

## A BAYESIAN APPROACH TO SHRINKAGE ESTIMATORS

Filipe das Neves RIZZO<sup>1</sup>  
Cristiane Alvarenga GAJO<sup>1</sup>  
Devanil Jaques de SOUZA<sup>1</sup>  
Lucas Monteiro CHAVES<sup>1</sup>

- *ABSTRACT: Estimators obtained by shrinking the least squares estimator are becoming widely used since the work of Stein, in the early 60's, where it is presented an estimator for the mean of multivariate normal that dominates the sample mean, and the work of Hoerl and Kennard, in the early 70's, on ridge estimators. In this work we present an approach using Bayesian and empirical Bayesian procedures to obtain some important shrinkage estimators.*
- *KEYWORDS: shrinkage estimators, empirical Bayes, minimum mean square estimator.*

### 1 Introduction

Two landmark papers, one by James-Stein (1961), introducing the James-Stein estimators, and another by Hoerl and Kennard (1970), introducing ridge estimators, explore the idea that shrinking the usual least squares estimator, besides the fact of introducing some bias, has the advantage of decreasing the mean square error (MSE). Both these works explore the idea is to shrinking the estimates. This is a conservative procedure used in many areas of statistic. Nowadays, that is known as biased estimation theory and is used when the least squared estimator can't be used in reason, for example, of multicollinearity. The theory of shrinkage estimators is a very active area of research and we can point, as examples, the LASSO theory (TIBSHIRANI, 1996) and Elastic net theory (ZOU-HASTIE, 2005).

In this short note, we present a systematic approach to shrinkage estimators using Bayesian and empirical Bayes procedures. Examples for ridge regression estimation and James-Stein estimator are present. These procedures are very well known and spread in the literature, but authors couldn't find any systematic similar approach. Since empirical Bayes procedure is not of common use, we review, in section 2, a classical example.

We end this work with a computational simulation showing the behavior of some shrinkage estimators.

### Shrinkage estimators

---

<sup>1</sup> University Federal of Lavras – UFLA , Department of Exact Sciences, CP 37 , CEP: 37200 - 000, Lavras, MG, Brazil.: rizzo.filipe@gmail.com; cristianegajo@yahoo.com.br; devaniljaques@dex.ufla.br; lucas@dex.ufla.br

Unbiased estimators of a vector of parameter have a problem that, in general, is not mentioned. Consider  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_k)$  an unbiased estimator for the vector  $\beta = (\beta_1, \dots, \beta_k)$ , then,

$$\begin{aligned} E\left[\|\hat{\beta}\|^2\right] &= E\left[\sum_{i=1}^k \beta_i^2\right] = \sum_{i=1}^k E[\beta_i^2] = \sum_{i=1}^k \left(\text{var}[\beta_i] + (E[\beta_i])^2\right) \\ &= \sum_{i=1}^k \text{var}[\beta_i] + \|\beta\|^2. \end{aligned}$$

We face here an unexpected fact: even though  $\hat{\beta}$  is an unbiased estimator of  $\beta$ ,  $\|\hat{\beta}\|^2$  is a biased estimator for  $\|\beta\|^2$ , tending to overestimate it. This fact has motivated James-Stein to obtain his famous estimator. The idea was to shrink the least squares estimator by a factor that depends of the value observed. Hoerl-Kennard, working with multiple linear regressions in the presence of multicollinearity, proposed an estimation denoted ridge regression estimator, that is a kind of shrinkage estimator. Nowadays there is a huge theory about these estimators.

## 2 Empirical Bayes method

When the knowledge of the prior distribution parameters is incomplete or unknown, some information may be obtained from the data itself and, in this way, it is possible to make inference about these parameters. This is the so-called Bayesian empirical method. To illustrate the empirical Bayes methodology, consider the classical example: Let  $X$  a Poisson random variable with parameter  $\mu$  and prior  $\pi(\mu)$  completely unknown. Given a size one sample  $X = x$ , the corresponding posterior distribution, according to Bayes theorem, is

$$\pi(\mu | x) = \frac{e^{-\mu} \mu^x \pi(\mu)}{g(x)} \text{ where } g(x) = \int_0^{\infty} \frac{e^{-\mu} \mu^x}{x!} \pi(\mu) d\mu.$$

In relation to the mean square error, the Bayes estimator  $\hat{\mu}_B$  for  $\mu$  is the mean of the posterior  $\pi(\mu | x)$ :

$$\begin{aligned} \hat{\mu}_B &= E[\mu | x] = \int_0^{\infty} \mu \frac{e^{-\mu} \mu^x}{x! g(x)} \pi(\mu) d\mu \\ &= \frac{(x+1)! g(x+1)}{x! g(x)} \underbrace{\int_0^{\infty} \frac{e^{-\mu} \mu^{x+1}}{(x+1)! g(x+1)} \pi(\mu) d\mu}_1 = \frac{(x+1)! g(x+1)}{x! g(x)}. \end{aligned}$$

Now, one can use the known past observations  $x_1, \dots, x_n$  to estimate  $g(x)$  and  $g(x+1)$ . Since  $g(x)$  is a discrete distribution it can be estimated by the relative frequencies

$$\hat{g}(x) = \frac{\text{card}\{i; x_i = x\}}{n}$$

$$\hat{g}(x+1) = \frac{\text{card}\{i; x_i = x+1\}}{n}$$

where  $\text{card}\{A\}$  is the number of elements of set A.

