

Problem 1

Question = Problem 1

Pg 1

$$P(x_i | \pi, r) = \binom{x_i + r - 1}{x_i} \pi^{x_i} (1 - \pi)^r$$

(a) Joint Likelihood:

$$L = \prod_{i=1}^n p(x_i) = \binom{x_i + r - 1}{x_i} \pi^{x_i} (1 - \pi)^r$$

$$= \binom{x_1 + r - 1}{x_1} \pi^{x_1} (1 - \pi)^r \cdot \binom{x_2 + r - 1}{x_2} \pi^{x_2} (1 - \pi)^r \cdot \dots \cdot \binom{x_n + r - 1}{x_n} \pi^{x_n} (1 - \pi)^r$$

$$= \left[\binom{x_1 + r - 1}{x_1} \binom{x_2 + r - 1}{x_2} \dots \binom{x_n + r - 1}{x_n} \right] \pi^{\sum x_i} (1 - \pi)^{nr}$$

(b) $\hat{\pi}_{ML} = \argmax_{\pi} L$

$$\argmax_{\pi} c \pi^{\sum x_i} (1 - \pi)^{nr} = \argmax_{\pi} \ln [c \pi^{\sum x_i} (1 - \pi)^{nr}]$$

$$L = \ln c + \sum x_i \ln \pi + nr \ln (1 - \pi)$$

$$\frac{\partial L}{\partial \pi} = \frac{\sum x_i}{\pi} - \frac{nr}{1 - \pi} = 0$$

$$\sum x_i (1 - \pi) - nr \pi = 0$$

$$\Rightarrow \sum x_i - \pi \sum x_i = nr \pi$$

$$\Rightarrow \sum x_i = \pi (\sum x_i + nr)$$

$$\Rightarrow \hat{\pi}_{ML} = \frac{\sum x_i}{\sum x_i + nr} = \frac{\bar{x}}{\bar{x} + r} \quad \bar{x} = \frac{\sum x_i}{n}$$

c) Prior $p(\pi) = \text{beta}(a, b)$

Pg 2

Derive MAP estimate $\hat{\pi}_{\text{MAP}}$ for π .

Find $p(\pi | x_1, x_2, \dots, x_n, r)$

Use

$$p(B|A) = \frac{p(A|B) \cdot p(B)}{p(A)} = \frac{p(A|B) \cdot p(B)}{\sum_B p(A|B) \cdot p(B)} \rightarrow \text{Baye's rule}$$

$$p(\pi | x_1, x_2, \dots, x_n, r) = \frac{p(x_1, x_2, \dots, x_n | \pi, r) \cdot p(\pi)}{\int_{\pi} p(x_1, x_2, \dots, x_n | \pi, r) \cdot p(\pi) d\pi}$$

$$\propto p(x_1, x_2, \dots, x_n | \pi, r) \cdot p(\pi)$$

$$\propto \prod_{i=1}^n \binom{x_i+r-1}{x_i} \pi^{x_i} (1-\pi)^{r-x_i} \frac{\pi^{a-1} (1-\pi)^{b-1}}{\Gamma(a)\Gamma(b)}$$

$$\propto \pi^{\sum x_i + a - 1} (1-\pi)^{nr + b - 1}$$

ignore the constant term.

$$\Rightarrow \propto \prod_{i=1}^n \binom{x_i+r-1}{x_i} \pi^{\sum x_i + a - 1} (1-\pi)^{nr + b - 1}$$

