

PARTH RATHOD

CSP 571 - DPA

A20458817

FALL 2021

Recitation Answers

Chapter 6

Question 1

a) The Best Subset selection will have the smallest training RSS as it will consider all the possible models unlike the others which has a greedy approach.

b) The Best subset selection model will have the high chances of choosing a model with less test RSS as it contains 2^p models whereas as the other two models will consider only $(1 + p(p+1))/2$ models.

c) i) TRUE :- because the $(k+1)$ variable model has one more predictor in addition to all of the predictors chosen for the k -variable model in the forward stepwise selection.

ii) TRUE :- because in the backward stepwise selection, k -variable model is obtained by removing one predictor from $(k+1)$ -variable model which will reduce the RSS of the model.

iii) FALSE :- because both the models follow different criteria.

iv) FALSE :- because both the models follow different criteria. Also, there is not link between the models obtained from forward and backward model.

v) FALSE :- because the best subset approach selects the model with $(k+1)$ predictors from among all feasible models with $(k+1)$ predictors. As a result, it does not ensure that the same predictors will be used for the k predictor model.

Question 2

a) Option (iii) is correct.

Because Regularization reduces the Test MSE by adding a penalty by decreasing variance and increasing bias. This penalty shrinks the coefficient and slope gets less steep.

b) Option (iii) is correct.

Ridge reduces predictors that do not have a significant link with the target variable, making them less flexible. It also decreases variance at the expense of increasing bias. To enhance prediction accuracy, the increase in bias should be smaller than the decrease in variance.

c) Option (ii) is correct.

Because non-linear techniques are more flexible than least squares, they may provide more accurate predictions.

Question 3

a) Option (iv) is correct.

As we increase s , the model becomes more and more flexible as the restriction on β is reducing, thus the coefficients increase from 0 to their least square estimate values. Thus, resulting in decreased RSS.

b) Option (ii) is correct.

As model is becoming more and more flexible the test RSS will reduce first and then start increasing when overfitting will start.

c) Option (iii) is correct.

Variance steadily increase with increase in model flexibility.

d) Option (iv) is correct.

Bias decreases with increase in model flexibility.

e) Option (v) is correct.

Irreducible error is model independent and does not depend on s .

Question 4

a) Option (iii) is correct.

As we increase λ , the model becomes less and less flexible as the restriction on β is increasing, thus the coefficients come close to 0 from their least square estimate values. Thus, resulting in increased RSS.

b) Option (ii) is correct.

As the model is becoming less and less flexible the test RSS will reduce first and then start increasing when overfitting will start.

c) Option (iv) is correct.

Variance steadily decreases with the decrease in model flexibility.

d) Option (iii) is correct.

Bias increases with decrease in model flexibility.

e) Option (v) is correct.

Irreducible error is model independent and does not depend on λ .

Question 5

a)

Question 5

a) Ridge Regression is given by:-

$$\text{Minimize:- } \sum_{i=1}^n (y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j x_{ij})^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2$$

here

$$n=p=2 \quad \& \quad \hat{\beta}_0 = 0$$

$$\therefore \min \rightarrow \left[(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2 \right] + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

$$=$$

b)

b) \therefore On expanding above Equation.

$$\rightarrow (y_1^2 + \hat{\beta}_1^2 x_{11}^2 + \hat{\beta}_2^2 x_{12}^2 - 2\hat{\beta}_1 x_{11} y_1 - 2\hat{\beta}_2 x_{12} y_1 + 2\hat{\beta}_1 \hat{\beta}_2 x_{11} x_{12})$$

$$+ (y_2^2 + \hat{\beta}_1^2 x_{21}^2 + \hat{\beta}_2^2 x_{22}^2 - 2\hat{\beta}_1 x_{21} y_2 - 2\hat{\beta}_2 x_{22} y_2 + 2\hat{\beta}_1 \hat{\beta}_2 x_{21} x_{22})$$

$$+ \lambda \hat{\beta}_1^2 + \lambda \hat{\beta}_2^2$$

Take derivative & equate it to zero.

$$\therefore \frac{\partial}{\partial \hat{\beta}_1} = 0$$

Also,

$$x_{11} = x_{12} = x_1 \quad \& \quad x_{21} = x_{22} = x_2 \quad \& \quad \text{divide by 2.}$$

$$\rightarrow (\hat{\beta}_1 x_1^2 - x_1 y_1 + \hat{\beta}_2 x_1^2) + (\hat{\beta}_1 x_2^2 - x_2 y_2 + \hat{\beta}_2 x_2^2) + \lambda \beta_1 = 0.$$

$$\rightarrow \hat{\beta}_1 (x_1^2 + x_2^2) + \hat{\beta}_2 (x_1^2 + x_2^2) + \lambda \beta_1 = x_1 y_1 + x_2 y_2$$

Add $2\hat{\beta}_1 x_1 x_2$ and $2\hat{\beta}_2 x_1 x_2$ on both sides.

$$\rightarrow \hat{\beta}_1 (x_1 + x_2)^2 + \hat{\beta}_2 (x_1 + x_2)^2 + \lambda \beta_1 = x_1 y_1 + x_2 y_2 \quad \text{--- (1)}$$

$\text{As } x_1 + x_2 = 0$

$$\therefore \lambda \hat{\beta}_1 = x_1 y_1 + x_2 y_2 + 2\hat{\beta}_1 x_1 x_2 + 2\hat{\beta}_2 x_1 x_2 \quad \text{--- (1)}$$

Similarly taking partial derivative wrt $\hat{\beta}_2$ we get.

