

ISLR Questions:

3) ISLR 6.8.3

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij}) \text{ subject to } \sum_{j=1}^p |\beta_j| \leq s$$

a) iv) is TRUE as 's' is increased, there is a less constraint on the model and it should always have a better training error (if 's' is increased, then the best model using a budget of 's' would be included when using a budget of 's').

If we consider $p=2$ for predictors, so the subject to constraint is $|\beta_1| + |\beta_2| \leq s$. When $s=0$ then $|\beta_1| + |\beta_2| \leq 0$ i.e. β_1, β_2 are 0 & the model's training RSS will be at its maximum. As 's' increases from 0, the co-efficients (β_1, β_2) will start to approach their least squares estimates values causing the training RSS to steadily decrease to the ordinary least square RSS. Hence iv) here is TRUE.

b) ii) is TRUE - test error will improve (decrease) to a point and then will worsen (increase) as constraints loosen and model overfits. If we see the same constraint then the coefficients are 0 when $s=0$, the model has a high test RSS. As we increase s from 0, the coefficients will assume non-zero values causing the

model to start to fit well to test data decreasing test RSS. However, there will reach a point when the coefficients approach ordinary least squares values and the model will start overfitting to the training data, which will increase test RSS. Therefore, as we increase s from 0, the test RSS will decrease initially and then eventually start increasing in a U-shape.

c) iii> is TRUE as variance always increases with fewer constraints. When $s=0$, the model basically predicts a constant leading to virtually no variance. As s increases from 0, the model starts including more coefficients and their values become highly dependent on training data. Therefore, as s increases from 0, the variance steadily increases.

d) iv> is TRUE - bias always decreases with more model flexibility. When $s=0$, the model is basically predicting a constant which is far from the actual value, thus the squared bias is high. As s increases, the coefficients become non-zero and fit the training data better which causes bias to decrease. Therefore, as s increases from 0 the squared bias steadily decreases.

c) \checkmark is TRUE as the irreducible error is a constant value, not related to model selection. Since the irreducible error is model independent, the choice of 's' does not affect the irreducible error associated with a model. Therefore, as 's' increases from 0, the irreducible error remains constant.

4) ISLR , 6.8.5

Ridge Regression : Minimize $\sum_{i=1}^n \left(Y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j X_{ij} \right)^2 + \lambda \sum_{i=1}^p (\hat{\beta}_i)^2$

a) by substituting $\hat{\beta}_0 = 0$ and $n=p=2$ specified in the problem into the general form of Ridge Regression shown above, we get

$$(Y_1 - \hat{\beta}_1 X_{11} - \hat{\beta}_2 X_{12})^2 + (Y_2 - \hat{\beta}_1 X_{21} - \hat{\beta}_2 X_{22})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

b) Expanding the equation from part a.

$$\begin{aligned} & (Y_1 - \hat{\beta}_1 X_{11} - \hat{\beta}_2 X_{12})^2 + (Y_2 - \hat{\beta}_1 X_{21} - \hat{\beta}_2 X_{22})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2) \\ &= (Y_1^2 + \hat{\beta}_1^2 X_{11}^2 + \hat{\beta}_2^2 X_{12}^2 - 2 Y_1 \hat{\beta}_1 X_{11} + 2 \hat{\beta}_1 \hat{\beta}_2 X_{11} X_{12} - 2 Y_1 \hat{\beta}_2 X_{12} \\ &+ (Y_2^2 + \hat{\beta}_1^2 X_{21}^2 + \hat{\beta}_2^2 X_{22}^2 - 2 Y_2 \hat{\beta}_1 X_{21} + -2 \hat{\beta}_1 \hat{\beta}_2 X_{21} X_{22} + 2 \hat{\beta}_1 \hat{\beta}_2 X_{21} X_{22} \\ &+ \lambda \hat{\beta}_1^2 + \lambda \hat{\beta}_2^2) \end{aligned}$$

