

Archit Matta
am5500

Homework 1
5/2/21

$$\textcircled{1} \quad p(X|\lambda) = \frac{\lambda^X e^{-\lambda}}{X!}$$

a) Joint likelihood of (x_1, x_2, \dots, x_N) would be equal to the product of the individual probabilities for x_i s.

$$\Rightarrow \left(\frac{\lambda^{x_1} e^{-\lambda}}{x_1!} \right) \cdot \left(\frac{\lambda^{x_2} e^{-\lambda}}{x_2!} \right) \cdot \dots \cdot \left(\frac{\lambda^{x_N} e^{-\lambda}}{x_N!} \right)$$

$$\Rightarrow \boxed{\frac{\lambda^{x_1+x_2+\dots+x_N}}{x_1! \cdot x_2! \cdot \dots \cdot x_N!} \cdot e^{-N\lambda}}$$

b) Maximum Likelihood Estimate for λ (λ_{ML})

$$\lambda_{ML} = \arg \max_{\lambda} (\text{Joint Likelihood})$$

We can take log of the function & maximise it.

$$x_{ML} = \arg \max_{\lambda} (\ln \left[\frac{\lambda^{x_1+x_2+\dots+x_N}}{x_1! \cdot x_2! \cdot \dots \cdot x_N!} \cdot e^{-N\lambda} \right])$$

Setting derivative w.r.t. λ as 0

$$\frac{d}{d\lambda} \left[(x_1+x_2+\dots+x_N) \lambda + \ln \left(\frac{1}{x_1! \cdot x_2! \cdot \dots \cdot x_N!} \right) - N\lambda \right] = 0$$

$$\Rightarrow \frac{(x_1+x_2+\dots+x_N)}{\lambda} - N = 0$$

$$\Rightarrow \boxed{\lambda_{ML} = \frac{x_1+x_2+\dots+x_N}{N}}$$

(C)

$$p(\lambda) = \text{gamma}(a, b)$$

$$p(\lambda) = \frac{b^a \lambda^{a-1} e^{-b\lambda}}{\Gamma(a)}$$

$$\lambda_{MAP} = ?$$

$$\lambda_{MAP} = \arg \max_{\lambda} (P[\lambda | x])$$

Taking log and maximising log since it is monotonic

$$\lambda_{MAP} = \arg \max_x (\ln (P[\lambda | x]))$$

$$P[\lambda | x] = P[x | \lambda] \cdot P[\lambda] \quad (\text{Baye's rule})$$

Constant

Setting the derivative equal to 0

$$\Rightarrow \frac{d}{d\lambda} [\ln(P[x | \lambda]) + \ln(P[\lambda])] = \ln(\text{constant}) = 0$$

$$\Rightarrow \frac{d}{d\lambda} \left[\ln \left(\frac{\lambda^{(x_1+x_2+\dots+x_n)} e^{-N\lambda}}{x_1! x_2! \dots x_n!} \right) + \ln \left[\frac{b^a \lambda^{a-1} e^{-b\lambda}}{\Gamma(a)} \right] - \ln(\text{constant}) \right] = 0$$

$$\Rightarrow \frac{d}{d\lambda} [N \ln \lambda - N\lambda + (a-1) \ln \lambda - b\lambda] = 0$$

(Ignoring Terms without λ)

$$\Rightarrow \frac{x + a - 1}{\lambda} - (b + N) = 0$$

\Rightarrow

$$\lambda_{MAP} = \frac{x + a - 1}{b + N}$$

where $x = x_1 + x_2 + \dots + x_n$

$$(d) P(\lambda|x) = \frac{P(x|\lambda) \cdot P(\lambda)}{\int P(x|\lambda) \cdot P(\lambda) d\lambda}$$

$$\Rightarrow P(\lambda|x) \propto P(x|\lambda) \cdot P(\lambda)$$

$$\Rightarrow P(\lambda|x) \propto \left(\frac{\lambda^{x_1+x_2+\dots+x_N}}{x_1! x_2! \dots x_N!} e^{-N\lambda} \right) \cdot \left(\frac{b^a \lambda^{a-1} e^{-b\lambda}}{\Gamma(a)} \right)$$

$$\Rightarrow P(\lambda|x) \propto \frac{b^a}{(x_1! x_2! \dots x_N!) \Gamma(a)} \cdot \lambda^{(x_1+x_2+\dots+x_N)+(a-1)} \cdot e^{-(N+b)\lambda}$$

We can see this looks like a Gamma Distribution

The gamma distribution can be obtained by multiplying constants.

