

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	4
2	Introduction: Predictive Inference	5
2.1	Why is predictive inference important?	5
2.2	Difference between parametric inference and predictive inference	5
2.2.1	When is predictive inference more useful?	5
2.2.2	When is parametric inference more useful?	5
2.3	The Bayesian Parametric Prediction Format	5
2.4	[Maybe] Example of Difference between results from Plug-in estimator and results using Predictive Inference	5
3	Chapter 1: Predictive Problems with Conjugate Priors	6
3.1	Prediction of Future Successes: Beta-Binomial (Geisser p. 73)	6
3.1.1	Derivation	6
3.1.2	R Implementation	7
3.1.3	Example	7
3.2	Survival Time: Exponential-Gamma (Geisser p. 74)	7
3.2.1	Derivation	7
3.2.2	R Implementation	9
3.2.3	Example	9
3.3	Poisson-Gamma Model (Hoff p. 43ff)	11
3.3.1	Derivation	11
3.3.2	R Implementation	15
3.3.3	Example	15
3.4	Normal Observation with Normal-Inverse Gamma Prior	17
3.4.1	One sample	17
3.4.1.1	Derivation	17
3.4.1.2	R Implementation	18
3.4.1.3	Example	19
3.4.2	Two samples	22
3.4.2.1	Derivation	22
3.4.2.2	R Implementation	23
3.4.2.3	Example	23
3.4.3	k samples: Comparing multiple groups	24
3.4.3.1	Derivation	25
3.4.3.2	R Implementation	25
3.4.3.3	Example	26
3.4.3.4	Ranking Treatments	28
4	Chapter 2: Normal Regression with Zellner's g-prior	29
4.0.0.1	Derivation	29
4.0.0.2	R Implementation	29
4.0.0.3	Example	29
5	Conclusion	30

1 Thesis Abstract

- (paragraph) Statement of the thesis topic and objectives
- (paragraph) Explanation of R package

2 Introduction: Predictive Inference

2.1 Why is predictive inference important?

2.2 Difference between parametric inference and predictive inference

2.2.1 When is predictive inference more useful?

2.2.2 When is parametric inference more useful?

[examples, comparisons]

2.3 The Bayesian Parametric Prediction Format

[Geisser p. 49]

Let

$$f(x^{(N)}, x_{(M)} | \theta) = f(x_{(M)} | x^{(N)}, \theta) f(x^{(N)} | \theta).$$

Here $x^{(N)}$ represents observed events and $x_{(M)}$ are future events. We calculate

$$f(x_{(M)}, x^{(N)}) = \int f(x^{(N)}, x_{(M)} | \theta) p(\theta) d\theta$$

where $p(\theta)$ is the prior density and

$$f(x_{(M)} | x^{(N)}) = \frac{f(x_{(M)}, x^{(N)})}{f(x^{(N)})} = \int f(x_{(M)} | \theta) p(\theta | x^{(N)}) d\theta$$

where

$$p(\theta | x^{(N)}) \propto f(x^{(N)} | \theta) p(\theta).$$

2.4 [Maybe] Example of Difference between results from Plug-in estimator and results using Predictive Inference

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]

