

Chapter 7. Biased Estimation

7.1 James - Stein Shrinkage Method

(1) Motivation

Let $\mathbf{z} \sim N_p(\boldsymbol{\mu}, \mathbf{I})$, then \mathbf{z} is a good estimator for $\boldsymbol{\mu}$ because it is MLE and UMVUE (uniformly minimum variance unbiased estimator). Even though \mathbf{z} is a good estimator for $\boldsymbol{\mu}$ in the sense of 1st moment (i.e., $E(\mathbf{z}) = \boldsymbol{\mu}$), however, James and Stein (1961), considered the 2nd moment (i.e., MSE : mean squared error);

$$E(\mathbf{z}'\mathbf{z}) = \sum_{i=1}^p E(z_i^2) = \sum (1 + \mu_i^2) = p + \boldsymbol{\mu}'\boldsymbol{\mu} > \boldsymbol{\mu}'\boldsymbol{\mu}.$$

Hence, $\mathbf{z}'\mathbf{z}$ is not an unbiased estimator of $\boldsymbol{\mu}'\boldsymbol{\mu}$. Based on this idea, they considered

$$\tilde{\boldsymbol{\mu}} = c\mathbf{z}, \quad 0 < c < 1$$

where c is called *shrinkage* constant.

(2) James-Stein estimator

James and Stein showed that $\exists \quad 0 < c < 1$ s.t. $\text{MSE}(\tilde{\boldsymbol{\mu}}) \leq \text{MSE}(\mathbf{z})$. We will

prove it. Note that

$$\begin{aligned}
\text{MSE}(\tilde{\boldsymbol{\mu}}) &= E[(\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu})'(\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu})] \\
&= E[(\tilde{\boldsymbol{\mu}} - E(\tilde{\boldsymbol{\mu}}) + E(\tilde{\boldsymbol{\mu}}) - \boldsymbol{\mu})'(\tilde{\boldsymbol{\mu}} - E(\tilde{\boldsymbol{\mu}}) + E(\tilde{\boldsymbol{\mu}}) - \boldsymbol{\mu})] \\
&= E[(\tilde{\boldsymbol{\mu}} - E(\tilde{\boldsymbol{\mu}}))'(\tilde{\boldsymbol{\mu}} - E(\tilde{\boldsymbol{\mu}}))] + [E(\tilde{\boldsymbol{\mu}}) - \boldsymbol{\mu}]'[E(\tilde{\boldsymbol{\mu}}) - \boldsymbol{\mu}] \\
&= \text{tr}[\text{Cov}(\tilde{\boldsymbol{\mu}})] + ||\text{Bias}(\tilde{\boldsymbol{\mu}})||^2
\end{aligned}$$

and therefore,

$$\begin{aligned}
\text{MSE}(\tilde{\boldsymbol{\mu}}) &= \text{tr}(c^2 \mathbf{I}) + (\boldsymbol{\mu} - E(\tilde{\boldsymbol{\mu}}))'(\boldsymbol{\mu} - E(\tilde{\boldsymbol{\mu}})) \\
&= pc^2 + (1 - c)^2 \boldsymbol{\mu}' \boldsymbol{\mu}.
\end{aligned}$$

On the other hand,

$$\text{MSE}(\mathbf{z}) = p$$

To find c minimizing $\text{MSE}(\tilde{\boldsymbol{\mu}})$, we take 1st derivative w.r.t. c , i.e.,

$$0 = \frac{\partial \text{MSE}(\tilde{\boldsymbol{\mu}})}{\partial c} = 2cp - 2(1 - c)\boldsymbol{\mu}' \boldsymbol{\mu} = 2c(p + \boldsymbol{\mu}' \boldsymbol{\mu}) - 2\boldsymbol{\mu}' \boldsymbol{\mu}$$

which gives

$$c = \frac{\boldsymbol{\mu}' \boldsymbol{\mu}}{p + \boldsymbol{\mu}' \boldsymbol{\mu}}$$

and then, we have

$$\tilde{\boldsymbol{\mu}} = \left(1 - \frac{p}{p + \boldsymbol{\mu}' \boldsymbol{\mu}}\right) \mathbf{z}$$

But, this estimator contains unknown parameter $\mu'\mu$, we replace it by its unbiased estimator $z'z - p$, i.e.,

$$\tilde{\mu} = \left(1 - \frac{p}{z'z}\right) z$$

Finally, James and Stein showed that

$$\tilde{\mu}_s = \left[1 - \frac{(p-2)}{z'z}\right] z$$

gives the minimum MSE, and it is called *Stein shrinkage* estimator.