Now, we want $\hat{\pi}_{\text{MAP}}$

$$= \underset{\pi}{\text{argmax}} \left(\binom{x_1+r-1}{x_1} \binom{x_2+r-1}{x_2} \dots \binom{x_n+r-1}{x_n} \pi^{\sum x_i + a - 1} (1-\pi)^{nr + b - 1} \right)$$

$$= \underset{\pi}{\text{argmax}} \sum_{i=1}^n \ln \binom{x_i+r-1}{x_i} + (\sum x_i + a - 1) \ln \pi + (nr + b - 1) \ln(1-\pi)$$

$$\frac{\partial}{\partial \pi} = \frac{\sum x_i + a - 1}{\pi} - \frac{(nr + b - 1)}{1-\pi} = 0$$

$$\Rightarrow (\sum x_i + a - 1)(1-\pi) = (nr + b - 1)\pi$$

Pg 3

$$\Rightarrow \left(\sum x_i + a - 1 \right) = \pi \left[(nr + b - 1) + \left(\sum x_i + a - 1 \right) \right]$$

$$\Rightarrow \sum x_i + a - 1 = \pi \left[\sum x_i + nr + a + b - 2 \right]$$

$$\Rightarrow \left| \hat{\pi}_{MAP} = \frac{\sum x_i + a - 1}{\sum x_i + nr + a + b - 2} \right|$$

d) Using Bayes's rule in previous step, we have already identified the kernel of the posterior distribution is beta ($\sum x_i + a$, $nr + b$)

$$c) E(\pi) = \frac{a}{a+b} \Rightarrow \frac{\sum x_i + a}{\sum x_i + a + nr + b}$$

$$Var(\pi) = \frac{ab}{(a+b)^2(a+b+1)} = \frac{(\sum x_i + a)(nr + b)}{(\sum x_i + a)^2 (\sum x_i + a + nr + b + 1)}$$

We observe the following:

- prior: $\pi \sim \text{beta}(a, b)$; posterior: $\pi \sim \text{beta}(\sum x_i + a, nr + b)$

So there is conjugacy, the posterior is in the same family.

We augment the parameters a and b with $\sum x_i$ and nr respectively.

- $\hat{\pi}_{MAP} = \frac{\sum x_i + a - 1}{\sum x_i + nr + a + b - 2} \neq E(\pi_{\text{posterior}}) = \frac{\sum x_i + a}{\sum x_i + a + nr + b}$

- As $n \rightarrow \infty$, $\sum x_i \rightarrow \infty$ (since $x_i \geq 0$), $\Rightarrow Var(\pi) \rightarrow 0$

Pg 4

$$\hat{\pi}_{\text{MAP}} = \frac{\sum x_i + a - 1}{\sum x_i + nr + a + b - 2} ; \hat{\pi}_{\text{ML}} = \frac{\sum x_i}{\sum x_i + nr}$$

$$\hat{\pi}_{\text{MAP}} = \hat{\pi}_{\text{ML}} \text{ when } a = b = 1$$

and $\pi \sim \text{bet}(a, b) \Rightarrow \pi \sim \text{Uniform}[0, 1]$ when $a = b = 1$