The parameters of the distribution would be -

$$\alpha = x_1 + x_2 + \dots + x_N + a$$

$$\beta = N + b$$

Thus, the posterior distribution of λ is gamma (α, β).

$$(e) \text{ Mean of Gamma} = \frac{\alpha}{\beta}, \text{ Variance of Gamma} = \frac{\alpha}{\beta^2}$$

$$\text{Mean} = \frac{x_1 + x_2 + \dots + x_N + a}{(N+b)}$$

$$\text{Variance} = \frac{(x_1 + x_2 + \dots + x_N + a)}{(N+b)^2}$$

We can see the following -

$$(2) \text{ Mean} = \lambda_{MAP} + \frac{1}{b+N}$$

$$(3) \text{ Variance} = (\lambda_{MAP} + \frac{1}{b+N})$$

$$(2) \text{ Mean} = ((\lambda_{ML} \times N) + a) / (N+b)$$

$$(4) \text{ Variance} = \frac{((\lambda_{ML} \times N) + a)}{(b+N)^2}$$

(2) @

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

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

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

$$E[y] = Xw \quad \text{since } y \sim N(Xw, \sigma^2 I)$$

 \Rightarrow

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

$$\text{Var}[w_{RK}] = E[w_{RK} w_{RK}^T] - E[w_{RK}] E[w_{RK}]^T$$

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

$$\Rightarrow \text{Var}[w_{RK}] = E[A X^T y y^T X A^T] - A X^T X w w^T X^T X A^T$$

$$= A X^T (\sigma^2 I + Xw w^T X^T) X A^T - A X^T X w w^T X^T X A^T$$

$$= A X^T \sigma^2 I X A^T + A X^T X w w^T X^T X A^T - A X^T w w^T X^T X A^T$$

$$\text{Var}[w_{RK}] = \sigma^2 (\lambda I + X^T X)^{-1} X^T X (\lambda I + X^T X)^{-1 T}$$

Right multiply both sides by $(\lambda I + X^T X)^T (\lambda I + X^T X)^{-1}$
 \hookrightarrow This is equal to $(\lambda (X^T X)^{-1} + I)$

$$\text{Var}[w_{RK}] (\lambda (X^T X)^{-1} + I)^T = \sigma^2 (\lambda I + X^T X)^{-1}$$

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

$$\text{Var}[w_{RK}] = \sigma^2 (\lambda I + X^T X)^{-1} ((I + \lambda (X^T X)^{-1})^{-1})^T$$

Left multiply both sides by $(\lambda I + X^T X) (\lambda (X^T X)^{-1})^{-1}$

$$(\lambda (X^T X)^{-1} + I) \text{Var}[w_{RK}] = \sigma^2 (X^T X)^{-1} (I + \lambda (X^T X)^{-1})^{-1}$$

Left multiplying both sides by $(\lambda(x^T x)^{-1} + I)^{-1}$

$$\text{Var}[w_{RE}] = \sigma^2 [\lambda(x^T x)^{-1} + I]^{-1} (x^T x) [\lambda(x^T x)^{-1} + I]^{-1 T}$$

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

$$\boxed{\text{Var}[w_{RE}] = \sigma^2 Z (x^T x)^{-1} Z^T}$$

(5)

$$w_{RE} = (X(I + x^T x)^{-1} x^T y$$

$$w_{LS} = (x^T x)^{-1} x^T y$$

$$X = V S V^T$$

$$w_{RE} = (\lambda I + x^T x)^{-1} (x^T x) \underbrace{(x^T x)^{-1} x^T y}_{w_{LS}}$$

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

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

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

Using $X = V S V^T$, we get $(x^T x)^{-1} = V S^{-2} V^T$

$$w_{RE} = (\lambda(V S^{-2} V^T) + I)^{-1} w_{LS}$$

$$w_{RE} = V (\lambda S^{-2} + I)^{-1} V^T w_{LS}$$

$$\boxed{w_{RE} = V M V^T w_{LS}}$$

where M is a diagonal matrix with $M_{ii} = \frac{s_{ii}^2}{\lambda + s_{ii}^2}$

$$\boxed{w_{RE} = V (\lambda S^{-2} + I)^{-1} V^T w_{LS}}$$

Problem 3

a) Here's a brief look at the obtained Wrr values.