7.2 Ridge Regression

7.2.1 Inference on the ridge regression

(1) Motivation and definition

The LSE $\hat{\beta} = (X'X)^{-1}X'y$ will be very unstable if there exists multicollinearity. Hoerl and Kennard (1970) suggested the *ridge estimator* defined as

$$\hat{\beta}(\theta) = (X'X + \theta I)^{-1}X'y,$$

where θ is called *biasing parameter* or *shrinkage parameter*

(2) Estimation of θ

(a) ridge trace

First, we draw a figure of $\hat{\beta}_j(\theta)$, $j = 0, \dots, p - 1$ for $\theta > 0$, and we estimate θ where $\hat{\beta}_j(\theta)$ is stabilized. This method is called *ridge trace*.

Fig. 7.1 (p. 280)

(b) GCV_θ

Wahba, Golub, and Heath (1979) suggested minimizing the GCV (generalized cross validation), defined as

$$GCV_\theta = \sum e_{i,\theta}^2 / \left[1 - \frac{1}{n} \text{tr}(\mathbf{H}_\theta) \right]^2,$$

where $e_{i,\theta}$ is residual and $\mathbf{H}_\theta = \mathbf{X}(\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\mathbf{X}'$ is a hat matrix.

Ex.7.1 (p.280)

7.2.2 Properties of ridge estimator

(1) Regularized LSE

We can show that

$$\hat{\beta}(\theta) = \text{Arg}_{\beta} \min (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \text{ s.t. } \beta'\beta = c^2$$

i.e.,

$$\hat{\beta}(\theta) = \text{Arg}_{\beta} \min (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \lambda(\beta'\beta - c^2),$$

where λ is a lagrangian multiplier.

Fig. 7.2 (p. 283)

(2) MSE

Theorem 7.1. $\text{MSE}(\hat{\beta}) - \text{MSE}(\hat{\beta}(\theta))$ is positive definite for some $\theta > 0$.

(Proof) Since

$$\begin{aligned} \text{MSE}(\hat{\beta}(\theta)) &= E[(\hat{\beta}(\theta) - \beta)(\hat{\beta}(\theta) - \beta)'] \\ &= \text{Cov}[(\hat{\beta}(\theta))] + [E(\hat{\beta}(\theta)) - \beta][E(\hat{\beta}(\theta)) - \beta]' \end{aligned}$$

and let $\mathbf{K} = (\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}$, then

$$\begin{aligned} \text{Cov}[\hat{\beta}(\theta)] &= \text{Cov}[(\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}] \\ &= (\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\sigma^2 \\ &= \mathbf{K}\mathbf{X}'\mathbf{X}\mathbf{K}\sigma^2 \end{aligned}$$

and

$$E(\hat{\beta}(\theta)) - \beta = E(\mathbf{K}\mathbf{X}'\mathbf{y}) - \beta = (\mathbf{K}\mathbf{X}'\mathbf{X} - \mathbf{K}\mathbf{K}^{-1})\beta = -\theta\mathbf{K}\beta$$

we have

$$\text{MSE}(\hat{\beta}(\theta)) = \mathbf{K}\mathbf{X}'\mathbf{X}\mathbf{K}\sigma^2 + \theta^2\mathbf{K}\beta\beta'\mathbf{K}$$

Here, $\text{Cov}(\hat{\beta}) = \text{MSE}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1}\sigma^2$, so that

$$\begin{aligned} \text{MSE}(\hat{\beta}) - \text{MSE}(\hat{\beta}(\theta)) &= (\mathbf{X}'\mathbf{X})^{-1}\sigma^2 - \mathbf{K}\mathbf{X}'\mathbf{X}\mathbf{K}\sigma^2 - \theta^2\mathbf{K}\beta\beta'\mathbf{K} \\ &= \mathbf{K}\mathbf{K}^{-1}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}^{-1}\mathbf{K}\sigma^2 - \mathbf{K}\mathbf{X}'\mathbf{X}\mathbf{K}\sigma^2 - \theta^2\mathbf{K}\beta\beta'\mathbf{K} \\ &= \mathbf{K}[\mathbf{K}^{-1}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}^{-1}\sigma^2 - \mathbf{X}'\mathbf{X}\sigma^2 - \theta^2\beta\beta']\mathbf{K} \end{aligned}$$