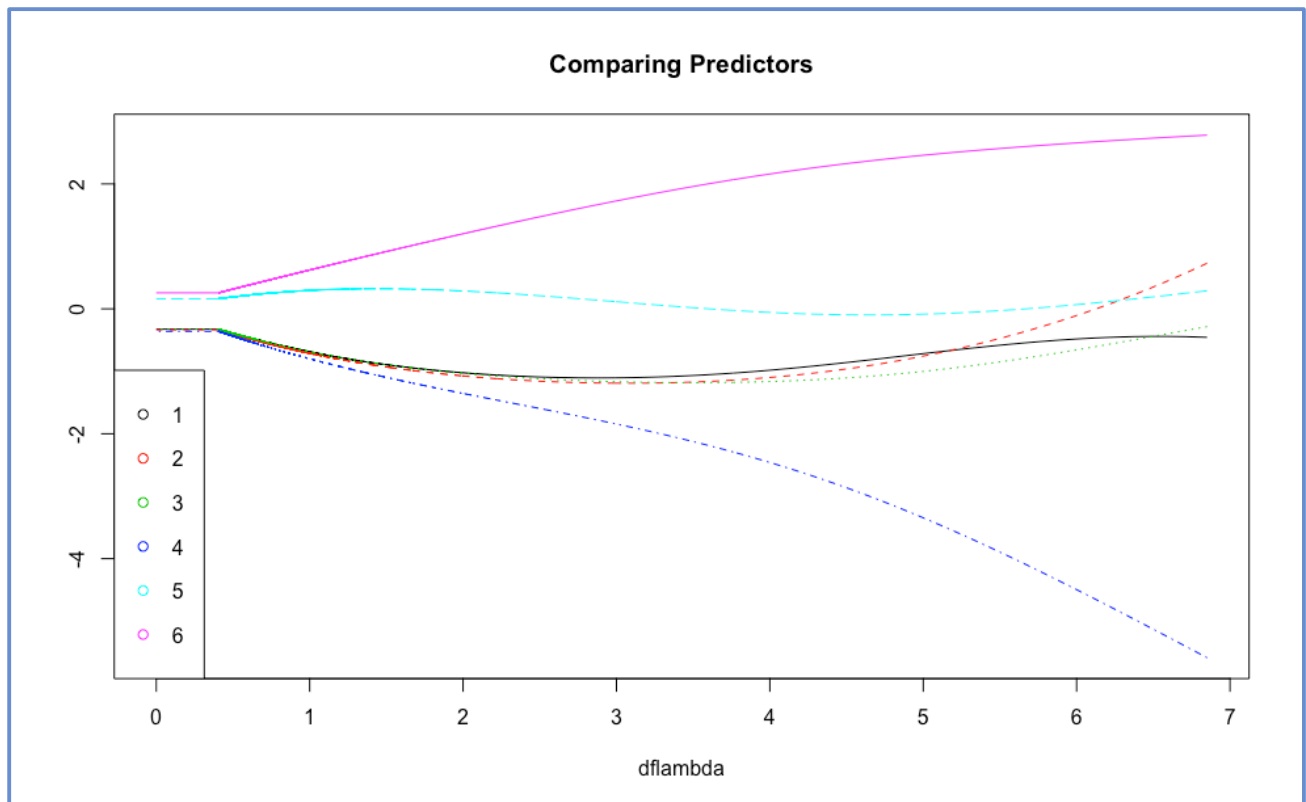
When we believe that our prior distribution of π is $U[0, 1]$, our estimates for π using ML and MAP are the same