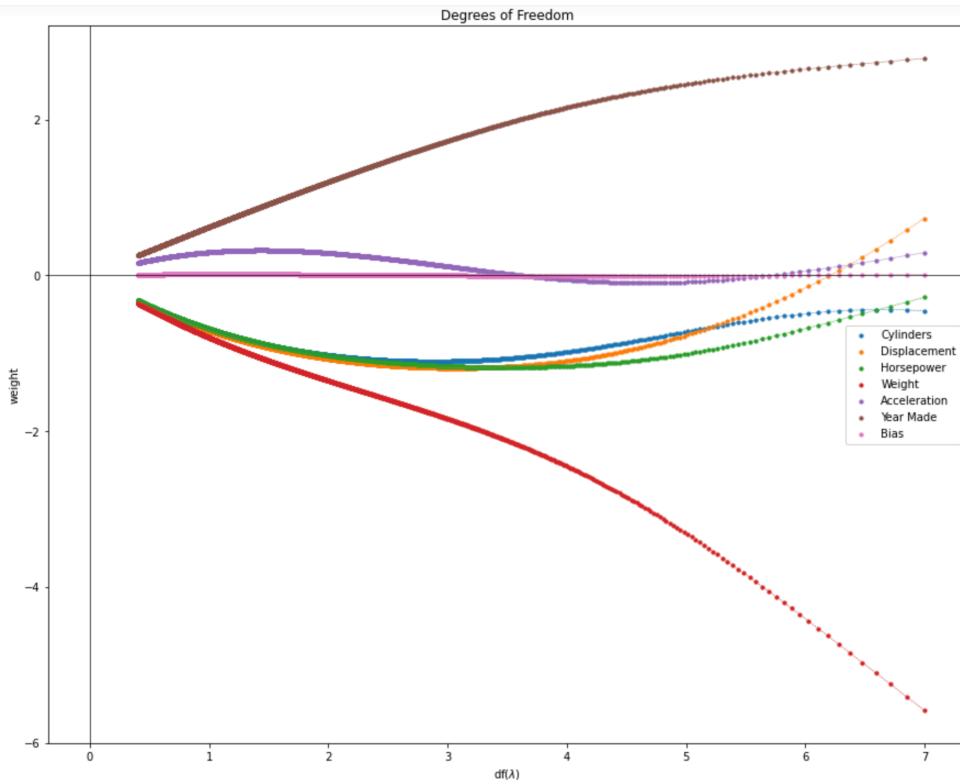
```
In [5]: #Printing the W for Ridge Regression for all Lambdas
pd.DataFrame(np.asarray(wrr).reshape(5001,7))
```

Out[5]:

	0	1	2	3	4	5	6
0	-0.456268	0.730185	-0.284646	-5.585580	0.289570	2.781397	0.010157
1	-0.445729	0.577783	-0.344522	-5.409677	0.251099	2.763334	0.008127
2	-0.441314	0.445755	-0.399202	-5.250282	0.216899	2.746404	0.006362
3	-0.441431	0.330230	-0.449221	-5.104973	0.186365	2.730448	0.004816
4	-0.444922	0.228265	-0.495063	-4.971812	0.159004	2.715338	0.003450
5	-0.450930	0.137581	-0.537159	-4.849218	0.134409	2.700969	0.002236
6	-0.458814	0.056385	-0.575889	-4.735887	0.112239	2.687253	0.001152
7	-0.468090	-0.016752	-0.611593	-4.630732	0.092206	2.674118	0.000179
8	-0.478388	-0.082984	-0.644570	-4.532839	0.074067	2.661502	-0.000698
9	-0.489422	-0.143252	-0.675083	-4.441430	0.057613	2.649351	-0.001492
10	-0.500972	-0.198333	-0.703368	-4.355839	0.042665	2.637622	-0.002212

Here's the graph between Degrees of Freedom and Weight

```
In [6]: #Plotting Degrees of Freedom vs Weight
wrr_np = np.asarray(wrr)
svd_np = np.asarray(svd)
plt.figure()
dimensions = ["Cylinders", "Displacement", "Horsepower", "Weight", "Acceleration", "Year Made", "Bias"]
for i in range(wrr_np[0].shape[0]):
    plt.plot(svd_np, wrr_np[:,i], linestyle = 'dashed', linewidth = 0.5)
    plt.scatter(svd_np, wrr_np[:,i], s = 10, label = dimensions[i])
plt.axhline(y = 0, color = 'black', linewidth = 0.7)
plt.axvline(x = 0, color = 'black', linewidth = 0.7)
plt.xlabel('df($\lambda$)')
plt.ylabel('weight')
plt.title('Degrees of Freedom')
plt.legend()
plt.show()
```



b)