we are given that : $X_{11} = X_{12} , X_{21} = X_{22} , Y_1 + Y_2 = 0$,
 $X_{11} + X_{21} = 0$ and $X_{12} + X_{22} = 0$

let $X_{11} = X_{12} = X_1$ & $X_{21} = X_{22} = X_2$

substituting values in equation :

$$= Y_1^2 - 2Y_1 X_1 (\hat{\beta}_1 + \hat{\beta}_2) + X_1^2 (\hat{\beta}_1 + \hat{\beta}_2)^2 + Y_2^2 - 2Y_2 X_2 (\hat{\beta}_1 + \hat{\beta}_2) + X_2^2 (\hat{\beta}_1 + \hat{\beta}_2)^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2) \rightarrow \text{equation } ①$$

taking the partial derivative to $\hat{\beta}_1$ and setting equation to 0 to minimize :

$$-2Y_1 X_1 + X_1^2 (2\hat{\beta}_1 + 2\hat{\beta}_2) - 2Y_2 X_2 + X_2^2 (2\hat{\beta}_1 + 2\hat{\beta}_2) + \lambda 2\hat{\beta}_1 = 0$$

$$-2Y_1 X_1 + 2X_1^2 (\hat{\beta}_1 + \hat{\beta}_2) - 2Y_2 X_2 + 2X_2^2 (\hat{\beta}_1 + \hat{\beta}_2) + 2\lambda \hat{\beta}_1 = 0$$

dividing equation by 2

$$-Y_1 X_1 + X_1^2 (\hat{\beta}_1 + \hat{\beta}_2) - \frac{1}{2} Y_2 X_2 + X_2^2 (\hat{\beta}_1 + \hat{\beta}_2) + \lambda \hat{\beta}_1 = 0$$

$$X_1 Y_1 - X_1^2 \hat{\beta}_2 + X_2 Y_2 + X_2^2 \hat{\beta}_2 = X_1^2 \hat{\beta}_1 + X_2^2 \hat{\beta}_1 + \lambda \hat{\beta}_1$$

$$\hat{\beta}_1 (X_1^2 + X_2^2 + \lambda) = X_1 Y_1 + X_2 Y_2 - \hat{\beta}_2 (X_1^2 + X_2^2)$$

$$\hat{\beta}_1 = \frac{X_1 Y_1 + X_2 Y_2 - \hat{\beta}_2 (X_1^2 + X_2^2)}{X_1^2 + X_2^2 + \lambda} \rightarrow ②$$

taking derivative of equation ① with respect to $\hat{\beta}_2$ and setting it to 0.

$$-2Y_1 X_1 + 2X_1^2 (\hat{\beta}_1 + \hat{\beta}_2) - 2Y_2 X_2 + 2X_2^2 (\hat{\beta}_1 + \hat{\beta}_2) + 2\lambda \hat{\beta}_2 = 0$$

$$x_1 y_1 - x_1^2 \hat{\beta}_1 + x_2 y_2 - x_2^2 \hat{\beta}_1 = \lambda \hat{\beta}_2 + x_1^2 \hat{\beta}_2 + x_2^2 \hat{\beta}_2$$

$$\hat{\beta}_2 (\lambda + x_1^2 + x_2^2) = x_1 y_1 + x_2 y_2 - \hat{\beta}_1 (x_1^2 + x_2^2)$$

$$\hat{\beta}_2 = \frac{x_1 y_1 + x_2 y_2 - \hat{\beta}_1 (x_1^2 + x_2^2)}{\lambda + x_1^2 + x_2^2} \rightarrow ③$$