Then, the empirical Bayes estimator is

$$\hat{\mu}_B = \frac{(x+1) \hat{g}(x+1)}{\hat{g}(x)}.$$

### 3 An approximation of the minimum mean squares error linear estimator for the mean

Consider independent observations  $Y_i, i = 1, \dots, n$ , with model

$$Y_i = \mu + \varepsilon_i,$$

where  $\mu$  is the parameter to be estimated and the  $\varepsilon_i$ 's are the associated independent errors. The least square estimator for  $\mu$  is

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}_n.$$

Let's consider all estimators of the form  $\delta_c(Y_1, \dots, Y_n) = c\bar{Y}_n$ . We want to determine the particular value of  $c$  that minimizes the mean squared error. To do that we take the function

$$\begin{aligned} MSE(c) &= E\left[(c\bar{Y}_n - \mu)^2\right] = c^2 E\left[\bar{Y}_n^2\right] - 2c\mu E\left[\bar{Y}_n\right] + \mu^2 \\ &= c^2 \left(\frac{\sigma^2}{n} + \mu^2\right) - 2c\mu^2 + \mu^2. \end{aligned}$$

The minimum value of  $MSE(c)$  defines the estimator

$$\hat{\mu} = \frac{\mu^2}{\frac{\sigma^2}{n} + \mu^2} \bar{Y}_n,$$

that depends on population parameters  $\mu$  and  $\sigma^2$ , both unknown. The natural way to bypass this problem is to replace them by their natural unbiased estimators  $\bar{Y}_n$  and  $S^2$ , resulting in the estimator

$$\hat{\mu}_1 = \frac{\bar{Y}_n^2}{\frac{S^2}{n} + \bar{Y}_n^2} \bar{Y}_n = \frac{1}{\frac{S^2}{n\bar{Y}_n^2} + 1} \bar{Y}_n \quad (1)$$

that is an approximation of minimum mean square error linear estimator for the mean.

In section 8 we compare, via simulation, the behavior of these estimators.

#### 4 Shrinkage estimators obtained from normality

The simplest example of a shrinkage estimator obtained from normality is: suppose  $Y \sim N(\theta, 1)$  and prior  $\theta \sim N(0, 1)$ . If a sample  $y = (y_1, \dots, y_n)$  is taken and  $\bar{y}$ , is its mean, the posterior distribution is

$$\begin{aligned} \pi(\theta | y) &\propto f(y | \theta) \pi(\theta) \\ &\propto \exp\left\{-\frac{1}{2} \sum_{i=1}^n (y_i - \theta)^2\right\} \exp\left\{-\frac{1}{2} \theta^2\right\} \\ &\propto \exp\left\{-\frac{1}{2} \frac{1}{n+1} \left(\theta - \frac{\bar{y}}{1+\frac{1}{n}}\right)^2\right\}. \end{aligned}$$

Then the Bayes estimator is  $\hat{\theta}_B = \frac{1}{1+\frac{1}{n}} \bar{y}$ , clearly a shrinkage estimator.

A more sophisticated argument is to obtain the estimator (1) using empirical Bayes approach.

If  $Y \sim \text{normal}(\mu, \sigma^2)$  with priori  $\mu \sim \text{normal}(0, \tau^2)$  and a sample  $y = (y_1, \dots, y_n)$  is observed, the posterior is normal (GRUBER, 2010 - p.50),

$$\pi(\mu | \bar{y}) \propto \exp\left\{-\frac{1}{2} \frac{\tau^2 \sigma^2 / n}{\tau^2 + \sigma^2 / n} \left(\mu - \frac{\tau^2 \bar{y}}{\tau^2 + \sigma^2 / n}\right)^2\right\}.$$

Then, using the mean of the posterior as the Bayes estimator for  $\mu$

$$\hat{\mu}_B = \frac{\tau^2 \bar{Y}_n}{\tau^2 + \sigma^2/n} = \frac{1}{\frac{\sigma^2}{n\tau^2} + 1} \bar{Y}_n.$$

We don't know  $\tau^2$  and the idea is to use empirical Bayes argument. The predictive distribution is a normal with mean zero and variance  $\frac{\sigma^2}{n} + \tau^2$  (appendix). Observing that if  $Z \sim normal(0, \sigma^2)$  as  $\sigma^2 = E[Z^2] - (E[Z])^2 = E[Z^2]$ , by the methods of moments,  $\sigma^2$  can be estimated by only one observation,  $\hat{\sigma}^2 = Z^2$ . Then we can estimate  $\frac{\sigma^2}{n} + \tau^2$  by  $\bar{Y}^2$ . It is reasonable to suppose  $\tau^2 \gg \frac{\sigma^2}{n}$ . Then we can use  $\bar{Y}^2$  as estimator of  $\tau^2$ . Replacing the parameters in  $\hat{\mu}_B$  by its respective estimators we get again the estimator (1).

The estimator (1) can also be obtained using mixture. Let's suppose that, in the model  $Y_i = \mu + \varepsilon_i$ ,  $Y_i \sim normal(\mu, \sigma^2)$  and  $\mu$  has a prior distribution  $normal(0, \tau^2)$ . We want to compute the mean square error  $E[(c\bar{Y}_n - \mu)^2] = E[c^2\bar{Y}_n^2 - 2c\bar{Y}_n\mu + \mu^2]$ , considering both  $\bar{Y}_n$  and  $\mu$  as random variables.

Computing these expected values:

$$f_{\bar{Y}_n}(\bar{y}) = \frac{1}{\sqrt{2\pi \frac{\sigma^2}{n}}} \exp\left\{-\frac{1}{2\frac{\sigma^2}{n}}(\bar{y} - \mu)^2\right\}$$

$$f_{\mu}(\mu) = \frac{1}{\sqrt{2\pi\tau^2}} \exp\left\{-\frac{1}{2\tau^2}\mu^2\right\}$$

$$E[c^2\bar{Y}_n^2 - 2c\bar{Y}_n\mu + \mu^2] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (c^2\bar{y}^2 - 2c\bar{y}\mu + \mu^2) f_{\bar{Y}_n}(\bar{y}) f_{\mu}(\mu) d\mu d\bar{y}.$$

