① 6.6.3     Lasso   as   s increases

estimate the regression coefficients by minimizing

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 \text{ subject to}$$

$$\sum_{j=1}^{p} |\beta_j| \leq s$$

for a particular value of s.

✳ As we increase s from 0

Train

(a) RSS :- (iv) steadily decrease

When s is sufficienty large, the β that minimizes RSS will be the least squares solution. Thus, training RSS will decrease monotonically.

(b) Test RSS :- (ii) Decrease initially, and then eventually start to increase in a U shape

When s=0     $y = \bar{y}$
As s increases, more β are included, flexibility of model ↑, ∴ test RSS at some point will start increasing leading to (overfitting) happens.

(c) Variance :- (iii) steadily increase

As s goes from 0 to some high value, the flexibility of model increases, & variance increases.

(d) Squared Bias :- (iv) Steadily decreases

→ Because, model flexibility increases with s.

(e) Irreducible error :- (v) Remain constant

Irreducible error $\varepsilon$ remains regardless of model complexity. Its not dependent on X.

ISLR

6.6.5

Ridge regression tends to give similar coefficient values to correlated variables, whereas lasso may give different coefficient values to correlated variables

Given

$$n=2, p=2, x_{11}=x_{12}, x_{21}=x_{22}.$$
$$y_1+y_2=0 \qquad x_{11}+x_{21}=0$$
$$x_{12}+x_{22}=0 \qquad \rightarrow \hat{\beta}_0=0$$

In Ridge - $\hat{\beta}_\lambda^R$ are values the minimizes

(a)

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

choose $\hat{\beta}_\lambda^R = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$ such that it minimizes :-

$$\sum_{i=1}^{n} \left( y_i - \hat{\beta}_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{2} \beta_j^2$$

$\Rightarrow$
$$\left( y_1 - \beta_0 - \left( \beta_1 x_{11} + \beta_2 x_{12} \right) \right)^2 +$$
$$\left( y_2 - \beta_0 - \left( \beta_1 x_{21} + \beta_2 x_{22} \right) \right)^2 +$$
$$\lambda \left( \beta_1^2 + \beta_2^2 \right)$$

$\hat{\beta}_0 = 0$

$\Rightarrow \left( y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12} \right)^2 + \left( y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22} \right)^2 +$
$$\lambda \left( \beta_1^2 + \beta_2^2 \right)$$

(b) $\left( y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12} \right)^2 + \left( y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22} \right)^2 +$
$$\lambda \left( \hat{\beta}_1^2 + \beta_2^2 \right)$$

$\Rightarrow \left( y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12} \right)^2 + \left( y_1 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22} \right)^2 +$
$$\lambda \left( \beta_1^2 + \beta_2^2 \right)$$

$\Rightarrow \left( y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12} \right)^2 * \left( \pm y_1 * \hat{\beta}_1 x_{11} * \hat{\beta}_2 x_{12} \right)^2 +$
$$\lambda \left( \beta_1^2 + \beta_2^2 \right)$$

$\Rightarrow$

$$\Rightarrow \quad (y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2 + \lambda(\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

$$\Rightarrow \quad (y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11})^2 + (-y_1 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{21})^2 + \lambda(\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

$$\Rightarrow \quad (y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11})^2 + (-y_1 + \hat{\beta}_1 x_{11} + \hat{\beta}_2 x_{11})^2 + \lambda(\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

$$\Rightarrow \quad (y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11})^2 + (-1)^2(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11})^2 + \lambda(\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