Now, to show $\text{MSE}(\hat{\beta}) - \text{MSE}(\hat{\beta}(\theta))$ is p.d., we are only to show

$$\Delta = \mathbf{K}^{-1}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}^{-1}\sigma^2 - \mathbf{X}'\mathbf{X}\sigma^2 - \theta^2\beta\beta'$$

is p.d. because \mathbf{K} is p.d. Note that

$$\begin{aligned} \mathbf{K}^{-1}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}^{-1} &= (\mathbf{X}'\mathbf{X} + \theta\mathbf{I})(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X} + \theta\mathbf{I}) \\ &= \mathbf{X}'\mathbf{X} + 2\theta\mathbf{I} + \theta^2(\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

implies that

$$\Delta = \theta[\{2\mathbf{I} + \theta(\mathbf{X}'\mathbf{X})^{-1}\}\sigma^2 - \theta\beta\beta']$$

Now, to show Δ is p.d., we need to show $2\sigma^2\mathbf{I} - \theta\beta\beta'$ is p.d., because

$(\mathbf{X}'\mathbf{X})^{-1}$ is p.d and $\theta > 0$. Let $\mathbf{a} \neq \mathbf{0}$, then

$$\mathbf{a}'(2\sigma^2\mathbf{I} - \theta\beta\beta')\mathbf{a} = 2\sigma^2\mathbf{a}'\mathbf{a} - \theta(\mathbf{a}'\beta)^2 \geq 2\sigma^2\mathbf{a}'\mathbf{a} - \theta(\mathbf{a}'\mathbf{a})(\beta'\beta) = (2\sigma^2 - \theta\beta'\beta)\mathbf{a}'\mathbf{a}$$

by the Cauchy-Schwarz inequality. Hence, if $2\sigma^2 - \theta\beta'\beta > 0$, then $2\sigma^2\mathbf{I} - \theta\beta\beta'$ is p.d., i.e., For $0 < \theta < 2\sigma^2 / \beta'\beta$, $\text{MSE}(\hat{\beta}) - \text{MSE}(\hat{\beta}(\theta))$ is p.d..

(3) Bayes estimator

Consider

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2\mathbf{I})$$

and let

$$\beta \sim N\left(\mathbf{0}, \frac{\sigma^2}{\theta}\mathbf{I}\right)$$

be the prior distribution of β . Then, we can show that the posterior distribution of β is

$$N_p((\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}, (\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\sigma^2)$$

Therefore, under the squared error loss, $\hat{\beta}(\theta) = (\mathbf{X}'\mathbf{X} + \theta\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$ is the Bayes estimator of β .

7.3 Principal Component Regression

(1) Motivation

When there are many covariates the PCR (principal component regression) uses a few principal components instead of the original covariates.

(1) Derivation of PCR

By the SVD (singular values decomposition) of the $n \times p$ matrix \mathbf{X} , we have

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}',$$

where \mathbf{U} and \mathbf{V} are $n \times p$ and $p \times p$ orthogonal matrix, respectively, and \mathbf{D} is $p \times p$ diagonal matrix with diagonal elements $d_1 \geq d_2 \geq \cdots \geq d_p \geq 0$ which are singular values of \mathbf{X} . Therefore,

$$\mathbf{X}'\mathbf{X} = \mathbf{U}\mathbf{D}^2\mathbf{V}'$$

which is just the spectral decomposition of $\mathbf{X}'\mathbf{X}$. Here, $d_1^2 \geq d_2^2 \geq \cdots \geq d_p^2 \geq 0$ are eigenvalues of $\mathbf{X}'\mathbf{X}$ and the corresponding eigenvectors are v_1, v_2, \cdots, v_p which are column vectors of \mathbf{V} , and they are orthogonal to each other. Let $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_p)$, where \mathbf{x}_i is n -vector for the i th co-variate, and let

$$\mathbf{z}_i = \mathbf{X}v_i = \mathbf{u}_i d_i, \quad i = 1, 2, \cdots, p,$$

where \mathbf{z}_i is called the i th principal component, and \mathbf{u}_i is the i th column vector of \mathbf{U} . Hence, the sample variance of \mathbf{z}_i is

$$\text{Var}(\mathbf{z}_i) = \text{Var}(\mathbf{X}v_i) = d_i^2/n, \quad i = 1, 2, \cdots, p.$$

Here, z_1 is called the 1st principal component, and the i -th principal component is $z_i = Xv_i$. Note that v_i is called the i th principal component direction, and it satisfies the following conditions;

$$\begin{aligned} & \max_{\|\alpha\|=1} \text{Var}(X\alpha), \\ & v_l' S \alpha = 0 \\ & l=1,2,\dots,i-1 \end{aligned}$$

where

$$S = X'X/n$$

is the sample variance-covariance matrix of covariates. Also, since $v_l' S \alpha = 0$, $z_i = X\alpha$ and $z_l = Xv_l$, $l = 1, 2, \dots, i-1$ are orthogonal to each other.