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

From (1) & (2)

$$\underline{\beta_1 = \beta_2}$$

c)

$$c) \min [(y_1 - \hat{\beta}_1 x_{11} - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_{21} - \hat{\beta}_2 x_{22})^2]$$

$$+ \lambda (|\hat{\beta}_1| + |\hat{\beta}_2|)$$

d)

d) Replacing penalty term from Ridge Regression
∴ the derivative term.

$$= \frac{\partial}{\partial \hat{\beta}} (1|\beta|) : \frac{1|\beta|}{\beta}$$

Same like ridge regression we get,

$$\frac{1|\beta_1|}{\beta_1} = \frac{1|\beta_2|}{\beta_2}$$

Provide β_1 & β_2 are both positive or both negative.

Chapter 7

Question 2

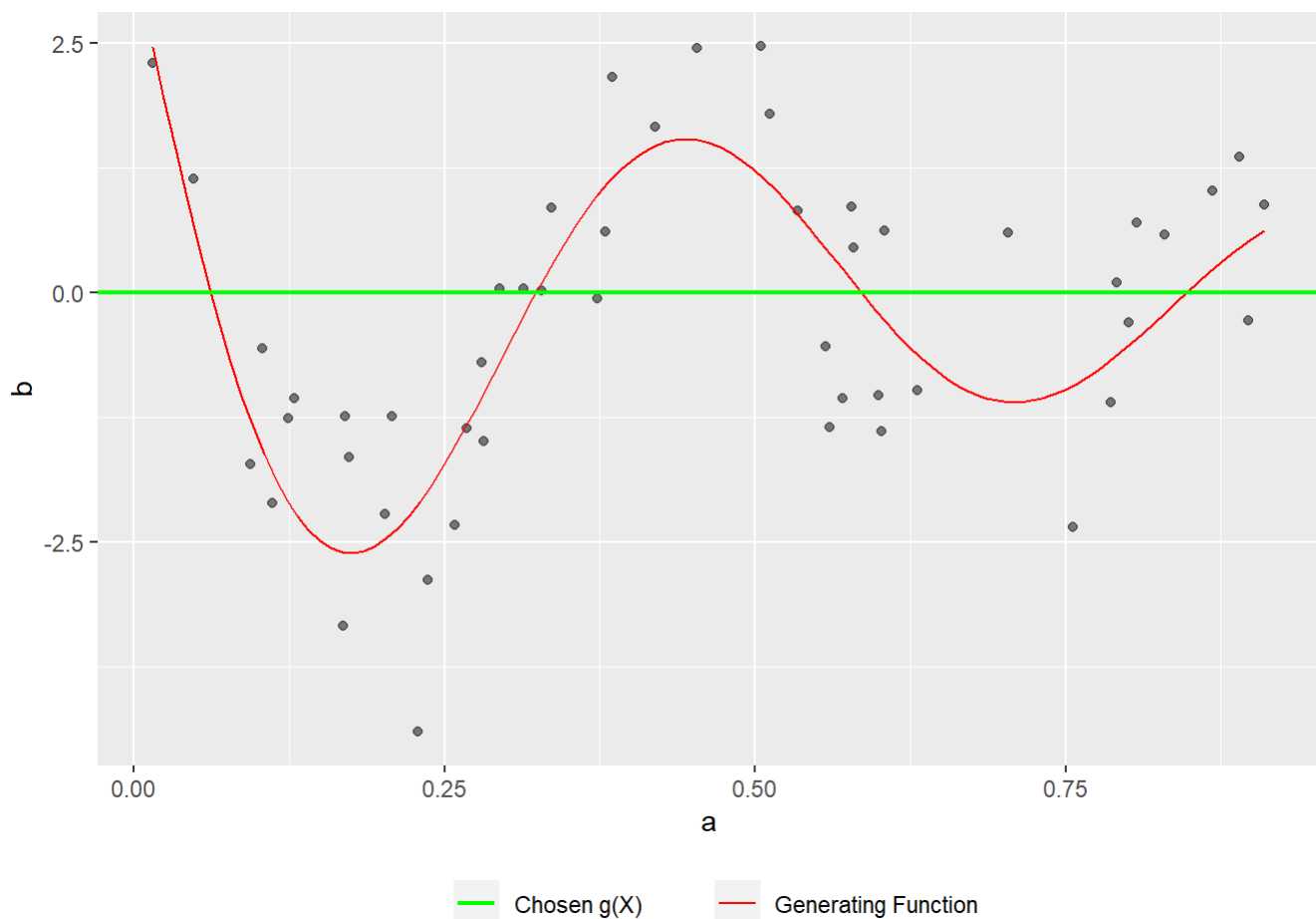
We need to generate some data before sketching \hat{g} under different conditions!

```
library(ggplot2)
set.seed(3)

a <- runif(50)
eps <- rnorm(50)
b <- sin(12*(a + 0.2)) / (a + 0.2) + eps
generating_fn <- function(a) {sin(12*(a + 0.2)) / (a + 0.2)}
df <- data.frame(a, b)
```

