

# **Elements Of Statistical Learning Solutions**

**Solutions to select exercises**

Krishan Bhalla

June 3, 2020

---

## CONTENTS

---

1	INTRODUCTION	3
2	CHAPTER 2 - OVERVIEW OF SUPERVISED LEARNING	4
3	CHAPTER 3 - LINEAR METHODS FOR REGRESSION	7

---

## INTRODUCTION

---

A selection of solutions to Elements of Statistical Learning by Hastie, Tibshirani, and Friedman. This will be continually updated as I work through the book.

I will not answer all exercises, but will endeavour to solve many.

This chapter largely exists to align latex chapter labels with those of the book.

---

## CHAPTER 2 - OVERVIEW OF SUPERVISED LEARNING

---

2.1 Suppose each of  $K$ -classes has an associated target  $t_k$  which is a vector of all zeros except one in the  $k$ th position. Show that classifying to the largest element of  $\hat{y}$  amounts to choosing the closest target  $\min_k \|t_k - \hat{y}\|$  if the elements of  $y$  sum to one.

Our norm here is the  $L^2$  norm.

$$\begin{aligned}
 \operatorname{argmin}_k \|t_k - \hat{y}\| &= \operatorname{argmin}_k \|t_k - \hat{y}\|^2 \\
 &= \operatorname{argmin}_k \sum_i (\hat{y}_i - (t_k)_i)^2 \\
 &= \operatorname{argmin}_k \sum_i (\hat{y}_i - \delta_{i,k})^2 \\
 &= \operatorname{argmin}_k \left( 1 - 2\hat{y}_k + \sum_i \hat{y}_i^2 \right) \\
 &= \operatorname{argmin}_k (-2\hat{y}_k)
 \end{aligned}$$

Where the last line follows as it is the only term dependant on  $k$ .  $\operatorname{argmin}$  is independent of scale, so the above states that the  $k$  corresponding to the minimum value of the norm is exactly the largest element of  $\hat{y}$

2.2 Show how to compute the Bayes decision boundary for the example in Figure 2.5.

Setup: There are 10 means  $m_k$  from a bivariate gaussian distribution  $N((1,0)^T, I)$  and we label this class blue. We take 10  $n_k$  from  $N((0,1)^T, I)$  and label this class orange. For each class there were 100 observations, where each observation was generated by picking one of the  $m_k$  ( $n_k$  resp) each with equal probability, and taking a point randomly from a  $N(m_k, I/5)$  distribution.

Let  $x$  be an observation. We equate posteriors for  $x$  being blue or yellow, and note that in our setup  $\mathbb{P}(\text{blue}) = \mathbb{P}(\text{orange})$  (priors are equal) to simplify to:

$$\sum_k \exp(-5 \|m_k - x\|^2) = \sum_k \exp(-5 \|n_k - x\|^2)$$

This defines a curve in the plane separating the two classes.

### 2.3 Derive equation 2.24.

Setup: Consider  $N$  data points uniformly distributed in a  $p$ -dimensional unit ball around the origin. Consider a nearest neighbour estimate at the origin. The median distance from the origin to the closest point is given by:

$$d(p_N) = \left(1 - \frac{1}{2}\right)^{1/p}$$

Solution: Let  $x$  be a point and let  $y = \|x\|$ .  $y$  has cdf equal to the ratio of a ball of radius  $y$  to a ball of radius 1, i.e.  $F(y) = y^p$ . The minimum over all  $x$  then has cdf  $F_{ymin}(y) = 1 - (1 - F(y))^N$  (a general fact for order statistics). Thus  $F_{ymin} = 1 - (1 - y^p)^N$ .

The median distance for  $ymin$  is when  $F_{ymin}(y) = 1/2$ . Solving for this yields the result.

### 2.4 Setup as in book. Projection in direction $a$ .

Pick an orthonormal basis of  $\mathbb{R}^p$  which includes the vector  $a$ , say  $a_1, \dots, a_p$  with  $a_1 = a$ . Then each  $x_i = \sum_j X_{i,j} a_j$ , and so  $z_i = X_{i,1}$  where  $X$  is the matrix with rows  $x_i$

The  $x_i$  have distribution  $N(0, I_p)$ , and under such a distribution each component of  $x_i$  has distribution  $N(0, 1)$ .