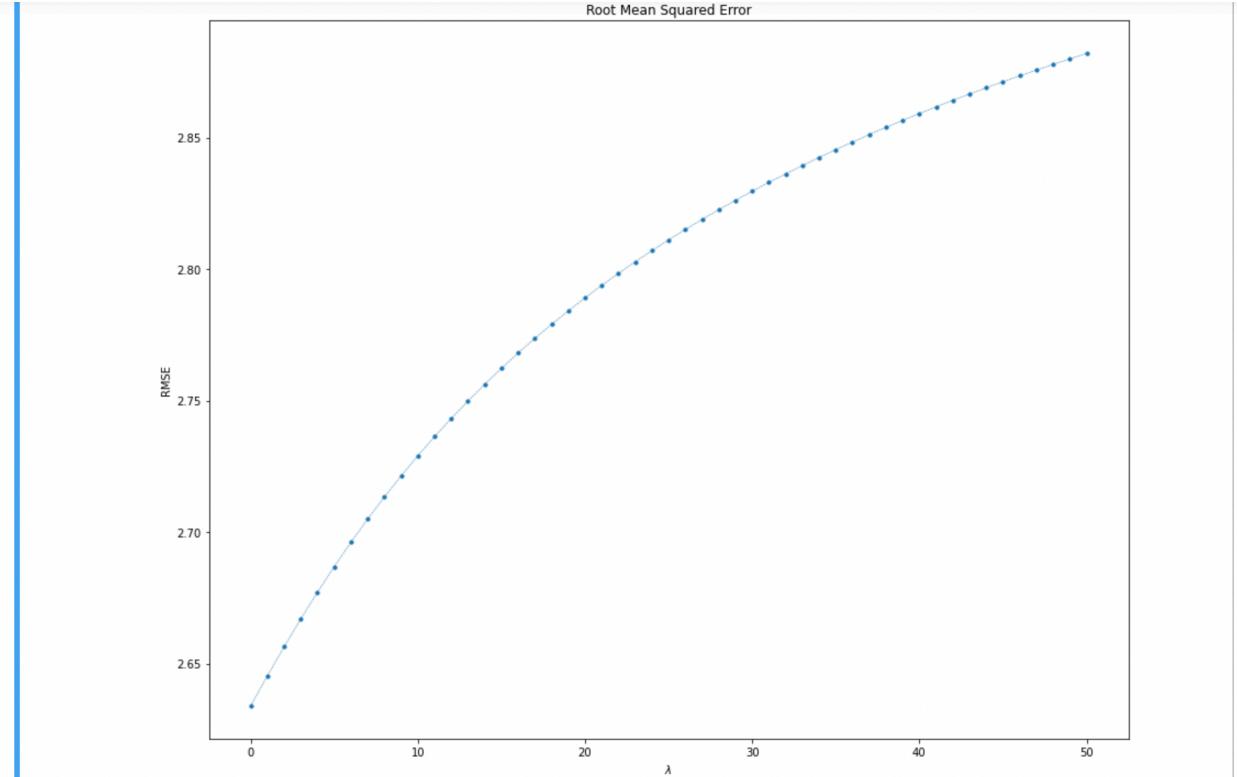
Here are the 2 dimensions that stand out

The 2 dimensions that stand out are 'Year Made' and 'Weight'. 'Year Made' has a high positive weight and thus indicates that newer cars have higher Miles per Gallon. 'Weight' has a high negative weight and thus indicates that lighter cars have higher Miles per Gallon.

c)

Here's the graph between RMSE and Lambda

```
In [9]: #Plotting RMSE vs Lambda
plt.plot(lambda[:51], rmse, linestyle = 'dashed', linewidth = 0.5)
plt.scatter(lambda[:51], rmse, s = 10)
plt.xlabel('$\lambda$')
plt.ylabel('RMSE')
plt.title('Root Mean Squared Error')
plt.show()
```