$f(\hat{\beta}_1, \hat{\beta}_2)$

$$2(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11})^2 + \lambda(\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

$$\frac{\partial(f)}{\partial\hat{\beta}_1} = 4(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11})(-x_{11})$$
$$+ 2\lambda\hat{\beta}_1 = 0$$

$$2\lambda\hat{\beta}_1 = 4y_1 x_{11} - 4\hat{\beta}_1 x_{11}^2 - 4\hat{\beta}_2 x_{11}^2$$

$$\lambda\hat{\beta}_1 = 2y_1 x_{11} - 2\hat{\beta}_1 x_{11}^2 - 2\hat{\beta}_2 x_{11}^2$$

$$\hat{\beta}_1(\lambda + 2x_{11}^2) = 2y_1 x_{11} - 2\hat{\beta}_2 x_{11}^2$$

$$\hat{\beta}_1 = \frac{2y_1 x_{11} - 2\hat{\beta}_2 x_{11}^2}{(\lambda + 2x_{11}^2)}$$

$$\frac{\partial(f)}{\partial(\hat{\beta}_2)} = 4\left(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11}\right)(-x_{11}) = 0$$
$$+ 2\lambda\hat{\beta}_2$$

$$\Rightarrow 2y_1 x_{11} - 2\hat{\beta}_1 x_{11}^2 - 2\hat{\beta}_2 x_{11}^2 = \lambda\hat{\beta}_2$$

$$\hat{\beta}_2 = \frac{2y_1 x_{11} - 2\hat{\beta}_1 x_{11}^2}{(\lambda + 2x_{11})^2}$$

Substituting

$$c = \frac{2y_1 x_{11}}{(\lambda + 2x_{11}^2)} \qquad \& \qquad \frac{c\lambda}{k} = \frac{-2x_{11}^2}{(\lambda + 2x_{11})^2}$$

$$\hat{\beta}_1 = c + k\hat{\beta}_2 \qquad —① \quad \Big\}$$
$$\hat{\beta}_2 = c + k\hat{\beta}_1 \qquad —② \quad \Big\}$$

$$\hat{\beta}_1 = c + k\left(c + k\hat{\beta}_1\right)$$
$$\hat{\beta}_1 = c + kc + k^2\hat{\beta}_1$$
$$\hat{\beta}_1 = \frac{c(1+k)}{(1-k^2)} \qquad\qquad —③$$

① → ②

$$\hat{\beta}_2 = c + k\left(c + k\hat{\beta}_2\right)$$
$$\hat{\beta}_2 = c + kc + k^2\hat{\beta}_1$$
$$\hat{\beta}_2 = \frac{c(1+k)}{(1-k^2)} \qquad\qquad —④$$

$$\therefore \quad \hat{\beta}_1 = \hat{\beta}_2$$

(c)

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$

$$\Rightarrow \left( y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12} \right)^2 +$$

$$\left( y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22} \right)^2 +$$

$$\lambda \left( |\beta_1| + |\beta_2| \right)$$

$$\Rightarrow 2 \left( y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{11} \right)^2 + \lambda \left( |\beta_1| + |\beta_2| \right)$$

(d) => Argue -non- unique coefficient estimates

Lasso will select $\hat{\beta}^L$ that minimizes.

$$\sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij}\right)^2 \text{ give to constraint}$$

$$\sum_{j=1}^{p} |\beta_j| \leq s$$

∴ Lasso coefficients must minimize

$$2 \left(y_1 - (\hat{\beta}_1 + \hat{\beta}_2) x_{11}\right)^2 \geq 0 \text{ subject to}$$

$$|\hat{\beta}_1| + |\hat{\beta}_2| \leq s$$

We can see that any $(\hat{\beta}_1, \hat{\beta}_2)$ that satisfy $\hat{\beta}_1 + \hat{\beta}_2 = \dfrac{y_1}{x_{11}}$ will have RSS of 0.

