} \right]$$

$$\gamma_n^2 = \left[\frac{1}{\gamma_0^2} + \frac{(n_1 + n_2)}{\sigma^2} \right]^{-1}$$

$$\{\delta | \mathbf{y}_1, \mathbf{y}_2, \mu, \sigma^2\} \sim \text{normal}(\delta_n, \tau_n^2), \text{ where}$$

$$\delta_n = \tau_n^2 \times \left[\frac{\delta_0}{\tau_0^2} + \frac{\sum_{i=1}^{n_1} (y_{i,1} - \mu) - \sum_{i=1}^{n_2} (y_{i,2} - \mu)}{\sigma^2} \right]$$

$$\tau_n^2 = \left[\frac{1}{\tau_0^2} + \frac{(n_1 + n_2)}{\sigma^2} \right]^{-1}$$

$$\{\sigma^2 | \mathbf{y}_1, \mathbf{y}_2, \mu, \delta\} \sim \text{inverse-gamma}\left(\frac{\nu_n}{2}, \frac{\nu_n \sigma_n^2}{2}\right), \text{ where}$$

$$\nu_n = \nu_0 + n_1 + n_2$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + \sum_{i=1}^{n_1} (y_{i,1} - [\mu + \delta])^2 + \sum_{i=1}^{n_2} (y_{i,2} - [\mu - \delta])^2$$

3.4.2.2 R Implementation The standard format R function `rpredNormIG2()` implements a Gibbs sampler to approximate the posterior distribution $p(\mu, \delta, \sigma^2 | \mathbf{y}_1, \mathbf{y}_2)$, from which to generate predictions for the two populations as follows:

1. Set initial values $\mu = \frac{\theta_1 + \theta_2}{2}$ and $\delta = \frac{\theta_1 - \theta_2}{2}$
2. Generate a single $\sigma^2 | \mathbf{y}_1, \mathbf{y}_2, \mu, \delta$
3. Generate a single $\mu | \mathbf{y}_1, \mathbf{y}_2, \delta, \sigma^2$
4. Generate a single $\delta | \mathbf{y}_1, \mathbf{y}_2, \mu, \sigma^2$
5. Predict $\tilde{y}_1 \sim \text{normal}(\mu + \delta, \sigma^2)$ and $\tilde{y}_2 \sim \text{normal}(\mu - \delta, \sigma^2)$

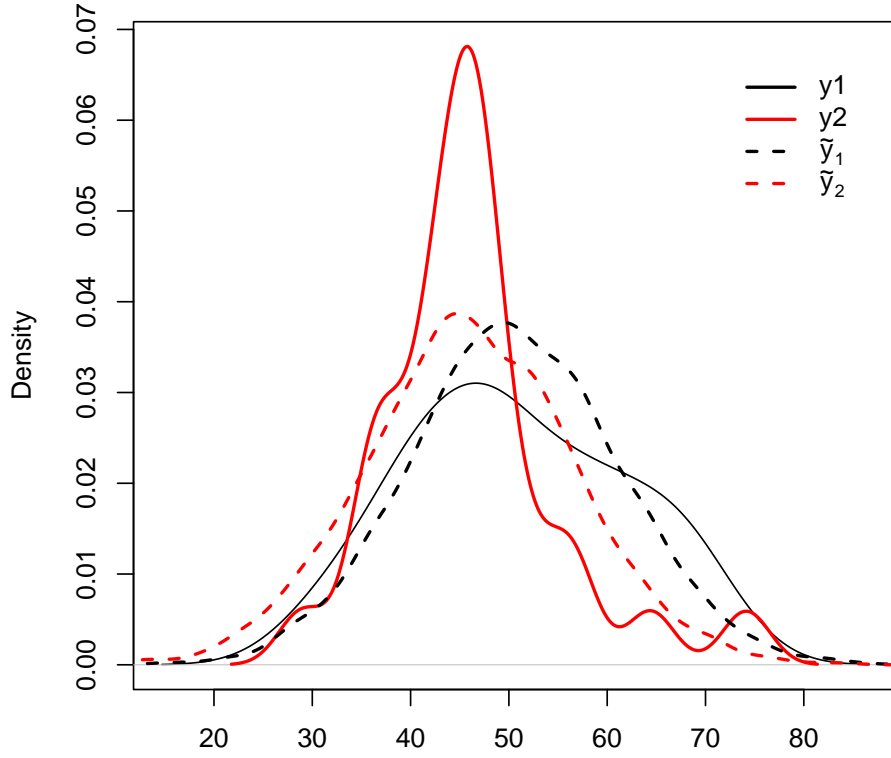
The user provides the two samples \mathbf{y}_1 and \mathbf{y}_2 along with values for $\mu_0, \sigma_0^2, \delta_0, \tau_0^2, \nu_0$, and desired prediction sample size N . The function returns N predictions for each population and the vectors of generated values for μ, δ , and σ^2 .

3.4.2.3 Example Hoff p. 128-129 *Analysis of math score data*

Math score data for two schools were based on results of a national exam in the United States, standardized to produce a nationwide mean of 50 and a standard deviation of 10. Unless the two schools were known in advance to be extremely exceptional, reasonable prior parameters can be based on this information. For the prior distributions of μ and σ^2 , we'll take $\mu_0 = 50$ and $\sigma_0^2 = 10^2 = 100$, although this latter value is likely to be an overestimate of the within-school sampling variability. We'll make these prior distributions somewhat diffuse, with $\gamma_0^2 = 25^2 = 625$ and $\nu_0 = 1$. For the prior distribution on δ , choosing $\delta_0 = 0$ represents the prior opinion that $\theta_1 > \theta_2$ and $\theta_2 > \theta_1$ are equally probable. Finally, since the scores are bounded between 0 and 100, half the difference between θ_1 and θ_2 must be less than 50 in absolute value, so a value of $\tau_0^2 = 25^2 = 625$ seems reasonably diffuse.

The results of a call to `rpredNormIG2(y1, y2, mu0, sigma0^2, delta0, tau0^2, N)` are summarized in the following plot.

2-samples: Density of Data and Predictions



3.4.3 k samples: Comparing multiple groups

For two-level data consisting of groups and units within groups, denote $y_{i,j}$ as the data on the i th unit in group j . We have the hierarchical normal model (Hoff p. 132ff):

$$\phi_j = \{\theta_j, \sigma^2\}, p(y|\phi_j) = \text{normal}(\theta_j, \sigma^2) \quad (\text{within-group model})$$

$$\psi_j = \{\mu, \tau^2\}, p(\theta_j|\psi) = \text{normal}(\mu, \tau^2) \quad (\text{between-group model})$$

We use standard semiconjugate normal and inverse-gamma prior distributions for the fixed but unknown parameters in the model:

$$\sigma^2 \sim \text{inverse-gamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$$

$$\tau^2 \sim \text{inverse-gamma}\left(\frac{\eta_0}{2}, \frac{\eta_0 \tau_0^2}{2}\right)$$

$$\mu \sim \text{normal}(\mu_0, \gamma_0^2)$$

3.4.3.1 Derivation As with the two-sample problem, joint posterior inferences for the unknown parameters can be made by constructing a Gibbs sampler to approximate the posterior distribution $p(\theta_1, \dots, \theta_m, \mu, \tau^2, \sigma^2 | \mathbf{y}_1, \dots, \mathbf{y}_m)$. For this we need the full conditional distribution of each parameter (Hoff pp. 134-135):

$$\{\mu | \theta_1, \dots, \theta_m, \tau^2\} \sim \text{normal} \left(\frac{\frac{m\bar{\theta}}{\tau^2} + \frac{\mu_0}{\gamma_0^2}}{\frac{m}{\tau^2} + \frac{1}{\gamma_0^2}}, \frac{1}{\frac{m}{\tau^2} + \frac{1}{\gamma_0^2}} \right)$$

$$\{\tau^2 | \theta_1, \dots, \theta_m, \mu\} \sim \text{inverse-gamma} \left(\frac{\eta_0 + m}{2}, \frac{\eta_0 \tau_0^2 + \sum (\theta_j - \mu)^2}{2} \right)$$

$$\{\theta_j | y_{1,j}, \dots, y_{n,j}, \sigma^2\} \sim \text{normal} \left(\frac{\frac{n_j \bar{y}_j}{\sigma^2} + \frac{1}{\tau^2}}{\frac{n_j}{\sigma^2} + \frac{1}{\tau^2}}, \frac{1}{\frac{n_j}{\sigma^2} + \frac{1}{\tau^2}} \right)$$

$$\{\sigma^2 | \theta, \mathbf{y}_1, \dots, \mathbf{y}_n\} \sim \text{inverse-gamma} \left(\frac{1}{2} \left[\nu_0 + \sum_{j=1}^m n_j \right], \frac{1}{2} \left[\nu_0 \sigma_0^2 + \sum_{j=1}^m \sum_{i=1}^{n_j} (y_{i,j} - \theta_j)^2 \right] \right).$$

Note that $\sum \sum (y_{i,j} - \theta_j)^2$ is the sum of squared residuals across all groups, conditional on the within-group means, and so the conditional distribution concentrates probability around a pooled-sample estimate of the variance.