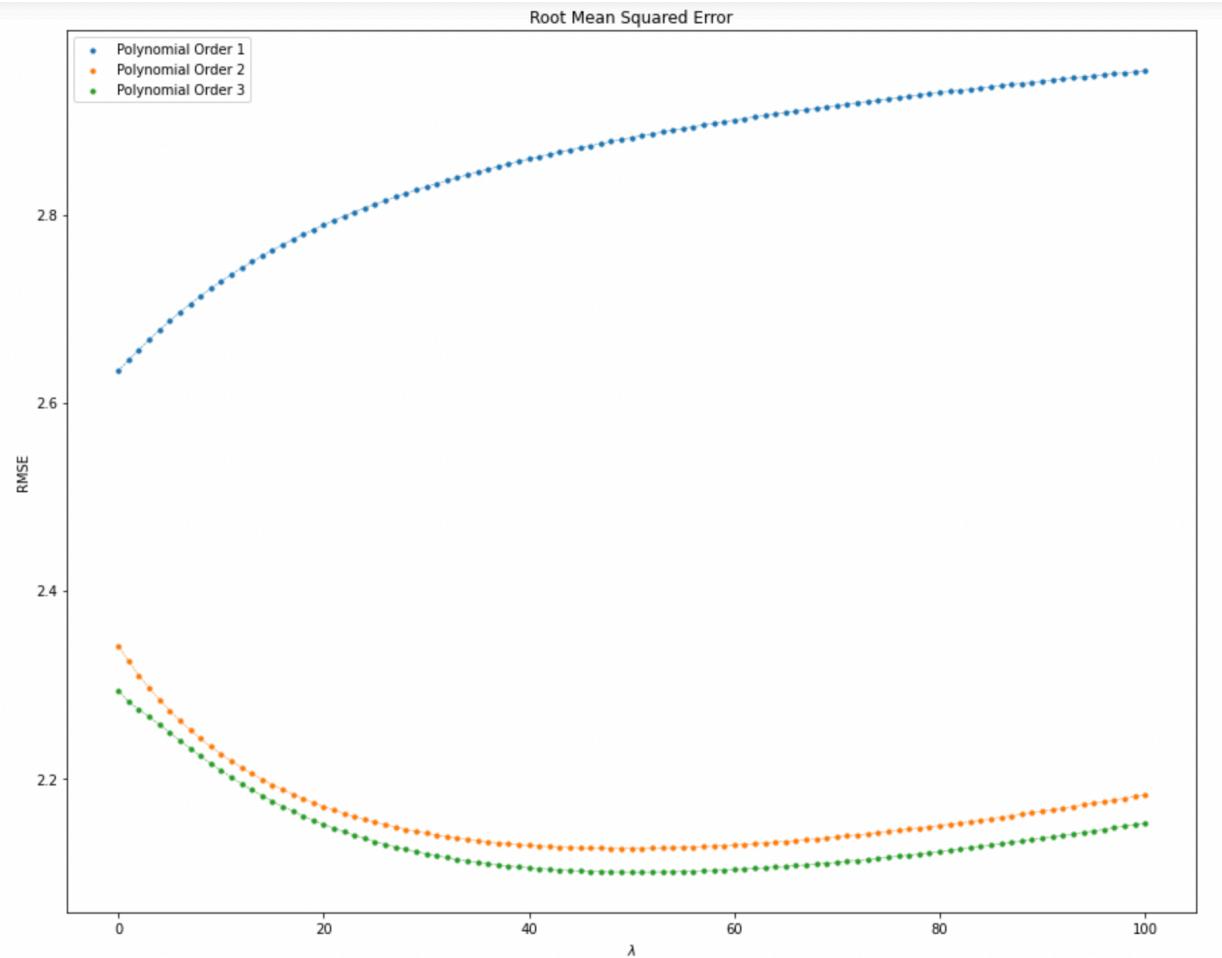
Here's what we observe by looking at this figure

Based on this figure, we can see that the optimal Lambda for this problem having the lowest RMSE is 0. We also know that for Lambda = 0, W for Ridge Regression is the same as W for Least Squares. Hence, based on this graph we should choose Least Squares for this problem.

d)

Here's the graph between RMSE and Lambda

```
#Plotting RMSE vs Lambda for different Polynomial Orders
plt.plot(lamda[:101], rmse_1, linestyle = 'dashed', linewidth = 0.5)
plt.scatter(lamda[:101], rmse_1, s = 10, label = 'Polynomial Order 1')
plt.plot(lamda[:101], rmse_2, linestyle = 'dashed', linewidth = 0.5)
plt.scatter(lamda[:101], rmse_2, s = 10, label = 'Polynomial Order 2')
plt.plot(lamda[:101], rmse_3, linestyle = 'dashed', linewidth = 0.5)
plt.scatter(lamda[:101], rmse_3, s = 10, label = 'Polynomial Order 3')
plt.xlabel('$\lambda$')
plt.ylabel('RMSE')
plt.title('Root Mean Squared Error')
plt.legend()
plt.show()
```



Here's what we observe by looking at this figure

Based on this plot, we should choose a Polynomial of Order 3 ($p = 3$) as it has the least RMSE overall. Our assessment of the ideal value of Lambda also changes to 51 from our initial value of 0 as the lowest RMSE value for $p = 2$ occurs at this point.