

Problem 1

(a) What is the joint likelihood of the data (x_1, \dots, x_n) ?

$$p(X|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$$

Joint Likelihood $\Rightarrow L(\lambda; x) = \prod_{i=1}^n p(X_i|\lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$

$$= \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!}$$

\nwarrow
 $x_1! \cdot x_2! \cdot x_3! \cdot \dots \cdot x_n!$

(b) MLE (λ_{MLE}) for λ

$$\log(L(\lambda; x)) = \sum_{i=1}^n x_i \log \lambda - n\lambda$$

$$\frac{d}{d\lambda} \left[\sum_{i=1}^n x_i \log \lambda - n\lambda \right] = 0$$

$$\frac{\sum_{i=1}^n x_i}{\lambda} - n = 0 \Rightarrow \frac{\sum_{i=1}^n x_i}{\lambda} = n$$

$$\therefore \hat{\lambda}_{MLE} = \frac{\sum_{i=1}^n x_i}{n} = \bar{X}$$

(c) Derive the maximum a posteriori (MAP) estimate λ_{MAP} for λ

$$P(\theta|X) \propto P(X|\theta) \cdot P(\theta)$$

posterior likelihood prior

* Assume that prior for λ is $p(\lambda) = \frac{\beta^a}{\Gamma(a)} \lambda^{a-1} e^{-\beta\lambda}$ which is a gamma(a, β) prior

* Likelihood = $\frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!}$
(from 1.a)

$$\lambda_{MAP} = \arg \max_{\lambda} \ln(p(\lambda|X))$$

$$= \arg \max_{\lambda} \ln \left[\frac{p(X|\lambda) \cdot p(\lambda)}{p(X)} \right] = \arg \max_{\lambda} \ln[p(X|\lambda)] + \ln[p(\lambda)] - \ln[p(X)]$$

$$= \arg \max_{\lambda} \sum_{i=1}^n x_i \log \lambda - n\lambda + \ln[\lambda^{a-1} e^{-\beta\lambda}]$$

$$= \sum_{i=1}^n x_i \log \lambda - n\lambda + \ln[\lambda^{a-1}] + \ln[e^{-\beta\lambda}]$$

$$= \sum_{i=1}^n x_i \log \lambda - n\lambda + (a-1) \ln \lambda - \beta\lambda$$

$$\frac{d}{d\lambda} \left[-(n+\beta)\lambda + \sum_{i=1}^n x_i \log \lambda + (a-1) \ln \lambda \right] = 0$$

$$-(n+\beta) + \frac{\sum_{i=1}^n x_i}{\lambda} + \frac{a-1}{\lambda} = 0$$

$$n+\beta = \frac{\sum_{i=1}^n x_i + a - 1}{\lambda}$$

$$\rightarrow \lambda_{MAP} = \frac{\sum_{i=1}^n x_i + a - 1}{n + \beta}$$

d) Derive the posterior distribution of λ and identify this distribution.

$$p(\lambda | X) \propto \underbrace{p(X | \lambda)}_{\text{Likelihood (joint)}} \cdot \underbrace{p(\lambda)}_{\text{prior} \rightarrow \text{gamma}}$$

$$\propto \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!} \cdot \frac{\beta^a \lambda^{a-1} e^{-\beta\lambda}}{\Gamma(a)}$$

$$\propto \frac{e^{-(n+\beta)\lambda} \lambda^{\sum_{i=1}^n x_i + a - 1}}{\prod_{i=1}^n x_i!} \left(\frac{\beta^a}{\Gamma(a)} \right)$$

but since posterior is a proportion of the product, we can get rid of constant + coefficients

$$p(\lambda | X) \propto e^{-(n+\beta)\lambda} \lambda^{\sum_{i=1}^n x_i + a - 1}$$

\therefore The distribution is Gamma $\left(\underbrace{\sum_{i=1}^n x_i + a}_{"a"}, \underbrace{n+\beta}_{"b"} \right)$

e) What is the mean and variance of λ under posterior distribution? Discuss how it relates to λ_{MLE} and λ_{MAP} .