3.4.3.2 R Implementation The standard format R function `rpredNormIGk()` implements a Gibbs sampler for posterior approximation of each unknown quantity by sampling from its full conditional distribution. From these posteriors, predictions are generated, as follows:

1. Set prior parameter values:

$$\begin{aligned} \nu_0, \sigma_0^2 & \text{ for } p(\sigma^2) \\ \eta_0, \tau_0^2 & \text{ for } p(\tau^2) \\ \mu_0, \gamma_0^2 & \text{ for } p(\mu). \end{aligned}$$

2. Set initial states for the unknown parameters:

$$\begin{aligned} \theta_1^{(1)} &= \bar{\mathbf{y}}_1, \dots, \theta_m^{(1)} = \bar{\mathbf{y}}_m \\ \mu^{(1)} &= \text{mean}(\theta_1^{(1)}, \dots, \theta_m^{(1)}) \\ \tau^{2(1)} &= \text{var}(\theta_1^{(1)}, \dots, \theta_m^{(1)}) \\ \sigma^{2(1)} &= \text{mean}(\text{var}(\mathbf{y}_1), \dots, \text{var}(\mathbf{y}_m)) \end{aligned}$$

3. For $s \in \{1, \dots, S\}$, sample

$$(a) \quad \mu^{(s+1)} \sim p(\mu | \theta_1^{(s)}, \dots, \theta_m^{(s)}, \tau^{2(s)})$$

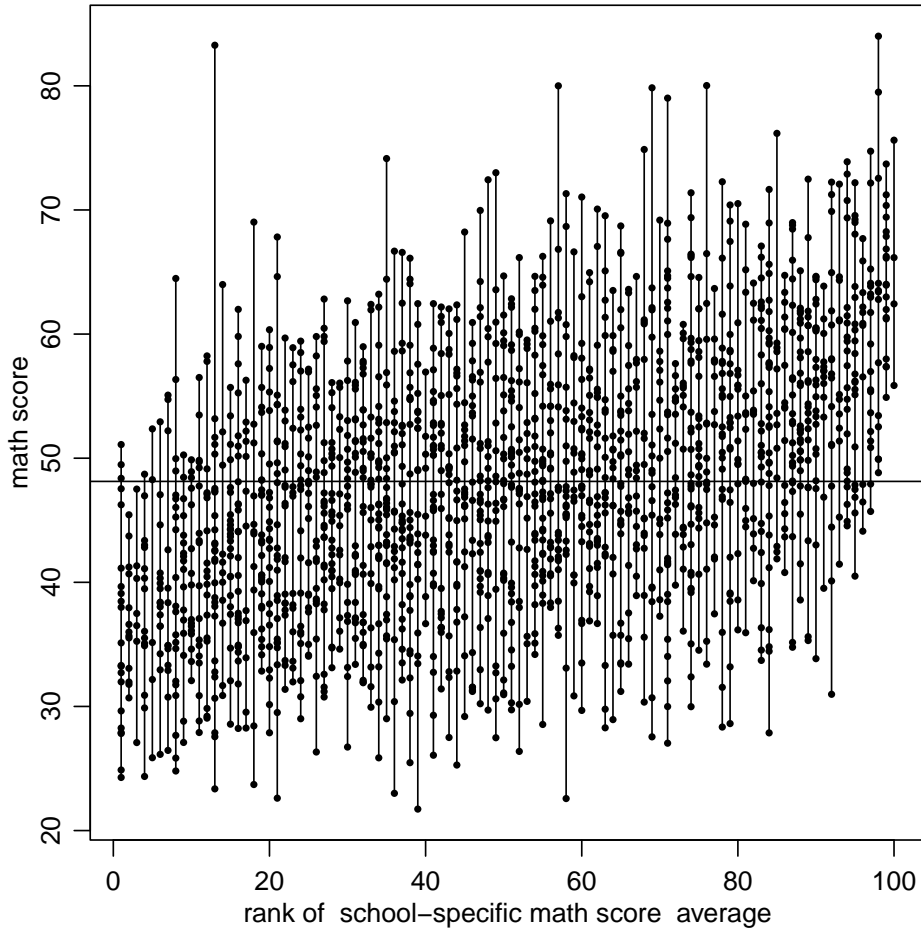
$$(b) \tau^{2(s+1)} \sim p\left(\tau^2 | \theta_1^{(s)}, \dots, \theta_m^{(s)}, \mu^{(s+1)}\right)$$

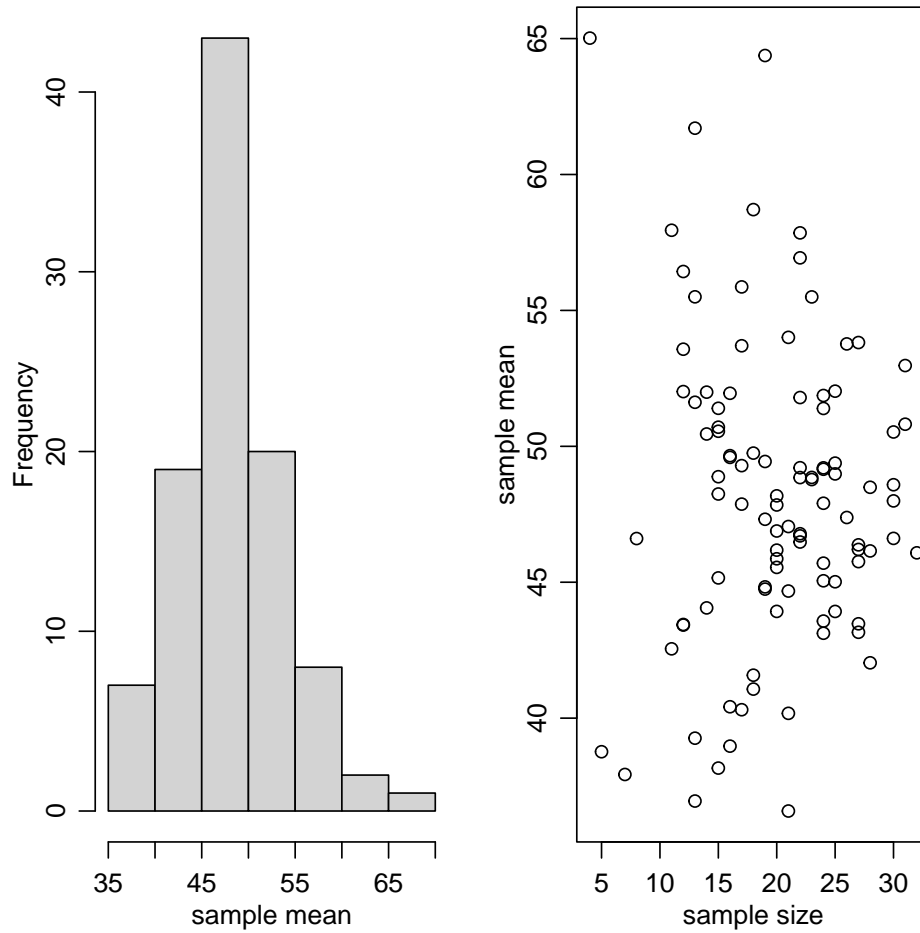
$$(c) \sigma^{2(s+1)} \sim p\left(\sigma^2 | \theta_1^{(s)}, \dots, \theta_m^{(s)}, \mathbf{y}_1, \dots, \mathbf{y}_m\right)$$

$$(d) \theta_j^{(s+1)} \sim p\left(\theta_j | \mu^{(s+1)}, \tau^{2(s+1)}, \sigma^{2(s+1)}, \mathbf{y}_j\right) \text{ for } j \in \{1, \dots, m\}$$

4. For $s \in \{1, \dots, S\}$, generate prediction $\tilde{y}_j^{(s)} \sim \text{normal}\left(\theta_j^{(s)}, \sigma^{2(s)}\right)$

3.4.3.3 Example Returning to the math scores example, data for 10th-grade students from 100 large urban schools (each having 10th-grade enrollment of at least 400) is summarized in the following plots.





For prediction, we'll use the following prior values (Hoff p. 137):

$\sigma_0^2 : 100$ (within-school variance)

$\nu_0 : 1$ (prior sample size)

$\tau_0^2 : 100$ (between-school variance)