2a) $\lambda \rightarrow \infty, m=0$

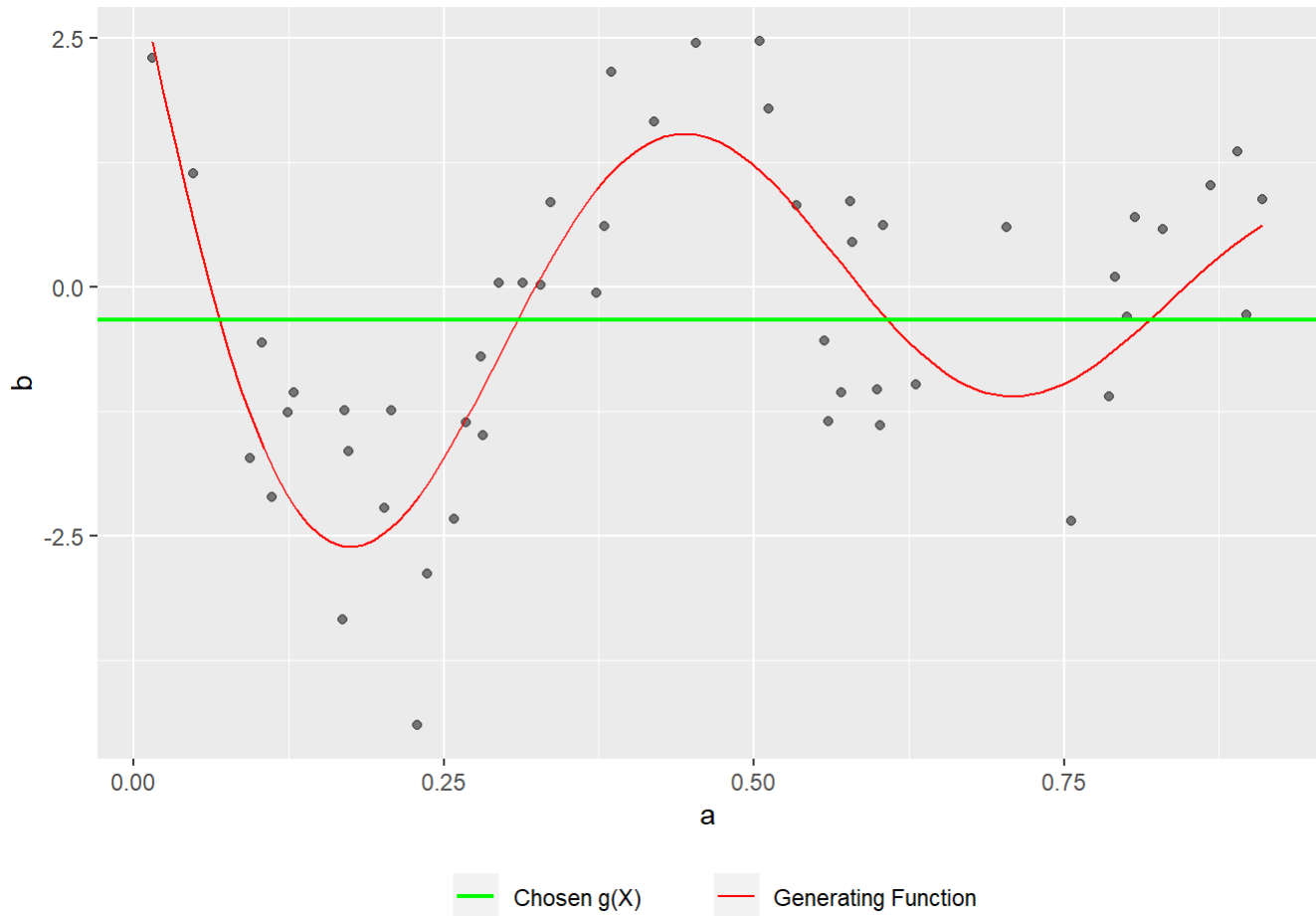
```
ggplot(df, aes(x = a, y = b)) +
  geom_point(alpha = 0.5) +
  stat_function(fun = generating_fn, aes(col = "Generating Function")) +
  geom_hline(aes(yintercept = 0, linetype = "Chosen g(X)", col = "green", size = 0.8) +
  scale_color_manual(values = "red") +
  theme(legend.position = "bottom", legend.title = element_blank())
```



