

Chapter 5. Advanced topics in modern linear models.

Linear models were largely developed in the pre-computer age of statistics, but even in today's computer era there are still good reasons to study and use them. They are simple and often provide an adequate and interpretable description of how the inputs (predictor) affect the output. For prediction purposes they can sometimes outperform fancier nonlinear models, especially in situations with small numbers of training cases, low signal-to-noise ratio or sparse data. In this chapter, we describe the linear methods for regression in penalized regularization framework, such as ridge regression, the least absolute shrinkage and selection operator (LASSO), in details. It is my firm belief that an understanding of linear methods is essential for understanding nonlinear ones. Please refer to the two books entitled "*The elements of statistical learning*" by Hastie, Tibshirani and Friedman (2011) and "*Statistical learning with sparsity: the Lasso and generalizations*" by Hastie, Tibshirani and Wainwright (2015) for more detailed discussions on the topics covered in this chapter.

5.1 The Bias-Variance Trade-off of estimating β

Generally, one may assume that the output/response Y and the input/predictor Z satisfy

$$y = f(\mathbf{Z}) + \varepsilon,$$

where $E(\varepsilon) = 0$ and $Var(\varepsilon) = \sigma^2$. When $f(\mathbf{z})$ is assumed to be $\mathbf{Z}\beta$, so far, we've been dealing with, $\hat{\beta}^{ls}$, the least square estimate, which has well-known properties (e.g. the Gauss-Markov, MLE in normal case). That is, the least squares estimate of the parameter β have the smallest variance among all linear unbiased estimates. But, can we do better?

Consider the mean squared error (MSE) of an estimator $\hat{\theta}$ in estimating θ :

$$MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2 = \underbrace{Var(\hat{\theta})}_{\text{BLUE}} + \underbrace{\{E(\hat{\theta}) - \theta\}^2}_{\text{squared bias}}.$$

The first term is the variance, while the second term is the squared bias. The Gauss-Markov theorem implies that the least squares estimator has the smallest mean squared error of all linear estimators with no bias. However, the restriction to unbiased estimates is not necessarily a wise one. There may well exist a biased estimator with smaller mean squared error. Such an estimator would trade a little bias for a larger reduction in variance. Biased estimators are commonly used. Any method that shrinks or sets to zero some of the least squares coefficients may result in a biased estimate.

For illustration, suppose we have an estimator $\hat{f}(\mathbf{z}) = \mathbf{Z}\hat{\beta}$. To see whether it is a good candidate, there are usually two questions we need to consider: (1) Is $\hat{\beta}$ close to the true β ? (*bias*). (2) Will $\hat{f}(\mathbf{z})$ fit future observations well? To answer the questions, we might consider the **mean squared error** of our estimate $\hat{\beta}$,

$$\begin{aligned} MSE(\hat{\beta}) &= E(\|\hat{\beta} - \beta\|^2) = E[(\hat{\beta} - \beta)^\top (\hat{\beta} - \beta)] \\ &= E(\hat{\beta}^\top \hat{\beta}) - \beta^\top \beta \\ &= E(Y^\top X (X^\top X)^{-1} (X^\top X)^{-1} X^\top Y) \\ &= \text{tr}(X (X^\top X)^{-1} (X^\top X)^{-1} X^\top Y Y^\top) \\ &= \text{tr}(X (X^\top X)^{-1} (X^\top X)^{-1} X^\top \sigma^2 I) \\ &= \sigma^2 \text{tr}(X (X^\top X)^{-1} X^\top) \end{aligned}$$

For example, in the least-square estimation, we have

$$MSE(\hat{\beta}^{ls}) = E(\|\hat{\beta}^{ls} - \beta\|^2) = \text{tr}(\text{Var}(\hat{\beta}^{ls})) = \sigma^2 \text{tr}((X^\top X)^{-1})$$

$\hat{f}(\mathbf{z})$ fits our data well doesn't mean that it will be a good fit to new data. In fact, suppose that we take new measurements y'_i at the same \mathbf{z}_i 's:

$$(\mathbf{z}_1, y'_1), (\mathbf{z}_2, y'_2), \dots, (\mathbf{z}_n, y'_n).$$

If $\hat{f}(\cdot)$ is a good model, then $\hat{f}(\mathbf{z})$ should also be close to the new target y'_i . This is the notion of **prediction error (PE)**. So good estimation should, on average, have small prediction error.