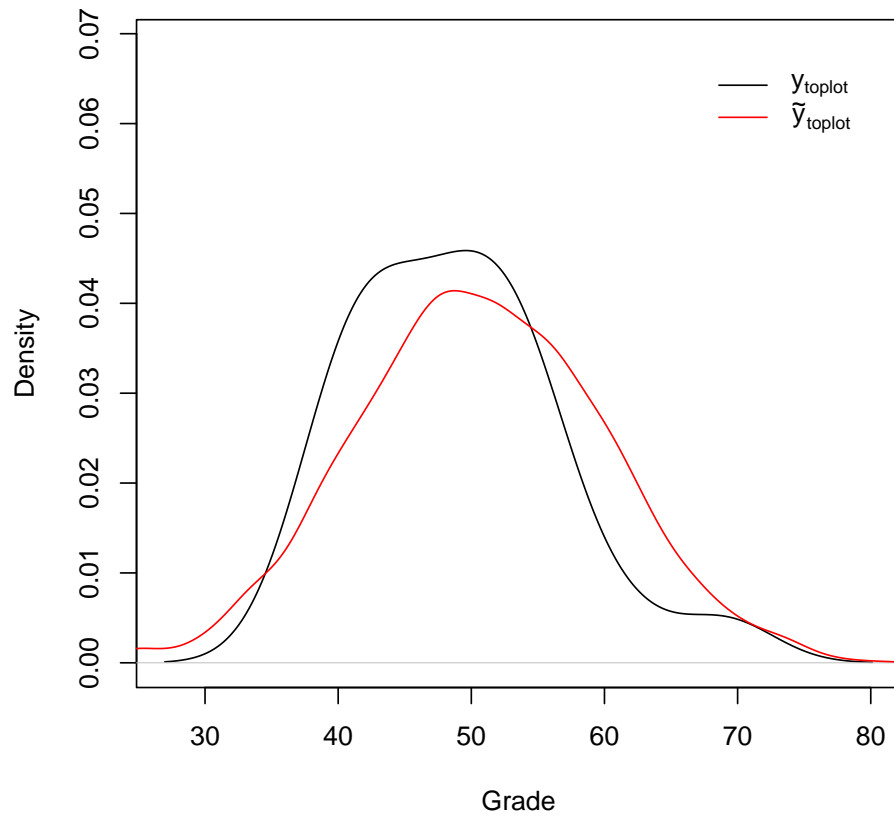
$\eta_0 : 1$ (prior sample size)

$\mu_0 : 50$ (prior mean of school means)

$\gamma_0^2 : 25$ (prior variance of school means)

Below: Pick a couple of schools that show different relationships between the data and the prediction

School 1 Data and Prediction



3.4.3.4 Ranking Treatments

4 Chapter 2: Normal Regression with Zellner's g -prior

4.1 Least Squares Estimation with Example (Hoff p. 149ff.)

Regression modeling is concerned with describing how the sampling distribution of one random variable Y varies with another variable or set of variables $\mathbf{x} = (x_1, \dots, x_p)$. Specifically, a regression model postulates a form for $p(y|\mathbf{x})$, the conditional distribution of Y given \mathbf{x} . Estimation of $p(y|\mathbf{x})$ is made using data y_1, \dots, y_n that are gathered under a variety of conditions $\mathbf{x}_1, \dots, \mathbf{x}_n$.

The normal linear regression model specifies that, in addition to $E[Y|\mathbf{x}]$ being linear, the sampling variability around the mean is i.i.d. normal:

$$\epsilon_1, \dots, \epsilon_n \stackrel{\text{i.i.d.}}{\sim} \text{normal}(0, \sigma^2)$$

$$Y_i = \boldsymbol{\beta}^T \mathbf{x}_i + \epsilon_i$$

This model provides a complete specification of the joint probability density of observed data y_1, \dots, y_n conditional upon $\mathbf{x}_1, \dots, \mathbf{x}_n$ and values of $\boldsymbol{\beta}$ and σ^2 :

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\beta}, \sigma^2) = \prod_{i=1}^n p(y_i | \mathbf{x}_i, \boldsymbol{\beta}, \sigma^2)$$

$$= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \boldsymbol{\beta}^T \mathbf{x}_i)^2 \right\} \quad (1)$$

Another way to write this joint probability density is in terms of the multivariate normal distribution: Let \mathbf{y} be the n -dimensional column vector $(y_1, \dots, y_n)^T$ and let \mathbf{X} be the $n \times p$ matrix whose i th row is $\mathbf{x}_i = \{x_{i,1}, x_{i,2}, \dots, x_{i,p}\}$. Then the normal regression model is

$$\{\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \sigma^2\} \sim \text{multivariate normal}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}),$$

where \mathbf{I} is the $p \times p$ identity matrix and

$$\mathbf{X}\boldsymbol{\beta} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} = \begin{pmatrix} \beta_1 x_{1,1} + \dots + \beta_p x_{1,p} \\ \vdots \\ \beta_1 x_{n,1} + \dots + \beta_p x_{n,p} \end{pmatrix} = \begin{pmatrix} E[Y_1 | \boldsymbol{\beta}, \mathbf{x}_1] \\ \vdots \\ E[Y_n | \boldsymbol{\beta}, \mathbf{x}_n] \end{pmatrix}$$

The density (1) depends on $\boldsymbol{\beta}$ through the residuals $(y_i - \boldsymbol{\beta}^T \mathbf{x}_i)$. We compute the ordinary least squares estimates

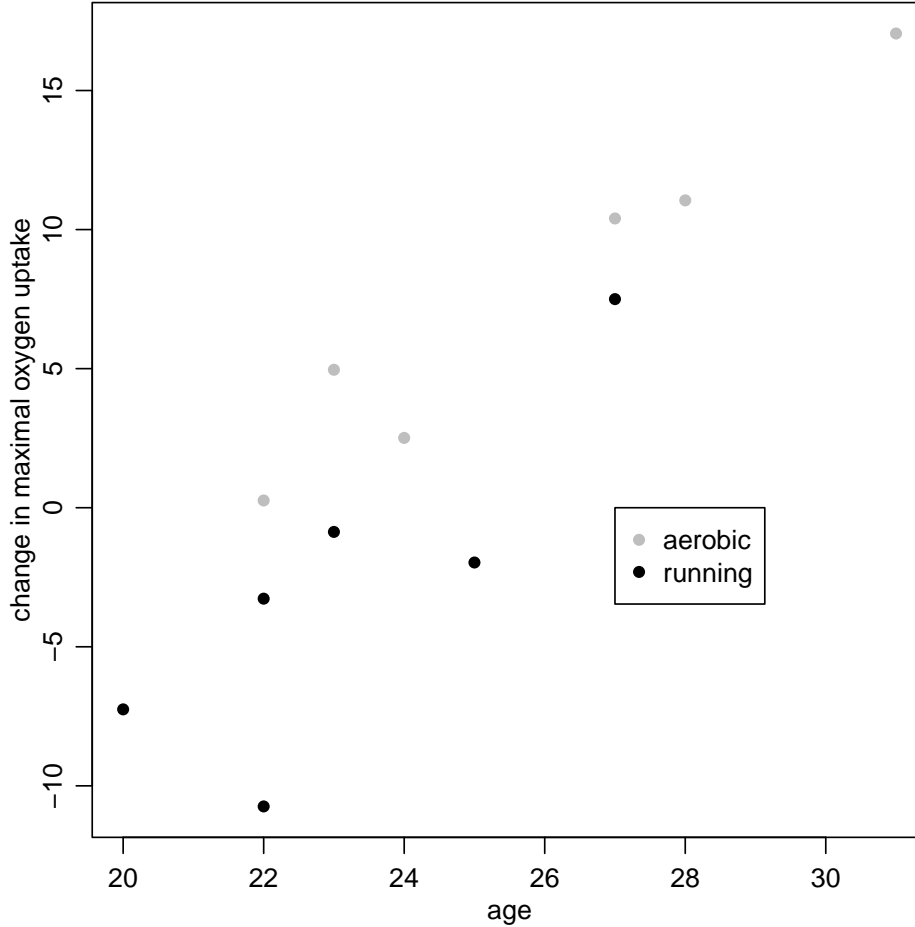
$$\hat{\boldsymbol{\beta}}_{ols} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

and

$$\hat{\sigma}_{ols}^2 = \frac{SSR(\hat{\boldsymbol{\beta}}_{ols})}{(n-p)} = \frac{\sum (y_i - \hat{\boldsymbol{\beta}}_{ols}^T \mathbf{x}_i)^2}{(n-p)}.$$

Example: Oxygen uptake (from Kuehl (2000), Hoff p. 149ff)

Twelve healthy men who did not exercise regularly were recruited to take part in a study of the effects of two different exercise regimens on oxygen uptake. Six of the twelve men were randomly assigned to a 12-week flat-terrain running program, and the remaining six were assigned to a 12-week step aerobics program. The maximum oxygen uptake of each subject was measured (in liters per minute) while running on an inclined treadmill, both before and after the 12-week program. Of interest is how a subject's change in maximal oxygen uptake may depend on which program they were assigned to. However, other factors, such as age, are expected to affect the change in maximal uptake as well. The results are shown here:



Hoff's regression model:

$$Y_i = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \beta_3 x_{i,3} + \beta_4 x_{i,4} + \epsilon_i, \text{ where} \quad (2)$$