The PCR uses just few principals z_1, z_2, \dots, z_M , $M \ll p$ out of z_1, z_2, \dots, z_p .

If the response and covariates are centered, we may write the fitted vector of the PCR as

$$\hat{y} = \sum_{m=1}^M \hat{\theta}_m z_m,$$

where $\hat{\theta}_m = z_m' y / z_m' z_m$ which is just the estimate of the slope for the simple linear regression of y on z_m . To estimate M , the number of principals to be used, is often estimated by CV.

7.4 PLS : Partial Least Squares

PLS is quite similar to the PCR, however, it uses the information of \mathbf{y} in addition to \mathbf{X} . Assume that the i th PLS direction vector \mathbf{v}_i satisfies the following condition;

$$\begin{aligned} \max_{\substack{||\boldsymbol{\alpha}||=1 \\ \mathbf{v}_l' \mathbf{S} \boldsymbol{\alpha} = 0 \\ l=1,2,\dots,i-1}} \quad & \text{Corr}^2(\mathbf{y}, \mathbf{X}\boldsymbol{\alpha}) \text{Var}(\mathbf{X}\boldsymbol{\alpha}) \end{aligned}$$

The PLS algorithm can be summarized as follows;

(algorithm)

- (i) By regressing \mathbf{y} on $\mathbf{x}_j, j = 1, 2, \dots, p$, compute the estimate of regression coefficients $\hat{v}_{1j} = \langle \mathbf{x}_j, \mathbf{y} \rangle$.
- (ii) Compute the 1st PLS direction vector $\mathbf{z}_1 = \sum \hat{v}_{1j} \mathbf{x}_j$.
- (iii) Regress \mathbf{y} on \mathbf{z}_1 , and compute the estimate of coefficient $\hat{\theta}_1 = \langle \mathbf{z}_1, \mathbf{y} \rangle$.
- (iv) Orthogonalize $\mathbf{x}_1, \dots, \mathbf{x}_p$ w.r.t. \mathbf{z}_1 , i.e, $\mathbf{x}_j - \frac{\langle \mathbf{z}_1, \mathbf{x}_j \rangle}{\langle \mathbf{z}_1, \mathbf{z}_1 \rangle} \mathbf{z}_1, j = 1, \dots, p$.
- (v) Repeat this process to the M th direction.

7.5 LASSO

7.5.1 Estimation of LASSO

(1) Definition

LASSO (least absolute shrinkage and selection operator) is suggested by Tibshirani (1996). The ridge regression is obtained by assigning L_2 restriction on β , however, the LASSO is obtained by assigning L_1 restriction on β , i.e.,

$$\begin{aligned}\hat{\beta}_L &= \arg \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \text{ s.t. } \sum_{j=1}^p |\beta_j| \leq s \\ &= \arg \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j|,\end{aligned}$$

where s and λ are parameters to be estimated. The solution of the above equation does not have a closed form solution, so that we are only to compute β numerically by packages like LARS (least angle regression).

We can obtain an explicit solution when we have one covariate. For simplicity, assume that $\sum y_i = 0$, $\sum x_i = 0$ and $\sum x_i^2 = 1$, and consider $y_i = \beta x_i + \epsilon_i$. Then, the LASSO estimate is given by

$$\hat{\beta}_L(\lambda) = \text{sgn}(\beta)(|\beta| - \lambda)_+,$$

where $(x)_+ = xI(x > 0)$. We can prove as follows; Note that

$$\begin{aligned}f(\beta) &= \frac{1}{2} \sum (y_i - x_i\beta)^2 + \lambda|\beta| \\ &= \frac{1}{2} \sum (y_i - x_i\hat{\beta} + x_i\hat{\beta} - x_i\beta)^2 + \lambda|\beta| \\ &= \frac{1}{2} \sum (y_i - x_i\hat{\beta})^2 + \frac{1}{2} (\hat{\beta} - \beta)^2 + \lambda|\beta|\end{aligned}$$

Hence,

$$\min_{\beta} f(\beta) \Leftrightarrow \min_{\beta} g(\beta) = \frac{1}{2}(\hat{\beta} - \beta)^2 + \lambda|\beta|$$

so that $\partial g(\beta)/\partial \beta = 0$ gives the LASSO estimate.