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

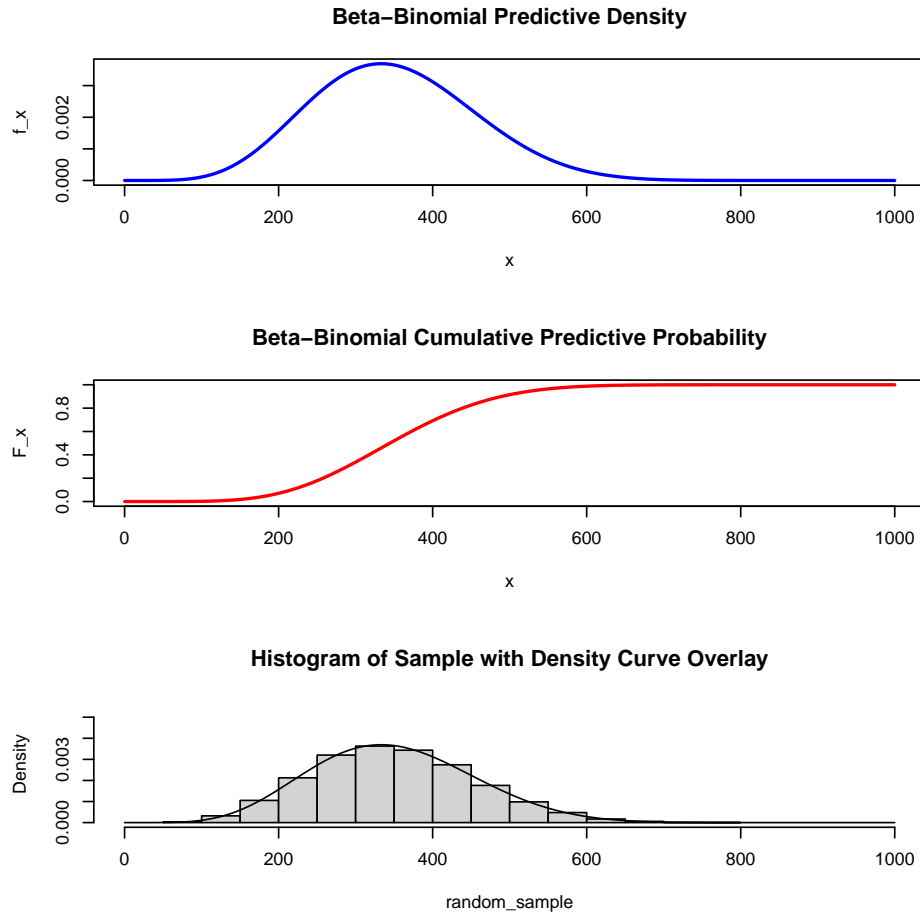
$$\begin{aligned} \Pr[R = r | t] &= \int \binom{M}{r} \theta^r (1 - \theta)^{M-r} p(\theta | X^{(N)}) d\theta \\ &= \binom{M}{r} \int \theta^r (1 - \theta)^{M-r} \frac{\Gamma(N + \alpha + \beta)}{\Gamma(t + \alpha) \Gamma(N - t + \beta)} \theta^{t+\alpha-1} (1 - \theta)^{N-t+\beta-1} d\theta \\ &= \frac{M!}{r!(M-r)!}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

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, 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.)

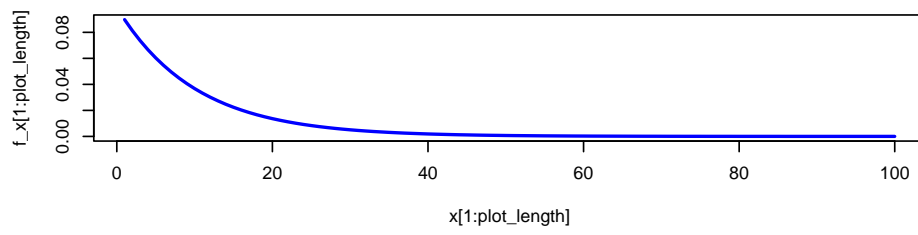
3.2.2 R Implementation

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.

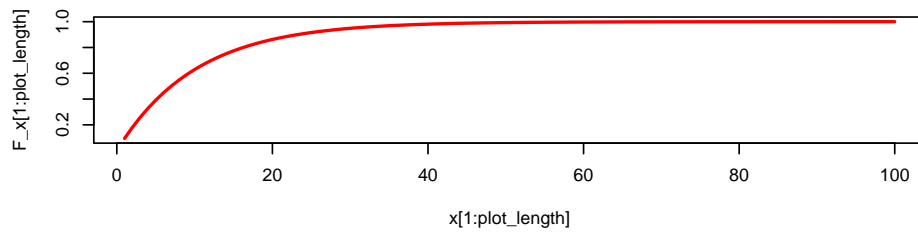
3.2.3 Example

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()`.

