

Name : PRASHAM D. SHETH

UNI : pds2136

### Problem - 1

(a)  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$   
 $\therefore p(x|\lambda) = \frac{\lambda^x \cdot e^{-\lambda}}{x!}$

$$\begin{aligned} L(\lambda; x_1, \dots, x_n) &= p(x_1, \dots, x_n | \lambda) \\ &= \prod_{i=1}^N p(x_i | \lambda) \\ &= \prod_{i=1}^N \frac{\lambda^{x_i} \cdot e^{-\lambda}}{x_i!} \end{aligned}$$

(b)  $\hat{\lambda}_{ML} = \underset{\lambda}{\operatorname{argmax}} \prod_{i=1}^N \frac{\lambda^{x_i} \cdot e^{-\lambda}}{x_i!}$   
 $= \underset{\lambda}{\operatorname{argmax}} \log \left( \prod_{i=1}^N \frac{\lambda^{x_i} \cdot e^{-\lambda}}{x_i!} \right)$   
 $= \underset{\lambda}{\operatorname{argmax}} \left[ \log \lambda \cdot \sum_{i=1}^N x_i - N \log \lambda - \sum_{i=1}^N \log(x_i!) \right]$   
 $= \underset{\lambda}{\operatorname{argmax}} \left[ \underbrace{(\log \lambda) \left( \sum_{i=1}^N x_i \right) - N \log \lambda}_{L_1} \right]$

$$L_1 = (\log \lambda) \cdot \left( \sum_{i=1}^N x_i \right) - N \log \lambda$$

$$\frac{\partial L_1}{\partial \lambda} = \frac{\sum_{i=1}^N x_i}{\lambda} - N$$

$$\frac{\partial L}{\partial \lambda} = 0$$

$$\therefore \frac{\sum_{i=1}^N x_i}{c} = N$$

$$\therefore \boxed{\hat{\lambda}_{ML} = \frac{\sum_{i=1}^N x_i}{N}}$$

(c) Prior distribution of gamma is assumed for  $\lambda$ .

$$\therefore p(\lambda) = \text{gamma}(a, b) = \frac{b^a \cdot \lambda^{a-1} \cdot e^{-b\lambda}}{\Gamma(a)}$$

$$\hat{\lambda}_{MAP} = \arg \max_{\lambda} \log(p(\lambda | x_1, \dots, x_N))$$

$$= \arg \max_{\lambda} [\log(p(x_1, \dots, x_N | \lambda)) + \log p(\lambda)]$$

$$= \arg \max_{\lambda} \left[ \log \left( \prod_{i=1}^N p(x_i | \lambda) \right) + \log p(\lambda) \right]$$

$$= \arg \max_{\lambda} \left[ \sum_{i=1}^N \log p(x_i | \lambda) + \log p(\lambda) \right]$$

$$= \arg \max_{\lambda} \left[ \sum_{i=1}^N \left( \frac{\lambda^{x_i} \cdot e^{-\lambda}}{x_i!} \right) + \log \left( \frac{b^a \cdot \lambda^{a-1} \cdot e^{-b\lambda}}{\Gamma(a)} \right) \right]$$

$$= \arg \max_{\lambda} \left[ \log \lambda \left( \sum_{i=1}^N x_i \right) - N \cdot \lambda + (\log \lambda)(a-1) - b\lambda \right]$$

$$\therefore \hat{\lambda}_{\text{MAP}} = \underset{\lambda}{\text{argmax}} \left[ \log \lambda \left[ \sum_{i=1}^N x_i + a - 1 \right] - \lambda (N+b) \right]$$

$$L_1 = \log \lambda \left( \sum_{i=1}^N x_i + a - 1 \right) - \lambda (N+b)$$

$$\frac{\partial L_1}{\partial \lambda} = 0$$

$$\therefore \frac{\sum_{i=1}^N x_i + a - 1}{\lambda} - (N+b) = 0$$

$$\therefore \boxed{\hat{\lambda}_{\text{MAP}} = \frac{\sum_{i=1}^N x_i + a - 1}{N+b}}$$

(d) for posterior distribution of  $\lambda$ ; we use Bayes' rule.

$$P(\lambda | x_1, \dots, x_n) = \frac{P(x_1, \dots, x_n | \lambda) \cdot P(\lambda)}{P(x_1, \dots, x_n)}$$

Here, the denominator normalizes the numerator & doesn't depend on  $\lambda$ .

$$\begin{aligned} \therefore P(\lambda | x_1, \dots, x_n) &\propto P(x_1, \dots, x_n | \lambda) \cdot P(\lambda) \\ &\propto \frac{\lambda^{\sum_{i=1}^N x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^N (x_i!)} \cdot \frac{b^a \cdot \lambda^{a-1} \cdot e^{-b\lambda}}{\Gamma(a)} \\ &\propto \lambda^{\left(\sum_{i=1}^N x_i + a - 1\right)} \cdot e^{-\lambda(n+b)} \end{aligned}$$