Consider the expected prediction error of an estimate $\hat{f}(\mathbf{z})$ at a particular point \mathbf{z}_0 ,

$$\begin{aligned} PE(\mathbf{z}_0) &= E_{Y|\mathbf{Z}=\mathbf{z}_0} [(Y - \hat{f}(\mathbf{Z}))^2 | \mathbf{Z} = \mathbf{z}_0] \\ &= \sigma_\epsilon^2 + \text{Bias}^2(\hat{f}(\mathbf{z}_0)) + \text{Var}(\hat{f}(\mathbf{z}_0)) \\ &= \sigma_\epsilon^2 + MSE(\hat{f}(\mathbf{z}_0)). \end{aligned} \quad (1)$$

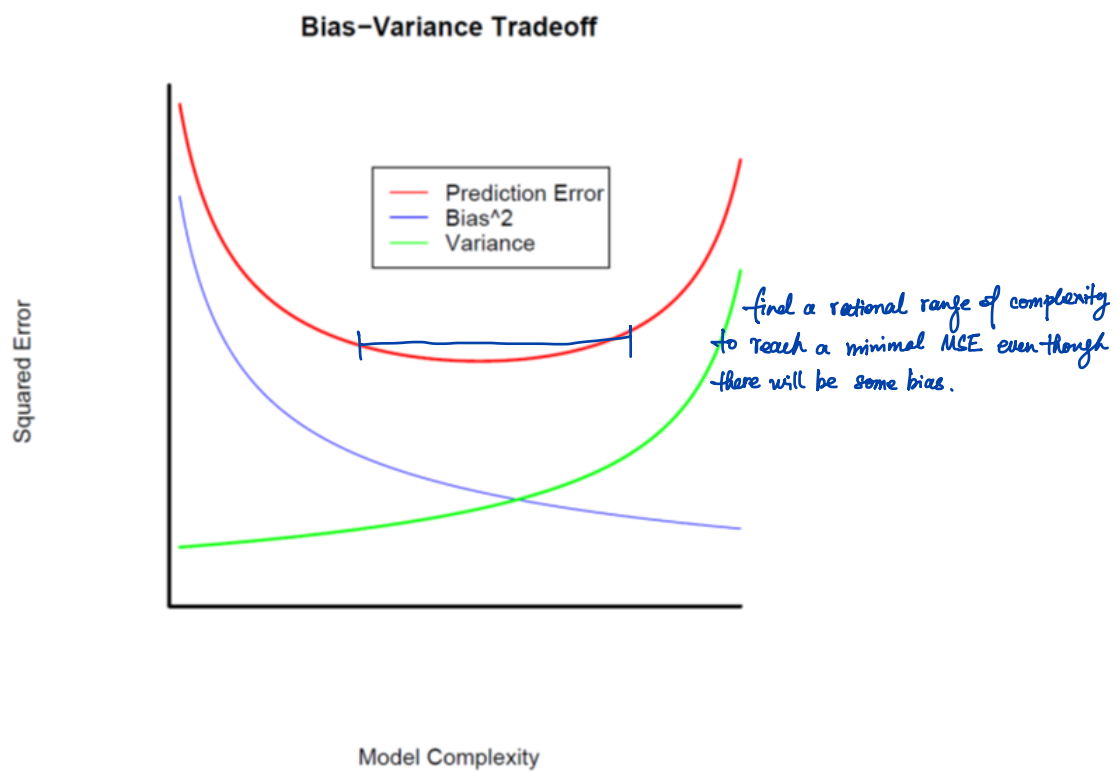
in paper, \mathbf{Z} is usually a random observation.

The relationship in (1) tells that the mean squared error is intimately related to prediction accuracy. In other words, the expected prediction error and mean squared error differ only by the constant σ_ϵ^2 , representing the variance of the new observation y_0 .