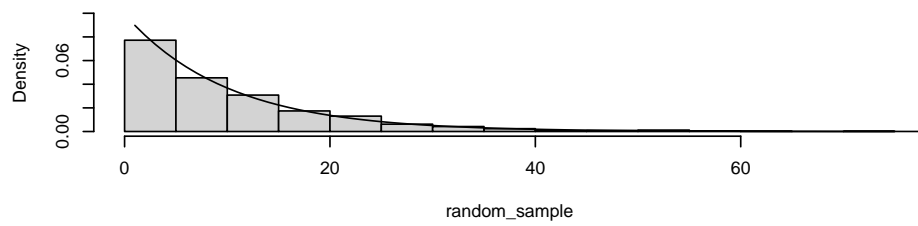
Exponential-Gamma Predictive Density



Exponential-Gamma Cumulative Predictive Probability



Histogram of Sample with Density Curve Overlay



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

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

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

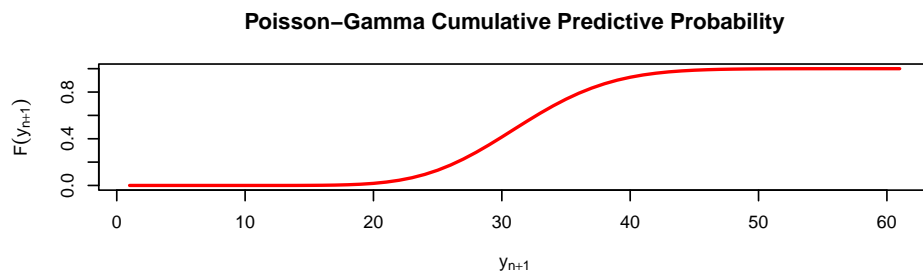
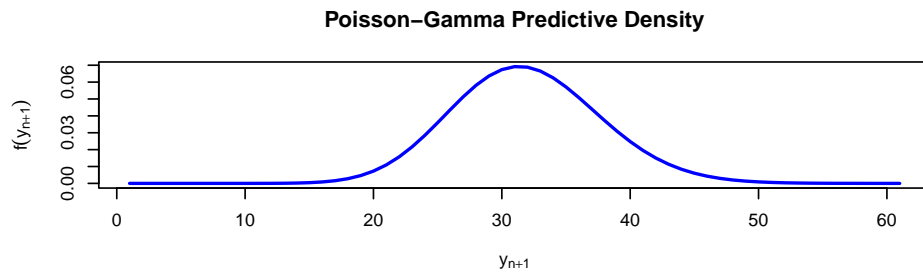
This result has been used to create standard format R functions `dpredPG()`, `ppredPG()`, and `rpredPG()` for the Poisson-Gamma distribution for density, cumulative probability, and random sampling, respectively (see appendix). These functions are exercised in the following example.

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

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

$$\bar{y} = 1.804$$

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

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

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

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

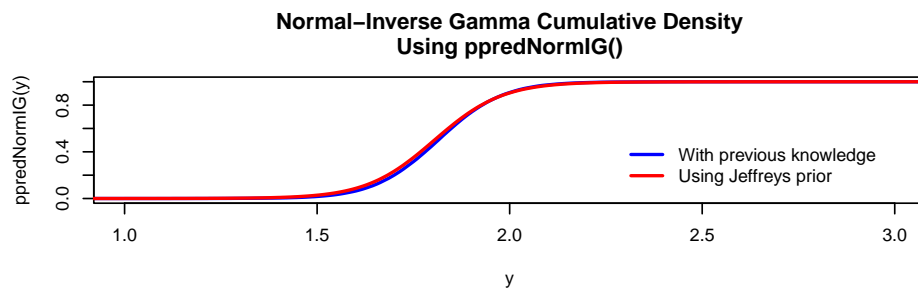
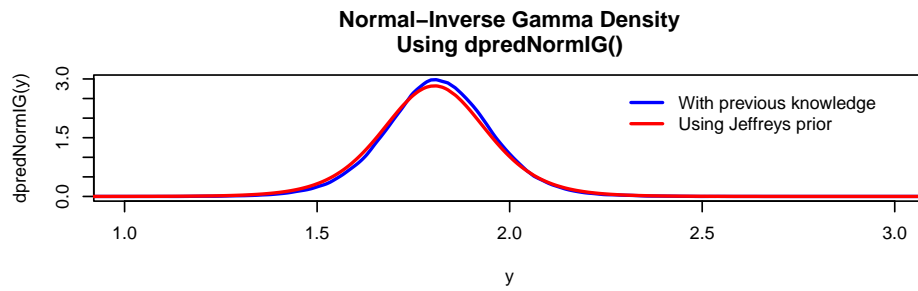
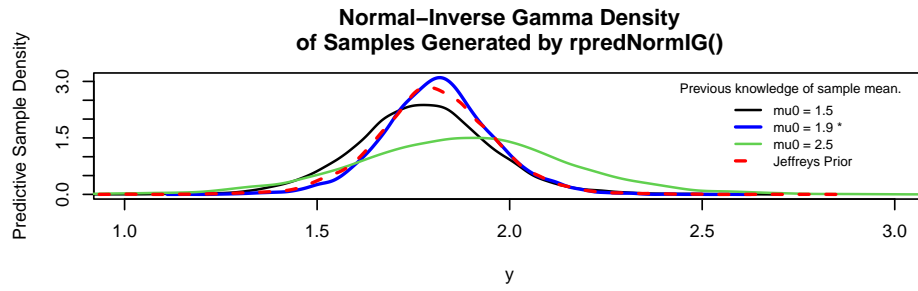
$$\sigma_n^2 = \frac{1}{10} \left[1 \cdot 0.01 + (9 - 1) \cdot 0.0169 + \frac{1 \cdot 9}{10} (1.804 - 1)^2 \right] = 0.0153$$