Problem 2

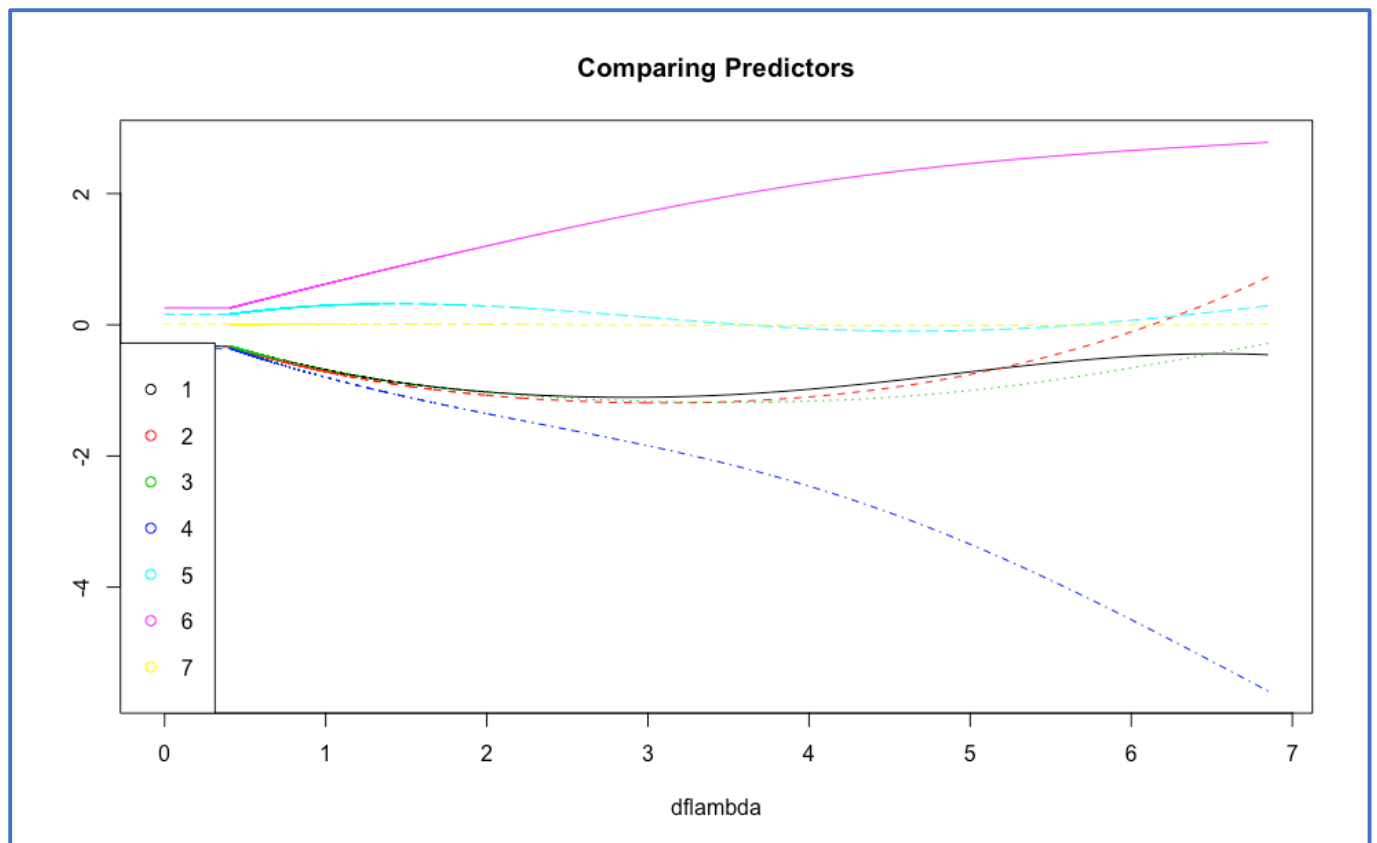
Part 1

a)

The following plot does NOT include the weight for the intercept



The next plot includes 7 predictors, including one for the intercept:



Part 1

b)

As lambda converges to 5000, all the coefficients converge to 0. However, we see that even for large values of lambda, the absolute value of the 4th and 6th predictors is farthest from 0. This means that these two predictors help predict the dependent variable, miles per gallon, the best.