Such a decomposition in (1) is known as the **bias-variance trade-off**. As model becomes more complex (more terms included), local structure/curvature can be picked up. But coefficient estimates suffer from high variance as more terms are included in the model. So introducing a little bias in our estimate for β might lead to a substantial decrease in variance, and hence to a substantial decrease in PE.

$$\begin{aligned} PE(\mathbf{z}_0) &= E[(Y - \hat{f}(\mathbf{z}_0))^2 | \mathbf{Z} = \mathbf{z}_0] \\ &= E[\epsilon + \hat{f}(\mathbf{z}_0) - \hat{f}(\mathbf{z}_0)]^2 \\ &= E[\epsilon^2 + \underbrace{(\hat{f}(\mathbf{z}_0) - \hat{f}(\mathbf{z}_0))^2}_{MSE} + \underbrace{2\epsilon(\hat{f}(\mathbf{z}_0) - \hat{f}(\mathbf{z}_0))}_0] \\ &= \sigma_\epsilon^2 + \underbrace{[f(\mathbf{z}_0) - E\hat{f}(\mathbf{z}_0)]^2}_{\text{bias}} + \underbrace{E[\hat{f}(\mathbf{z}_0) - E\hat{f}(\mathbf{z}_0)]^2}_{\text{variance}} \end{aligned}$$

over the training process.



5.2 Shrinkage methods

As introduced before, we have an input vector $Z = (Z_1, Z_2, \dots, Z_p)$ and an output or response variable Y . The linear regression model has the form

$$Y = \beta_0 + \sum_{j=1}^p Z_j \beta_j + \varepsilon = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

As explained in section 5.1, there are two reasons why we are often not satisfied with the least squares estimates.

(1) The first is *prediction accuracy*: the least square estimates often have low bias but large variance. Prediction accuracy can sometimes be improved by shrinking or setting some coefficients to zero. By doing so, we sacrifice a little bit bias to reduce the variance of the predicted values, and hence may improve the overall prediction accuracy.

(2) The second reason is *interpretation*. With a large number of predictors, we often would like to determine a smaller subset that exhibit the strongest effects. In order to get the “big picture”, we are willing to sacrifice some of the small details.

There are a number of approaches to variable subset selection with linear regression, e.g. best-subset selection, forward and backward-Stepwise selection, etc. Due to limitation of time, I will skip the introduction of the subset selection methods and focus on the shrinkage methods.

By retaining a subset of the predictors and discarding the rest, subset selection produces a model that is interpretable and has possibly lower prediction error than the full model. However, because it is a discrete process—variables are either retained or discarded—it often exhibits high variance, and so doesn’t reduce the prediction error of the full model. Shrinkage methods are more continuous and stable, and don’t suffer as much from high variability.

5.2.1 Ridge regression

Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge regression is to minimize

$$\sum_{i=1}^n (y_i - \mathbf{Z}_i \boldsymbol{\beta})^2 \quad \text{s.t.} \quad \sum_{j=1}^p \beta_j^2 \leq t$$

i.e. $(\mathbf{Y} - \mathbf{Z}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{Z}\boldsymbol{\beta}) \quad \text{s.t.} \quad \sum_{j=1}^p \beta_j^2 \leq t$

★ optimization
Lagrange & KKT.

By convention,

(1) \mathbf{Z} is assumed to be standardized (mean=0, variance=1)

(2) \mathbf{Y} is assumed to be centered

no penalty on Intercept.

Equivalently, we can write the ridge coefficients as the minimizer of the following penalized residual sum of square (PRSS):

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} PRSS(\beta)_{\ell_2} = \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - \mathbf{z}_i \beta)^2 + \lambda \left(\sum_{j=1}^p \beta_j^2 - 1 \right) \right\}$$

$$= \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - \mathbf{z}_i \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

$$= \arg \min_{\beta} \|\mathbf{Y} - \mathbf{Z}\beta\|_2^2 + \lambda \|\beta\|_2^2 \quad (2)$$

Its solution may have smaller average PE than $\hat{\beta}^{ls}$. Note that $PRSS(\beta)_{\ell_2}$ is convex, and hence has a unique solution. Taking derivatives, we obtain:

$$\frac{\partial PRSS(\beta)_{\ell_2}}{\partial \beta} = -2\mathbf{Z}^T(\mathbf{Y} - \mathbf{Z}\beta) + 2\lambda\beta$$

$$\hat{\beta}_{\lambda}^{\text{ridge}} = (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^T \mathbf{Y}$$

$\lambda > 0$ ($\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I}_p$) is always invertible

where \mathbf{Z} is standardized and \mathbf{Y} is centered. The solution is indexed by a tuning parameter λ , which makes the problem non-singular even if $\mathbf{Z}^T \mathbf{Z}$ is not invertible. This was the original motivation for ridge regression (Hoerl and Kennard, 1970).