In particular this means that each  $X_{i,j}$  has distribution  $N(0, 1)$  and so the  $z_i$  do.

The squared distance from the origin is just  $z_i^2$ , with distribution  $\chi_1^2$ , and this has mean 1.

2.6 Consider a regression problem with inputs  $x_i$  and outputs  $y_i$ , and a parameters model  $f_\theta(x)$  to be fit by least squares. Show that if there are observations with tied or identical values of  $x$  then the fir can be obtained from a reduced weighted least squares problem.

The problem can of finding  $\theta$  amounts to solving the following:

$$\operatorname{argmin}_\theta (y - f_\theta(x))^T (y - f_\theta(x))$$

Denote by  $z_1, \dots, z_M$  the unique values of  $x$  in our training set, denote by  $n_j$  the number of occurrences of value  $z_j$ . Then let  $t_j = \frac{1}{n_j} \sum_{i=1}^{n_j} y_i$ . If we can get to the following, we're done (as it is in the form of weighted regression).

$$\operatorname{argmin}_\theta \sum_j n_j (t_j - f_\theta(z_j))^2$$

Expanding the initial expression we get (denoting by  $y_{i,j}$  the  $i$ th value of  $y$  corresponding to input  $z_j$ ):

$$\begin{aligned}
(y - f_\theta(x))^T (y - f_\theta(x)) &= \sum_i (y_i - f_\theta(x_i))^2 \\
&= \sum_i (y_i^2 + f_\theta(x_i)^2 - 2y_i f_\theta(x_i)) \\
&= \sum_i (y_i^2 + f_\theta(x_i)^2 - 2y_i f_\theta(x_i)) \\
&= \sum_{j=1}^M \sum_{i=1}^{n_j} y_{i,j}^2 - 2f_\theta(z_j) y_{i,j} + f_\theta(z_j)^2 \\
&= \sum_{j=1}^M \left( \sum_i y_{i,j}^2 \right) - 2n_j f_\theta(z_j) t_j + n_j f_\theta(z_j)^2 \\
&= \sum_{j=1}^M n_j (t_j - f_\theta(z_j))^2 - \sum_{j=1}^M n_j t_j^2 + \sum_{j=1}^M \left( \sum_i y_{i,j}^2 \right)
\end{aligned}$$

This last trick of adding 0 leaves us with an expression where only the first sum is dependant on  $\theta$ . Thus when taking  $\operatorname{argmin}_\theta$  we can ignore the last two terms.

This leaves us with the required equivalence.:

$$\operatorname{argmin}_\theta (y - f_\theta(x))^T (y - f_\theta(x)) = \operatorname{argmin}_\theta \sum_j n_j (t_j - f_\theta(z_j))^2$$

---

CHAPTER 3 - LINEAR METHODS FOR REGRESSION

---

3.1 Show that the F-statistic for dropping a single coefficient of a model is equivalent to the square of the corresponding z score

Let  $X$  be our data, and let  $v_{j,j}$  be the  $j$ th diagonal element of  $V = (X^T X)^{-1}$ .  $z_j = \hat{\beta}_j / \hat{\sigma} \sqrt{v_{j,j}}$  is the z score.

The F statistic is

$$F = \frac{(RSS_0 - RSS_1) / (p_1 - p_0)}{RSS_1 / (N - p_1 - 1)}$$

Where the regression models are have  $p_1 + 1$  and  $p_0 + 1$  degrees of freedom respectively. We also know that  $\hat{\sigma}^2$  is equivalent to the denominator. In the case of dropping a single variable, this simplifies to:

$$F = \frac{RSS_0 - RSS_1}{\hat{\sigma}^2}$$

Where  $\hat{\sigma}$  is derived from the bigger model.

Thus our question can be simplified to showing that:

$$RSS_0 - RSS_1 = \hat{\beta}_j^2 / v_{j,j}$$

We know  $\hat{\beta} \sim N(\beta, \sigma^2 V)$  and so  $\hat{\beta}_j \sim N(\beta_j, \sigma^2 v_{j,j})$ .