Considering the lasso constraint, the solution to the optimization problem will be where the contour of $2\left(y_1 - (\hat{\beta}_1 + \hat{\beta}_2) x_{11}\right)^2$ touch the lasso diamond.
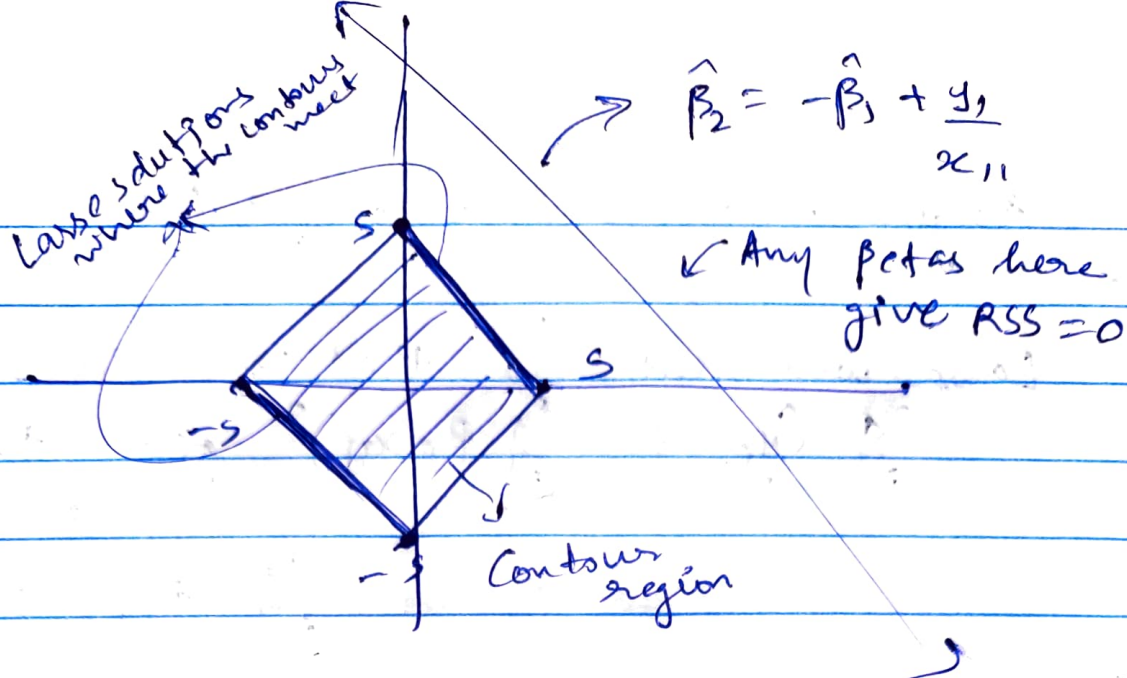
As $(\hat{\beta}_1, \hat{\beta}_2)$ vary along $\hat{\beta}_2 = -\hat{\beta}_1 + \dfrac{y_1}{x_{11}}$ the contour will touch lasso diamond at many points instead of one.

$$\left( \underbrace{|\hat{\beta}_1| + |\hat{\beta}_2| \leq s}_{} \right) \qquad \underbrace{\hat{\beta}_1 + \hat{\beta}_2 = s}_{}$$

$$\hat{\beta}_1 + \hat{\beta}_2 = -s$$

$$\hat{\beta_2} = -\hat{\beta_1} + \frac{y_1}{x_{11}}$$

Lasso solutions where the contour meet

Any betas here give RSS = 0

S

S

-S

-1

Contour region

This means there are $\infty$ solutions

$\hat{\beta}_\lambda^L$

$$\hat{\beta}_\lambda^L \in \left\{ (\hat{\beta_1}, \hat{\beta_2}) : \hat{\beta_1} + \hat{\beta_2} = s, \right.$$
$$\left. \hat{\beta_1} \in [0, s], \hat{\beta_2} \in [0, s] \right\}$$

$\cup$

$$\left\{ (\hat{\beta_1}, \hat{\beta_2}) : \hat{\beta_1} + \hat{\beta_2} = -s, \right.$$
$$\left. \hat{\beta_1} \in [-s, 0], \hat{\beta_2} \in [-s, 0] \right\}$$