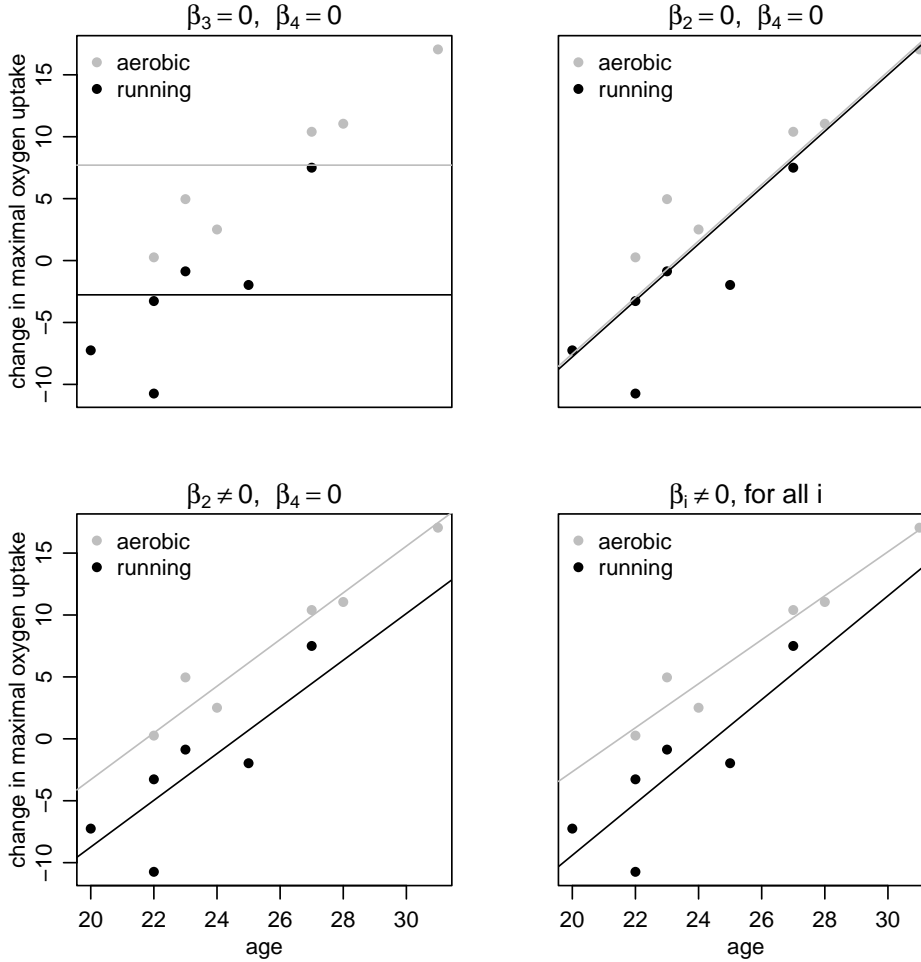
$x_{i,1} = 1$ for each subject i
 $x_{i,2} = 0$ if subject j is on the running program, 1 if on aerobic
 $x_{i,3} = \text{age of subject } i$
 $x_{i,4} = x_{i,2} \times x_{i,3}$

Under this model the conditional expectations of Y for the two different levels of $x_{i,1}$ are

$$E[Y|\mathbf{x}] = \beta_1 + \beta_3 \times (age) \text{ if } x_1 = 0, \text{ and}$$

$$E[Y|\mathbf{x}] = (\beta_1 + \beta_2) + (\beta_3 + \beta_4) \times (age) \text{ if } x_1 = 1$$

In other words, the model assumes that the relationship is linear in age for both exercise groups, with the difference in intercepts given by β_2 and the difference in slopes given by β_4 . If we assumed that $\beta_2 = \beta_4 = 0$, then we would have identical lines for both groups. If we assumed $\beta_2 \neq 0$ and $\beta_4 = 0$ then we would have a different line for each group but they would be parallel. Allowing all coefficients to be non-zero gives us two unrelated lines. Some different possibilities are depicted graphically below:



Let's find the least squares regression estimates for the model (2), and use the results to evaluate the differences between the two exercise groups. The ages of the 12 subjects, along with their observed changes in maximal oxygen uptake, are

$$\mathbf{x}_3 = (23, 22, 22, 25, 27, 20, 31, 23, 27, 28, 22, 24)$$

$$\mathbf{y} = (-0.87, -10.74, -3.27, -1.97, 7.50, -7.25, 17.05, 4.96, 10.40, 11.05, 0.26, 2.51),$$

with the first six elements of each vector corresponding to the subjects in the running group and the latter six corresponding to subjects in the aerobics group. After construct-

ing the 12×4 matrix $\mathbf{X} = (\mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \mathbf{x}_4)$, the matrices $\mathbf{X}^T \mathbf{X}$ and $\mathbf{X}^T \mathbf{y}$ can be computed, from which we get $\hat{\beta}_{ols} = (-51.29, 13.11, 2.09, -0.32)^T$:

This means that the estimated linear relationship between uptake and age has an intercept and slope of -51.29 and 2.09 for the running group, and $-51.29 + 13.11 = -38.18$ and $2.09 - 0.32 = 1.77$ for the aerobics group. These two lines are plotted in the fourth panel of Figure XX. We obtain unbiased estimate $\sigma^2 = SSR(\hat{\beta}_{ols})/(n - p) = 8.54$, and use this to compute the standard error of the components of $\hat{\beta}_{ols}$, which are 12.25, 15.76, 0.53, and 0.65, respectively. Comparing the values of $\hat{\beta}_{ols}$ to their standard errors suggests that the evidence for differences between the two exercise regimens is not very strong.

“Comparing the values of $\hat{\beta}_{ols}$ to their standard errors:”

Difference in Intercept:

$$H_0 : Intercept_{running} - Intercept_{aerobic} = 0; H_A : Intercept_{running} - Intercept_{aerobic} \neq 0$$

$$H_0 : \beta_1 - (\beta_1 + \beta_2) = -\beta_2 = 0 \text{ (that is } \beta_2 = 0); H_A : \beta_2 \neq 0$$

$$T = \frac{\beta_2 - 0}{SE_{\beta_2}} = \frac{13.11}{15.76} = 0.49$$

→ $p = 0.79$ → fail to reject H_0 and conclude no significant difference in intercept

Difference in Slope:

$$H_0 : Slope_{running} - Slope_{aerobic} = 0; H_A : Slope_{running} - Slope_{aerobic} \neq 0$$

$$H_0 : \beta_3 - (\beta_3 + \beta_4) = 0 \text{ (that is } \beta_4 = 0); H_A : \beta_4 \neq 0$$

$$T = \frac{\beta_4 - 0}{SE_{\beta_4}} = \frac{-0.32}{0.65} = 0.83$$

→ $p = 0.68$ → fail to reject H_0 and conclude no significant difference in slope

[1] 8.542477

	beta.ols	SE.ols	CIL	CIU
x1	-51.2939459	12.2522126	-78.5935768	-23.994315
x2	13.1070904	15.7619762	-22.0127811	48.226962
x3	2.0947027	0.5263585	0.9219028	3.267503
x4	-0.3182438	0.6498086	-1.7661075	1.129620

4.2 Bayesian Estimation for a Regression Model (Hoff p. 154ff)

4.2.1 Derivation

4.2.1.1 A semiconjugate prior distribution Hoff proposes a semiconjugate prior distribution for β and σ^2 to be used when there is information available about the parameters. The sampling density of the data (Equation 1) is

$$p(\mathbf{y}|\mathbf{X}, \beta, \sigma^2) \propto \exp\left\{-\frac{1}{2\sigma^2}\text{SSR}(\beta)\right\} = \exp\left\{-\frac{1}{2\sigma^2}[\mathbf{y}^T\mathbf{y} - 2\beta^T\mathbf{X}^T\mathbf{y} + \beta^T\mathbf{X}^T\mathbf{X}\beta]\right\}.$$

The role that β plays in the exponent looks very similar to that played by \mathbf{y} , and the distribution of \mathbf{y} is multivariate normal. This suggests that a multivariate normal prior distribution for β is conjugate: if $\beta \sim \text{multivariate normal}(\beta_0, \Sigma_0)$, then

$$\begin{aligned} p(\beta|\mathbf{y}, \mathbf{X}, \sigma^2) &\propto p(\mathbf{y}|\mathbf{X}, \beta, \sigma^2) \times p(\beta) \\ &\propto \exp\left\{-\frac{1}{2}(-2\beta^T\mathbf{X}^T\mathbf{y}/\sigma^2 + \beta^T\mathbf{X}^T\mathbf{X}\beta/\sigma^2) - \frac{1}{2}(-2\beta^T\Sigma_0^{-1}\beta_0 + \beta^T\Sigma_0^{-1}\beta)\right\} \\ &= \exp\left\{\beta^T(\Sigma_0^{-1}\beta_0 + \mathbf{X}^T\mathbf{y}/\sigma^2) - \frac{1}{2}\beta^T(\Sigma_0^{-1} + \mathbf{X}^T\mathbf{X}/\sigma^2)\beta\right\} \end{aligned}$$

This is proportional to a multivariate normal density, with

$$\text{Var}[\beta|\mathbf{y}, \mathbf{X}, \sigma^2] = (\Sigma_0^{-1} + \mathbf{X}^T\mathbf{X}/\sigma^2)^{-1} \quad (3)$$

$$\text{E}[\beta|\mathbf{y}, \mathbf{X}, \sigma^2] = (\Sigma_0^{-1} + \mathbf{X}^T\mathbf{X}/\sigma^2)^{-1}(\Sigma_0^{-1}\beta_0 + \mathbf{X}^T\mathbf{y}/\sigma^2). \quad (4)$$

As usual, we can gain some understanding of these formulae by considering some limiting cases. If the elements of the prior precision matrix Σ_0^{-1} are small in magnitude, then the conditional expectation $\text{E}[\beta|\mathbf{y}, \mathbf{X}, \sigma^2]$ is approximately equal to $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$, the least squares estimate. On the other hand, if the measurement precision is very small (σ^2 is very large), then the expectation is approximately β_0 , the prior expectation.