Since equation ② = equation ③ we can say that
 $\hat{\beta}_1 = \hat{\beta}_2$

c) Lasso : Minimize $\frac{1}{2} (y_1 - \hat{\beta}_1 x_{11})^2 + (y_2 - \hat{\beta}_1 x_{21})^2 + \lambda |\hat{\beta}_1|$
 Similar to Ridge Regression, we can write the optimization problem for this case as :

$$\text{Minimize} : (Y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (Y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{12})^2 + \lambda (|\hat{\beta}_1| + |\hat{\beta}_2|)$$

d) Consider the lasso constraints as :

$$i) |\hat{\beta}_1| + |\hat{\beta}_2| \leq s \rightarrow ①$$

$$ii) (Y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (Y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{12})^2 \rightarrow ②$$

equation ② can be simplified by using the assumptions

$$x_{11} = x_{12}, x_{21} = x_{22}, x_{11} + x_{21} = 0, x_{12} + x_{22} = 0$$

$$Y_1 + Y_2 = 0$$

which minimizes to

$$(Y_1 - (\hat{\beta}_1 + \hat{\beta}_2) x_{11})^2 \rightarrow ③$$

from the 1st constraint

$|\hat{\beta}_1| + |\hat{\beta}_2| \leq s \rightarrow$ this corresponds to the lasso diamond

$\hat{\beta}_1 + \hat{\beta}_2 = s \rightarrow$ this corresponds to the edge on the lasso diamond

solution to the optimization problem:

$$\hat{\beta}_1 + \hat{\beta}_2 = \frac{y_1}{x_{11}}$$

as values for $\hat{\beta}_1$ & $\hat{\beta}_2$ vary along the line $\hat{\beta}_1 + \hat{\beta}_2 = \frac{y_1}{x_{11}}$,

the contours related to equation ③ will touch the lasso diamond edge at different points

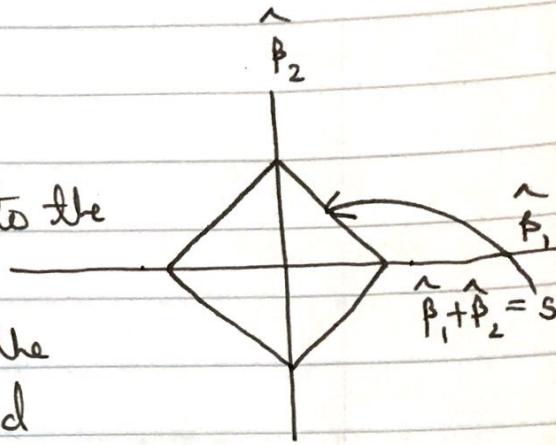
for this reason, the entire edge $\hat{\beta}_1 + \hat{\beta}_2 = s$ can be a solution to the lasso optimization problem as well as edge at $\hat{\beta}_1 + \hat{\beta}_2 = -s$.

Hence, the lasso optimization problem does not have a unique solution.

The general forms of solutions are the two line segments as follows:

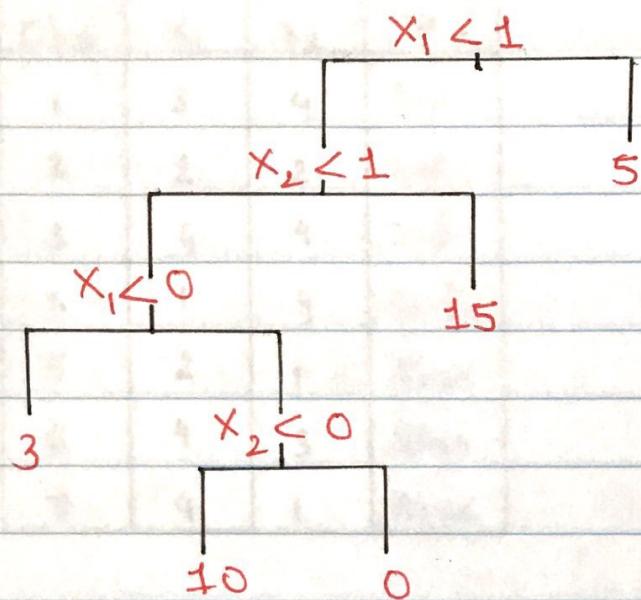
① $\hat{\beta}_1 + \hat{\beta}_2 = s$ such that $\hat{\beta}_1 \geq 0$ and $\hat{\beta}_2 \geq 0$