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 $\tilde{y}_2 \sim \text{normal}(\mu - \delta, \sigma^2)$

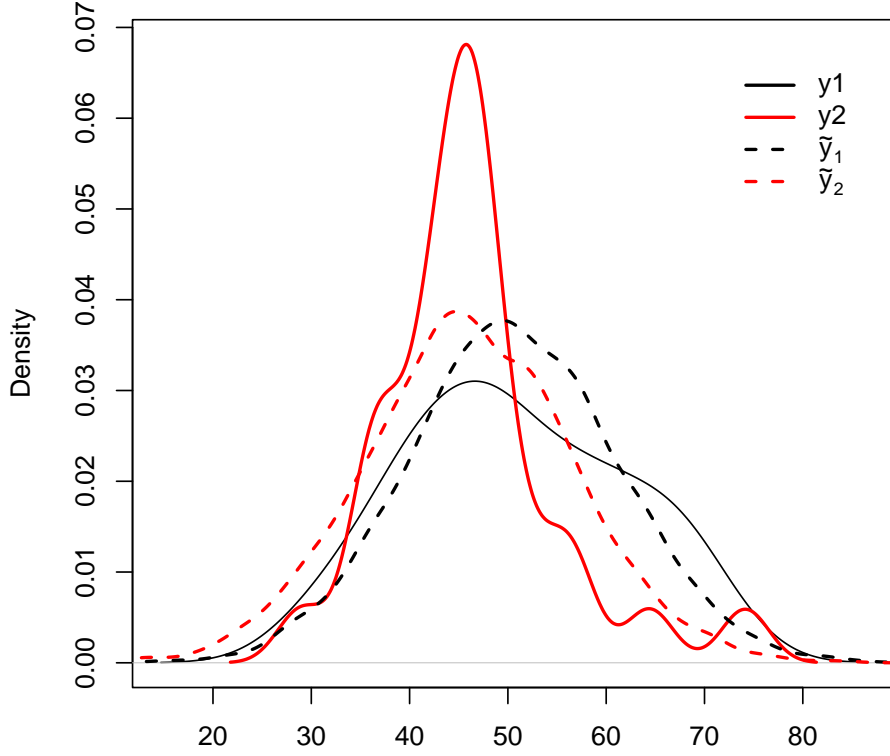
The user provides the two samples \mathbf{y}_1 and \mathbf{y}_2 along with values for $\mu_0, \sigma_0^2, \delta_0, \tau_0^2, \nu_0$, and desired prediction sample size N . The function returns N predictions for each population and the vectors of generated values for μ, δ , and σ^2 .