Computing each integral:

$$E[c^2\bar{Y}_n^2] = c^2 \int_{-\infty}^{\infty} \bar{y}^2 \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi \frac{\sigma^2}{n}}} \exp\left\{-\frac{1}{2\frac{\sigma^2}{n}}(\bar{y} - \mu)^2\right\} \frac{1}{\sqrt{2\pi\tau^2}} \exp\left\{-\frac{1}{2\tau^2}\mu^2\right\} d\mu d\bar{y}$$

$$= c^2 \int_{-\infty}^{\infty} \bar{y}^2 \frac{1}{\sqrt{2\pi\left(\frac{\sigma^2}{n} + \tau^2\right)}} \exp\left\{-\frac{1}{2\left(\frac{\sigma^2}{n} + \tau^2\right)} \bar{y}^2\right\} d\bar{y}$$

$$= c^2 \left(\frac{\sigma^2}{n} + \tau^2\right).$$

$$E[\mu \bar{Y}_n] = \int_{-\infty}^{\infty} \mu \left[ \int_{-\infty}^{\infty} \bar{y} \frac{1}{\sqrt{2\pi\frac{\sigma^2}{n}}} \exp\left\{-\frac{1}{2\frac{\sigma^2}{n}}(\bar{y} - \mu)^2\right\} \frac{1}{\sqrt{2\pi\tau^2}} \exp\left\{-\frac{1}{2\tau^2}\mu^2\right\} d\bar{y} \right] d\mu$$

$$= \int_{-\infty}^{\infty} \mu \frac{1}{\sqrt{2\pi\tau^2}} \exp\left\{-\frac{1}{2\tau^2}\mu^2\right\} \underbrace{E[\bar{y}]}_{\mu} d\mu.$$

$$E[\mu^2] = Var[\mu] + \underbrace{(E[\mu])^2}_0 = \tau^2.$$

All together:

$$E\left[c^2 \bar{Y}^2 - 2c\bar{Y}\mu + \mu^2\right] = c^2 \left(\frac{\sigma^2}{n} + \tau^2\right) - 2c\tau^2 + \tau^2.$$

So, the minimum value of the mean squared error  $E\left[(c\bar{Y} - \mu)^2\right]$  is obtained when

$c_0^* = \frac{\tau^2}{\frac{\sigma^2}{n} + \tau^2}$ . Then, the MSE estimator for  $\mu$  is:

$$\hat{\mu}_2 = \frac{\tau^2}{\frac{\sigma^2}{n} + \tau^2} \bar{Y}_n.$$

Replacing the parameters by its estimators we, again, get the estimator (1).

## 5 A Bayesian interpretation of ridge regression estimators

Consider the multiple linear regression

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \quad \boldsymbol{\varepsilon}_{n \times 1} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_{n \times n}).$$

We will suppose  $\sigma^2$  known and a priori distribution for the parameter vector  $\boldsymbol{\beta}_{p \times 1}$  given by a  $p$ -variate normal

$$\boldsymbol{\beta}_{p \times 1} | \sigma^2 \sim N_p\left(\mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}\right).$$

The posterior distribution given data  $\mathbf{Y}$  and design matrix  $\mathbf{X}$  is

$$\begin{aligned} f_{\boldsymbol{\beta}}(\boldsymbol{\beta}, \sigma^2 | \mathbf{Y}, \mathbf{X}, \lambda) &\propto f_{\mathbf{Y}}(\mathbf{Y} | \mathbf{X}, \boldsymbol{\beta}, \sigma^2) f_{\boldsymbol{\beta}}(\boldsymbol{\beta} | \sigma^2, \lambda) \\ &\propto \sigma^{-n} \exp\left[-\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\right] \\ &\propto \sigma^{-n} \exp\left[-\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\right] \sigma^{-p} \exp\left[-\frac{\lambda}{2\sigma^2} \boldsymbol{\beta}' \boldsymbol{\beta}\right] \\ &\propto \sigma^{-(n+p)} \exp\left[-\frac{1}{2\sigma^2} \left[(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) - \lambda \boldsymbol{\beta}' \boldsymbol{\beta}\right]\right] \\ &\propto \sigma^{-(n+p)} \exp\left[-\frac{1}{2\sigma^2} \left[\mathbf{Y}' \mathbf{Y} - \mathbf{Y}' \mathbf{X} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{X}' \mathbf{Y} + \boldsymbol{\beta}' \mathbf{X}' \mathbf{X} \boldsymbol{\beta} + \lambda \boldsymbol{\beta}' \boldsymbol{\beta}\right]\right] \\ &\propto \sigma^{-(n+p)} \exp\left[-\frac{1}{2\sigma^2} \left[\mathbf{Y}' \mathbf{Y} - \mathbf{Y}' \mathbf{X} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{X}' \mathbf{Y} + \boldsymbol{\beta}' [\mathbf{X}' \mathbf{X} \boldsymbol{\beta} + \lambda \mathbf{I}] \boldsymbol{\beta}\right]\right]. \end{aligned}$$

To explicit the distribution it is necessary to complete the square in

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}' \boldsymbol{\beta} &= \mathbf{Y}' \mathbf{Y} - \mathbf{Y}' \mathbf{X} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{X}' \mathbf{Y} + \boldsymbol{\beta}' \mathbf{X}' \mathbf{X} \boldsymbol{\beta} + \lambda \boldsymbol{\beta}' \boldsymbol{\beta} \\ &= \mathbf{Y}' \mathbf{Y} - 2\mathbf{Y}' \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\beta}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta}. \end{aligned}$$

So, we need to determine  $\mathbf{W}$  such that

$$\begin{aligned} (\boldsymbol{\beta} - \mathbf{W})' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) (\boldsymbol{\beta} - \mathbf{W}) &= \mathbf{Y}' \mathbf{Y} - 2\mathbf{Y}' \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\beta}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} \\ &= \boldsymbol{\beta}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} - 2\mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} + \mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \mathbf{W}. \end{aligned}$$

Then  $\mathbf{W} = (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}' \mathbf{Y}$  and  $\mathbf{Y}' \mathbf{X} \boldsymbol{\beta} = \mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta}$ .

