

Introductory Material and Linear Methods

Robert Schmidt

Abstract

The following are notes on the key results from the **Elements of Statistical Learning** text. They were primarily derived from course notes and readings in the Stanford STATS 315: *Modern Applied Statistics* series.

Contents

1	Modeling basics	2
1.1	Linear model bias-variance decomposition	2
1.2	K-NN bias-variance decomposition	3
1.3	Variance of $\hat{\beta}$	3
1.4	Relationship between t and F statistic for dropping coefficients	4
2	Ridge regression	5
2.1	Derivation using the SVD	5
2.2	Eigenvalue interpretation	5
2.3	Relationship between λ and the coefficients	6
2.4	Effect of collinearity on ridge regression	7
2.5	Ridge when $p \gg N$	8
2.6	Bayesian theory: ridge estimate is the mean of the posterior	9
3	Degrees of freedom	10
3.1	Degrees of freedom: ridge	10
3.2	Degrees of freedom: LAR and lasso	10
3.2.1	Derivation: linear regression	10
3.2.2	Derivation: ridge	11
4	Methods using derived input directions	12
4.1	Principal component regression up to j coefficients	12
4.2	PCR connection to OLS/ridge	12
4.3	PCR vs. ridge	12
5	Zero-residual fits	13
5.1	There are infinitely many OLS solutions with zero residuals	13
5.2	Ridge solution has minimum norm among all zero-residual solutions as $\lambda \rightarrow 0$	13

1 Modeling basics

1.1 Linear model bias-variance decomposition

1) Overview

Conceptually, the bias-variance tradeoff is the problem of how closely one should fit a model to presented data. The tradeoff stipulates that there exists some minimum to the total error such that the model has optimal bias (how different the constructed function is from the true function) and variance (how much, on average, the estimated function differs from its expected value – essentially controlling the closeness of the function fit to the training data).

2) Mathematical specification

Mathematically, this can be derived from the squared-error loss formula, and similarly for other losses.

We consider squared error for illustrative purposes in this case.

Suppose $Y = f(X) + \varepsilon$, for noise ε where $\mathbb{E}[\varepsilon] = 0$ and $\text{Var}[\varepsilon] = \sigma_\varepsilon^2$, and errors independent of Y .

Thus, the *MSE* at a point x_0 is as follows:

$$\begin{aligned} \text{MSE}(x_0) &= \mathbb{E}[(Y - \hat{f}(x_0))^2 \mid X = x_0] \\ &= \mathbb{E}[(f(x_0) + \varepsilon - \hat{f}(x_0))^2] \\ &= \mathbb{E}[(\{f(x_0) - \hat{f}(x_0)\} + \varepsilon)^2] \end{aligned}$$

Now, we expand this expression.

$$\begin{aligned} \text{MSE}(x_0) &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2 + 2\varepsilon(f(x_0) - \hat{f}(x_0)) + \varepsilon^2] \\ &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + 2\mathbb{E}[\varepsilon(f(x_0) - \hat{f}(x_0))] + \mathbb{E}[\varepsilon^2] \quad \text{by linearity of expectation} \end{aligned}$$

First, consider the middle term. By error independence, we can simplify this term significantly.

$$2\mathbb{E}[\varepsilon(f(x_0) - \hat{f}(x_0))] = 2\mathbb{E}[\varepsilon]\mathbb{E}[f(x_0) - \hat{f}(x_0)] = 0 \quad \text{since } \mathbb{E}[\varepsilon] = 0$$

Since the middle term cancels completely, we see that

$$\begin{aligned} \text{MSE}(x_0) &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + \mathbb{E}[\varepsilon^2] \\ &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + \sigma_\varepsilon^2 \\ &= \text{Avg}(\text{model error}) + (\text{irreducible error}) \end{aligned}$$

Let us now expand the model error via the bias-variance decomposition.

$$\begin{aligned} \text{Model error} &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] \\ &= \mathbb{E}[(f(x_0) - \mathbb{E}\hat{f}(x_0) + \mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2] \\ &= \mathbb{E}[(f(x_0) - \mathbb{E}\hat{f}(x_0))^2 + 2(f(x_0) - \mathbb{E}\hat{f}(x_0))(\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0)) + (\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2] \end{aligned}$$