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

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

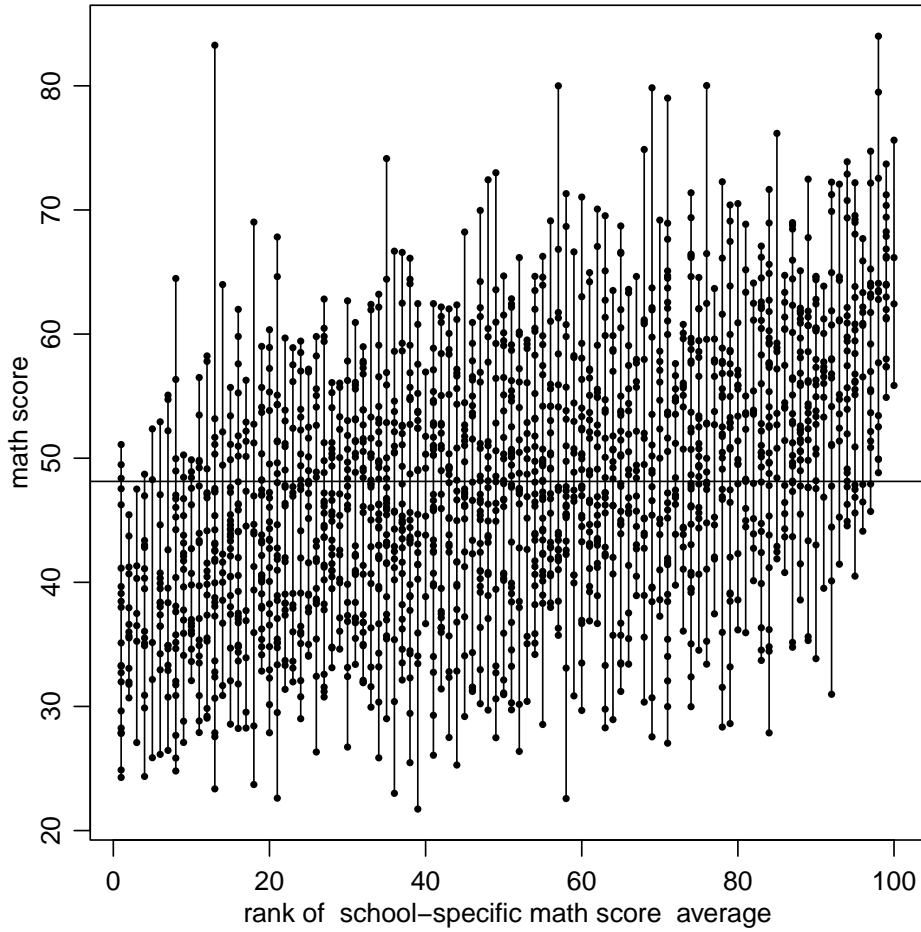
$$(b) \tau^{2(s+1)} \sim p\left(\tau^2 | \theta_1^{(s)}, \dots, \theta_m^{(s)}, \mu^{(s+1)}\right)$$