As λ increases, the penalty term becomes more and more important in the equation. As $\lambda \rightarrow \infty$, this forces $g(x) \rightarrow 0$. We therefore get $\hat{g}(x) = 0$

2b) $\lambda=\infty, m=1$

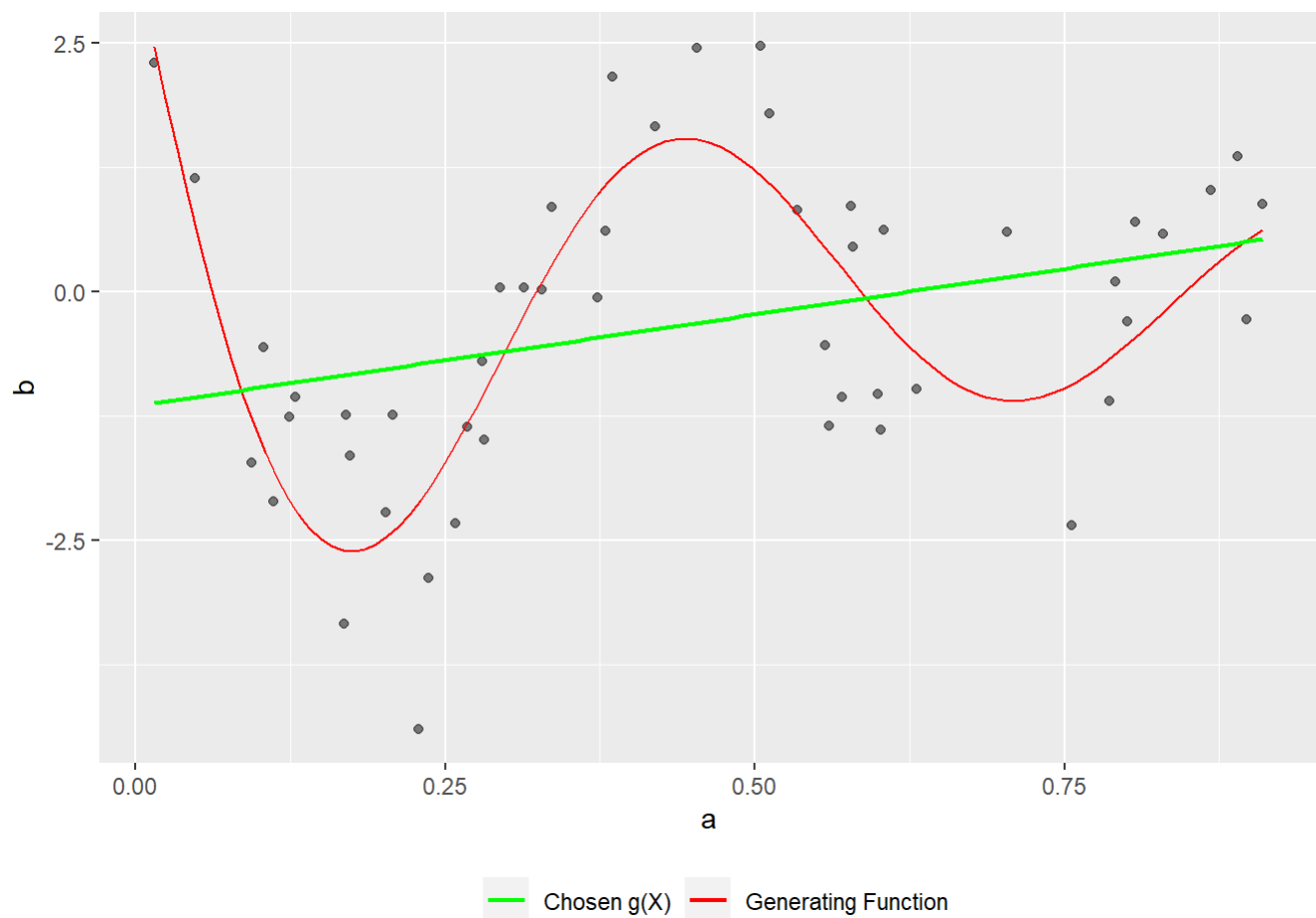
```
ggplot(df, aes(x = a, y = b)) +
  geom_point(alpha = 0.5) +
  stat_function(fun = generating_fn, aes(col = "Generating Function")) +
  geom_hline(aes(yintercept = mean(b)), linetype = "Chosen g(X)", col = "green", size = 0.8) +
  scale_color_manual(values = "red") +
  theme(legend.position = "bottom", legend.title = element_blank())
```



As $\lambda \rightarrow \infty$, this forces $g'(x) \rightarrow 0$. This means we would get $\hat{g}(x) = c$.