② $\hat{\beta}_1 + \hat{\beta}_2 = -s$ such that $\hat{\beta}_1 \leq 0$ & $\hat{\beta}_2 \leq 0$

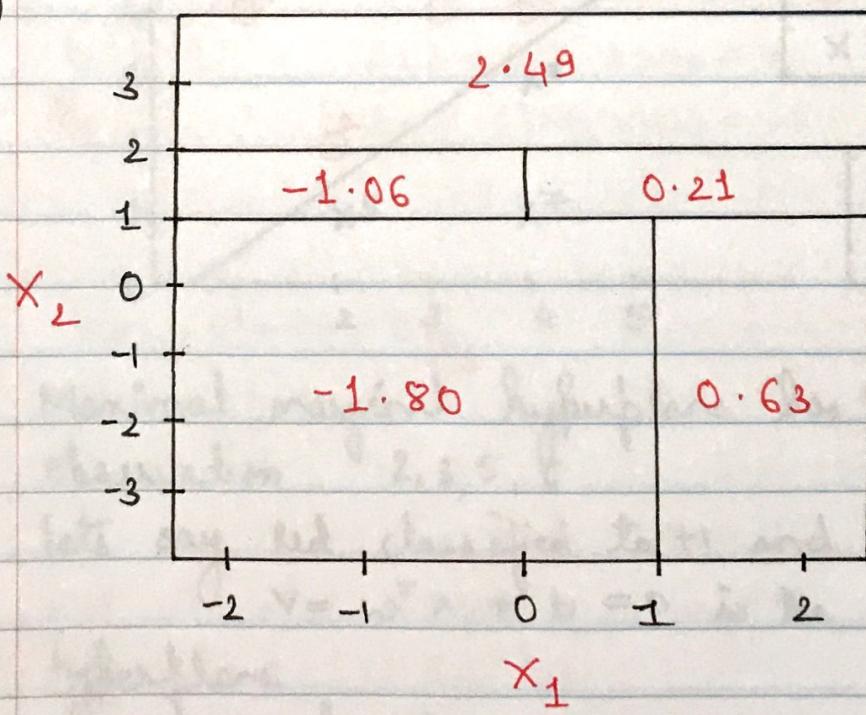


5) ISLR 8.4.5

a)



b)

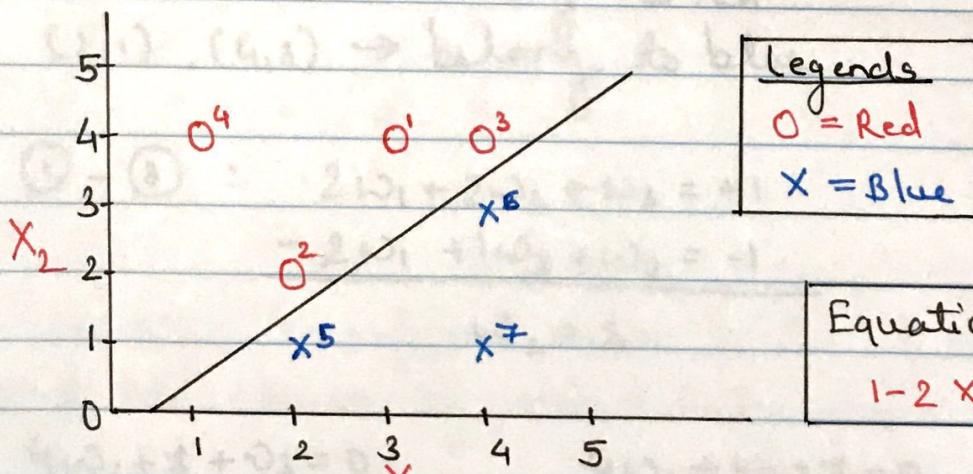


6) ISLR 9.7.3 (2.0, (4.3))

a)