Remark 1. It is worthwhile to point out that, the ridge solutions are not equivariant under scaling of the inputs, and so this is why one normally standardizes the inputs before solving (2). In addition, notice that the intercept β_0 has been left out of the penalty term. Penalization of the intercept would make the procedure depend on the origin chosen for \mathbf{Y} ; that is, adding a constant c to each of the targets y_i would not simply result in a shift of the predictions by the same amount c . Henceforth we assume that this centering has been done, so that the input/design matrix has p (rather than $p + 1$) columns.

Note that the solution is indexed by λ , a complexity parameter that controls the amount of shrinkage: the larger the value of λ , the greater the amount of shrinkage:

- (1) λ controls the size of the coefficients
- (2) λ controls the amount of regulation
- (3) As $\lambda \downarrow 0$, we obtain the least square solutions
- (4) As $\lambda \uparrow \infty$, we have $\hat{\beta}_{\lambda}^{\text{ridge}} = 0$ (intercept-only model)

Remark 2. We next give a simple proof that $\hat{\beta}_\lambda^{ridge}$ is biased. To this end, let $\mathbf{R} = \mathbf{Z}^\top \mathbf{Z}$ and assume \mathbf{Z} has full column rank. Write

$$A \in \mathbb{R}^{p \times p} \quad A^{-1} \sim O(p^3).$$

$$\begin{aligned} \hat{\beta}_\lambda^{ridge} &= (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{Y} \\ &= (\mathbf{R} + \lambda \mathbf{I}_p)^{-1} \mathbf{R} (\mathbf{R}^{-1} \mathbf{Z}^\top \mathbf{Y}) \\ &= [\mathbf{R}(\mathbf{I}_p + \lambda \mathbf{R}^{-1})]^{-1} \mathbf{R} \hat{\beta}^{ls} \\ &= (\mathbf{I}_p + \lambda \mathbf{R}^{-1})^{-1} \hat{\beta}^{ls}. \end{aligned}$$

Thus,

$$\begin{aligned} E(\hat{\beta}_\lambda^{ridge}) &= E[(\mathbf{I}_p + \lambda \mathbf{R}^{-1})^{-1} \hat{\beta}^{ls}] \\ &= (\mathbf{I}_p + \lambda \mathbf{R}^{-1})^{-1} \beta \\ &\neq \beta \quad (\forall \beta \neq 0, \lambda \neq 0) \end{aligned}$$

Note that 1 is not an eigenvalue of the matrix $(\mathbf{I}_p + \lambda \mathbf{R}^{-1})^{-1}$, as $\lambda \neq 0$ and \mathbf{R} is positive definite.

On the other hand, ridge regression can be translated into the least square problem alike with the **data augmentation** approach. The ℓ_2 -PRSS can be written as:

$$\begin{aligned} PRSS(\beta)_{\ell_2} &= \sum_{i=1}^n (y_i - \mathbf{Z}_i \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \\ &= \sum_{i=1}^n (y_i - \mathbf{Z}_i \beta)^2 + \sum_{j=1}^p (0 - \sqrt{\lambda} \beta_j)^2. \end{aligned}$$

Hence, the ℓ_2 criterion can be recast as another least squares problem for another data set.

The ℓ_2 criterion is simply the RSS for the augmented data set:

$$\mathbf{Z}_\lambda = \begin{pmatrix} z_{1,1} & z_{1,2} & z_{1,3} & \cdots & z_{1,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ z_{n,1} & z_{n,2} & z_{n,3} & \cdots & z_{n,p} \\ \hline \sqrt{\lambda} & 0 & 0 & \cdots & 0 \\ 0 & \sqrt{\lambda} & 0 & \cdots & 0 \\ 0 & 0 & \sqrt{\lambda} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \sqrt{\lambda} \end{pmatrix}, \quad \mathbf{Y}_\lambda = \begin{pmatrix} y_1 \\ \vdots \\ y_n \\ \hline 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

augmented data
(added p data points).