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}' \boldsymbol{\beta} &= \mathbf{Y}' \mathbf{Y} + \boldsymbol{\beta}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} - 2\mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} \\ &= \mathbf{Y}' \mathbf{Y} + \boldsymbol{\beta}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} - 2\mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \boldsymbol{\beta} \\ &\quad + \mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \mathbf{W} - \mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \mathbf{W} \\ &= (\boldsymbol{\beta} - \mathbf{W})' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) (\boldsymbol{\beta} - \mathbf{W}) + \mathbf{Y}' \mathbf{Y} - \mathbf{W}' (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}) \mathbf{W}. \end{aligned}$$

Therefore

$$f_{\beta}(\beta, \sigma^2 | \mathbf{Y}, X) \propto \exp\left[(\beta - W)'(X'X + \lambda I)(\beta - W) + Y'Y - W'(X'X + \lambda I)W\right] \\ \propto (\beta - W)'(X'X + \lambda I)(\beta - W).$$

The posterior is a multivariate normal with mean  $W = (X'X + \lambda I)^{-1} X'Y$ , therefore the Bayes estimator is

$$\hat{\beta}(\lambda) = (X'X + \lambda I)^{-1} X'Y.$$

This is the ridge regression estimator with ridge parameter  $\lambda$  (Hoerl-Kennard, 1970). For  $\lambda = 0$  we have the ordinary least square estimator and  $\|\hat{\beta}(\lambda)\| < \|\hat{\beta}(0)\|$  and, for  $\lambda > 0$   $\hat{\beta}(\lambda)$  is a shrinkage estimator.

## 6 The James-Stein estimator as an empirical Bayes estimator

Consider the random normal vector  $X = (X_1, \dots, X_p)$ ,  $X \sim N_p(\theta, I)$  and suppose independent normal priors distributions for the parameters:  $\theta_i \sim N(0, \sigma^2)$ . If one observation  $X = (X_1, \dots, X_p)$  is taken, the Bayes estimator is

$$\hat{\theta}_i = \frac{\sigma^2}{\sigma^2 + 1} X_i = \left(1 - \frac{1}{\sigma^2 + 1}\right) X_i.$$

As  $\sigma^2$  is unknown, using an empirical Bayes procedure it can be estimate from the data  $X = (X_1, \dots, X_p)$ . The predictive distribution is a multivariate independent normal with variance  $\sigma^2 + 1$ ,  $N_p(0, (\sigma^2 + 1)I)$ . Consider the random variables

$$Z_i = \frac{X_i}{\sqrt{\sigma^2 + 1}} \Rightarrow Z_i \sim N(0, 1)$$

$$\sum_{i=1}^p Z_i^2 = \frac{\sum_{i=1}^p X_i^2}{(\sqrt{\sigma^2 + 1})^2} = \frac{X'X}{\sigma^2 + 1}.$$

As sum of square  $p$  independent standard normal, it has chi-square distribution with parameter  $p$ , that is  $Y = \frac{X'X}{\sigma^2 + 1} \sim \chi^2(p)$ .

As  $Y = \frac{X'X}{\sigma^2 + 1} \Leftrightarrow \frac{1}{X'X} = \frac{1}{Y} \frac{1}{\sigma^2 + 1}$ , taking expectation,

$$E\left[\frac{1}{X'X}\right] = E\left[\frac{1}{(\sigma^2 + 1)Y}\right] = \frac{1}{\sigma^2 + 1} E\left[\frac{1}{Y}\right].$$

As



$$Y \sim \text{Gama}\left(\alpha = \frac{p}{2}, \beta = \frac{1}{2}\right)$$

$$\begin{aligned} E\left[\frac{1}{Y}\right] &= \int_0^{\infty} \frac{1}{y} \underbrace{\frac{\left(\frac{1}{2}\right)^{\frac{p}{2}}}{\Gamma\left(\frac{p}{2}\right)} y^{\left(\frac{p}{2}-1\right)} e^{-\frac{1}{2}y}}_{f_Y(y)} dy \\ &= \frac{\Gamma\left(\frac{p}{2}-1\right)}{\Gamma\left(\frac{p}{2}\right)\left(\frac{1}{2}\right)^{-1}} \underbrace{\int_0^{\infty} \frac{\left(\frac{1}{2}\right)^{\frac{p-1}{2}}}{\Gamma\left(\frac{p}{2}-1\right)} y^{\left(\frac{p}{2}-1\right)-1} e^{-\frac{1}{2}y} dy}_{\text{Gama}\left(\alpha=\frac{p-1}{2}, \beta=\frac{1}{2}\right)} \\ &= \frac{\Gamma\left(\frac{p}{2}-1\right)}{\Gamma\left(\frac{p}{2}\right)\left(\frac{1}{2}\right)^{-1}} = \frac{\Gamma\left(\frac{p}{2}-1\right)}{2\left(\frac{p}{2}-1\right)\Gamma\left(\frac{p}{2}-1\right)} = \frac{1}{2\left(\frac{p}{2}-1\right)} = \frac{1}{p-2}. \end{aligned}$$

It follows that,

$$E\left[\frac{1}{XX}\right] = \frac{1}{\sigma^2+1} \frac{1}{p-2}$$

and, therefore,

$$E\left[\frac{p-2}{XX}\right] = \frac{1}{\sigma^2+1}.$$

By the methods of moments we can estimate  $\frac{1}{\sigma^2+1}$  by  $\frac{p-2}{XX}$ . As

$$\hat{\theta}_i = \frac{\sigma^2}{\sigma^2+1} x_i = \left(1 - \frac{1}{\sigma^2+1}\right) x_i,$$

replacing  $\frac{1}{\sigma^2+1}$  by  $\frac{p-2}{XX}$ , the result is the famous James - Stein estimator (James-Stein, 1961)

$$\left(\hat{\theta}_{JS}\right)_i = \left(1 - \frac{p-2}{XX}\right) X_i$$

in vectorial form

$$\hat{\theta}_{JS} = \left(1 - \frac{p-2}{XX}\right) X.$$

Observe that the James - Stein estimator is a shrinkage estimator for  $p > 2$ .