Obs	x_1	x_2	Y
1	3	4	Red
2	2	2	Red
3	4	4	Red
4	1	4	Red
5	2	1	Blue
6	4	3	Blue
7	4	1	Blue

b)



Equation of hyperplane
 $1 - 2X_1 + 2X_2 = 0$

Maximal marginal hyperplane has to be in between observation 2, 3, 5, 6

lets say Red classified to +1 and blue classified to -1
 $y = w^T x_i + b = 0$ is the equation for a hyperplane

for this problem it takes like

$$0 = w_1 x_1 + w_2 x_2 + w_3$$

positive points $[(2, 2), (4, 4)]$

negative points $[(2,1), (4,3)]$

Hyperplane linear equations:

$$\textcircled{1} \text{ Obs 2 : } 2\omega_1 + 2\omega_2 + \omega_3 = 1$$

$$\textcircled{2} \text{ Obs 3 : } 4\omega_1 + 4\omega_2 + \omega_3 = 1$$

$$\textcircled{3} \text{ Obs 5 : } 2\omega_1 + 1\omega_2 + \omega_3 = -1$$

$$\textcircled{4} \text{ Obs 6 : } 4\omega_1 + 3\omega_2 + \omega_3 = -1$$

Co-ordinates: $(2,2), (4,4), (2,1), (4,3)$

$(2,2), (4,4) \rightarrow$ belong to red

$(2,1), (4,3) \rightarrow$ belong to blue

$$\begin{aligned}\textcircled{1} - \textcircled{3} : \quad & 2\omega_1 + 2\omega_2 + \omega_3 = +1 \\ & \underline{- 2\omega_1 + 1\omega_2 + \omega_3 = -1} \\ & \omega_2 = 2\end{aligned}$$

$$4\omega_1 + ? + \omega_3 = 0$$

$$2\omega_1 + 2(2) + \omega_3 = 1$$

$$4\omega_1 + \omega_3 + ? = 0$$

$$\underline{- 2\omega_1 + \omega_3 + 3 = 0}$$

$$2\omega_1 + 4 = 0$$

$$\therefore \omega_1 = -2$$

put $\omega_1 = -2$ & $\omega_2 = 2$ in eqⁿ ①

$$2(-2) + 2(2) + \omega_3 = 1$$

$$\therefore \omega_3 = 1$$

hence equation of hyperplane :

$$1 - 2x_1 + 2x_2 = 0$$

c) From the equation of the hyperplane

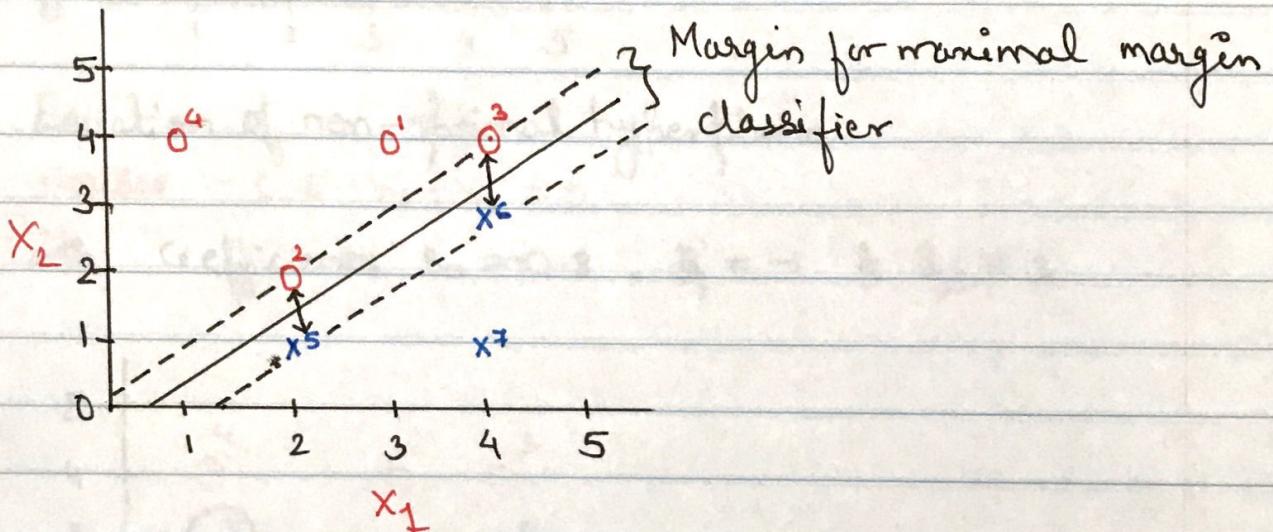
$$1 - 2x_1 + 2x_2 = 0 \quad \text{Red} = +1 \quad \text{Blue} = -1$$