2c) $\lambda=\infty, m=2$

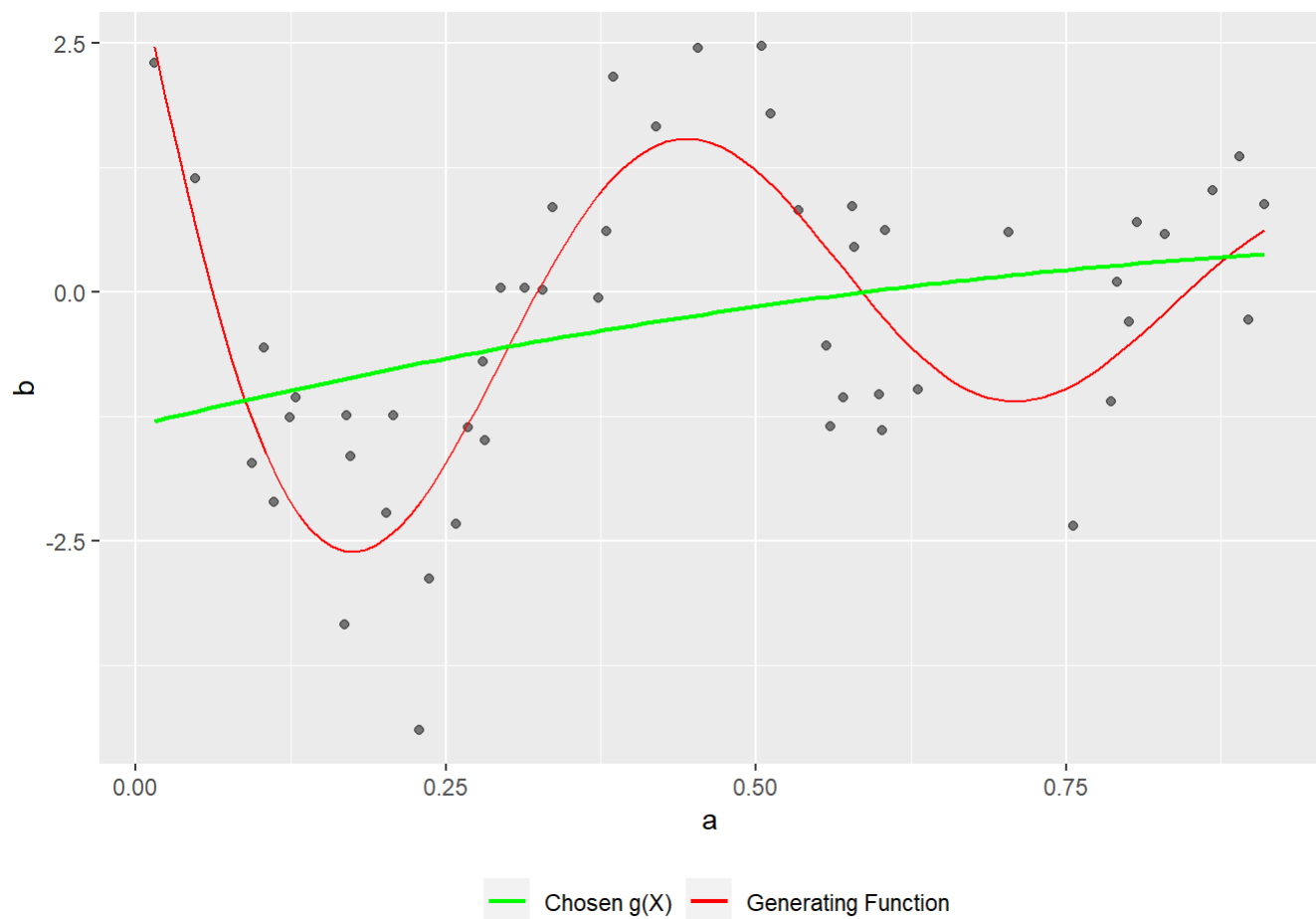
```
ggplot(df, aes(x = a, y = b)) +
  geom_point(alpha = 0.5) +
  stat_function(fun = generating_fn, aes(col = "Generating Function")) +
  geom_smooth(method = "lm", formula = "y ~ x", se = F, size = 0.8, aes(col = "Chosen g(X)")) +
  scale_color_manual(values = c("green", "red")) +
  theme(legend.position = "bottom", legend.title = element_blank())
```



As $\lambda \rightarrow \infty$, this forces $g''(x) \rightarrow 0$. This means we would get $\hat{g}(x) = ax + b$

2d) $\lambda = \infty, m = 3$

```
ggplot(df, aes(x = a, y = b)) +
  geom_point(alpha = 0.5) +
  stat_function(fun = generating_fn, aes(col = "Generating Function")) +
  geom_smooth(method = "lm", formula = "y ~ x + I(x^2)", se = F, size = 0.8, aes(col = "Chosen g
(X)")) +
  scale_color_manual(values = c("green", "red")) +
  theme(legend.position = "bottom", legend.title = element_blank())
```