Again, we will focus on the middle term. If we note that $f(x_0)$ and $\mathbb{E}\hat{f}(x_0)$ are not random in regards to the expectation over the data, we can cancel this term:

$$\begin{aligned} 2\mathbb{E}[(f(x_0) - \mathbb{E}\hat{f}(x_0))(\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))] &= 2(f(x_0) - \mathbb{E}\hat{f}(x_0))\mathbb{E}[\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0)] \\ &= 2(f(x_0) - \mathbb{E}\hat{f}(x_0))[\mathbb{E}\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0)] \\ &= 0 \end{aligned}$$

So, we finally get:

$$\begin{aligned} \text{Model error} &= (\mathbb{E}\hat{f}(x_0) - f(x_0))^2 + \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0))^2] \\ &= \text{Avg}(\text{bias}^2) + \text{Avg}(\text{variance}) \end{aligned}$$

And, returning to the overall average error, we see that:

$$\begin{aligned} \text{MSE}(x_0) &= (\mathbb{E}\hat{f}(x_0) - f(x_0))^2 + \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0))^2] + \sigma_\varepsilon^2 \\ &= (\text{bias}^2) + (\text{variance}) + (\text{irreducible error}) \\ &= \text{bias-variance decomposition} \end{aligned}$$

1.2 K-NN bias-variance decomposition

Consider squared error: $\text{Err}(x_0) = \mathbb{E}[(y - \hat{f}_k(x_0))^2 \mid \mathbf{X} = x_0]$.
 Note that \hat{f}_k is the average of the K nearest neighbors:

$$\hat{f}_k(x_0) = \frac{1}{k} \sum_{\ell=1}^k y_{(\ell)} = \frac{1}{k} \sum_{\ell=1}^k (f(x_{(\ell)}) + \varepsilon_{(\ell)})$$

Assume that f is continuous and $f(x_{(\ell)}) \approx f(x_0) + \frac{1}{k} \sum_{\ell=1}^k \varepsilon_{(\ell)}$.
 This implies that $\frac{1}{k} \sum_{\ell=1}^k (f(x_{(\ell)}) + \varepsilon_{(\ell)}) \approx f(x_0) + \frac{1}{k} \sum_{\ell=1}^k \varepsilon_{(\ell)}$.
 Hence,

$$\begin{aligned} \text{Var}[\hat{f}_k(x_0)] &= \frac{1}{k^2} \sum_{\ell=1}^k \text{Var}(\varepsilon_{(\ell)}) \\ &= \frac{1}{k^2} \cdot k \cdot \sum_{\ell=1}^k \text{Var}(\varepsilon) = \frac{1}{k} \sigma_{\varepsilon}^2 \end{aligned}$$

In conclusion:

$$\boxed{\text{Err}(x_0) = \sigma_{\varepsilon}^2 + \left[f(x_0) - \frac{1}{k} \sum_{\ell=1}^k f(x_{(\ell)}) \right]^2 + \frac{\sigma_{\varepsilon}^2}{k}}$$

1.3 Variance of $\hat{\beta}$

1) Assumptions

1. y_i uncorrelated with constant variance σ^2
2. x_i fixed/not random

2) Derivation

Consider the linear model $y = \mathbf{X}\beta + \varepsilon$.

- Let \mathbf{X} be fixed. This implies $\mathbb{E}(y) = \mathbf{X}\beta$.
- $\mathbb{E}(\hat{\beta}) = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta$.
- So,

$$\begin{aligned} \hat{\beta} - \mathbb{E}\hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (y - \mathbf{X} \beta) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon \end{aligned}$$

$$\begin{aligned} \implies \text{Var}(\hat{\beta}) &= \mathbb{E}(\hat{\beta} - \mathbb{E}\hat{\beta})(\hat{\beta} - \mathbb{E}\hat{\beta})^T \\ &= \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon)(\varepsilon^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\varepsilon \varepsilon^T) \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\varepsilon) \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \quad \text{since } \mathbb{E}(\varepsilon) = 0 \end{aligned}$$

- Recall the assumption that all entries of y are uncorrelated with variance σ^2 . Hence,

$$\implies \text{Var}(\varepsilon) = \sigma^2 \mathbf{I}_N$$

$$\begin{aligned} \implies \text{Var}(\hat{\beta}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \end{aligned}$$

This matches equation (3.8) on page 47 of *ESL*.