Under the null-hypothesis  $\beta_j = 0$  and so:  $\hat{\beta}_j = \sigma \sqrt{v_{j,j}} Z$  where  $Z \sim N(0, 1)$  Thus  $\hat{\beta}_j^2 = \sigma^2 v_{j,j} Z^2 = \sigma^2 v_{j,j} Q$  where  $Q \sim \chi_1^2$

Similarly  $RSS_0, RSS_1$  have distribution  $\sigma^2 \chi_{N_i}^2$  where  $N_i$  is the number of degrees of freedom. Thus  $RSS_0 - RSS_1 \sim \sigma^2 \chi_1^2$

Hence  $\hat{\beta}_j^2 / v_{j,j}$  and  $RSS_0 - RSS_1$  have the same distribution. They further test the same hypothesis and thus must be identical.

3.3 a. Prove the Gauss-Markov theorem: the least squares estimate of a parameter  $a^T \beta$  has a variance no bigger than that of any other linear unbiased estimate of  $a^T \beta$ .

b. Secondly, show that if  $\hat{V}$  is the variance-covariance matrix of the least squares estimate of  $\beta$  and  $\tilde{V}$  is the variance covariance matrix of any other linear unbiased estimate, then  $\hat{V} \leq \tilde{V}$ , where  $B \leq A$  if  $A - B$  is positive semidefinite.

First note that part b implies part a. If  $\beta$  has dimension 1, then  $V$  is just the variance of beta, and  $\leq$  is equivalent to the normal  $\leq$  operator. Taking the inner product with  $a$  is just a linear operation. We thus only need to show b.

Suppose  $\hat{\beta}$  is the OLS estimate of  $\beta$  and that  $\tilde{\beta}$  is another linear unbiased estimate. The the variance-covariance matrices be  $\hat{V}$  and  $\tilde{V}$  resp.  $\hat{\beta} = (X^T X)^{-1} X^T y$  so  $\hat{\beta} = Cy$  say, and write  $\tilde{\beta} = (C + D)y$  for some non-zero  $m \times n$  matrix  $D$ .

$$\begin{aligned} \mathbb{E}(\tilde{\beta}) &= \mathbb{E}((C + D)(X\beta + \epsilon)) \\ &= \mathbb{E}\left(\left(X^T X\right)^{-1} X^T + D\right)(X\beta + \epsilon) \\ &= \mathbb{E}((\beta + DX\beta)) + \zeta \mathbb{E}(\epsilon) \\ &= \mathbb{E}((\beta + DX\beta)) \\ &= \beta + DX\beta \\ &= \beta \end{aligned}$$

Where the last lines follow as  $\tilde{\beta}$  is unbiased. In particular  $DX\beta = 0$  and so  $DX = 0$  as beta is an unobserved parameter to be estimated.

$$\begin{aligned} \tilde{V} &= \text{Var}(\tilde{\beta}) \\ &= \text{Var}((C + D)y) \\ &= (C + D)(C + D)^T \text{Var}(y) \\ &= \sigma^2 (C + D)(C + D)^T \\ &= \sigma^2 (CC^T + CD^T + DC^T + DD^T) \\ &= \sigma^2 (CC^T + CD^T + DC^T + DD^T) \\ &= \hat{V} + \sigma^2 (+CD^T + DC^T + DD^T) \\ &= \hat{V} + \sigma^2 \left( \left(X^T X\right)^{-1} X^T D^T + DX \left(X^T X\right)^{-1} + D^T D \right) \\ &= \hat{V} + \sigma^2 DD^T \end{aligned}$$

Where we know that  $DD^T$  is positive semi-definite<sup>1</sup> and so are done.

<sup>1</sup>  $v^T DD^T v = \|D^T v\|^2 \geq 0 \forall v$



3.4 Show how the vector of least squares coefficients can be obtained from a single pass of the Gram–Schmidt procedure (Algorithm 3.1). Represent your solution in terms of the QR decomposition of  $X$ .

Let  $X = QR$  be the QR decomposition of  $X$ . This can be attained via Gram-Schmidt. We assume that  $R$  has no zeros on the diagonal (i.e. the variables are linearly independent). Then

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T y \\ &= (R^T R)^{-1} R^T Q^T y \\ &= R^{-1} Q^T y\end{aligned}$$

We can invert  $R$  via backpropagation, and  $Q^T y$  is easily calculable.