## 7 A non-normal example

In the example (CASELLA-BERGER, 2010, pp. 295), suppose  $X_1, \dots, X_n$  i.i.d. *Bernoulli*( $p$ ). It is well known that the maximum likelihood and unbiased estimator for the parameter  $p$  is

$$\hat{p} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i .$$

The MSE for  $\hat{p}$  is

$$\begin{aligned} E\left[(\hat{p} - p)^2\right] &= \text{var}[\hat{p}] - \left(E[(\hat{p} - p)]\right)^2 \\ &= \text{var}\left[\bar{X}\right] = \frac{p(1-p)}{n} . \end{aligned}$$

To compute the Bayes estimator we consider that  $Y = \sum_{i=1}^n X_i$  is *Binomial*( $n, p$ )

, and use a *Beta*( $\alpha, \beta$ ) as the priori for  $p$ . Then, the posteriori is given by:

$$\pi(p | y) = \frac{1}{B(y + \alpha, n - y + \beta)} p^{y + \alpha - 1} (1 - p)^{n - y + \beta - 1} .$$

So, the posterior is a *Beta*( $y + \alpha, n - y + \beta$ ) and the Bayes estimator for the parameter  $p$  is its mean:

$$\hat{p}_B = \frac{y + \alpha}{\alpha + \beta + n} .$$

The mean squared error for this estimator is:

$$E\left[(\hat{p}_B - p)^2\right] = \text{Var}[\hat{p}_B] + \left(E[\hat{p}_B - p]\right)^2 = \frac{np(1-p)}{(\alpha + \beta + n)^2} + \left(\frac{np + \alpha}{\alpha + \beta + n} - p\right)^2 .$$

Observe that, taking  $\alpha = \beta = \sqrt{\frac{n}{4}}$ , we get that the MSE is constant, that is, it

doesn't depend of the parameter  $p$ . In this case the estimator  $\hat{p}_B = \frac{y + \sqrt{\frac{n}{4}}}{\sqrt{n} + n}$  is minimax (MOOD et al, 1963, theorem 14, page 350) and

$$E\left[(\hat{p}_B - p)^2\right] = \frac{np(1-p)}{(\sqrt{n}+n)^2} + \left(\frac{np + \sqrt{\frac{n}{4}}}{\sqrt{n}+n} - p\right)^2$$

$$= \frac{np(1-p)}{(\sqrt{n}+n)^2} + \frac{\left[\left(np + \sqrt{\frac{n}{4}}\right) - 2p\sqrt{\frac{n}{4}} - np\right]^2}{(\sqrt{n}+n)^2} = \frac{n}{4(n+\sqrt{n})^2}.$$

Now, as  $E\left[(\hat{p}_B - p)^2\right] = \frac{n}{4(n+\sqrt{n})^2}$  and  $E\left[(\hat{p} - p)^2\right] = \frac{p(1-p)}{n}$ .

Figure 1 helps to choose between these two estimators.

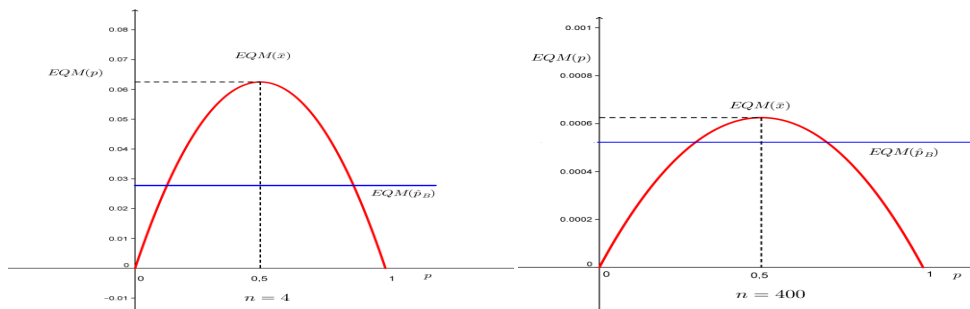


Figure 1 - Comparison of MSE of  $\hat{p}$  (blue) and  $\hat{p}_B$  (red) for different sample sizes

For small values of  $n$ , clearly  $\hat{p}_B$  is a better choice, except if there is a strong reason to believe that  $p$  is near 0 or 1. For high values of  $n$ , the choice is the estimator  $\hat{p}$ , except for a strong belief that  $p$  is near  $1/2$ .

### 8 Comparing the estimator $\hat{\mu}_1$ and the pseudo-estimator $\hat{\mu}$

To emphasize the level of complexity implied by the shrinking procedure we compare the behavior of the pseudo estimator  $\hat{\mu} = \frac{\mu^2}{\frac{\sigma^2}{n} + \mu^2} \bar{Y}_n$  and the shrinkage estimator

$$\hat{\mu}_1 = \frac{1}{\frac{S^2}{n\bar{Y}_n} + 1} \bar{Y}_n .$$

Considering the distributions Normal, Uniform and Double Exponential, we compute the mean squared error using 10.000 samples of sizes 20, 40 and 100 with variances 4, 16 and 36. The results are shown in Figures 2, 3 and 4.