1.4 Relationship between t and F statistic for dropping coefficients

Let RSS_0 be the result of the regression $Y \sim X_1 + X_2 + \dots + X_{p-1}$ (for simplicity, I will denote $X_{(p)} = (X_1, \dots, X_p)$), and let RSS_1 be the residual sum of squares from the regression of Y on the whole set of predictors, $Y \sim X_1 + \dots + X_p$. For this problem, the F -statistic for dropping the p^{th} predictor is

$$F = \frac{(RSS_0 - RSS_1)/(p_1 - p_0)}{RSS_1/(N - p_1 - 1)} = \frac{RSS_0 - RSS_1}{RSS_1/(N - p - 1)} = \frac{RSS_0 - RSS_1}{\hat{\sigma}^2}$$

For $\hat{\beta}_{(p)}$ as the set of coefficients for the predictor group $X_{(p)}$, we can explicitly define RSS_0 :

$$RSS_0 = (Y - X_{(p)}\hat{\beta}_{(p)})^T(Y - X_{(p)}\hat{\beta}_{(p)})$$

To analyze the RSS formulas, we will regress X_p on $X_{(p)}$ to get the residual Z_p .

By regressing $Y \sim X_{(p)} + Z_p$, we get the coefficients $(\hat{\beta}_p, \hat{\beta}_{(p)})$, such that $\hat{\beta}_p = \frac{\langle Y, Z_p \rangle}{\|Z_p\|^2}$.

Thus,

$$RSS_1 = (Y - X_{(p)}\hat{\beta}_{(p)} - Z_p\hat{\beta}_p)^T(Y - X_{(p)}\hat{\beta}_{(p)} - Z_p\hat{\beta}_p)$$

We can now find an expression for the numerator of the F -statistic:

$$\begin{aligned} RSS_0 - RSS_1 &= 2(Y - X_{(p)}\hat{\beta}_{(p)})^T Z_p \hat{\beta}_p - \|Z_p\|^2 \hat{\beta}_p^2 \\ &= 2Y^T Z_p \hat{\beta}_p - 2X_{(p)}^T Z_p \hat{\beta}_p - \|Z_p\|^2 \hat{\beta}_p^2 \\ &= 2Y^T Z_p \hat{\beta}_p - \|Z_p\|^2 \hat{\beta}_p^2 \quad \text{since } Z_p, \hat{\beta}_p \text{ are orthogonal} \\ &= \frac{\langle Y, Z_p \rangle^2}{\|Z_p\|^2} = \hat{\beta}_p^2 \|Z_p\|^2 \end{aligned}$$

In summary, we have found that $F = \frac{\hat{\beta}_p^2}{\hat{\sigma}^2} \|Z_p\|^2$.

From ESL, we see that $t = \frac{\hat{\beta}_p}{\hat{\sigma}\sqrt{v_p}}$, so $t^2 = \frac{\hat{\beta}_p^2}{\hat{\sigma}^2 v_p}$. Thus, to prove that $t^2 = F$, we need to focus on v_p .

By the definition of v_p , we know that $v_p = [(\mathbf{X}^T \mathbf{X})^{-1}]_{pp}$, where $\mathbf{X} = [X_{(p)} \quad X_p]$. So, we see that

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} X_{(p)}^T X_{(p)} & X_{(p)}^T X_p \\ X_p^T X_{(p)} & X_p^T X_p \end{bmatrix}$$

To invert this matrix, we can employ Schur's complement; for $\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$, the entry we seek in $\mathbf{M}^{-1} = (\mathbf{X}^T \mathbf{X})^{-1}$ is of the form $(\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B})^{-1}$. First, we will find $\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B}$:

$$\begin{aligned} \mathbf{D} - \mathbf{CA}^{-1}\mathbf{B} &= X_p^T X_p - (X_p^T X_{(p)})(X_{(p)}^T X_{(p)})^{-1} X_{(p)}^T X_p \\ &= \|X_p - X_p^T (X_{(p)}^T X_{(p)})^{-1} X_{(p)}\|^2 \\ &= \|X_p - \frac{\langle X_p, X_{(p)} \rangle}{\|X_{(p)}\|^2} X_{(p)}\|^2 \end{aligned}$$

Recall that Z_p is the residual from regressing X_p on $X_{(p)}$, so it is of the form $Z_p = X_p - \frac{\langle X_p, X_{(p)} \rangle}{\|X_{(p)}\|^2} X_{(p)}$, which is the expression in the norm above.

Hence, $\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B} = \|Z_p\|^2 \implies v_p = (\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B})^{-1} = (\|Z_p\|^2)^{-1}$.

Plugging this expression into t :

$$t^2 = \frac{\hat{\beta}_p^2}{\hat{\sigma}^2 v_p} = \frac{\hat{\beta}_p^2}{\hat{\sigma}^2} \|Z_p\|^2 = F$$

In conclusion, we see that, when dropping the last of p variables, $t^2 = F$.