The 4th predictor is weight

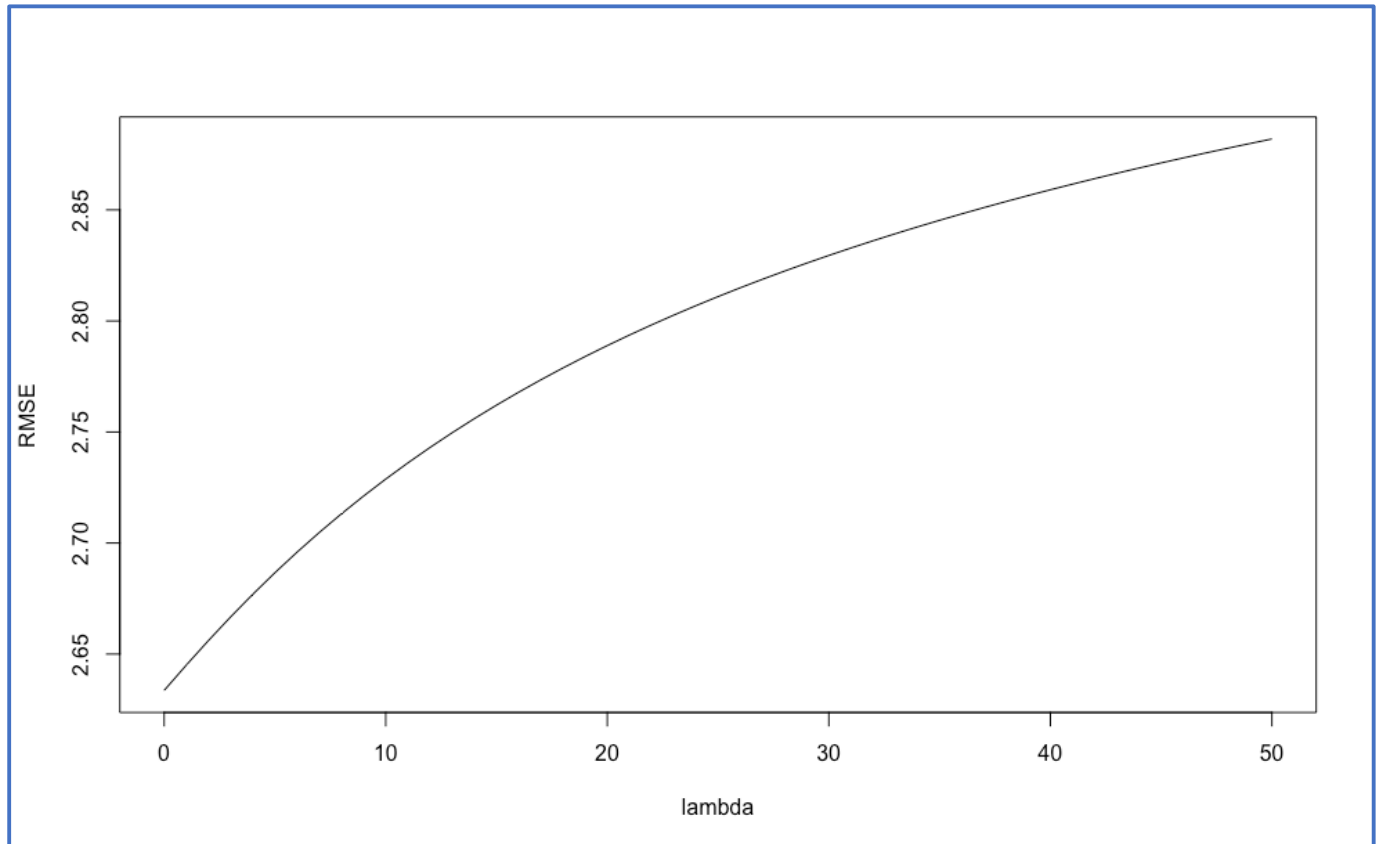
The 6th predictor is model year

According to this plot, weight has a negative relationship with mpg. The higher the weight of a car, the lower its mileage.

Model year has a positive relationship with mpg. The newer the car, the higher its mpg.

Part 1

c)