From the above form we can identify that posterior distribution of  $c$  must be a gamma with parameters  $(\sum_{i=1}^N x_i + a, N+b)$

$$\therefore P(c | x_1, \dots, x_N) = \text{gamma} \left( \sum_{i=1}^N x_i + a, N+b \right)$$

(e) ~~As~~ we know,  $E[X] = \frac{a}{b}$  and  $\text{Var}[X] = \frac{a}{b^2}$  where  $X \sim \text{gamma}(a, b)$

Now, as

$$P(c | x_1, \dots, x_N) \sim \text{gamma} \left( \sum_{i=1}^N x_i + a, N+b \right)$$

Mean value of  $c$  =  $E[c] = \frac{\sum_{i=1}^N x_i + a}{N+b}$

Variance of  $c$  =  $\text{Var}[c] = \frac{\sum_{i=1}^N x_i + a}{(N+b)^2}$

→ This clearly shows that the mean value of  $c$  under the posterior is greater than  $\hat{\lambda}_{MAP}$ .

→  $\hat{\lambda}_{ML}$  and this mean value of  $c$  under the posterior would be equal when  $a$  and  $b$  would be zero.

## Problem - 2 :-

$$\rightarrow y_i \stackrel{\text{iid}}{\sim} N(x_i^T \cdot w, \sigma^2)$$

Using data we already have approximated  $w_{RR}$  as

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

$\rightarrow$  We also know that

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

$$E[w_{LS}] = w$$

$$\text{Var}[w_{LS}] = \sigma^2 (X^T X)^{-1}$$

$$\rightarrow \text{As, } w_{RR} = (\lambda I + X^T X)^{-1} X^T y$$

we can write it as

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

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

From this we can calculate the expected value for  $w_{RR}$

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

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



Now, for variance of  $w_{RR}$ ,

$$w_{RR} = (\lambda I + X^T X)^{-1} (X^T X) \cdot w_{LS}$$

$$= \left( (X^T X) (\lambda \cdot (X^T X)^{-1} + I) \right)^{-1} \cdot (X^T X) w_{LS}$$

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

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

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

$$\therefore w_{RR} = Z \cdot w_{LS}$$

$$\text{Var} [w_{RR}] = Z \cdot \text{Var} [w_{LS}] \cdot Z^T$$

$$= Z \cdot \sigma^2 (X^T X)^{-1} \cdot Z^T$$

$$= \sigma^2 \cdot Z (X^T X)^{-1} \cdot Z^T$$

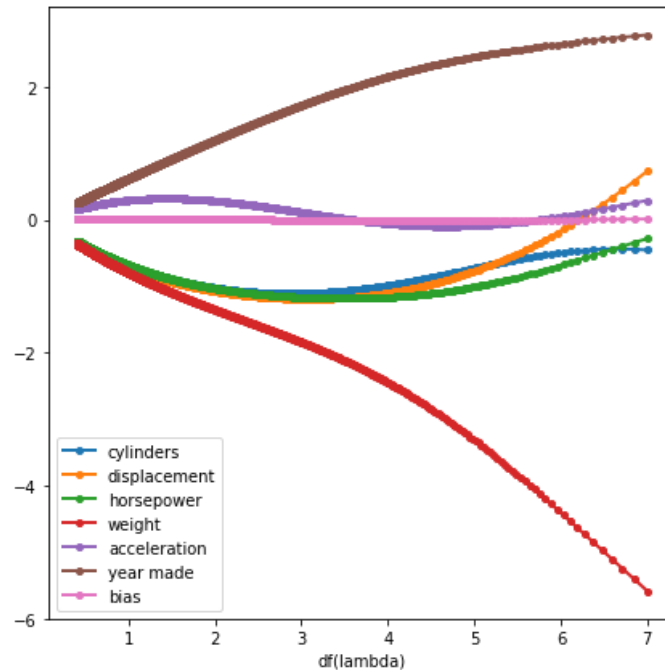
$$\therefore E[w_{RR}] = (\lambda I + X^T X)^{-1} X^T X \omega$$

$$\text{Var} [w_{RR}] = \sigma^2 Z (X^T X)^{-1} Z^T$$

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

### Problem – 3:

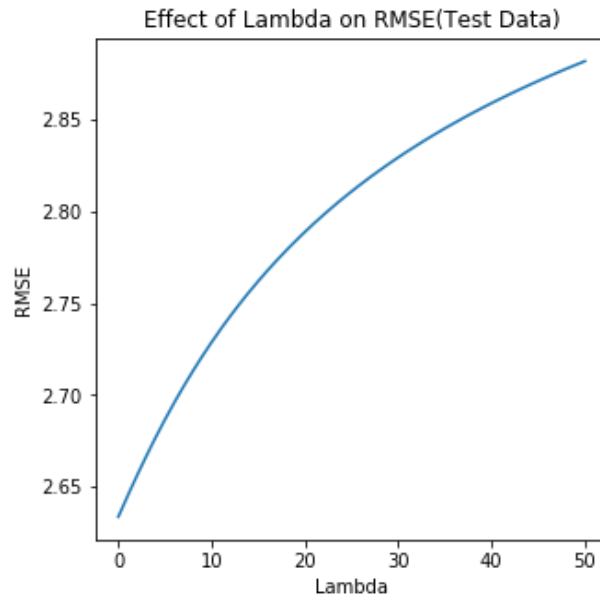
(a) For  $\lambda = 0; 1; 2; 3; \dots; 5000$ , solve for WRR. (Notice that when  $\lambda = 0$ ,  $WRR = WLS$ .) In one figure, plot the 7 values in WRR as a function of  $df(\lambda)$ . You will need to call a built in SVD function to do this as discussed in the slides. Be sure to label your 7 curves by their dimension in  $X$ .