So,

$$\mathbf{Z}_\lambda = \begin{pmatrix} \mathbf{Z} \\ \sqrt{\lambda} \mathbf{I}_p \end{pmatrix}, \quad \mathbf{Y}_\lambda = \begin{pmatrix} \mathbf{Y} \\ 0 \end{pmatrix}$$

The “least squares” solution for the augmented data set is:

$$\begin{aligned} (\mathbf{Z}_\lambda^\top \mathbf{Z}_\lambda)^{-1} \mathbf{Z}_\lambda^\top \mathbf{Y}_\lambda &= \left((\mathbf{Z}^\top, \sqrt{\lambda} \mathbf{I}_p) \begin{pmatrix} \mathbf{Z} \\ \sqrt{\lambda} \mathbf{I}_p \end{pmatrix} \right)^{-1} (\mathbf{Z}^\top, \sqrt{\lambda} \mathbf{I}_p) \begin{pmatrix} \mathbf{Y} \\ 0 \end{pmatrix} \\ &= (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{Y} \end{aligned}$$

which is simply the ridge solution.

Computing the ridge solutions via the SVD.

When computing $\hat{\beta}_\lambda^{\text{ridge}}$ numerically, matrix inversion should be avoided, since inverting $\mathbf{Z}^\top \mathbf{Z}$ can be computationally expensive: $O(p^3)$. Instead, the **singular value decomposition** is utilized, that is,

$$\mathbf{Z} = \mathbf{U} \mathbf{D} \mathbf{V}^\top,$$

$n \times p \quad p \times p \quad p \times p$

Where

$\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p)$ is an $n \times p$ orthogonal matrix.

$\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_p)$ is a $p \times p$ diagonal matrix consisting of the singular values $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$.

$$\begin{aligned} \hat{\beta}_{ls} &= (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y} \\ \hat{\beta}_r &= (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{Y} \\ &= (\mathbf{I}_p + \lambda (\mathbf{Z}^\top \mathbf{Z})^{-1})^{-1} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y} \\ &= (\mathbf{I}_p + \lambda (\mathbf{Z}^\top \mathbf{Z})^{-1})^{-1} \hat{\beta}_{ls} \end{aligned}$$

$$V^\top = \begin{pmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \\ \vdots \\ \mathbf{v}_p^\top \end{pmatrix} \text{ is a } p \times p \text{ orthogonal matrix.}$$

$$\begin{aligned} &= (VDU^\top UDV^\top + \lambda I_p)^{-1} VDU^\top Y \\ &= (VD^\top V^\top + V\lambda I_p V^\top)^{-1} VDU^\top Y \\ &= V \text{diag}(d_i^2 + \lambda)^{-1} \underbrace{V^\top V}_{I} \text{diag}(d_i) U^\top Y \\ &= V \text{diag} \frac{d_i}{d_i^2 + \lambda} U^\top Y \end{aligned}$$

It can be shown that

$$\begin{aligned} &V \text{diag}(d_j^2) V^\top + V \text{diag}(\lambda) V^\top \\ \hat{\beta}_\lambda^{\text{ridge}} &= (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{Y} \\ &= V \text{diag} \left(\frac{d_j}{d_j^2 + \lambda} \right) U^\top \mathbf{Y} \end{aligned}$$

by using the following **eigen (or spectral) decomposition** of $\mathbf{Z}^\top \mathbf{Z}$:

$$\begin{aligned} \mathbf{Z}^\top \mathbf{Z} &= (UDV^\top)^\top (UDV^\top) \\ &= VD^\top U^\top UDV^\top \\ &= VD^2V^\top \end{aligned}$$

Proof:

A consequence is that:

$$\begin{aligned}\hat{\mathbf{Y}}^{ridge} &= \mathbf{Z}\hat{\boldsymbol{\beta}}_{\lambda}^{ridge} \\ &= \sum_{j=1}^p \left(\mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^{\top} \right) \mathbf{Y}\end{aligned}$$

Remark 3. (Optional) Ridge regression also has a relationship with principal components analysis (PCA). The fact is that, the derived variable $\gamma_j = \mathbf{Z}\mathbf{v}_j = d_j\mathbf{u}_j$ is the j th principal component (PC) of \mathbf{Z} , hence, ridge regression projects \mathbf{Y} onto these components with large d_j . Ridge regression shrinks the coefficients of low-variance (small singular value) components.