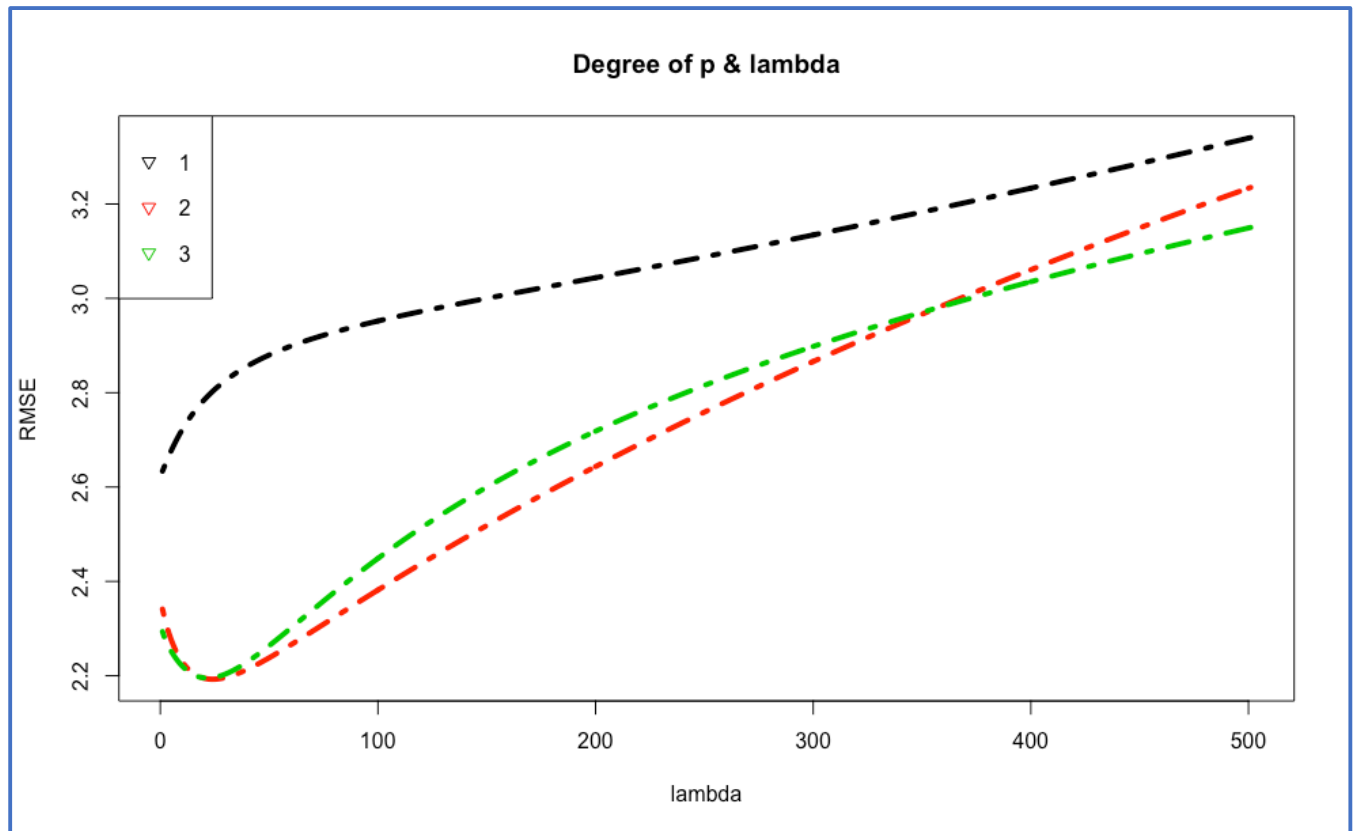


From this graph, we can infer the following:

1. The optimal value of lambda is $\lambda = 0$
2. $\lambda = 0$ corresponds to the least squares solution. Thus, least squares is a better fit for this data set, as compared to ridge regression.
3. When $\lambda > 0$, our estimates for the true parameters are no longer unbiased. However, although bias increases, ridge regression reduces the variance in the model. From the above plot, we infer that the reduction in the variance is not greater than the increase in the bias. This is why RMSE is increasing as lambda increases.
4. This makes sense intuitively. There are only 7 predictors (including intercept) in this data set. And $n = 350$, which is very large. So $n \gg p$. Hence, ridge does not achieve a big decrease in the variance which is low to start off with. This is why OLS is a better estimator.

Part 2

d)



The numbers 1, 2 and 3 in the legend denote the degree of the polynomial used in the fitted model.

a) We see that $p=1$ yields much higher RMSE for any λ

b) We see that $p=2$ and $p=3$ yield very similar RMSE values

It seems like we could choose either $p=2$ or 3, since they both yield similar results.

We should go ahead **and choose $p=2$** because:

- we see that it outperforms $p=3$ for low values of λ .

- Also, increasing predictors in the regression will increase the variance. So unless there is a significant reduction in the bias, we should stick to the simpler model. Hence, we should choose $p=2$

The ideal value of lambda is **lambda = 23**. We found this using the following code:

```
> which.min(testErrorMat[,2])  
[1] 23
```