As in most normal sampling problems, the semiconjugate prior distribution for σ^2 is an inverse-gamma distribution. Letting $\gamma = 1/\sigma^2$ be the measurement precision, if $\gamma \sim \text{gamma}(\nu_0/2, \nu_0\sigma_0^2/2)$, then

$$\begin{aligned} p(\gamma|\mathbf{y}, \mathbf{X}, \beta) &\propto p(\gamma)p(\mathbf{y}|\mathbf{X}, \beta, \gamma) \\ &\propto [\gamma^{\nu_0/2-1}\exp(-\gamma \times \nu_0\sigma_0^2/2)] \times [\gamma^{n/2}\exp(-\gamma \times \text{SSR}(\beta)/2)] \\ &= \gamma^{(\nu_0+n)/2-1}\exp(-\gamma[\nu_0\sigma_0^2 + \text{SSR}(\beta)]/2), \end{aligned}$$

which we recognize as a gamma density, so that

$$\{\sigma^2|\mathbf{y}, \mathbf{X}, \beta\} \sim \text{inverse-gamma}([\nu_0 + n]/2, [\nu_0\sigma_0^2 + \text{SSR}(\beta)]/2).$$

Constructing a Gibbs sampler to approximate the joint posterior distribution $p(\beta, \sigma^2|\mathbf{y}, \mathbf{X})$ is then straightforward: given current values $\{\beta^{(s)}, \sigma^{2(s)}\}$, new values can be generated by

1. updating β :

- (a) compute $\mathbf{V} = \text{Var}[\beta|\mathbf{y}, \mathbf{X}, \sigma^{2(s)}]$ and $\mathbf{m} = \text{E}[\beta|\mathbf{y}, \mathbf{X}, \sigma^{2(s)}]$
- (b) sample $\beta^{(s+1)} \sim \text{multivariate normal}(\mathbf{m}, \mathbf{V})$

2. updating σ^2 :

- (a) compute $\text{SSR}(\beta^{(s+1)})$
- (b) sample $\sigma^{2(s+1)} \sim \text{inverse-gamma}([\nu_0 + n]/2, [\nu_0\sigma_0^2 + \text{SSR}(\beta^{(s+1)})]/2)$.

To create a sample from the predictive distribution of responses: for each $s \in \{1, \dots, S\}$, draw $\epsilon^{(s)} \sim N(0, \sigma^{2(s)})$. Then compute

$$y^{(s)} = \beta^{(s)T} \mathbf{X} + \epsilon.$$

4.2.1.2 Default and weakly informative prior distributions In situations where prior information is unavailable or difficult to quantify, an alternative “default” class of prior distributions is given. Specification of the prior parameters (β_0, Σ_0) and (ν_0, σ_0^2) that represent actual prior information for a Bayesian analysis can be difficult. For a prior distribution that is not going to represent real prior information about the parameters, we choose one that is as minimally informative as possible. The resulting posterior distribution, then, will represent the posterior information of someone who began with little knowledge of the population being studied. Here we will employ Zellner’s “ g -prior” (Zellner, 1986). We choose $\beta_0 = \mathbf{0}$ and $\Sigma_0 = k(\mathbf{X}^T \mathbf{X})^{-1}$, $k = g\sigma^2$, $g > 0$, which satisfies a desired condition that the regression parameter estimation be invariant to changes in the scale of the regressors. With this, equations 3 and 4 reduce to

$$\text{Var}[\beta|\mathbf{y}, \mathbf{X}, \sigma^2] = [\mathbf{X}^T \mathbf{X}/(g\sigma^2) + \mathbf{X}^T \mathbf{X}/\sigma^2]^{-1} = \frac{g}{g+1} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \quad (5)$$

$$\text{E}[\beta|\mathbf{y}, \mathbf{X}, \sigma^2] = [\mathbf{X}^T \mathbf{X}/(g\sigma^2) + \mathbf{X}^T \mathbf{X}/\sigma^2]^{-1} \mathbf{X}^T \mathbf{y}/\sigma^2 = \frac{g}{g+1} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (6)$$

Letting

$$\mathbf{V} = \frac{g}{g+1} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \text{ and } \mathbf{m} = \frac{g}{g+1} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

we arrive at posteriors

$$\{\sigma^2|\mathbf{y}, \mathbf{X}\} \sim \text{inverse-gamma}([\nu_0 + n]/2, [\nu_0\sigma_0^2 + \text{SSR}_g]/2) \quad (7)$$

$$\{\beta|\mathbf{y}, \mathbf{X}, \sigma^2\} \sim \text{multivariate normal} \left(\frac{g}{g+1} \hat{\beta}_{ols}, \frac{g}{g+1} \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1} \right). \quad (8)$$

$$\text{Here } \text{SSR}_g = \mathbf{y}^T \mathbf{y} - \mathbf{m}^T \mathbf{V}^{-1} \mathbf{m} = \mathbf{y}^T (\mathbf{I} - \frac{g}{g+1} \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \mathbf{y}.$$

Simple Monte Carlo approximation can be used to sample from the joint posterior density $p(\sigma^2, \beta|\mathbf{y}, \mathbf{X})$ as follows. Here g is typically set to the number of prior observations. Then:

1. sample $\sigma^2 \sim \text{inverse-gamma}([\nu_0 + n]/2, [\nu_0\sigma_0^2 + \text{SSR}_g]/2)$
2. sample $\beta \sim \text{multivariate normal}\left(\frac{g}{g+1}\hat{\beta}_{ols}, \frac{g}{g+1}\sigma^2[\mathbf{X}^T\mathbf{X}]^{-1}\right)$.

To create a sample from the predictive distribution of responses, draw $\epsilon \sim N(0, \sigma^2)$. Then for each triplet $(\beta, \sigma^2, \epsilon)$ we have

$$y = \beta^T \mathbf{X} + \epsilon.$$

4.2.2 R Implementation

The standard format R function

```
rpredNormReg(S=1,Xpred,X,y,beta0,Sigma0,nu0=1,s20=1,gprior = TRUE)
```

approximates the joint posterior density $p(\sigma^2, \beta | \mathbf{y}, \mathbf{X})$ using one of the two methods described above, generates S triplets $(\beta^{(s)}, \sigma^{2(s)}, \epsilon^{(s)} \sim N(0, \sigma^{2(s)})$, and returns S predictions $y = X_{pred}\beta^{(s)} + \epsilon^{(s)}$.

The function defaults to Zellner's location-invariant g-prior, in which case input values for `beta0`, `Sigma0`, `nu0`, and `s20` are ignored. If the user wants to employ Hoff's semi-conjugate prior as defined in section 4.2.1.1 above, all input variables must be specified, with `gprior = FALSE`.

4.2.3 Example

In the example below (Hoff data and code found [here](#)) to employ Hoff's semi-conjugate prior we use

$\beta_0 = \hat{\beta}_{ols} = (-51.29, -51.29, -51.29, -51.29)$ (ordinary least squares estimator of β)

$$\Sigma_0 = (X^T X)^{-1} \sigma^2 n = \begin{pmatrix} 1801.4 & -1801.4 & -77.02 & 77.02 \\ -1801.4 & 2981.28 & 77.02 & -122.03 \\ -77.02 & 77.02 & 3.32 & -3.32 \\ 77.02 & -122.03 & -3.32 & 5.07 \end{pmatrix} \text{ (sampling variance of } \hat{\beta}_{ols} \text{)}$$

$\nu_0 = 1$ (prior sample size)

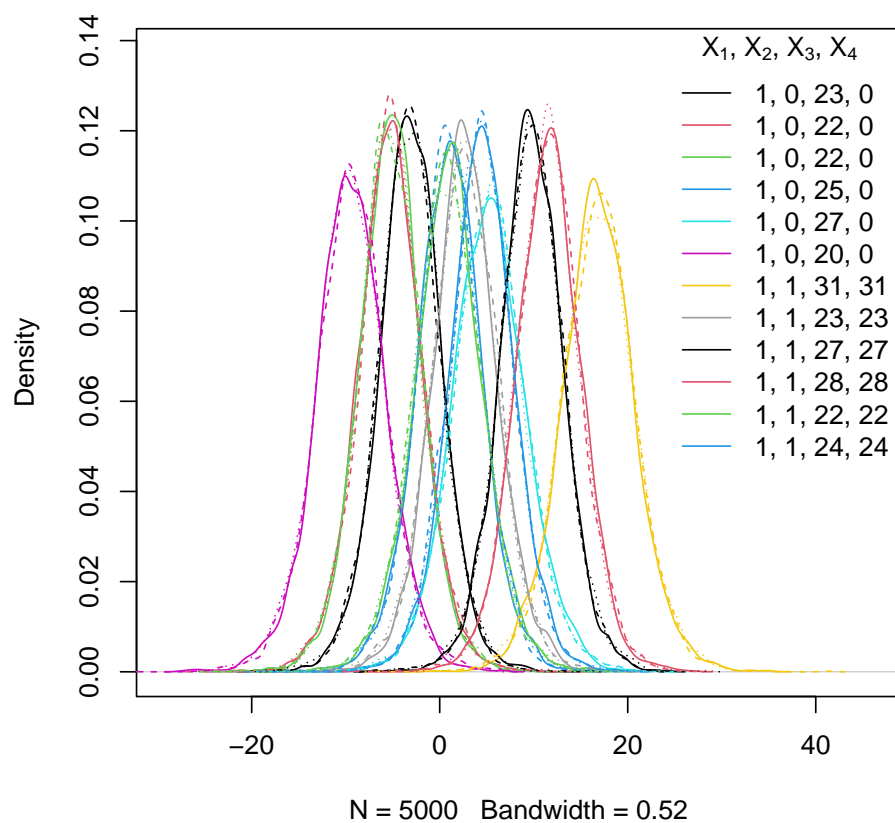
$$\sigma_0^2 = \frac{\sum e_i}{n-1} = 6.21 \text{ (variance of the residuals)}$$