2 Ridge regression

2.1 Derivation using the SVD

As derived in ESL, $\hat{\beta}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.

We will re-write this estimate in terms of the SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$.

Here, I use the convention that:

$$\begin{cases} \mathbf{U} \text{ is } n \times p \text{ orthonormal} \\ \mathbf{D} \text{ is } p \times p \text{ diagonal with singular values } d_i \text{ on the diagonal} \\ \mathbf{V} \text{ is } p \times p \text{ orthogonal} \end{cases}$$

$$\begin{aligned} \hat{\beta}(\lambda) &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T + \lambda \mathbf{I})^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{D}^2 \mathbf{V}^T + \lambda \mathbf{V} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \quad \text{since } \mathbf{U}^T \mathbf{U} = \mathbf{I}_{p \times p} = \mathbf{V} \mathbf{V}^T \\ &= (\mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}) \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{V}^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \quad \text{since } (\mathbf{V}^T)^{-1} = \mathbf{V} \text{ given that } \mathbf{V} \text{ is orthogonal} \\ &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y} \end{aligned}$$

Hence, in terms of the SVD, $\boxed{\hat{\beta}(\lambda) = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y}}$.

2.2 Eigenvalue interpretation

Consider $\hat{\beta}^{\text{ridge}}$ derived from the SVD:

$$\hat{\beta} = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y}$$

We can use this estimate to find \hat{y}^{ridge} :

$$\begin{aligned} \hat{y}^{\text{ridge}} &= \mathbf{X} \hat{\beta}^{\text{ridge}} \\ &= (\mathbf{U} \mathbf{D} \mathbf{V}^T) \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y} \quad \text{by the SVD of } \mathbf{X} \\ &= \mathbf{U} \mathbf{D} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y} \end{aligned}$$

Consider further the innermost matrix calculation $\mathbf{D} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D}$. This will be a diagonal matrix with elements $\frac{d_j^2}{d_j^2 + \lambda}$:

$$\mathbf{D} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} = \begin{bmatrix} \frac{d_1^2}{d_1^2 + \lambda} & & & \\ & \ddots & & \\ & & \frac{d_j^2}{d_j^2 + \lambda} & \\ & & & \ddots \end{bmatrix}$$

Next, we note that $\mathbf{U}^T \mathbf{y}$ represents the coordinates of \mathbf{y} in the basis spanned by the p columns of \mathbf{U} .

So, summing over the columns, we see that:

$$\begin{aligned} \hat{y}^{\text{ridge}} &= \mathbf{X} \hat{\beta}^{\text{ridge}} \\ &= \sum_{j=1}^p u_j \left(\frac{d_j^2}{d_j^2 + \lambda} \right) u_j^T \mathbf{y} \\ &= \sum_{j=1}^p u_j \left(\frac{d_j^2}{d_j^2 + \lambda} \right) \langle u_j, \mathbf{y} \rangle \end{aligned}$$

Thus, we see that ridge shrinks and regularizes according to the eigenvalues. Specifically, ridge has the greatest shrinkage effect on coefficients with small d_j corresponding to the directions in the column space of \mathbf{X} having small variance.

2.3 Relationship between λ and the coefficients

Again, recall that the ridge fit vector is of the form

$$\hat{\beta}(\lambda) = \mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y}$$

Using this decomposition, we now consider $\|\hat{\beta}_\lambda\|_2^2 = \|\mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y}\|_2^2$.

First, note that this expression is of the form $\|\hat{\beta}_\lambda\|^2 = \|\mathbf{V}\theta\|^2$, for $\theta = (\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y}$.

So, our first simplification is that

$$\begin{aligned}\|\hat{\beta}_\lambda\|^2 &= \|\mathbf{V}\theta\|^2 = (\mathbf{V}\theta)^T(\mathbf{V}\theta) \\ &= \theta^T\mathbf{V}^T\mathbf{V}\theta \\ &= \theta^T\theta = \|\theta\|^2 \text{ since } \mathbf{V} \text{ is orthogonal}\end{aligned}$$

Continuing with the expansion of the norm expression, $\|\hat{\beta}_\lambda\|^2 = \|\theta\|^2 = \|(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y}\|^2$.