3.5 Show that the ridge regression problem (using  $\alpha$  to denote the constant intercept vector)

$$\hat{\beta} = \operatorname{argmin}_{\beta, \alpha} (y - \alpha - X\beta)^T (y - \alpha - X\beta) + \lambda \|\beta\|^2$$

is equivalent to the problem:

$$\hat{\beta}^c = \operatorname{argmin}_{\beta^c, \alpha^c} (y - \alpha^c - \tilde{X}\beta^c)^T (y - \alpha^c - \tilde{X}\beta^c) + \lambda \|\beta^c\|^2$$

Where  $\tilde{X} = X - \bar{X}$ , and bar represents the  $N \times p$  matrix where each value in the  $j$ th column is the mean  $x_j$ . Give the correspondence between  $\beta^c$  and the original  $\beta$ . Do the same for the lasso.<sup>2</sup>

This problem is easier with summation

$$\begin{aligned}& (y - \alpha - X\beta)^T (y - \alpha - X\beta) + \lambda \|\beta\|^2 \\ &= (y - \alpha - \bar{X}\beta - (X - \bar{X})\beta)^T (y - \alpha - \bar{X}\beta - (X - \bar{X})\beta) + \lambda \|\beta\|^2 \\ &= (y - (\alpha + \bar{X}\beta) - (X - \bar{X})\beta)^T (y - (\alpha + \bar{X}\beta) - (X - \bar{X})\beta) + \lambda \|\beta\|^2 \\ &= (y - \alpha^c - \tilde{X}\beta^c)^T (y - \alpha^c - \tilde{X}\beta^c) + \lambda \|\beta^c\|^2\end{aligned}$$

Where  $\alpha_i^c = \alpha_i + \sum_{j=1}^p \bar{x}_j \beta_j$  in all coordinates, and  $\beta^c = \beta$ . This is an expression of our desired form. This problem is equivalent to demeaning the data and adjusting the intercept. One can do exactly the same for the lasso.

<sup>2</sup> This is all a tad odd as regression with intercept is desired, but we only penalise the non-intercept terms.

3.6 Show that the ridge regression estimate is the mean (and mode) of the posterior distribution, under a Gaussian prior  $\beta \sim N(0, \tau I)$  and Gaussian sampling model  $y \sim N(X\beta, \sigma^2 I)$ . Find the relationship between the ridge parameter  $\lambda$  and the variances  $\tau$  and  $\sigma^2$

This question really states that the pdf of the posterior is proportional to the pdfs of  $y$  given  $\beta$  and  $\beta$ . Hence:

$$f(\beta|y) \propto f(\beta)f(y|\beta) \quad (1)$$

$$\log(f(\beta|y)) = C + \log(f(\beta)) + \log(f(y|\beta)) \quad (2)$$

$$= C + -\frac{1}{2\tau}\beta^T\beta * \log\left(\frac{1}{\sqrt{2 * \pi * \tau}}\right) \quad (3)$$

$$- \frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) * \log\left(\frac{1}{\sqrt{2 * \pi * \sigma^2}}\right) \quad (4)$$

$$= C + -\frac{1}{2\tau}\beta^T\beta + -\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) \quad (5)$$

Where we absorb terms into the constant as needed. We have recovered that:

$$f(\beta|y) = C_1 e^{-\frac{1}{2\sigma^2} (\frac{\sigma^2}{\tau}\beta^T\beta + (y - X\beta)^T (y - X\beta))} \quad (6)$$

$$= C_1 e^{-\frac{1}{2\sigma^2} (\beta^T (X^T X + \frac{\sigma^2}{\tau} I) \beta + y^T y - y^T X\beta - \beta^T X^T y)} \quad (7)$$

$$= C_1 e^{-\frac{1}{2\sigma^2} (\beta^T \Sigma \beta + y^T y - y^T X\beta - \beta^T X^T y)} \quad (8)$$

$$= C_1 e^{-\frac{1}{2\sigma^2} ((\Sigma\beta - X^T y)^T \Sigma^{-1} (\Sigma\beta - X^T y) + y^T y - y^T X^T X y)} \quad (9)$$

$$= C_2 e^{-\frac{1}{2\sigma^2} (\Sigma\beta - X^T y)^T \Sigma^{-1} (\Sigma\beta - X^T y)} \quad (10)$$

$$(11)$$

Where  $\Sigma = (X^T X + \frac{\sigma^2}{\tau} I)$  is a  $(p + 1) * (p + 1)$  matrix, and for (11) we absorbed into the constant any terms not dependant on  $\beta$

Thus the posterior has a multivariate gaussian pdf (up to scaling), so the mean and mode of the distribution are identical, and can be found by maximising (5) over  $\beta$ , which due to the minus sign is sufficient.