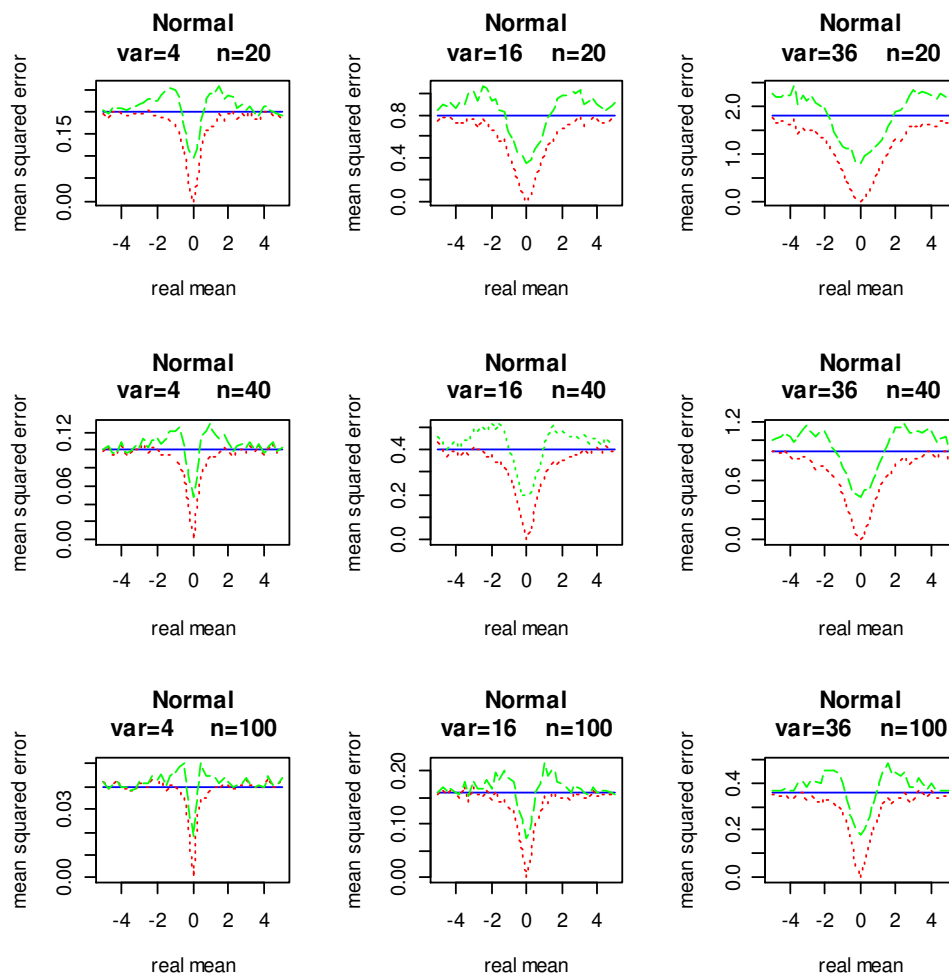


Figure 2 - MSE for  $\bar{Y}_n$  (blue – solid),  $\hat{\mu}$  (red – dotted) and  $\hat{\mu}_1$  (green – long dash) for Normal distribution with different variances and sample sizes.

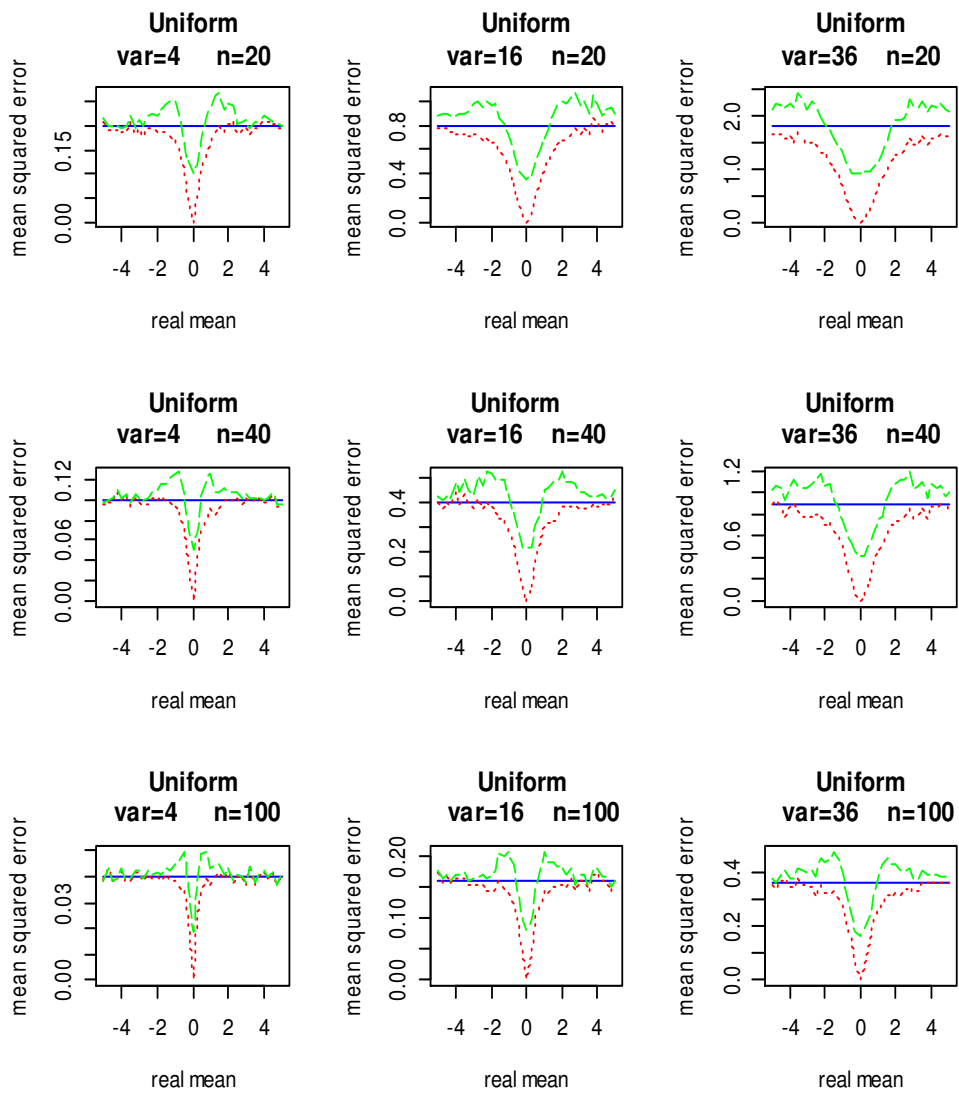


Figure 3 - MSE for  $\bar{Y}_n$  (blue – solid),  $\hat{\mu}$  (red – dotted) and  $\hat{\mu}_1$  (green – long dash) for Uniform distribution with different variances and sample sizes.