Again taking particular note of the expressions involving \mathbf{D} , we recall that:

$$(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D} = \begin{bmatrix} \frac{d_1}{d_1^2 + \lambda} & & & \\ & \ddots & & \\ & & \frac{d_j}{d_j^2 + \lambda} & \\ & & & \ddots \end{bmatrix}$$

So, the j^{th} component of θ is

$$[(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y}]_j = \frac{d_j}{d_j^2 + \lambda} \langle \mathbf{U}, \mathbf{y} \rangle_j$$

Hence,

$$\|\hat{\beta}_\lambda\|^2 = \sum_j \left(\frac{d_j}{d_j^2 + \lambda} \right)^2 \left(\langle \mathbf{U}, \mathbf{y} \rangle_j \right)^2$$

We therefore see that as $\lambda \downarrow$, $\boxed{\|\hat{\beta}_\lambda\|^2 \uparrow}$

2.4 Effect of collinearity on ridge regression

Let $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_p]$.

Suppose two of the predictors are perfectly collinear and orthogonal to the remaining predictors.

More concretely, consider $\mathbf{x}_1 = \mathbf{x}_2$ and $\mathbf{x}_1^T \mathbf{x}_j = 0 \ \forall j \in \{3, 4, \dots, p\}$.

Recall that the ridge solution is $\hat{\beta}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.

We will first consider the term inside the parentheses:

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \mathbf{x}_1^T \mathbf{x}_3 & \cdots \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \mathbf{x}_2^T \mathbf{x}_3 & \cdots \\ \mathbf{x}_3^T \mathbf{x}_1 & \mathbf{x}_3^T \mathbf{x}_2 & \mathbf{x}_3^T \mathbf{x}_3 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} = \begin{bmatrix} \|\mathbf{x}_1\|^2 & \|\mathbf{x}_1\|^2 & 0 & \cdots & 0 \\ \|\mathbf{x}_1\|^2 & \|\mathbf{x}_1\|^2 & 0 & \cdots & 0 \\ 0 & 0 & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \cdots & \mathbf{B} \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} = \begin{bmatrix} \|\mathbf{x}_1\|^2 + \lambda & \|\mathbf{x}_1\|^2 & \mathbf{0} \\ \|\mathbf{x}_1\|^2 & \|\mathbf{x}_1\|^2 + \lambda & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{B} + \lambda \mathbf{I} \end{bmatrix}$$

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} = \begin{bmatrix} \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} \begin{bmatrix} \|\mathbf{x}_1\|^2 + \lambda & -\|\mathbf{x}_1\|^2 \\ -\|\mathbf{x}_1\|^2 & \|\mathbf{x}_1\|^2 + \lambda \end{bmatrix} & \mathbf{0} \\ \mathbf{0} & (\mathbf{B} + \lambda \mathbf{I})^{-1} \end{bmatrix}$$

Given that the variables are identical, $\hat{\beta}_1(\lambda) = \hat{\beta}_2(\lambda)$. We will consider $\hat{\beta}_1(\lambda)$.

To find $\hat{\beta}_1(\lambda)$, we multiply the first row of $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}$ by $\mathbf{X}^T \mathbf{y}$:

$$\begin{aligned} \hat{\beta}_1(\lambda) &= \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} [\|\mathbf{x}_1\|^2 + \lambda, \ -\|\mathbf{x}_1\|^2, \ 0, \ \dots, \ 0] \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} [\|\mathbf{x}_1\|^2 + \lambda, \ -\|\mathbf{x}_1\|^2, \ 0, \ \dots, \ 0] \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_p^T \end{bmatrix} \mathbf{y} \\ &= \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} [\|\mathbf{x}_1\|^2 + \lambda, \ -\|\mathbf{x}_1\|^2] \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \end{bmatrix} \mathbf{y} \end{aligned}$$

Noting that $\mathbf{x}_1 = \mathbf{x}_2$:

$$\begin{aligned} \hat{\beta}_1(\lambda) &= \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} [\|\mathbf{x}_1\|^2 + \lambda, \ -\|\mathbf{x}_1\|^2] \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_1^T \end{bmatrix} \mathbf{y} \\ &= \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} ((\|\mathbf{x}_1\|^2 + \lambda)(\mathbf{x}_1^T \mathbf{y}) - \|\mathbf{x}_1\|^2 \mathbf{x}_1^T \mathbf{y}) \\ &= \frac{1}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} (\|\mathbf{x}_1\|^2 \mathbf{x}_1^T \mathbf{y} + \lambda \mathbf{x}_1^T \mathbf{y} - \|\mathbf{x}_1\|^2 \mathbf{x}_1^T \mathbf{y}) \\ &= \frac{\lambda \mathbf{x}_1^T \mathbf{y}}{2\lambda\|\mathbf{x}_1\|^2 + \lambda^2} \\ &= \frac{\mathbf{x}_1^T \mathbf{y}}{2\|\mathbf{x}_1\|^2 + \lambda} \end{aligned}$$

In conclusion, we see that $\hat{\beta}_1(\lambda) = \hat{\beta}_2(\lambda) = \boxed{\frac{\mathbf{x}_1^T \mathbf{y}}{2\|\mathbf{x}_1\|^2 + \lambda}}$.