This amounts to

$$\operatorname{argmin}_{\beta} \left( \frac{1}{\tau}\beta^T\beta + \frac{1}{\sigma^2} (y - X\beta)^T (y - X\beta) \right)$$

And hence is equivalent to solving:

$$\operatorname{argmin}_{\beta} \left( \frac{\sigma^2}{\tau}\beta^T\beta + (y - X\beta)^T (y - X\beta) \right)$$

Hence the ridge regression parameter is  $\sigma^2/\tau$

3.9 *Forward stepwise regression:* Suppose we have the QR decomposition for the  $N \times q$  matrix  $X_1$  in a multiple regression problem with response  $y$ , and suppose we have an additional  $p - q$  predictors in the matrix  $X_2$ . Denote the residual by  $r$ . Describe an efficient procedure for establishing which additional variable will reduce the residual sum of squares the most.

Intuition - pick the column  $\hat{v}$  such that  $v$  has the least angle with  $r$ . i.e.

$$\hat{v} = \operatorname{argmax}_{v \in \{\text{columns } X_2\}} \frac{|r^T v|}{\|v\|}$$

With this in mind, let  $u_j = x_j - \frac{|r^T x_j|}{\|x_j\|} \frac{x_j}{\|x_j\|}$  where the  $x_j$  are the columns of  $X_2$ , and let  $v_j = \frac{u_j}{\|u_j\|}$

$$RSS = r^T r$$

$$r = y - \hat{y}$$

$$= y - R^{-1} Q^T y$$

$$\text{let } r_j = r - (r^T u_j) u_j$$

$$\text{and } RSS_j = RSS - 2(r^T u_j)^2 + (r^T u_j)^2$$

$$\text{then } RSS_j = RSS - (r^T u_j)^2$$

Where  $RSS_j$  is the residual sum of squares for our new model. This verifies our intuition,  $RSS_j$  is minimised when we pick the column of  $X_2$  with least angle to  $r$ . I assume the efficiency in the question comes from the ease of inverting  $R$  compared to  $X^T X$  (backpropagation will do), and that we can extend our QR decomposition to include the new variable easily via Gram-Schmidt. In particular most of what is needed for Gram-Schmidt is already computed when taking inner product with the residual.

3.10 *Backward stepwise regression.* Suppose we have the multiple regression fit of  $y$  on  $X$  along with the standard errors and Z-scores. We wish to establish which variable, when dropped, will increase the RSS the least. How would you do this?

From question 3.1 we know that the F-statistic for dropping a single coefficient of a model is equivalent to the square of the corresponding Z score. Further, we know the F statistic in the case of dropping 1 variable is:

$$\frac{(RSS_0 - RSS_1)}{RSS_1 / (N - p_1 - 1)}$$

where  $N$ ,  $p_1$  and  $RSS_1$  are constant. In particular, the change in RSS is proportional to the F-score with a constant that does not depend

on choice of variable to be dropped, and so the change in RSS is proportional to the square of the z-score in the same manner. Thus the difference will be smallest (smallest increase in RSS) if the variable has minimal z-score in our model.

3.12 Show the ridge regression estimates can be obtained by OLS regression on an augmented data set. Add  $p$  rows to the centered matrix  $X$ ,  $\sqrt{\lambda}I_p$ , and augment  $y$  with  $p$  zeros.

Let  $X_2 = [X^T, \sqrt{\lambda}I_p]^T$  be our augmented matrix, and let  $y_2 = [y^T, 0]^T$  be the augmented response. Under OLS, we have  $\beta_2 = (X_2^T X_2)^{-1} X_2^T y_2$ .

$$\begin{aligned} X_2^T X_2 &= [X^T, \sqrt{\lambda}I_p][X^T, \sqrt{\lambda}I_p]^T \\ &= X^T X + \lambda I_p \\ X_2^T y_2 &= [X^T, \sqrt{\lambda}I_p][y^T, 0]^T \\ &= X^T y \end{aligned}$$

Thus  $\beta = (X^T X + \lambda I_p)^{-1} X^T y$  which is exactly the ridge regression beta for our non-augmented dataset.