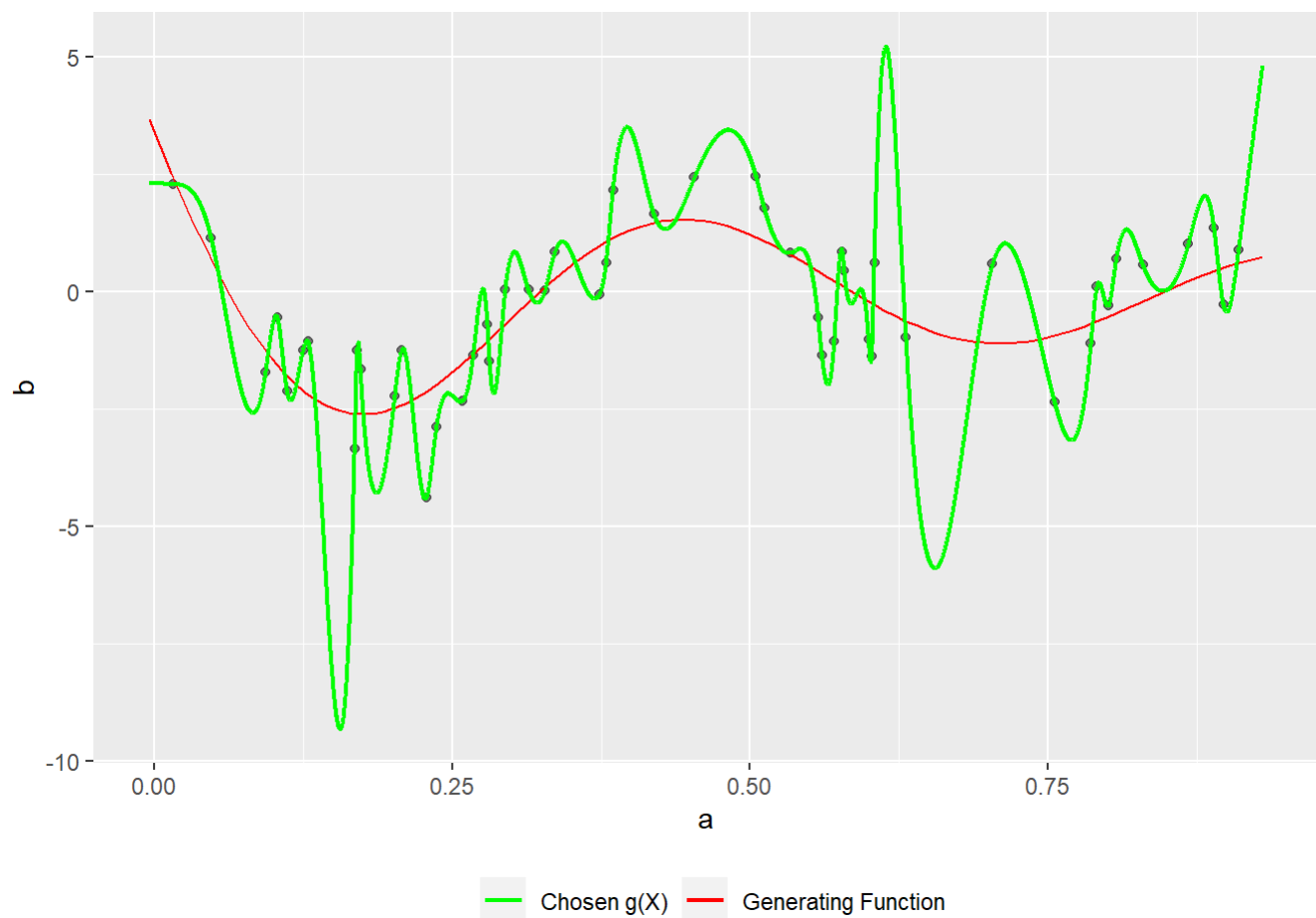



As $\lambda \rightarrow \infty$, this forces $g(3)(x) \rightarrow 0$. This means we would get $\hat{g}(x) = ax^2 + bx + c$.

2e) $\lambda=0, m=3$

```
interp_spline <- smooth.spline(x = df$a, y = df$b, all.knots = T, lambda = 0.0000000000001)
fitted <- predict(interp_spline, x = seq(min(a) - 0.02, max(a) + 0.02, by = 0.0001))
fitted <- data.frame(x = fitted$x, fitted_y = fitted$y)

ggplot(df, aes(x = a, y = b)) +
  geom_point(alpha = 0.5) +
  stat_function(fun = generating_fn, aes(col = "Generating Function")) +
  geom_line(data = fitted,
            aes(x = x, y = fitted_y, col = "Chosen g(X)"), size = 0.8) +
  scale_color_manual(values = c("green", "red")) +
  theme(legend.position = "bottom", legend.title = element_blank())
```