When $\lambda = 0$, we return to OLS estimation.

With multicollinearity present, $(\mathbf{X}^T \mathbf{X})^{-1}$ does not exist and the problem is intractable.

However, in the limiting case where $\lambda \rightarrow 0$, we see that $\hat{\beta}_1(\lambda) \rightarrow \frac{\mathbf{x}_1^T \mathbf{y}}{2\|\mathbf{x}_1\|^2}$.

2.5 Ridge when $p \gg N$

Let p be the number of predictors for N rows in matrix \mathbf{X} .

For $p \gg N$, consider the SVD of $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$. In this case, \mathbf{D} will look like:

$$\mathbf{D} = \begin{bmatrix} d_1 & & & & & \\ & d_2 & & & & \\ & & \ddots & & & \\ & & & d_n & & \\ & & & & 0 & \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix}$$

Note that the off-diagonal entries of \mathbf{D} are also 0.

Recalling the result from the previously-derived ridge solution, we know that $\hat{\beta}(\lambda) = \mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{Y}$.

Consider the limiting ridge solution as $\lambda \rightarrow 0$. We will focus on the innermost term in the expression, $(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}$.

$$(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D} = \begin{bmatrix} \frac{d_1}{d_1^2 + \lambda} & & & & & \\ & \frac{d_2}{d_2^2 + \lambda} & & & & \\ & & \ddots & & & \\ & & & \frac{d_n}{d_n^2 + \lambda} & & \\ & & & & 0 & \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix}$$

$$\lim_{\lambda \rightarrow 0} ((\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}) = \begin{bmatrix} \frac{1}{d_1} & & & & & \\ & \frac{1}{d_2} & & & & \\ & & \ddots & & & \\ & & & \frac{1}{d_n} & & \\ & & & & 0 & \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix} = \begin{cases} 1/d_i, & d_i \neq 0 \\ 0, & d_i = 0 \end{cases}$$

$$\text{Let } \mathbf{D}^\dagger = \begin{bmatrix} d_1^\dagger & & \\ & d_2^\dagger & \\ & & \ddots \end{bmatrix} = \begin{cases} 1/d_i, & d_i \neq 0 \\ 0, & d_i = 0 \end{cases}$$

This limiting value of \mathbf{D} is the generalized inverse of \mathbf{D} , and can be plugged into the overall limiting ridge solution:

$$\boxed{\lim_{\lambda \rightarrow 0} \hat{\beta}(\lambda) = \mathbf{V}\mathbf{D}^\dagger\mathbf{U}^T\mathbf{Y}}$$

2.6 Bayesian theory: ridge estimate is the mean of the posterior

Assume that $\mathbf{Y} \mid \beta \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I})$ and $\beta \sim N(0, \tau^2 \mathbf{I})$. Thus:

$$f(\mathbf{Y} \mid \beta) \propto \exp\left(-\frac{\|\mathbf{Y} - \mathbf{X}\beta\|^2}{2\sigma^2}\right) \quad \text{and} \quad f(\beta) \propto \exp\left(-\frac{\|\beta\|^2}{2\tau^2}\right)$$

Note that $f(\beta \mid \mathbf{Y}) = \frac{f(\beta, \mathbf{Y})}{f(\mathbf{Y})} = \frac{f(\beta)f(\mathbf{Y} \mid \beta)}{f(\mathbf{Y})} \propto f(\beta)f(\mathbf{Y} \mid \beta)$ since $f(\mathbf{Y})$ is a constant in terms of the function. So, plugging in the expressions we found for these functions,

$$f(\beta \mid \mathbf{Y}) \propto \exp\left(-\frac{\|\mathbf{Y} - \mathbf{X}\beta\|^2}{2\sigma^2} - \frac{\|\beta\|^2}{2\tau^2}\right) = \exp\left(-\left(\frac{(\mathbf{Y} - \mathbf{X}\beta)^T(\mathbf{Y} - \mathbf{X}\beta)}{2\sigma^2} + \frac{\beta^T \beta}{2\tau^2}\right)\right)$$