Orthonormal \mathbf{Z} in ridge regression

In the special case of an orthonormal design matrix, that is \mathbf{Z} is orthonormal, then $\mathbf{Z}^{\top}\mathbf{Z} = I_p$. In this case, a couple of closed form properties exist. Let $\hat{\boldsymbol{\beta}}^{ls}$ denote the LS solution for our orthonormal \mathbf{Z} , then

$$\hat{\boldsymbol{\beta}}_{\lambda}^{ridge} = \frac{1}{1 + \lambda} \hat{\boldsymbol{\beta}}^{ls}$$

and

$$\hat{\beta}_j^{ridge} = \frac{\hat{\beta}_j^{ols}}{1 + \lambda}.$$

The optimal choice of λ minimizing the expected prediction error is:

$$\lambda^* = \frac{p\sigma^2}{\sum_{j=1}^p \beta_j^2}$$

where $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)$ is the true coefficient vector.

This illustrates the essential “shrinkage” feature of ridge regression. Applying the ridge regression penalty has the effect of shrinking the estimates towards zero. The penalty introduces bias but reducing the variance of the estimate. Ridge estimator does not threshold, since the shrinkage is smooth (proportional to the original coefficient).

Furthermore, it can be checked that, the variance of the ridge regression estimate is

$$\text{Var}(\hat{\boldsymbol{\beta}}^{ridge}) = \sigma^2 \mathbf{W}_{\lambda} (\mathbf{Z}^{\top} \mathbf{Z})^{-1} \mathbf{W}_{\lambda} = \sigma^2 (\mathbf{Z}^{\top} \mathbf{Z} + \lambda I_p)^{-1} \mathbf{Z}^{\top} \mathbf{Z} (\mathbf{Z}^{\top} \mathbf{Z} + \lambda I_p)^{-1},$$

where $\mathbf{W}_\lambda = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{Z}$. Recall that $\hat{\boldsymbol{\beta}}^{ridge} = \mathbf{W}_\lambda \hat{\boldsymbol{\beta}}^{ls}$. The bias of the ridge regression estimate is

$$\text{bias}(\hat{\boldsymbol{\beta}}^{ridge}) = -\lambda \mathbf{W}_\lambda (\mathbf{Z}^\top \mathbf{Z})^{-1} \boldsymbol{\beta} = -\lambda (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \boldsymbol{\beta}$$

It can be shown that

try to proof it.

- (1) the total variance $\sum_{j=1}^p \text{Var}(\hat{\beta}_j^{ridge})$ is a monotone decreasing sequence with respect to λ .
- (2) the total squared bias $\sum_{j=1}^p \text{Bias}^2(\hat{\beta}_j^{ridge})$ is a monotone increasing sequence with respect to λ .

Existence Theorem

There always exists a λ such that the MSE of $\hat{\boldsymbol{\beta}}^{ridge}$ is less than the MSE of $\hat{\boldsymbol{\beta}}^{OLS}$.

□. find λ^ such that $\text{MSE}(\hat{\boldsymbol{\beta}}_{ridge}) \leq \text{MSE}(\hat{\boldsymbol{\beta}}_{ls})$.*

The proof of this existence theorem is left for students' take-home exercise. A rather surprising result with somewhat radical implications: even if the model we fit ^{*function of λ*} is exactly correct and follows the exact distribution we specify, we can always obtain a better (in terms of MSE) estimator by shrinking towards zero.

Smoother matrix and effective degrees of freedom

A **Smoother matrix** S is a linear operator satisfying:

$$\hat{\mathbf{Y}} = \mathbf{S} \mathbf{Y}$$

Example: In ordinary least squares, recall the hat matrix

$$\mathbf{H} = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$$

Recall that $\text{rank}(\mathbf{Z}) = \text{tr}(\mathbf{H}) = p$, which is how many degrees of freedom are used in the model. Analogously, we can define the **effective degrees of freedom** (or effective number of parameters) for a smoother to be:

$$\text{df}(\mathbf{S}) \triangleq \text{tr}(\mathbf{S})$$

In ridge regression, the fits are given by:

$$\hat{\mathbf{Y}} = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top \mathbf{Y}$$

So the smoother or “hat” matrix in ridge takes the form:

$$\mathbf{S}_\lambda = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top$$

So the effective degrees of freedom in ridge regression are given by:

$$\text{df}(\lambda) = \text{tr}(\mathbf{S}_\lambda) = \text{tr}(\mathbf{Z}(\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_p)^{-1} \mathbf{Z}^\top) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

Note that $\text{df}(\lambda)$ is monotone decreasing in λ .