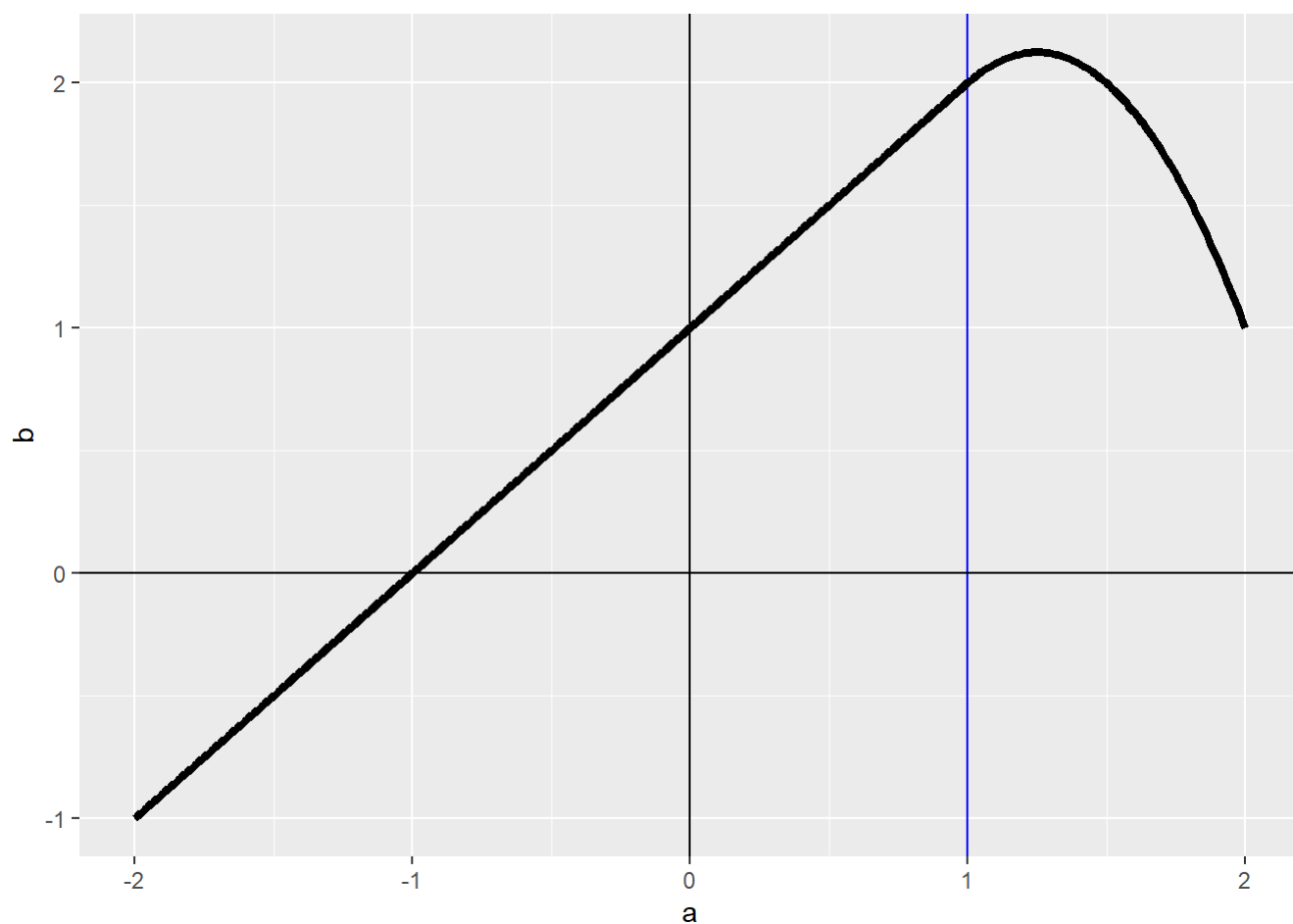


However, since $\lambda=0$, the penalty term no longer plays any role in the selection of $\hat{g}(x)$. For this reason, we can achieve $RSS = 0$

Question 3

```
a = seq(-2, 2, 0.01)
b = 1 + a + -2 * (a - 1)^2 * (a >= 1)
df <- data.frame(a, b)

ggplot(df, aes(x = a, y = b)) +
  geom_vline(xintercept = 0) +
  geom_vline(xintercept = 1, col = "blue") +
  geom_hline(yintercept = 0) +
  geom_line(size = 1.5)
```