$S = 5000$ (sample size for predictive distribution random draw)

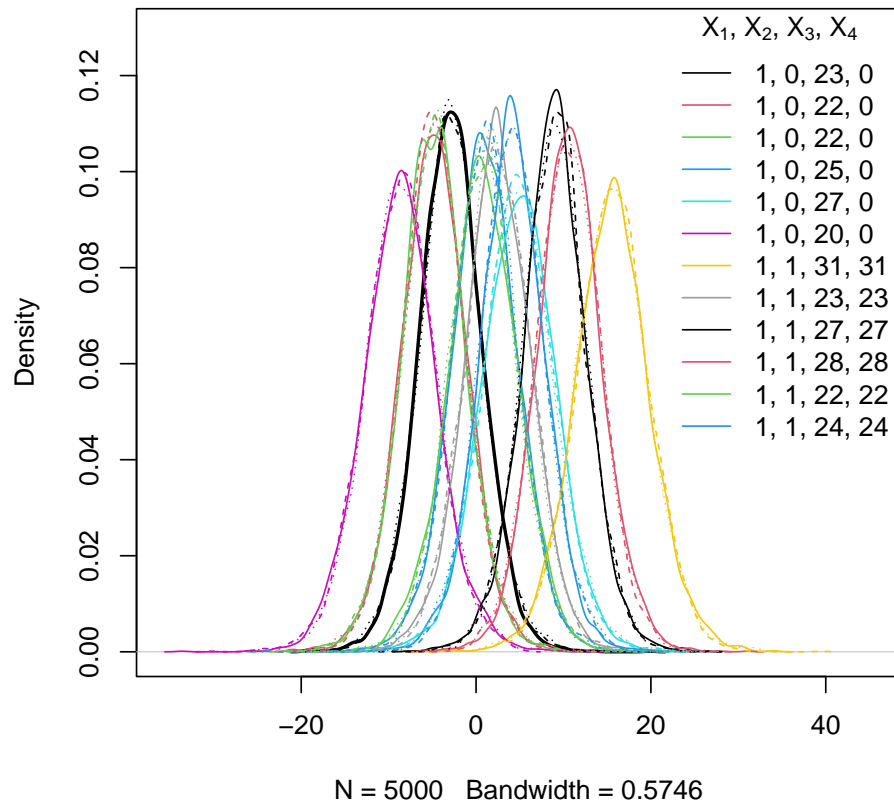
To do: clean up write-up; make plots comparing g-prior and non-g-prior results; circle back with Dean about appropriate priors for non-g-prior case; why is Hoff looking at the difference between the two cases with prior info very close to the sample characteristics?

ALSO INCLUDE IN WRITE-UP EXPLICIT PREDICTION STEP, E.G. "Y = ..."

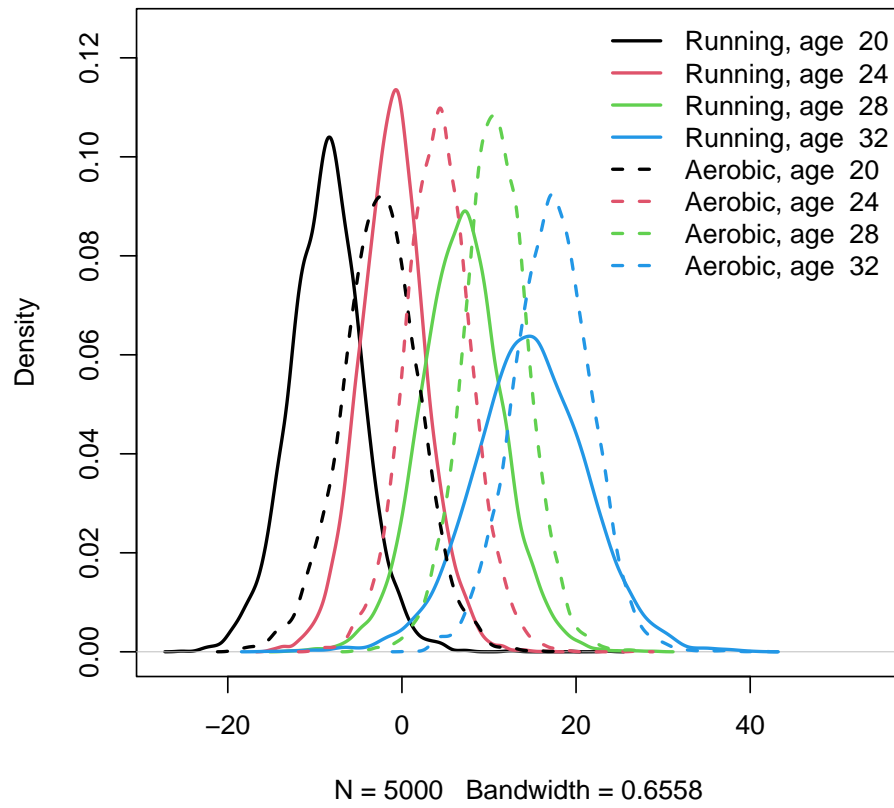
Predictive Density Using Semi-Conjugate Prior



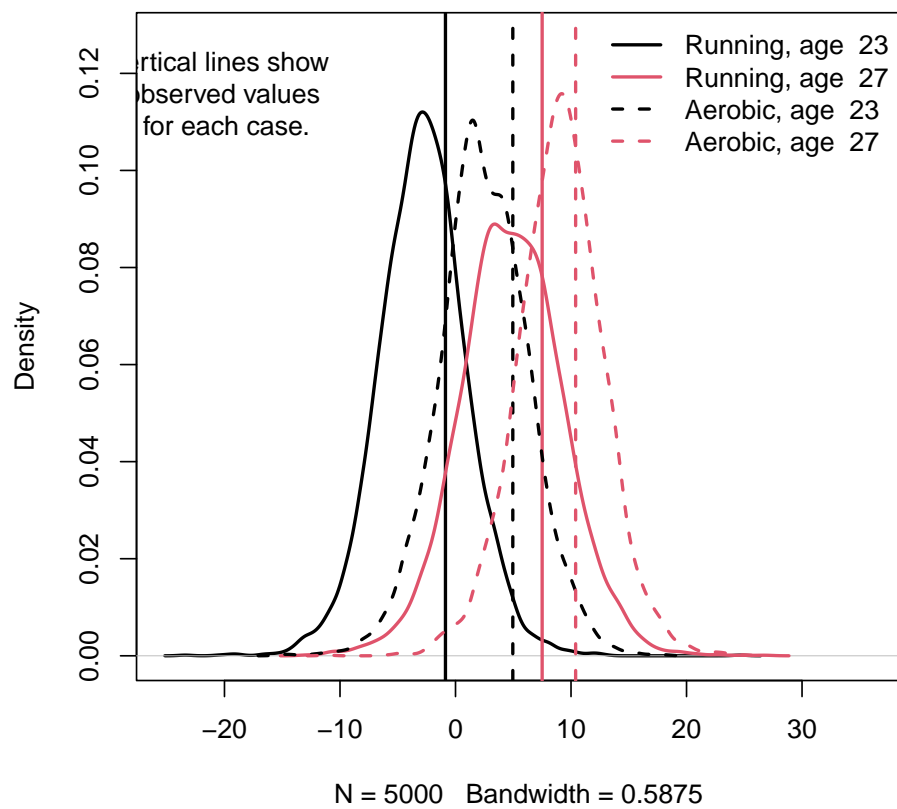
Predictive Density Using Zellner's g-prior