① Majority vote :-

Accross 10 estimates

$G, G, G, G, R, R, R, R, R$

O/P prediction → Red

② Average probability

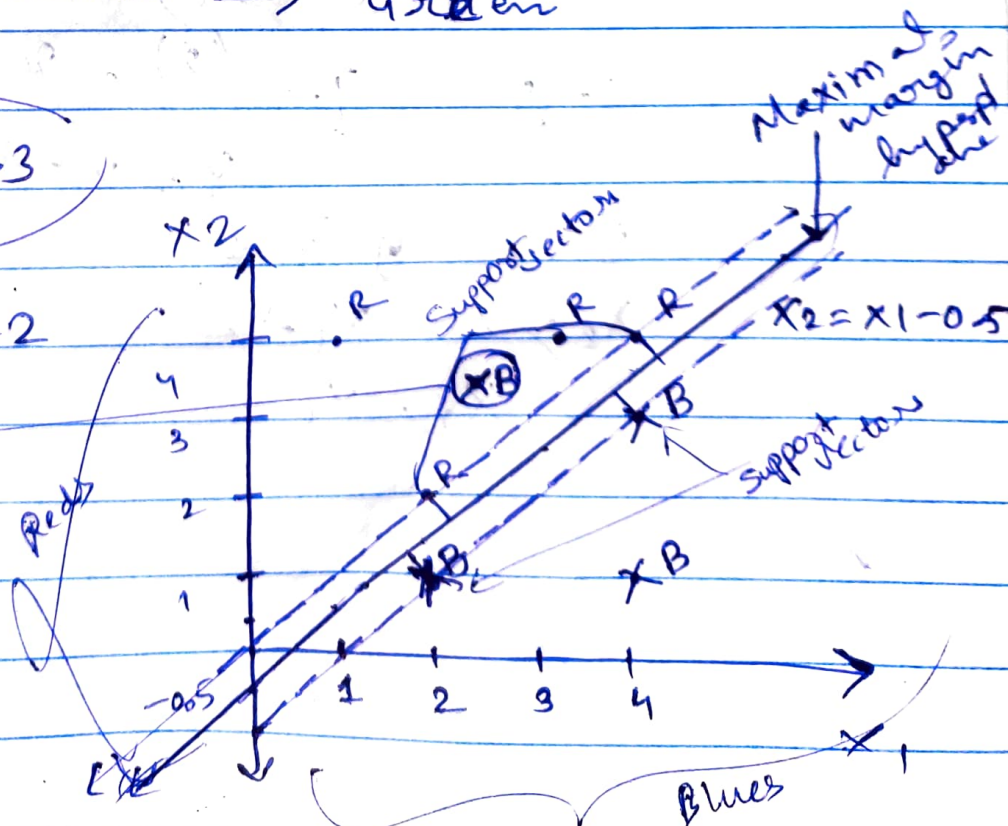Ave $P(C = "Red" | x) = 0.45$

Ave $P(C = "Green" | x) = 0.55$

O/P prediction → Green

ISLR  9.7.3

(a') $n = 7$, $p = 2$

-- Additional point

Support vectors
are on dashed
lines.



Maximal margin hyperplane

$x_2 = x_1 - 0.5$

Support vector

Support vector

Red

Blues

(b)     Visually  it  looks  like  the  optimal
        seperating   hyperplane  is

$$x_2 - x_1 + 0.5 = 0$$

Hyper planes ⇒

$(2,1)$ $d_1$ → $\dfrac{|2-1+0.5|}{\sqrt{2}}$ $=$ $\dfrac{1.5}{\sqrt{2}}$

$\begin{matrix} d_2 \\ (2,2) \end{matrix}$ → $\dfrac{|0.5|}{\sqrt{2}}$ $=$ $\dfrac{0.5}{\sqrt{2}}$

① 
$(2,2)$
$y - y_1 = 1(x - x_1)$
$y - 2 = 1(x - 2)$
$\boxed{y = x}$ —① $\qquad$ $\boxed{x_2 = x_1}$

② 
$(2,1)$
$y - y_1 = 1(x - x_1)$
$y - y_1 = x - 2$
$\boxed{x_2 = x_1 - 1}$

Margin $= \dfrac{\sqrt{2}}{4}$

(c) Classification Rule for Maximal margin classifier.