Gamma distribution $\left(\sum_{i=1}^n x_i + a, n+\beta \right)$

$$\text{Mean} = E[\lambda] = \frac{\sum_{i=1}^n x_i + a}{n+\beta}$$

$$\text{variance} = \text{var}[\lambda] = \frac{\sum_{i=1}^n x_i + a}{(n+\beta)^2}$$

The mean and variance of λ can help us to describe a distribution in which our point estimates λ_{MLE} and λ_{MAP} fall within.

$$\lambda_{MAP} = \frac{\sum_{i=1}^n x_i + a - 1}{n+\beta} \text{ is the mode of the posterior distribution (Gamma),}$$

where the mode takes on the form of $E[\lambda] - \frac{1}{\beta}$.

whereas λ_{MLE} is the value of the random variable parameter that maximizes the entire distribution.

Problem 2

Given $y_i \stackrel{iid}{\sim} N(x_i^T w, \sigma^2)$ and $w_{RR} = (\lambda I + X^T X)^{-1} X^T y$

Calculate $E[w_{RR}]$, $\text{var}[w_{RR}]$

We know

$E[w_{LS}]$

$(X^T X)^{-1} X^T y$

$$E[w_{RR}] = E[(\lambda I + X^T X)^{-1} X^T y]$$

$$= E[(\lambda I + X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T y]$$

$$= E[(\lambda I + X^T X)^{-1} X^T X w_{LS}] \quad \underbrace{(X^T X)^{-1} X^T X = I}_{=1} \quad (\lambda I + X^T X)^{-1} X^T X w$$

$$= E[(X^T X) (\lambda (X^T X)^{-1} + I)^{-1} (X^T X) w_{LS}] = w_{ML}$$

$$= E[(\lambda (X^T X)^{-1} + I)^{-1} (X^T X)^{-1} X^T X w_{LS}]$$

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

$$\text{var}[w_{RR}] = \text{var}[(\lambda (X^T X)^{-1} + I)^{-1} w_{LS}]$$

$$\text{Let } Z = (\lambda (X^T X)^{-1} + I)^{-1}$$

$$\text{var}(Z w_{LS}) = Z \text{var}(w_{LS}) Z^T$$

$$= Z \sigma^2 (X^T X)^{-1} Z^T = \sigma^2 (I + \lambda (X^T X)^{-1})^{-1} (X^T X)^{-1} (I + \lambda (X^T X)^{-1})^{-1}$$

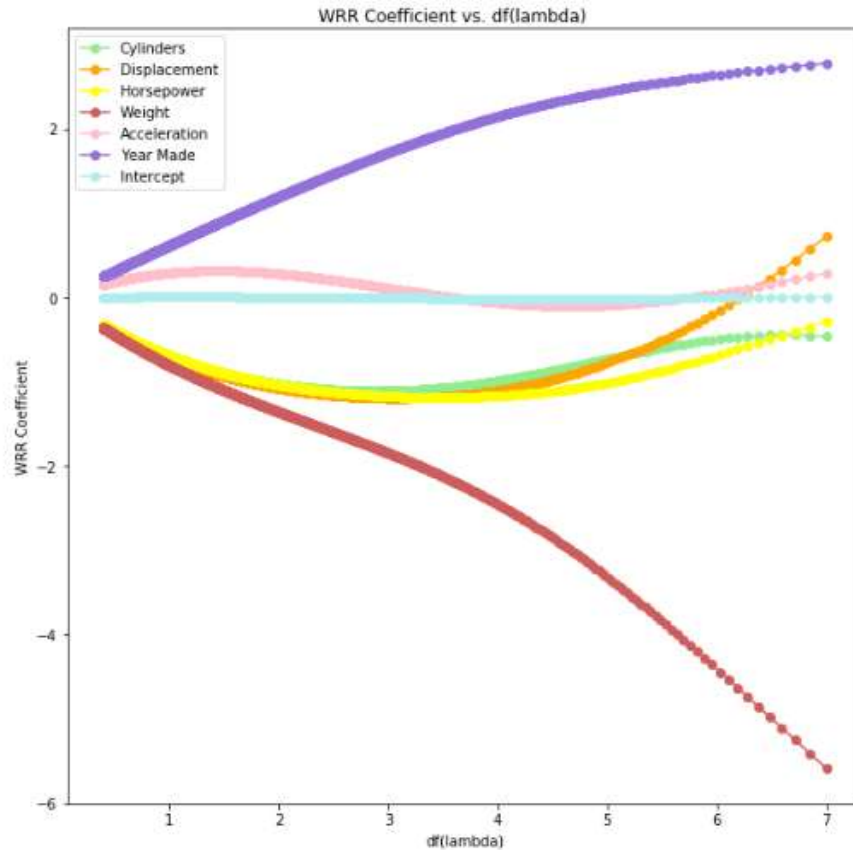
$$= (\lambda I + X^T X)^{-1} X^T X (X^T X)^{-1} X^T X (\lambda I + X^T X)^{-1}$$