(2) Estimation of s .

Note that $\sum |\beta_j| \leq s \sum \beta_j^2 / |\beta_j| \leq s$, which is similar form to the ridge regression case. Hence, under this restriction, the ridge regression estimator is given by

$$\tilde{\beta} = (X'X + \lambda W^-)^{-1} Xy,$$

where $W = \text{diag}(|\tilde{\beta}_1|, \dots, |\tilde{\beta}_p|)$ and λ satisfies $\sum |\tilde{\beta}_j| = s$. Here, W^- denotes the generalized inverse of W which might contain $\tilde{\beta}_j = 0$. Now, the hat matrix is

$$H_s = X(X'X + \lambda W^-)^{-1} X'$$

and to estimate s , we minimize the generalized cross validation, i.e.,

$$GCV(s) = \frac{1}{n} \sum (y_i - \hat{y}_i)^2 / \{1 - \frac{1}{n} \text{tr}(H_s)\}^2$$

7.5.2 Properties of LASSO

(1) Geometric interpretation

Fig. 7.3 (p. 290)

(2) Bayesian interpretation

ridge regression estimator = posterior mean when the prior for β is Gaussian.

LASSO estimator = posterior mode when the prior for β is double exponential (Laplace).

(3) Comparisons

Under the orthonormal design, i.e., $X'X = I$,

(i) variable selection : $\hat{\beta}_{S,j} = \hat{\beta}_j I(|\hat{\beta}_j| > \lambda)$

(ii) ridge regression : $\hat{\beta}_{R,j} = \hat{\beta}_j / (1 + \lambda)$

(iii) LASSO : $\hat{\beta}_{L,j} = \text{sgn}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+$

Fig. 7.4 (p. 290)

7.6 Example (Prostate Cancer Data)

(1) Data description

response : lpsa (log of the level of prostate specific antigen)

covariates : X_1, \dots, X_8

LSE, variable selection, ridge reg., LASSO, PCR, PLS : Table 7.4 (p. 293)

Fig. 7.7 and 7.8 (p. 294)

R code (Fig. 7.9, p. 294)

7.7 Bayes Estimator and Biased Estimators

7.7.1 Bayes Estimator

(1) Posterior distribution

Let $f(x; \theta)$ be the pdf of the r.v. X , i.e.,

$$X \sim f(x; \theta),$$

where θ is a unknown parameter. In the non-Bayesian point of view, θ is regarded as a unknown constant. On the other hand, in the Bayesian point of view, θ is regarded as a r.v. with pdf. Hence, in the Bayesian point of view, $f(x; \theta)$ is regarded as a conditional pdf of X when θ is given, i.e.,

$$X \sim f(x|\theta)$$

Now, let the pdf of θ be $g(\theta; \gamma)$, i.e.,

$$\theta \sim g(\theta; \gamma),$$

where γ is also a parameter. We call $g(\theta; \gamma)$, the pdf of θ , as a *prior distribution*.

Also, the product of the conditional pdf $f(x|\theta)$ and $g(\theta; \gamma)$, the prior pdf of θ , is called a *posterior distribution*, i.e.,

$$g(\theta|x; \gamma) = \frac{f(x|\theta)g(\theta; \gamma)}{\int f(x|\theta)g(\theta; \gamma)d\theta}$$

(2) Bayes estimator

Let $\delta \equiv \delta(x)$ be an estimator of θ , and let $L(\theta, \delta)$ be a loss function. Then, $L(\theta, \delta) = (\theta - \delta)^2$ is called a squared error loss function, and $L(\theta, \delta) = |\theta - \delta|$ is called an absolute error loss function. Further, the expectation of the loss function $R(\theta, \delta) = E[L(\theta, \delta)]$ is called a risk function, i.e.,

$$R(\theta, \delta) = \int_{\chi} L(\theta, \delta) f(x|\theta) dx$$

Now, the Bayes risk is defined as

$$r(\theta, \delta) = \int_{\Theta} R(\theta, \delta) g(\theta; \gamma) d\theta$$