"Red" if $\beta_0 + \beta_1 X_1 + \beta_2 X_2 > 0$ else "Blue"

$\longrightarrow (0.5, -1, 1)$

∴ "Red" when $X_2 > X_1 - 0.5$

& "Blue" when $X_2 \leq X_1 - 0.5$

The maximal marginal classifier thus becomes

$$f(x) = X_2 - X_1 + 0.5$$
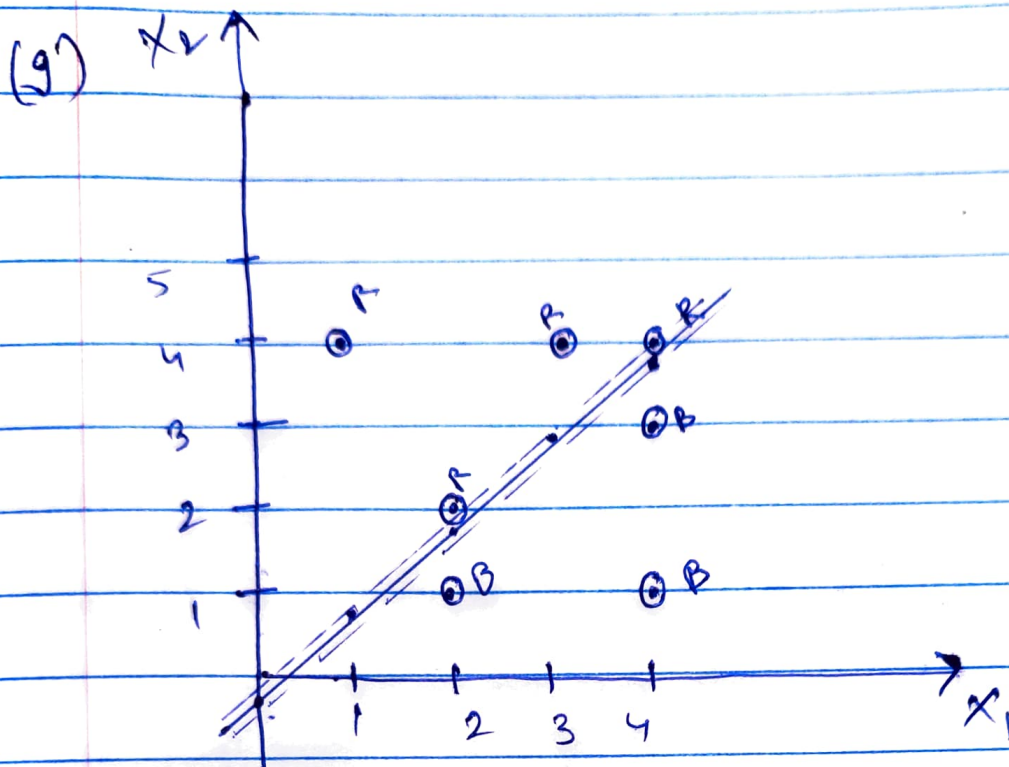
here if $f(x') > 0 \longrightarrow$ "Blue"

else "Red"

(d)
⟩ Check answer (a)
(e)

(f) Observation ⑦ is not a support vector. Any small movement of this point won't change the fact that $X_2 - X_1 + 0.5 = 0$.

If observation 7 moves inside of the margin then it will start influencing the position of maximal hyperplane.

(g)



$$x_2 = x_1 - 0.25$$

The Margin M in this case is extremely
small.

Not an optimal hyperplance

(h)   Check answer (a)