(b) Two dimensions clearly stand out over the others. Which ones are they and what information can we get from this?

- We can clearly see that “weight” and “year made” clearly stand out over others. The dominance of these 2 show the importance of both the values in determining the dependent variable (i.e.) the miles per gallon for that car. If we reduce the degrees of freedom for the model, by increasing the value of  $\lambda$ , at a certain point when  $df(\lambda)$  is nearly 4, we see that the coefficient for acceleration becomes almost zero. The X-axis of the above plot shows the Degree of freedom which essentially is the function of  $\lambda$ .
- Another key thing is that the increase in value of “year made” increases the miles per gallon which follows the intuitive way of understanding that the newer cars would have higher mileage than the older ones where as with increase in “weight” of the car, the mileage always decreases. These 2 things remain the same at any value of  $\lambda$  (as the curves for both the features don’t cross the line for 0 at any point in the plot)

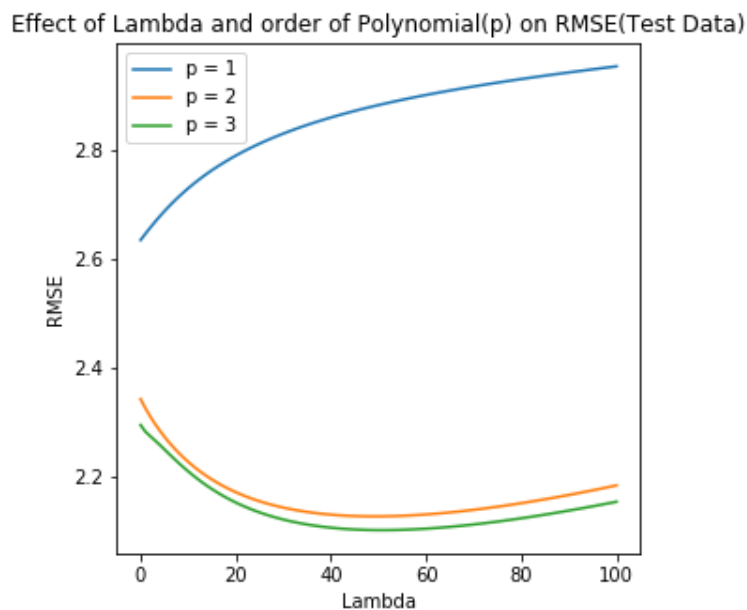
(c) For  $\lambda = 0; \dots; 50$ , predict all 42 test cases. Plot the root mean squared error (RMSE) on the test set as a function of  $\lambda$  —not as a function of  $df(\lambda)$ . What does this figure tell you when choosing for this problem (and when choosing between ridge regression and least squares)?



- From this plot we can see the RMSE value increases with increase in  $\lambda$  in this case. So, the RMSE value for  $\lambda = 1$  is less than RMSE value for  $\lambda = (i+1)$ .
- Also, from the above statement in this case, we can say that the RMSE for  $\lambda = 0$  would be the minimum.
- Thus, in this case we would prefer least square (same as ridge regression with  $\lambda = 0$ ) over the ridge regression.

(d) In one figure, plot the test RMSE as a function of  $\lambda = 0; \dots; 100$  for  $p = 1; 2; 3$ . Based on this plot, which value of  $p$  should you choose and why? How does your assessment of the ideal value of  $\lambda$  change for this problem?

- Below figure shows the plot for RMSE as a function of  $\lambda$  when  $p = 1, 2, 3$ .





- From the plot we can see that
  - RMSE values for  $p = 1$  increases with increase in value of  $\lambda$
  - RMSE values for  $p = 2, 3$  decrease with increase in  $\lambda$  (initially) while they after a certain  $\lambda$  starts on increasing with increase in  $\lambda$ .
- The RMSE values for  $p = 3$  are the minimum and hence we can say that the model with  $p = 3$  perform the best on the test set and hence, we will select the value of  $p$  as 3.
- In the case of  $p = 1$ , as the minimum RMSE was obtained at the point where  $\lambda = 0$ , we would have selected that value of  $\lambda$  had we selected  $p$  to be 1. For  $p = 2$  the minimum is no longer corresponding to the value of  $\lambda = 0$ . Now the minimum value of RMSE can be approximately seen to be obtained at  $\lambda = 42$ . Similar is the case for  $p = 3$ . Hence, with introduction of  $p^{\text{th}}$  order polynomial terms we would require to have need for regularization and hence the ideal value for  $\lambda$  no longer remains 0.