Hence, classification rule is:

$$0.5 - x_1 + x_2 > 0$$

$$\text{with } \beta_0 = 0.5, \beta_1 = -1, \beta_2 = 1$$

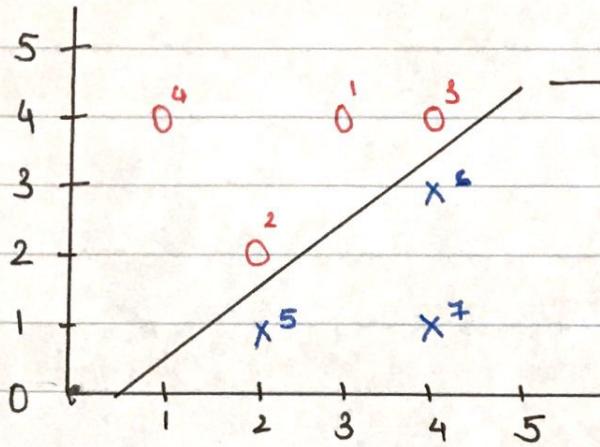
d)



e) There are 4 support vectors located at observations 2, 3, 5 & 6 and are shown by arrows in the graph shown in d)

f) A slight movement in the 7th observation would not affect the maximal marginal hyperplane because its position after a slight movement would be outside the margin

g)



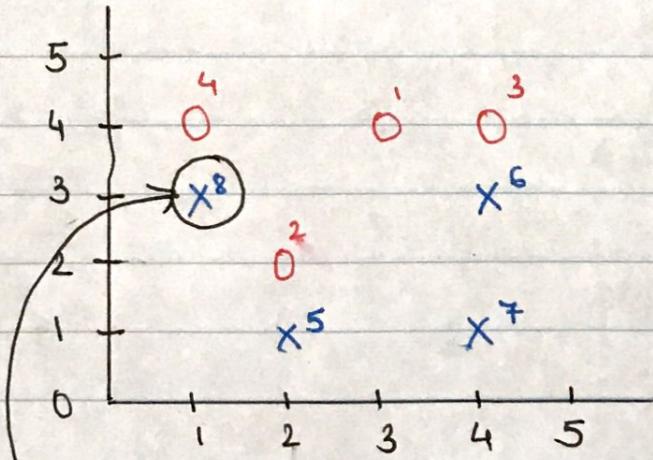
→ Non-optimal separating hyperplane because it is closer to the ~~blue~~ data points than the red.

Equation of non-optimal hyperplane

$$\text{---} -0.8 - x_1 + x_2 = 0$$

with coefficients $B_0 = -0.8$, $B_1 = -1$ & $B_2 = 1$

h)



The addition of this observation makes it so the two classes are no longer separable by a hyperplane.