3.13 Derive expression 3.62 and show that  $\hat{\beta}^{pcr}(p) = \hat{\beta}^{ls}$

Given an SVD of  $X$ ,  $X = UDV^T$  say, with  $V = [v_1, \dots, v_p]$ , the principal components  $z_m$  in equation 3.61 are defined as  $Xv_m$  for all  $m$ . The principal component regression is, for any  $M \leq p$ :

$$\begin{aligned} \hat{y}^{pcr} &= \bar{y} + \sum_{m=1}^M \hat{\theta}_m z_m \\ &= \bar{y} + X \sum_{m=1}^M \hat{\theta}_m v_m \end{aligned}$$

So we can just set  $\hat{\beta}^{pcr}(M) = \sum_{m=1}^M \hat{\theta}_m v_m$  and we're done. Now it remains to show that  $\hat{\beta}^{pcr}(p) = \hat{\beta}^{ls}$ .

$$\begin{aligned}
\hat{\beta}^{pcr}(p) &= \sum_{m=1}^p \hat{\theta}_m v_m \\
&= V \left[ \frac{z_1^T y}{z_1^T z_1}, \dots, \frac{z_p^T y}{z_p^T z_p} \right]^T \\
&= V \left[ \frac{z_1^T y}{d_1^2}, \dots, \frac{z_p^T y}{d_p^2} \right]^T \\
&= V D^{-2} [z_1^T y, \dots, z_p^T y]^T \\
&= V D^{-2} [u_1^T d_1^T y, \dots, u_p^T d_p^T y]^T \\
&= V D^{-2} D [u_1^T y, \dots, u_p^T y]^T \\
&= V D^{-1} U^T y
\end{aligned}$$

Where we made use of the SVD. Now:

$$\begin{aligned}
\hat{\beta}^{ls} &= (X^T X)^{-1} X^T y \\
&= (V D U^T U D V^T)^{-1} U D V^T y \\
&= (V D D V^T)^{-1} V D U^T y \\
&= D^{-2} (V V^T)^{-1} V D U^T y \\
&= D^{-2} V D U^T y \\
&= V D^{-1} U^T y
\end{aligned}$$

Using the orthonormality of  $U$  and  $V$ .

3.14 Show that in the orthogonal case, partial least squares stops after  $m = 1$  steps.

Assume  $X$  is such that each column has mean zero, unit variance, and the columns are orthogonal. let  $z = \sum_i (x_i^T y) x_i$  and  $\hat{\theta} = \frac{z^T y}{z^T z}$ .

In this case

$$\begin{aligned}
z^T z &= \sum_i \sum_j (x_i^T y) (x_j^T y) x_i^T x_j \\
&= \sum_i (x_i^T y)^2 (x_i^T x_i) \\
&= \sum_i (x_i^T y)^2
\end{aligned}$$

Similarly  $z^T y = \sum_i (x_i^T y)^2$ . Now let  $x_j^{(1)} = x_j - \frac{z^T x_j}{z^T z} z$ . Then let:

$$x_j^{(1)} = x_j - \frac{z^T x_j}{z^T z} z$$

In the next iteration, we have:

$$\begin{aligned}\langle x_j^{(1)}, y \rangle &= x_j^T y - \frac{(x_j^T y) (x_j^T x_j)}{\sum_i (x_i^T y)^2} z^T y \\ &= x_j^T y - \frac{x_j^T y}{\sum_i (x_i^T y)^2} \sum_i (x_i^T y)^2 \\ &= x_j^T y - x_j^T y \\ &= 0\end{aligned}$$

So algorithm 3.3 (page 81 in my edition) terminates after 1 step.

3.19 Show that  $\|\hat{\beta}^{ridge}\|$  increases as the tuning parameter  $\lambda \rightarrow 0$ . Does the same property hold for the Lasso and PLS?