We will focus our simplification efforts on the function inside the exponential.

$$\begin{aligned} \frac{(\mathbf{Y} - \mathbf{X}\beta)^T(\mathbf{Y} - \mathbf{X}\beta)}{2\sigma^2} + \frac{\beta^T \beta}{2\tau^2} &= \frac{\mathbf{Y}^T \mathbf{Y} - 2\mathbf{Y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta}{2\sigma^2} + \frac{\beta^T \beta}{2\tau^2} \\ &= \beta^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) \beta - \left(\frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2}\right) \beta + c \end{aligned}$$

This expression is quadratic in β ; we will factor it into the form $(\beta - \mu)^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) (\beta - \mu) + c$.

Such a factorization would show that $f(\beta \mid \mathbf{Y}) \propto \exp((\beta - \mu)^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) (\beta - \mu))$, which would imply that $E[f(\beta \mid \mathbf{Y})] = \mu$ since $f(\beta \mid \mathbf{Y})$ would be distributed multivariate Gaussian.

Since $(\beta - \mu)^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) (\beta - \mu) + c = \beta^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) \beta - \left(\frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2}\right) \beta + c'$, we need $-2\mu^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) \beta = -\left(\frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2}\right) \beta$.

Focusing on this expression to solve for μ :

$$\begin{aligned} -2\mu^T \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right) \beta &= -\left(\frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2}\right) \beta \\ \implies \mu &= \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right)^{-1} \left(\frac{\mathbf{Y}^T \mathbf{X}}{2\sigma^2}\right)^T \\ \implies \mu &= \left(\frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2}\right)^{-1} \left(\frac{\mathbf{X}^T \mathbf{Y}}{2\sigma^2}\right) = \left(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I}\right)^{-1} \mathbf{X}^T \mathbf{Y} \end{aligned}$$

Thus, $\mu = (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$ is the mean of the posterior distribution, $f(\beta \mid \mathbf{Y})$.

Recall that the ridge regression solution is $\hat{\beta}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$. This solution is identical to our result with the stipulation that $\lambda = \sigma^2 / \tau^2$. Hence, the ridge estimate is indeed the mean of the posterior distribution.

3 Degrees of freedom

3.1 Degrees of freedom: ridge

The effective degrees of freedom for ridge is the trace of the fitted "hat" matrix. Mathematically:

$$\text{df}(\lambda) = \text{Tr}(\mathbf{H}_\lambda) = \text{Tr}[\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T]$$

Using the SVD of \mathbf{X} , we can gain insight into how the degrees of freedom expression relates to the singular values.

$$\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T = \mathbf{U} \mathbf{D}(\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T$$

In section 2.2, we showed that the eigenvalues of the above expression are given by $\frac{d_j^2}{d_j^2 + \lambda}$.

Hence, the trace of the above expression is the sum of these eigenvalues:

$$\boxed{\text{df}(\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}}$$

Hence, we arrive at expression (3.50) in ESL.

As a footnote to this exercise, note that:

- $\text{df}(\lambda) = p$ when $\lambda = 0$
- $\text{df}(\lambda) \rightarrow 0$ when $\lambda \rightarrow \infty$

3.2 Degrees of freedom: LAR and lasso

In section 3.4.4 of ESL, the authors delve into the details of least angle regression (LAR) and its connection to the lasso. We need a different notion of effective degrees of freedom to understand how fitting to the data influences the number of predictors that enter the final set in these selection-oriented algorithms. As presented in ESL, we consider the degrees of freedom of the fitted vector $\hat{\mathbf{y}}$:

$$\text{df}(\hat{\mathbf{y}}) = \frac{1}{\sigma^2} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i)$$

This definition also leads to the known ridge result derived in the above section 3.1; this shows that considering degrees of freedom this way induces an "apples to apples" comparison between these adaptive methods and the more straightforward linear methods previously under consideration.

3.2.1 Derivation: linear regression

For illustrative purposes, I will show how the above definition for degrees of freedom leads to the known linear regression result, where the degrees of freedom of the fit equals the number of predictors in the model ($\text{df} = p$).

Let e_i be the unit vector with i^{th} element = 1 and the rest = 0.