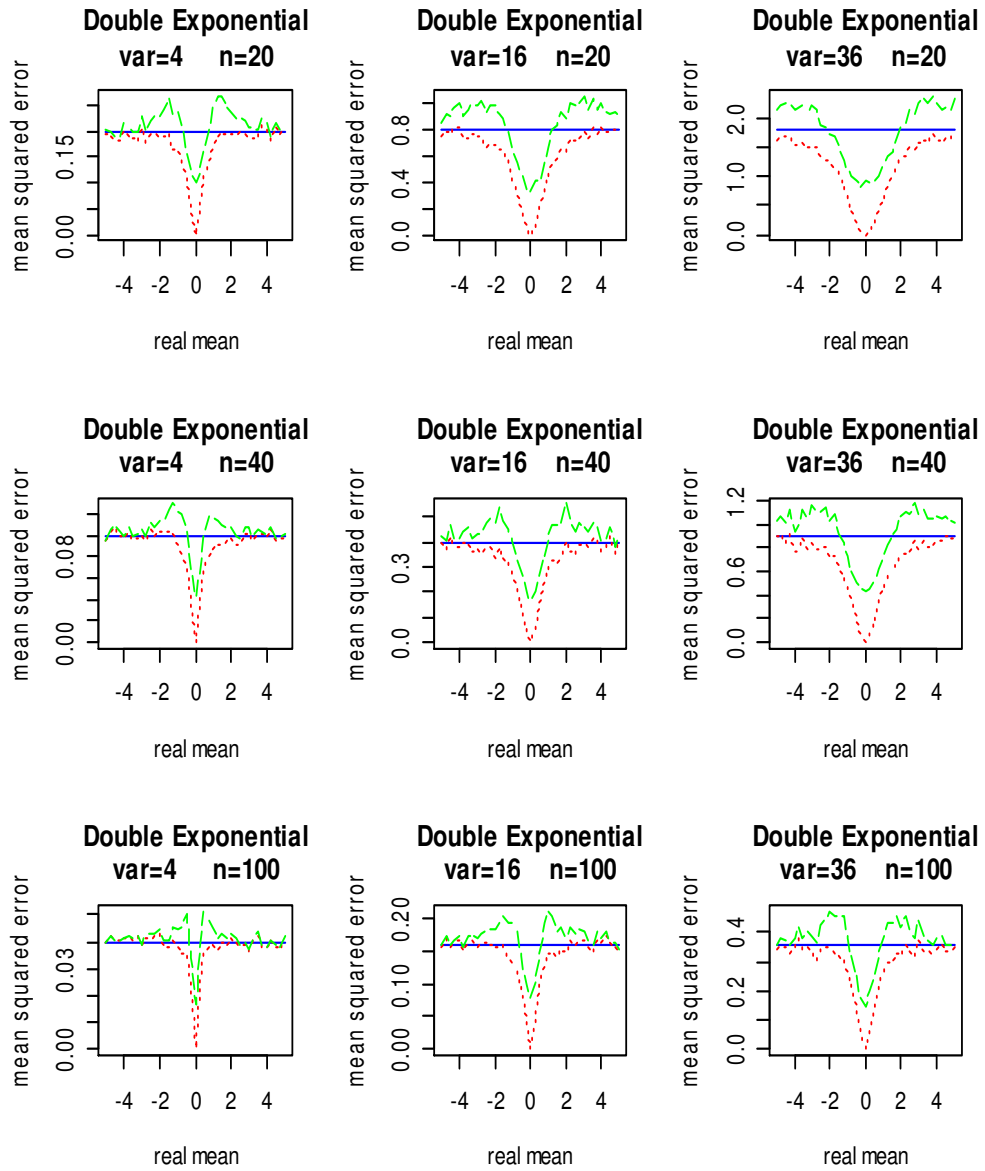


Figure 4 - MSE for  $\bar{Y}_n$  (blue – solid),  $\hat{\mu}$  (red – dotted) and  $\hat{\mu}_1$  (green – long dash) for Double Exponential distribution with different variances and sample sizes.

About figures 2, 3 and 4, observe that the pseudo-estimator  $\hat{\mu}$  behaves as it is supposed to, that is, except for some fluctuation when the real mean is far from the origin, it has MSE lower than the estimator  $\bar{Y}_n$ . On the other hand, the estimator  $\hat{\mu}_1$  has a complex behavior and only for the real mean near the origin it has MSE lower than the estimator  $\bar{Y}_n$ . We couldn't find or guess a good reason for that strange behavior.

## Conclusion

The shrinkage procedure shows to be fully compatible with Bayesian and empirical Bayesian method. The major problem, when shrinking an estimator, is to decide the amount of reduction to be applied. The Bayesian method can help, given an interpretation of the necessary reduction to be used. The complexity introduced by shrinkage is far from been completely understood.

RIZZO, F. N.; GAJO, C. A.; SOUZA, D. J.; CHAVES, L. M. Uma abordagem Bayesiana do estimador de encolhimento. *Rev. Bras. Biom.* São Paulo, v.33, n.4, p.485-602, 2015.