Throughout this question I use regression without intercept as the intercept is not included in the authors formulation of the penalty in the Lagrangian form.

a) Ridge

For data  $X$  with mean 0 and unit variance, we have

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

Then:

$$\begin{aligned}\|\hat{\beta}^{ridge}\|^2 &= \left( (X^T X + \lambda I)^{-1} X^T y \right)^T (X^T X + \lambda I)^{-1} X^T y \\ &= \left( V(D^T D + \lambda I)^{-1} D U^T y \right)^T V(D^T D + \lambda I)^{-1} D U^T y \\ &= y^T U D^2 (D^T D + \lambda I)^{-2} U^T y\end{aligned}$$

Using the SVD, the commutativity of  $D$  and  $V$ , and the orthonormality of  $V$ .  $D^T D + \lambda I$  is diagonal, and so  $D^2 (D^T D + \lambda I)^{-2}$  is too with entries  $\frac{d_j^2}{d_j^2 + \lambda}$  on the diagonal. Let  $z_j = (U^T y)_j$  Then we see:

$$\begin{aligned}\|\hat{\beta}^{ridge}\|^2 &= \sum_i \sum_j \frac{d_j^2 z_i^T z_j}{d_j^2 + \lambda} \\ &= \sum_j \frac{d_j^2 z_j^2}{d_j^2 + \lambda}\end{aligned}$$

All terms are non-negative, and increase with decreasing  $\lambda$  so  $\hat{\beta}^{ridge}$  must too. Recall that we can view the ridge regression estimate as the solution to

$$\hat{\beta}^{ridge} = \arg \min_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_2^2$$

b) Lasso

Under a similar formulation, we have

$$\hat{\beta}^{lasso} = \operatorname{argmin}_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_1$$

In a similar way one can see that if  $\lambda$  is sufficiently large, it dominates this expression (for fixed  $X$  and  $y$ ), and so  $\hat{\beta}^{lasso}$  will decrease in norm with increasing  $\lambda$ .

c) PLS

?

3.23. Please refer to the book. LAR

a) Correlation with residuals remains constant in absolute value.

$$\begin{aligned} \frac{1}{N} X^T (y - u(\alpha)) &= \frac{1}{N} (X^T y - \alpha X^T X \hat{\beta}) \\ &= \frac{1}{N} \left( X^T y - \alpha X^T X (X^T X)^{-1} X^T y \right) \\ &= \frac{1}{N} (X^T y - \alpha X^T y) \\ &= \frac{1 - \alpha}{N} X^T y \end{aligned}$$

$$\text{Thus } \left| \frac{1}{N} X^T (y - u(\alpha)) \right| = (1 - \alpha) [\lambda, \dots, \lambda]^T$$

Which is exactly the required result in vector notation.

b) Explicit form of correlation

The question as stated ignores the need for an absolute value, so we shall assume that  $\langle x_j, y - u(\alpha) \rangle \geq 0$  for every  $j$ , else replace  $x_j$  by  $-x_j$  in our data.

Correlations are given by covariance divided by the product of the standard deviations. as everything (data and response) is assumed to be standardised, we have:

$$\begin{aligned} (y - u(\alpha))^T (y - u(\alpha)) &= y^T y - 2y^T u(\alpha) + u(\alpha)^T u(\alpha) \\ &= y^T y - 2\alpha y^T X \hat{\beta} + \alpha^2 \hat{\beta}^T X^T X \hat{\beta} \\ &= y^T y - 2\alpha y^T X \hat{\beta} + \alpha^2 \hat{\beta}^T X^T X (X^T X)^{-1} X^T y \\ &= y^T y - 2\alpha y^T X \hat{\beta} + \alpha^2 \hat{\beta}^T X^T y \\ &= y^T y + \alpha(\alpha - 2) y^T X \hat{\beta} \end{aligned}$$