which can be rewritten as

$$\begin{aligned} r(\theta, \delta) &= \int_{\Theta} R(\theta, \delta) g(\theta; \gamma) d\theta \\ &= \int_{\Theta} \int_{\chi} L(\theta, \delta) f(x|\theta) g(\theta; \gamma) dx d\theta \\ &= \int_{\chi} \int_{\Theta} L(\theta, \delta) f(x|\theta) g(\theta; \gamma) d\theta dx \\ &= \int_{\chi} E[L(\Theta, \delta) | X = x] dx, \end{aligned}$$

where $E[L(\Theta, \delta)|X = x]$ is called *posterior risk*. The *Bayes estimator* minimizes the Bayes risk or the posterior risk, i.e.,

$$\hat{\delta}_{Bayes} = \operatorname{argmin}_{\delta} r(\theta, \delta).$$

We can show that the Bayes estimator is a posterior mean under squared error loss, and a posterior mode under absolute error loss.

7.7.2 Stein shrinkage estimator and the empirical Bayes estimator

If the parameter γ in the prior pdf $g(\theta; \gamma)$ is unknown, it should be estimated. To estimate γ , we use the marginal pdf of the posterior pdf by integrating it w.r.t. θ . If we replace the estimate of γ in the Bayes estimator, the resulting estimator is called *empirical Bayes estimator*.

Theorem 7.2 Assume that $\mathbf{z} \sim N_p(\boldsymbol{\mu}, \mathbf{I})$ and the prior for $\boldsymbol{\mu}$ is $\boldsymbol{\mu} \sim N_p(\mathbf{0}, \sigma^2 \mathbf{I})$. Then, the Stein shrinkage estimator is the empirical bayes estimator under the squared error loss.

(Proof) First, the posterior becomes

$$\begin{aligned}
& f(\mathbf{z}|\boldsymbol{\mu})g(\boldsymbol{\mu}) \\
& \propto \exp \left[-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})'(\mathbf{z} - \boldsymbol{\mu}) \right] \exp \left[-\frac{1}{2\sigma^2}\boldsymbol{\mu}'\boldsymbol{\mu} \right] \\
& \propto \exp \left[-\frac{1}{2}(\mathbf{z}'\mathbf{z} - 2\boldsymbol{\mu}'\mathbf{z} + \boldsymbol{\mu}'\boldsymbol{\mu} + \frac{1}{\sigma^2}\boldsymbol{\mu}'\boldsymbol{\mu}) \right] \\
& \propto \exp \left[-\frac{1}{2} \left\{ (1 + \frac{1}{\sigma^2})[\boldsymbol{\mu}'\boldsymbol{\mu} - 2\frac{\sigma^2}{1 + \sigma^2}\boldsymbol{\mu}'\mathbf{z} + \frac{\sigma^2}{1 + \sigma^2}\mathbf{z}'\mathbf{z}] \right\} \right].
\end{aligned}$$

If we let

$$w = \frac{\sigma^2}{1 + \sigma^2}$$

then, the posterior can be written as

$$\begin{aligned}
& \propto \exp \left[-\frac{1}{2w}(\boldsymbol{\mu}'\boldsymbol{\mu} - 2w\boldsymbol{\mu}'\mathbf{z} + w^2\mathbf{z}'\mathbf{z}) \right] \exp \left[-\frac{1}{2}(1 - w)\mathbf{z}'\mathbf{z} \right] \\
& \propto \exp \left[-\frac{1}{2w}(\boldsymbol{\mu} - w\mathbf{z})'(\boldsymbol{\mu} - w\mathbf{z}) \right] \exp \left[-\frac{1}{2}(1 - w)\mathbf{z}'\mathbf{z} \right].
\end{aligned}$$

Therefore, the posterior becomes

$$\boldsymbol{\mu}|\mathbf{z} \sim N_p(w\mathbf{z}, w\mathbf{I}).$$

Since the Bayes estimator is the posterior mean under squared error loss,

we have

$$w\mathbf{z} = \frac{\sigma^2}{1 + \sigma^2}\mathbf{z} = \left(1 - \frac{1}{1 + \sigma^2}\right)\mathbf{z},$$

however, it contains unknown parameter σ^2 , we obtain the marginal distribution of \mathbf{z} to compute the estimator of σ^2 .