- RESUMO: Estimadores obtidos por encolhimento do estimador de mínimos quadrados usual têm sido amplamente utilizados, desde os trabalhos de James-Stein (1961) e Hoerl e Kennard (1970), sobre estimadores de cumeieira. Neste trabalho, apresentamos uma abordagem para alguns importantes estimadores de encolhimento, utilizando procedimentos Bayesianos e Bayesianos empíricos.
- PALAVRAS-CHAVE: Estimador de encolhimento, Bayes empírico, estimador de quadrados mínimos.

## References

- CARLIN, B. P.; LOUIS, T. A.. *Bayes and empirical Bayes methods for data analysis*. New York:Chapman & Hall/CRC, 2000.
- CASELLA, G.; BERGER, R. L. *Inferência estatística*. 2. ed. São Paulo: Cengage Learning, 2010. 573p.
- GRUBER, M. H. J. *Improving efficiency by shrinkage: the James-Stein and ridge regression estimators*. New York: CRC, 1998. 648 p.
- GRUBER, M. H. J. *Regression estimators: a comparative study*. Baltimore: The Johns Hopkins Press/CRC, 2010. 412p.
- HOERL, A. E.; KENNARD, R. W. , Biased Estimation for Nonorthogonal Problems, *Technometrics*, v. 12, n. 1, p. 55-67, 1970.
- MOOD, A. M.; BOES, F. A.; GRAYBILL, F. A. *Introduction to the theory of statistics*. 3rd ed. Columbus: McGraw-Hill, 1974. 564p.
- JAMES, W.; STEIN, C. Estimation with Quadratic Loss. Proc. 4<sup>th</sup> Berkeley Sympos. Math. Statist. And Prob., v. 1, 1961, Univ. California Press, Berkeley, Calif., p.361-379, 1961.

TIBSHIRANI, R. Regression shrinkage and selection via the Lasso. *J.R. Statis. Soc.*, v.58, n.1, p.267-288, 1996.

ZOU, H.; HASTIE, T. Regularization and variable selection via the elastic net. *J.R. Statis. Soc.*, v. 68, n.2, p.301-320, 2005.

Received in 11.05.2015

Approved after revised in 12.08.2015



## Appendix

The predictive for normal and mean prior normal

$$\begin{aligned}
 f(\bar{y}) &= \int_{-\infty}^{\infty} (\bar{y} | \mu) \pi(\mu) d\mu \\
 &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2/n}} \exp\left\{-\frac{1}{2\sigma^2/n}(\bar{y}-\mu)^2\right\} \frac{1}{\sqrt{2\pi\tau^2}} \exp\left\{-\frac{1}{2\tau^2}\mu^2\right\} d\mu \\
 &= \frac{1}{\sqrt{2\pi/n\sigma\tau}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{n}{2\sigma^2}(\bar{y}^2 - 2\bar{y}\mu + \mu^2) - \frac{1}{2\tau^2}\mu^2\right\} d\mu \\
 &= \frac{1}{\sqrt{2\pi/n\sigma\tau}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{\mu^2 \left[ \left(-\frac{1}{2}\right) \left(\frac{n}{\sigma^2} + \frac{1}{\tau^2}\right) \right] + \frac{2n\bar{y}\mu}{2\sigma^2} - \frac{n\bar{y}^2}{2\sigma^2}\right\} d\mu \\
 &= \frac{1}{\sqrt{2\pi/n\sigma\tau}} \exp\left\{-\frac{n\bar{y}^2}{2\sigma^2}\right\} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{1}{\left(\frac{\sigma^2\tau^2}{\sigma^2+n\tau^2}\right)} \left[\mu^2 - \frac{2n\bar{y}\mu\tau^2}{(\sigma^2+n\tau^2)}\right]\right\} d\mu \\
 &= \frac{1}{\sqrt{2\pi/n\sigma\tau}} \exp\left\{-\frac{n\bar{y}^2}{2\sigma^2}\right\} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{1}{\left(\frac{\sigma^2\tau^2}{\sigma^2+n\tau^2}\right)} \left\{ \left[\mu - \frac{n\bar{y}\tau^2}{(\sigma^2+n\tau^2)}\right]^2 - \left(\frac{n\bar{y}\tau^2}{(\sigma^2+n\tau^2)}\right)^2 \right\}\right\} d\mu \\
 &= \frac{1}{\sqrt{\frac{2\pi\sigma^2\tau^2}{n}}} \exp\left\{-\frac{n\bar{y}^2}{2\sigma^2} + \frac{n^2\bar{y}^2\tau^4}{(\sigma^2+n\tau^2)2\sigma^2\tau^2}\right\} \sqrt{\frac{\sigma^2\tau^2}{\sigma^2+n\tau^2}} \\
 &\quad \int \frac{1}{\sqrt{2\pi\left(\frac{\sigma^2\tau^2}{\sigma^2+n\tau^2}\right)}} \exp\left\{-\frac{1}{2} \frac{1}{\left(\frac{\sigma^2\tau^2}{\sigma^2+n\tau^2}\right)} \left[\mu - \frac{n\bar{y}\tau^2}{(\sigma^2+n\tau^2)}\right]^2\right\} d\mu
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\sqrt{2\pi\left(\frac{\sigma^2+n\tau^2}{n}\right)}} \exp\left\{-\frac{1}{2}\left[\frac{n\sigma^2\tau^2+n^2\tau^4-n^2\tau^4}{(\sigma^2+n\tau^2)\sigma^2\tau^2}\right]\bar{y}^2\right\} \\
&= \frac{1}{\sqrt{2\pi\left(\frac{\sigma^2+n\tau^2}{n}\right)}} \exp\left\{-\frac{1}{2}\left[\frac{1}{\left(\frac{\sigma^2+n\tau^2}{n}\right)}\right]\bar{y}^2\right\} \\
&= \frac{1}{\sqrt{2\pi\left(\frac{\sigma^2}{n}+\tau^2\right)}} \exp\left\{-\frac{1}{2\left(\frac{\sigma^2}{n}+\tau^2\right)}\bar{y}^2\right\}.
\end{aligned}$$