Setting  $\alpha = 1$  we get  $RSS = y^T y - y^T X \hat{\beta}$ . Thus:

$$(y - u(\alpha))^T (y - u(\alpha)) = y^T y + \alpha(\alpha - 2)(y^T y - RSS)$$

And so:

$$\begin{aligned} (y - u(\alpha))^T (y - u(\alpha)) &= (1 - \alpha)^2 y^T y + \alpha(2 - \alpha)RSS \\ &= N(1 - \alpha)^2 + \alpha(2 - \alpha)RSS \end{aligned}$$

$$\begin{aligned} \text{Corr}(x_j, y - u(\alpha)) &= \frac{\langle x_j, y - u(\alpha) \rangle / N}{\sqrt{\langle x_j, x_j \rangle / N} \sqrt{\langle y - u(\alpha), y - u(\alpha) \rangle / N}} \\ &= \frac{\lambda(1 - \alpha)}{1 \cdot \sqrt{(1 - \alpha)^2 + \alpha(2 - \alpha)RSS / N}} \\ &= \frac{(1 - \alpha)}{\sqrt{(1 - \alpha)^2 + \alpha(2 - \alpha)RSS / N}} \lambda \end{aligned}$$

As required.

c) Show the LAR algorithm keeps correlations tied and monotonically decreasing.

Part a) showed that the LAR algorithm keeps correlations tied and this is exactly the step taken for an active set of variables in the  $k$ th step. Part b) shows that correlations are monotonically decreasing with  $\alpha$  as numerator falls faster than the denominator for  $\alpha$  in  $[0, 1]^3$ . In particular  $\lambda(0) = \lambda$  and  $\lambda(1) = 0$ .

3.29 Please see the book. Ridge regression with duplicate variables

Setup: The data  $X$  has identical columns  $x$ .

Let  $X = [x, \dots, x]$  say, for  $x$  some column vector, where  $X$  has dimension  $n \times m$ . The ridge regression beta is given by

$$(X^T X + \lambda I)^{-1} X^T y$$

Note  $X^T y = [x^T y, \dots, x^T y]^T$   $X^T X$  is a constant matrix with every value  $x^T x$ .

Let  $M = X^T X / (x^T x)$  be the matrix of ones.

Looking for an inverse of the form  $sM + \lambda^{-1}I$  yields

$$s = \frac{-x^T x}{\lambda(\lambda + mx^T x)}$$

<sup>3</sup> One can and should compute the gradient and check that it is negative on paper.



Where one uses that  $M^T M = mM$ . We have

$$(sM + \lambda^{-1}I)[x^T y, \dots, x^T y]^T = [c, \dots, c]^T$$

Where

$$\begin{aligned} c &= \frac{-m \cdot x^T x \cdot x^T y}{\lambda(\lambda + m \cdot x^T x)} + \frac{x^T y}{\lambda} \\ &= \frac{(\lambda + m \cdot x^T x)x^T y - m \cdot x^T x \cdot x^T y}{\lambda(\lambda + m \cdot x^T x)} \\ &= \frac{\lambda x^T y}{\lambda(\lambda + m \cdot x^T x)} \\ &= \frac{x^T y}{\lambda + m \cdot x^T x} \end{aligned}$$

3.30 Consider the elastic net optimisation problem:

$$\min_{\beta} \|y - X\beta\|_2^2 + \lambda \left[ \alpha \|\beta\|_2^2 + (1 - \alpha) \|\beta\|_1 \right]$$

Show how one can turn this into a lasso problem using an augmented version of  $X$  and  $y$

Let  $X_2 = [X, \sqrt{\alpha\lambda}I_p]^T$  where  $X$  is  $n \times p$  and we assume  $X$  has been standardised while  $y$  has mean 0 (this is an easy augmentation to remove the constant in regression). Let  $y_2 = [y, 0]$  where we have augmented  $y$  with  $p$  additional zeros. This is the setup in 3.12 except here we use  $\alpha\lambda$  in place of  $\lambda$

$$\begin{aligned} \|y_2 - X_2\beta\|_2^2 &= \left\| \begin{bmatrix} y - X\beta \\ \sqrt{\alpha\lambda}\beta \end{bmatrix} \right\|_2^2 \\ &= \|y - X\beta\|_2^2 + \|\sqrt{\alpha\lambda}\beta\|_2^2 \\ &= \|y - X\beta\|_2^2 + \alpha\lambda \|\beta\|_2^2 \end{aligned}$$

With these data, the original problem becomes a lasso optimisation problem with parameter  $\lambda(1 - \alpha)$  :

$$\min_{\beta} \|y_2 - X_2\beta\|_2^2 + \lambda(1 - \alpha) \|\beta\|_1$$