Predictive Density Using Zellner's g-prior

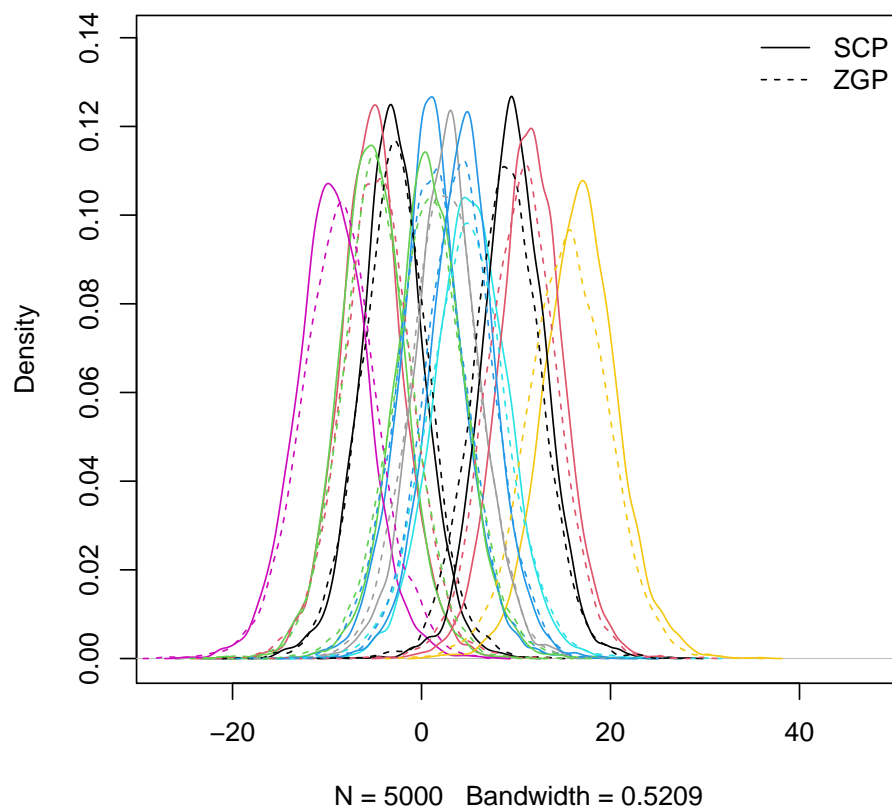


Predictive Density Using Zellner's g-prior



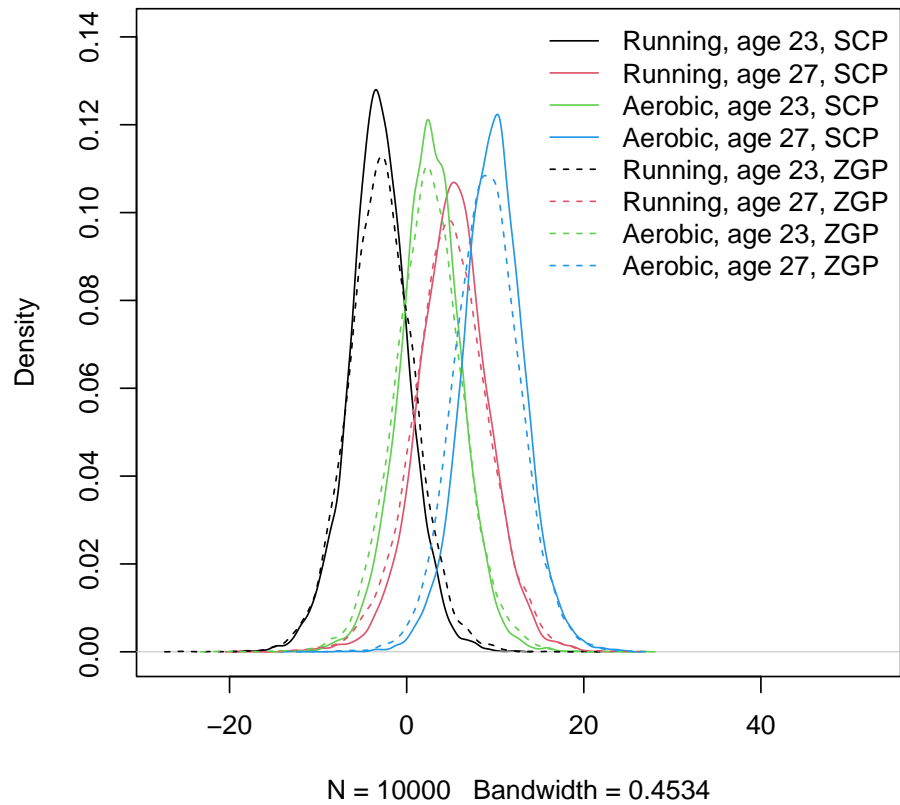
Comparing predictions using semi-conjugate prior vs. Zellner's g-prior:

Semi-conjugate prior vs. Zellner's g-prior



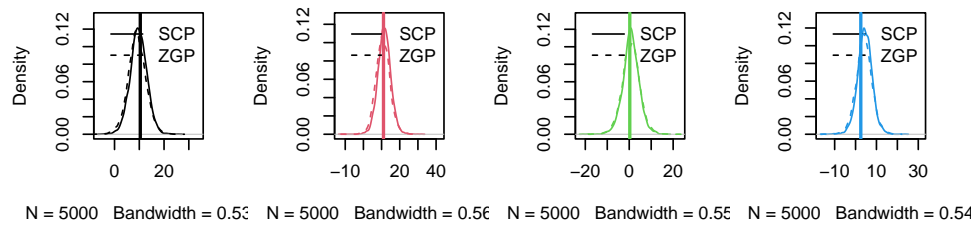
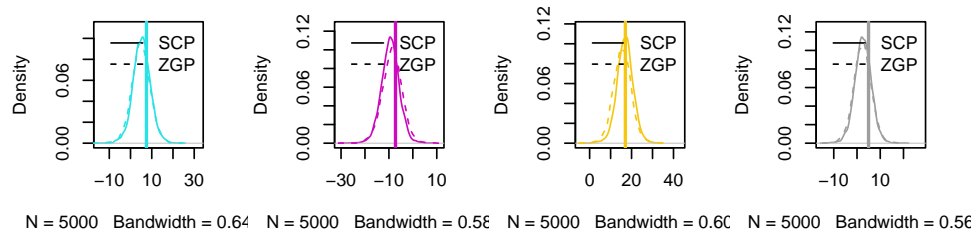
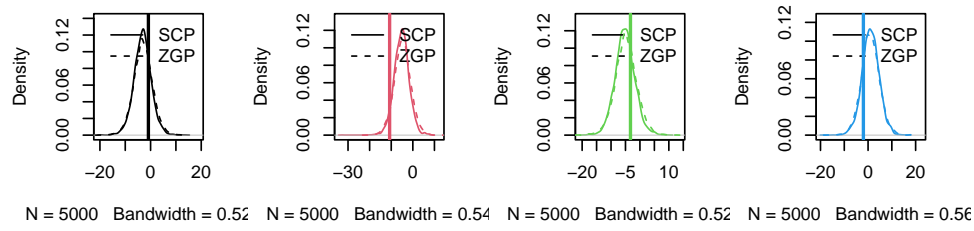
Inspection of the graphs shows the predicted distributions using Zellner's g-prior shrink toward 0, and have greater variance than those predicted using Hoff's semi-conjugate prior.

Semi-conjugate prior vs. Zellner's g-prior



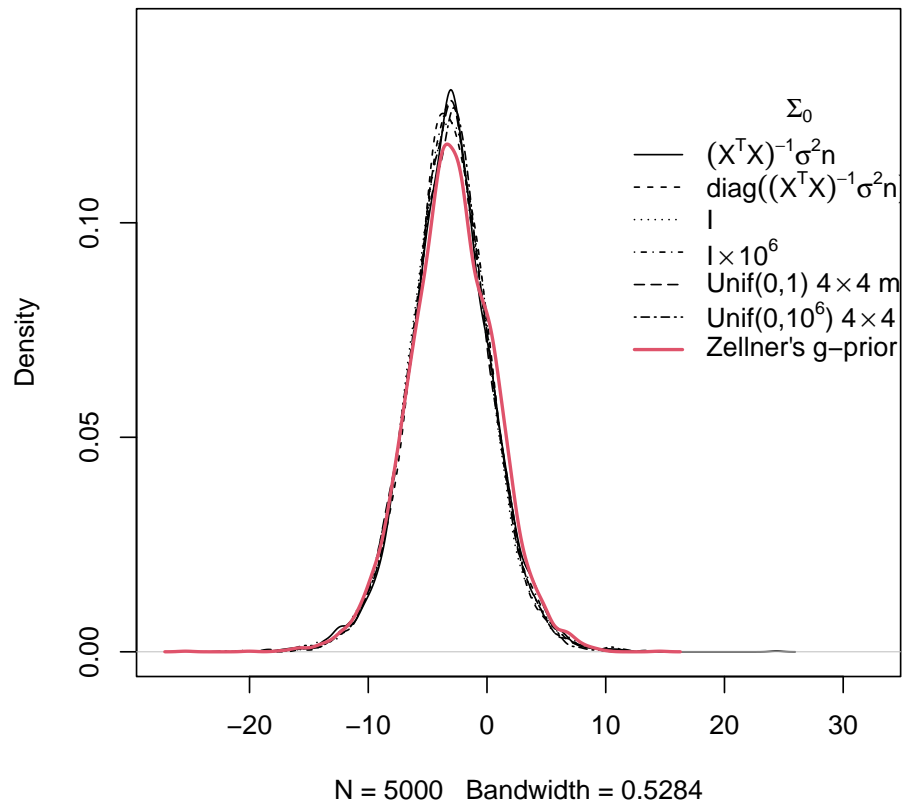
EXPLAIN WHY SCP PREDICTIONS HAVE TALLER DISTRIBUTIONS THAN ZGP PREDICTIONS. ALSO WHY ZGP PREDICTIONS SHRINK TOWARD 0

Comparing observed values to predictive distributions:

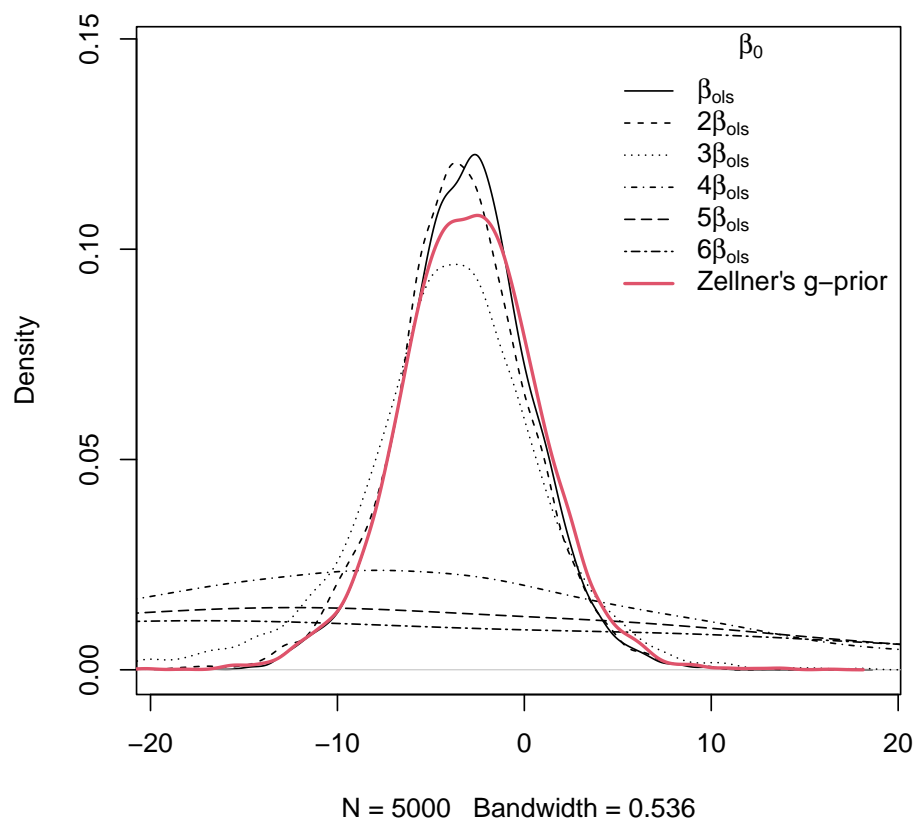


The following plots exhibit the influence of varying the prior info with the semi-conjugate prior.

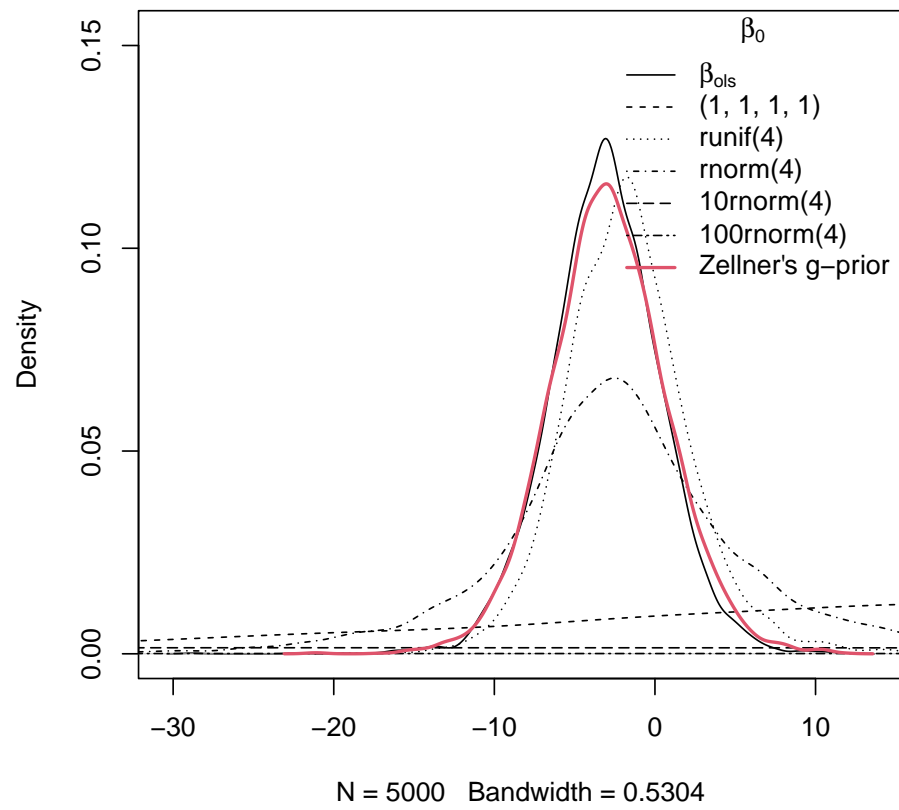
Varying Prior Information: Σ_0



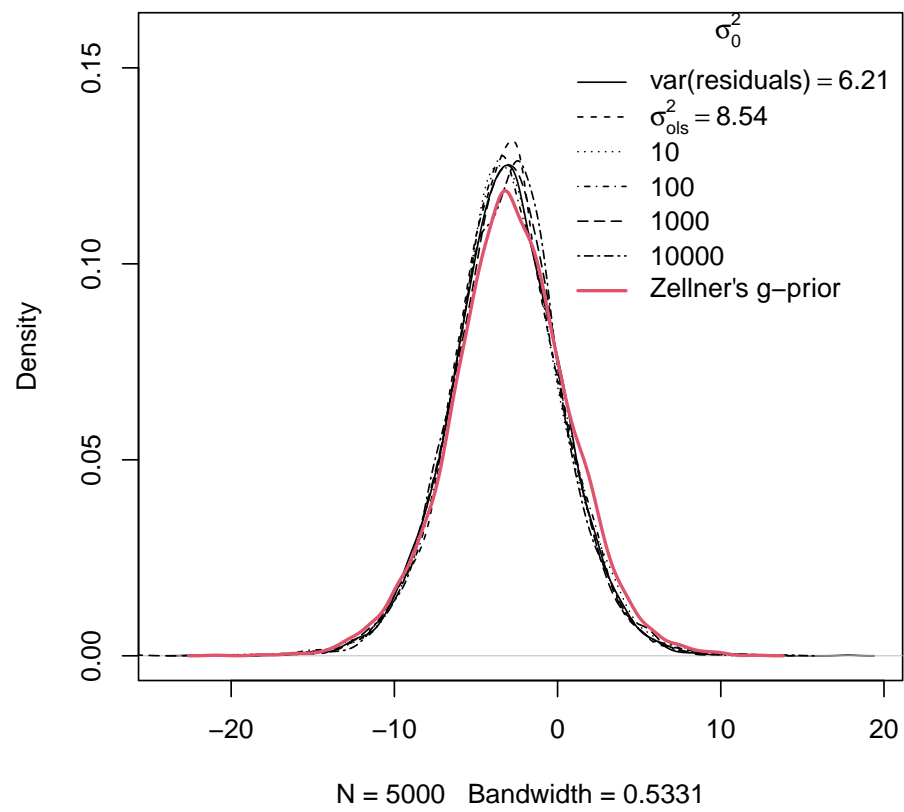
Varying Prior Information: Scaling β_0



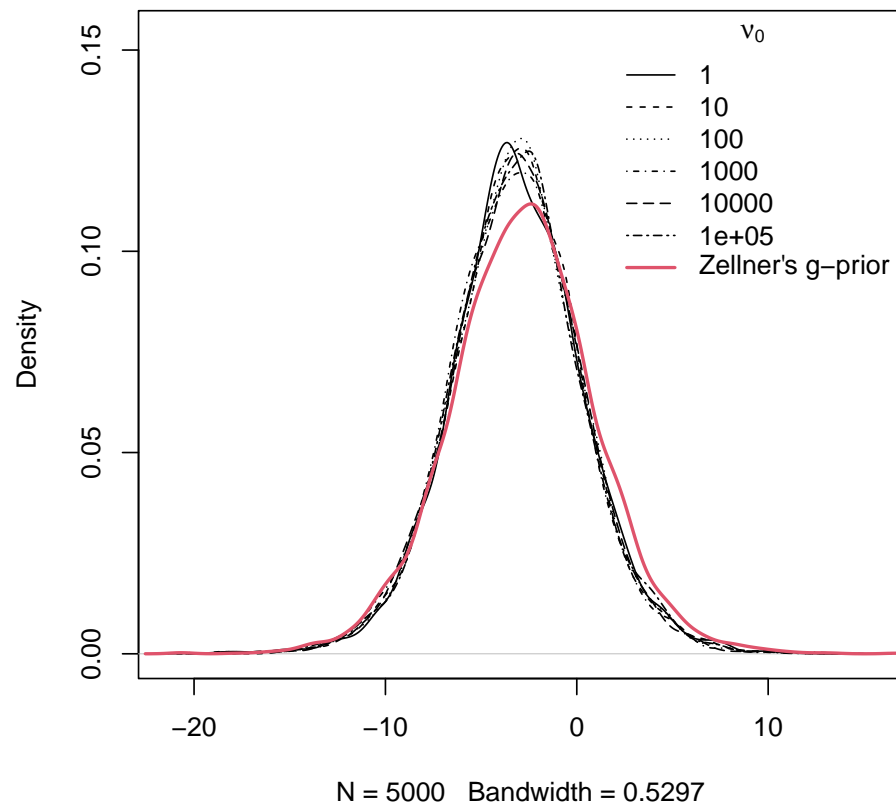
Varying Prior Information: Various β_0



Varying Prior Information: Scaling σ_0^2



Varying Prior Information: Scaling v_0



5 Conclusion

6 References