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

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

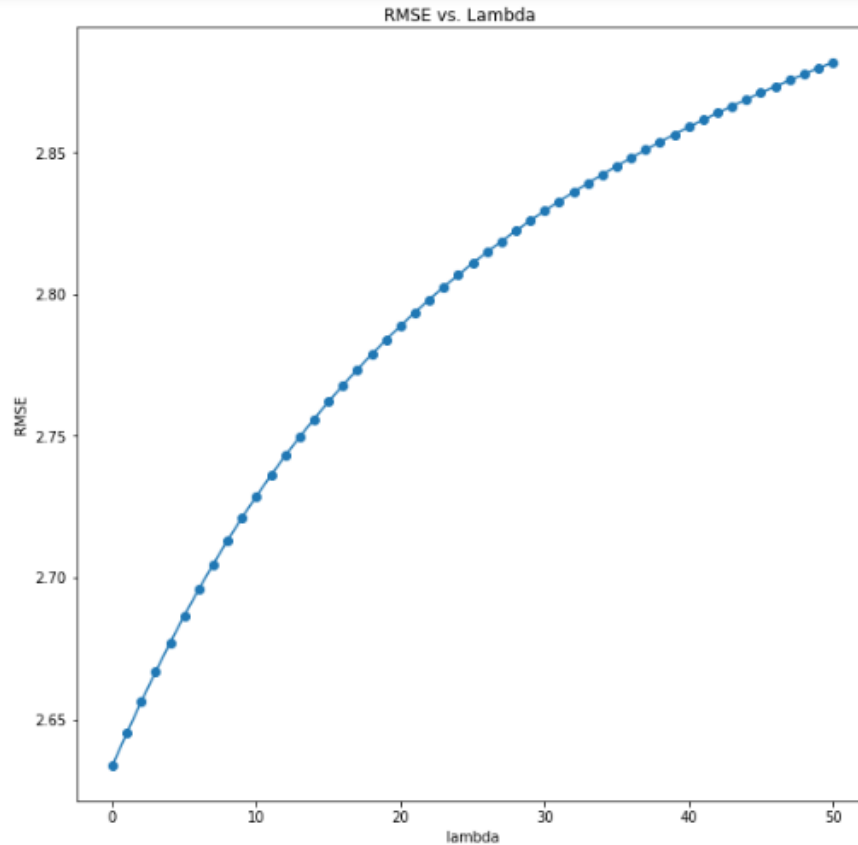
$$(\lambda I + X^T X)^{-1}$$

$$\lambda I$$



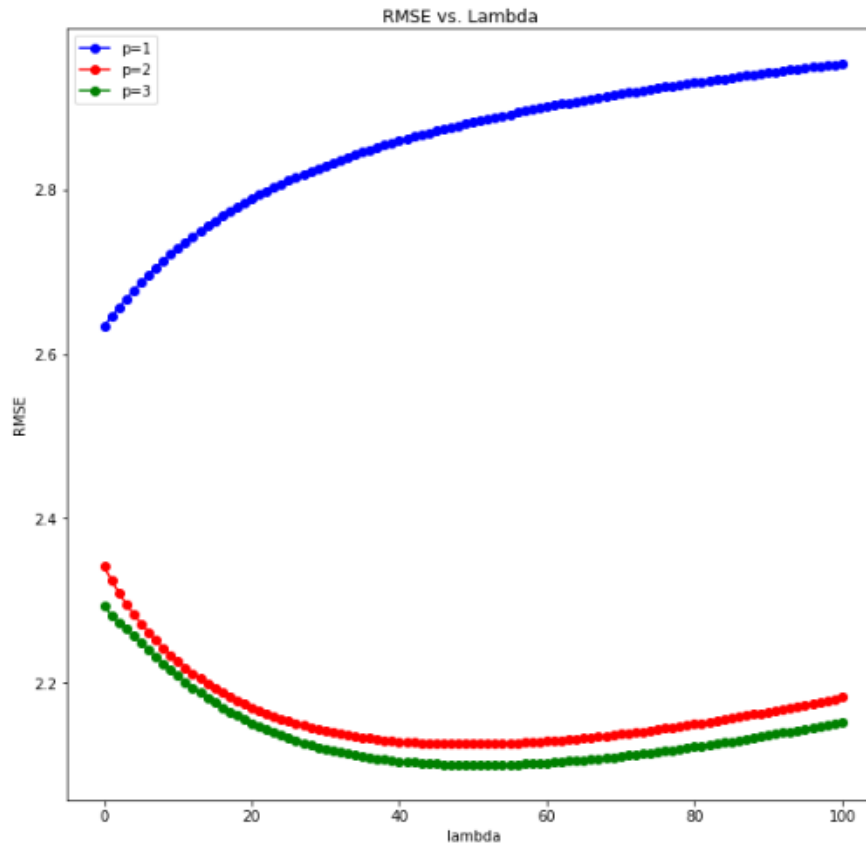
Part b)

The two dimensions that stand out the most over the others are 'Weight' and 'Year Made'. We can glean that these two variables, as the $df(\lambda)$ decreases as a result of our λ increasing and causing the denominator of $df(\lambda)$ to get very large, this causes our wrr coefficients to go toward 0. It seems that Year Made and Weight are the two variables whose wrr coefficient does not go to 0 as quickly as the other variables. As a result, the two dimensions above that stand out have the most effect on the dependent variable, mpg for the car.



Insights:

We can see from the graph that at lower values of lambda, we experience low RMSE's. However, as the lambda increases, the RMSE increases, albeit at a decreasing rate. Since the RMSE is minimized when $\lambda = 0$, we are effectively solving a least squares problem. Therefore, we might as well just use a least squares regression method to solve for/estimate w .



Part 2d)

We should choose the $p = 3$ technique, which involves the transformation and standardization of X to include squared and cubic terms. The general trend of this line decreases when lambda is in between $[0,40]$ and then slowly increases after that, which may signal potential overfitting when lambda is large and the model is more complex. Since $p=3$ line overall results in a RMSE measure of lambda = 0 through 100, we want to choose this one since it has the least amount of model error. The lambda for which the RMSE is minimized is when lambda = 40.

Observations about the other lines include that the $p=2$ line also performs surprisingly well compared to $p=3$, but has a slightly higher RMSE across all lambdas. In stark contrast to these two, $p=1$ line performs the worse with an overall increasing RMSE across all lambdas, with the minimum RMSE occurring at lambda = 0. This implies that this model does not do terribly well at predicting our DV values. ¶

As for the optimal lambda value for the $p=3$ lines, which is 40, this value differs from our previous lambda choice, which was lambda = 0 for the $p = 1$ plot from earlier.

The moral of the story is that the transformation and standardization of our data can help us to achieve a much better model with less error.