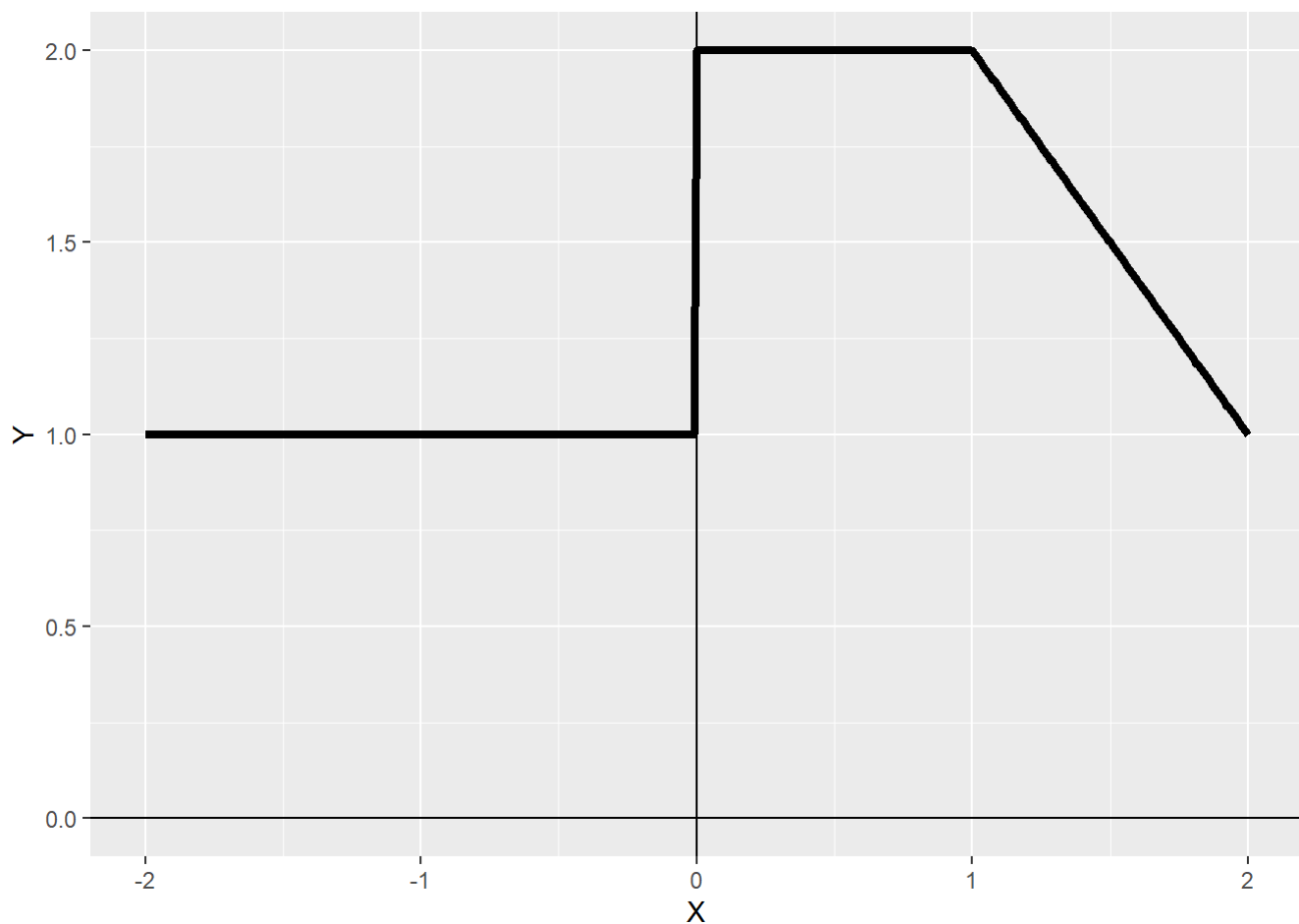


The curve is linear between -2 and 1 with $y=1+x$ Quadratic between 1 , and 2 with $y=1+x-2(x-1)^2$

Question 4

```
X = seq(-2, 2, 0.01)
Y = 1 + (X >= 0 & X <= 2) - (X - 1)*(X >= 1 & X <= 2) + 3*(X - 3)*(X >= 3 & X <= 4) + 3*(X > 4 &
X <= 5)
df <- data.frame(X, Y)

ggplot(df, aes(x = X, y = Y)) +
  geom_vline(xintercept = 0) +
  geom_hline(yintercept = 0) +
  geom_line(size = 1.5)
```



The curve is constant between -2 and 0 with $y=1$ Constant between 0 and 1 with $y=2$ Linear between 1 and 2 with $y=3-x$.

Question 5

- As $\lambda \rightarrow \infty$, will g_1 or g_2 have the smaller training RSS? Answer:- The smoothing spline g_2 will most likely have the smaller training RSS since it is a higher order polynomial owing to the penalty term's order (it will be more flexible).
- As $\lambda \rightarrow \infty$, will g_1 or g_2 have the smaller test RSS? Answer:- The test RSS will depend on the distribution of test data. If we have to provide the behavior of test RSS based on the nature of curve, g_2 will have more test RSS as it is more flexible and hence may overfit the data.
- For $\lambda = 0$, will g_1 or g_2 have the smaller training and test RSS? Answer:- If $\lambda=0$, we have $g_1=g_2$, so they will have the same training and test RSS.