Hence, we can rewrite the i^{th} element of the covariance expression as:

$$\text{Cov}(\hat{y}_i, y_i) = \text{Cov}(e_i^T \hat{\mathbf{y}}, e_i^T \mathbf{y}) = e_i^T \text{Cov}(\hat{\mathbf{y}}, \mathbf{y}) e_i$$

Recall the standard OLS result that $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Thus,

$$\begin{aligned} \text{Cov}(\hat{\mathbf{y}}, \mathbf{y}) &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Cov}(\mathbf{y}, \mathbf{y}) \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \text{ since } \text{Cov}(\mathbf{y}, \mathbf{y}) = \sigma^2 \mathbf{I} \end{aligned}$$

We can now reduce this expression to the i^{th} element:

$$\begin{aligned} \text{Cov}(\hat{y}_i, y_i) &= e_i^T \text{Cov}(\hat{\mathbf{y}}, \mathbf{y}) e_i \\ &= \sigma^2 e_i^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T e_i \\ &= \sigma^2 (\mathbf{X}^T e_i)^T (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T e_i) \end{aligned}$$

Note here that $(\mathbf{X}^T e_i)^T = x_i$ for $1 \leq i \leq N$. Hence, we see that

$$\text{Cov}(\hat{y}_i, y_i) = \sigma^2 x_i^T (\mathbf{X}^T \mathbf{X})^{-1} x_i$$

To arrive at our final expression, we sum from $i = 1$ to N and divide by σ^2 .

$$\begin{aligned}
\text{df}(\hat{\mathbf{y}}) &= \frac{1}{\sigma^2} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) \\
&= \frac{1}{\sigma^2} \sum_{i=1}^N \sigma^2 x_i^T (\mathbf{X}^T \mathbf{X})^{-1} x_i \\
&= \sum_{i=1}^N x_i^T (\mathbf{X}^T \mathbf{X})^{-1} x_i \\
&= \sum_{i=1}^N \text{Tr} \left[x_i^T (\mathbf{X}^T \mathbf{X})^{-1} x_i \right] \\
&= \sum_{i=1}^N \text{Tr} \left[x_i x_i^T (\mathbf{X}^T \mathbf{X})^{-1} \right] \\
&= \text{Tr} \left[\left(\sum_{i=1}^N x_i x_i^T \right) (\mathbf{X}^T \mathbf{X})^{-1} \right] \\
&= \text{Tr} \left[(\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \right] \\
&= \text{Tr}(\mathbf{I}_{p \times p}) \\
&= p
\end{aligned}$$

In conclusion, we see that the new definition of degrees of freedom matches the existing result for linear regression. When the linear model has $k < p$ predictors, we see also that $\text{df} \rightarrow k$.

3.2.2 Derivation: ridge

Let us also check that this definition for the degrees of freedom leads to the previously derived ridge result. Recall again for ridge that the fitted vector has the form:

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta}^{\text{ridge}} = \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

First, let us find $\text{Cov}(\hat{\mathbf{y}}, \mathbf{y})$:

$$\begin{aligned}
\text{Cov}(\hat{\mathbf{y}}, \mathbf{y}) &= \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \text{Cov}(\mathbf{y}, \mathbf{y}) \\
&= \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \sigma^2
\end{aligned}$$

Adapting this to the result for the i^{th} element:

$$\begin{aligned}
\text{Cov}(\hat{y}_i, y_i) &= e_i^T \text{Cov}(\hat{\mathbf{y}}, \mathbf{y}) e_i \\
&= \sigma^2 (\mathbf{X}^T e_i)^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T e_i) \\
&= \sigma^2 x_i^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} x_i
\end{aligned}$$

Finally, as in the OLS derivation, we sum these from $i = 1$ to N and divide by σ^2 .

$$\begin{aligned}
\text{df}(\hat{\mathbf{y}}) &= \frac{1}{\sigma^2} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) \\
&= \sum_{i=1}^N \text{Tr} \left[x_i x_i^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \right] \\
&= \text{Tr} \left[\left(\sum_{i=1}^N x_i x_i^T \right) (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \right] \\
&= \text{Tr} \left[(\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \right] \\
&= \text{Tr} \left[\mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \right] \\
&= \text{Tr}(\mathbf{H}_\lambda)
\end{aligned}$$

This is the same starting point as in the above section 3.1, which was shown to lead to expression (3.50) in ESL.

4 Methods using derived input directions

4.1 Principal component regression up to j coefficients

4.2 PCR connection to OLS/ridge

4.3 PCR vs. ridge

5 Zero-residual fits

5.1 There are infinitely many OLS solutions with zero residuals

5.2 Ridge solution has minimum norm among all zero-residual solutions as $\lambda \rightarrow 0$