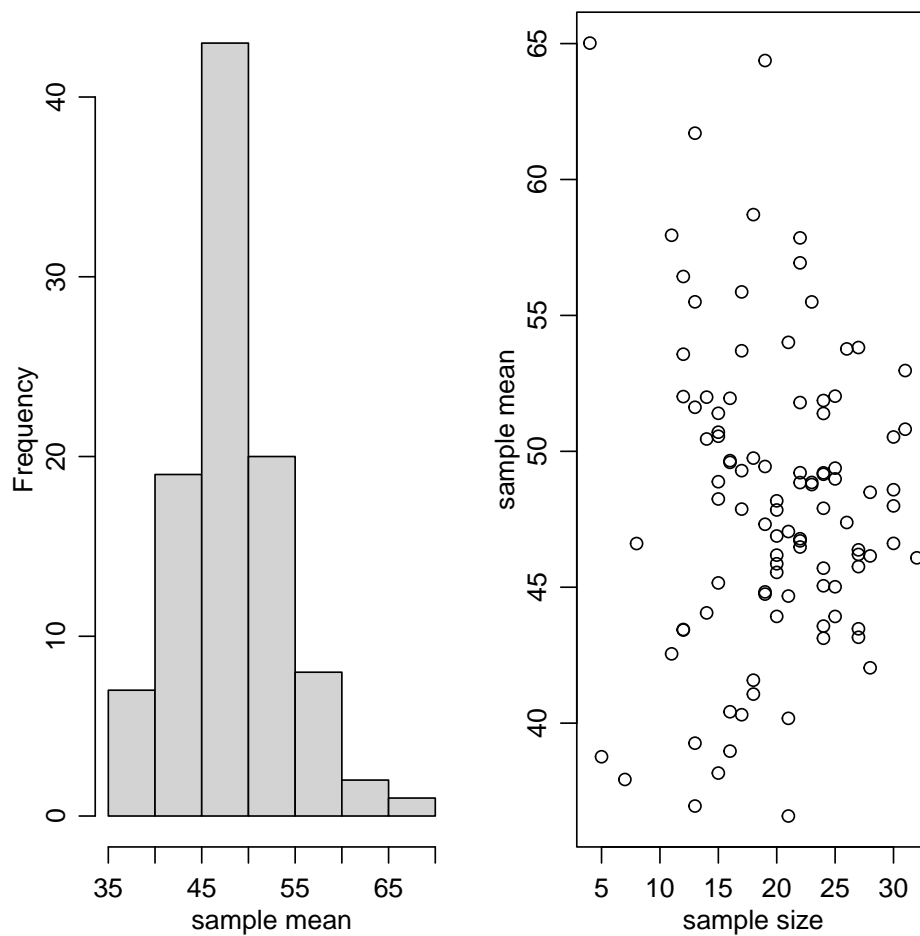
$$(c) \sigma^{2(s+1)} \sim p\left(\sigma^2 | \theta_1^{(s)}, \dots, \theta_m^{(s)}, \mathbf{y}_1, \dots, \mathbf{y}_m\right)$$

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

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.



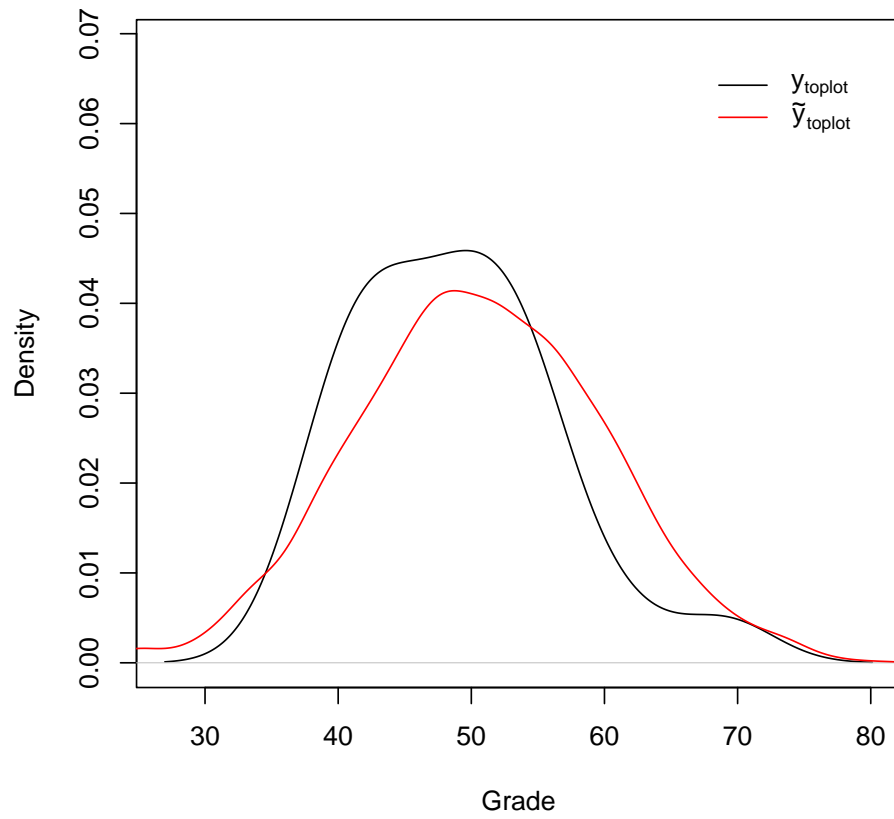


For prediction, we'll use the following prior values (Hoff p. 137):

- $\sigma_0^2 : 100$ (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

4.0.0.1 Derivation

4.0.0.2 R Implementation

4.0.0.3 Example

5 Conclusion