$$\int f(\mathbf{z}|\boldsymbol{\mu})g(\boldsymbol{\mu})d\boldsymbol{\mu} \propto \exp[-\frac{1}{2}(1-w)\mathbf{z}'\mathbf{z}], \mathbf{z} \sim N_p(\mathbf{0}, (1+\sigma^2)\mathbf{I})$$

Now, consider

$$V \equiv \frac{\mathbf{z}'\mathbf{z}}{1+\sigma^2} \sim \chi^2(p)$$

then

$$E\left(\frac{1}{V}\right) = \int_0^\infty \frac{1}{v} \frac{v^{\frac{p}{2}-1} e^{-v/2}}{\Gamma(\frac{p}{2}) 2^{\frac{p}{2}}} dv = \frac{\Gamma(\frac{p}{2}-1) 2^{\frac{p}{2}-1}}{\Gamma(\frac{p}{2}) 2^{\frac{p}{2}}} = \frac{\Gamma(\frac{p}{2}-1)}{2(\frac{p}{2}-1)\Gamma(\frac{p}{2}-1)} = \frac{1}{p-2}$$

hence,

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

Therefore, the unbiased estimator of $(1+\sigma^2)^{-1}$ is $(p-2)/\mathbf{z}'\mathbf{z}$, and finally the empirical Bayes estimator is

$$\left(1 - \frac{p-2}{\mathbf{z}'\mathbf{z}}\right)\mathbf{z}$$

which is just the Stein shrinkage estimator.

7.7.3 Bayes estimator in regression

Consider a multiple linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

which is equivalent to $\mathbf{y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, and assume that the prior to $\boldsymbol{\beta}$ be

$$\boldsymbol{\beta} \sim N_p(\mathbf{m}, \sigma^2\mathbf{V})$$

Now, the posterior of $\boldsymbol{\beta}$ becomes

$$\begin{aligned} & \exp \left[-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \exp \left[-\frac{1}{2\sigma^2}(\boldsymbol{\beta} - \mathbf{m})'\mathbf{V}^{-1}(\boldsymbol{\beta} - \mathbf{m}) \right] \\ & \propto \exp \left[-\frac{1}{2\sigma^2}(\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{V}^{-1}\boldsymbol{\beta} - 2\boldsymbol{\beta}'\mathbf{X}'\mathbf{y} - 2\boldsymbol{\beta}'\mathbf{V}^{-1}\mathbf{m} + \dots) \right] \\ & \propto \exp \left[-\frac{1}{2\sigma^2}\{\boldsymbol{\beta}'(\mathbf{X}'\mathbf{X} + \mathbf{V}^{-1})\boldsymbol{\beta} - 2\boldsymbol{\beta}'(\mathbf{V}^{-1}\mathbf{m} + \mathbf{X}'\mathbf{y}) + \dots\} \right] \\ & \equiv \exp \left[-\frac{1}{2\sigma^2}\{(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)'(\mathbf{X}'\mathbf{X} + \mathbf{V}^{-1})(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)\} \right], \end{aligned}$$

where $\boldsymbol{\mu}_\beta$ is the mean of the posterior, then we let

$$-2\boldsymbol{\beta}'(\mathbf{X}'\mathbf{X} + \mathbf{V}^{-1})\boldsymbol{\mu}_\beta \equiv -2\boldsymbol{\beta}'(\mathbf{V}^{-1}\mathbf{m} + \mathbf{X}'\mathbf{y})$$

Hence,

$$\boldsymbol{\mu}_\beta = (\mathbf{X}'\mathbf{X} + \mathbf{V}^{-1})^{-1}(\mathbf{V}^{-1}\mathbf{m} + \mathbf{X}'\mathbf{y}) \equiv \tilde{\boldsymbol{\beta}}_{Bayes}$$

which is the Bayes estimator under the squared error loss.

As a special case, let $\mathbf{V} = \lambda^{-1}\mathbf{I}_p$ for some $\lambda > 0$ and $\mathbf{m} = \mathbf{0}$, then

$$\tilde{\boldsymbol{\beta}}_{Bayes} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$$

which is just the ridge regression estimator. Next, if $\mathbf{V} = \lambda^{-1}(\mathbf{X}'\mathbf{X})^{-1}$ and $\mathbf{m} = \mathbf{0}$, then

$$\tilde{\boldsymbol{\beta}}_{Bayes} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (1 + \lambda)^{-1}\hat{\boldsymbol{\beta}}$$

which is called *James – Stein regression estimator*.