

# ELEN 4720: Machine learning for Signals, Information & Data

## Homework-1

Date - 10<sup>th</sup> October 2021

Name - Rakshith Kamath

UNI - 9K3165

### Problem - 1 -

a)  $(x_1, x_2, x_3, \dots, x_N) \stackrel{\text{iid}}{\sim} \frac{\lambda^x e^{-\lambda}}{x!}$

since each of the observations are independent of each other, the joint probability will be product of the individual probabilities.

$$\therefore p(x_1, x_2, \dots, x_N | \lambda) = \prod_{i=1}^N p(x_i | \lambda)$$

$$p(x_1, x_2, x_3, \dots, x_N | \lambda) = \prod_{i=1}^N \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

b) Maximum likelihood will be to find the value of  $\lambda$  for which the joint probability we derived would be maximum.

This can be defined by a function as follows-

$$\lambda_{ML} = \arg \max_{\lambda} \prod_{i=1}^N \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

Taking logarithm doesn't change the value of  $\lambda_{ML}$  in this case.

$$\begin{aligned} \lambda_{ML} &= \arg \max_{\lambda} \left[ \log \left( \prod_{i=1}^N \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \right) \right] \\ &= \arg \max_{\lambda} \left[ \sum_{i=1}^N x_i \log \lambda - \sum_{i=1}^N \lambda \log e - \sum_{i=1}^N \log(x_i!) \right] \end{aligned}$$

The last term doesn't have a  $\lambda$  term & can be ignored for arg max calculation.

$$\lambda_{ML} = \arg \max_{\lambda} \left[ \log \lambda \underbrace{\sum_{i=1}^N x_i - N\lambda}_{f} \right]$$

consider  $f = \log \lambda \sum_{i=1}^N x_i - N\lambda$

To get  $\lambda_{ML}$ ,  $\frac{\partial f}{\partial \lambda} = 0$

Applying the same,

$$\frac{1}{\lambda_{ML}} \sum_{i=1}^N x_i - N = 0$$

$$\lambda_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$

c)  $\lambda_{MAP}$  uses posterior distribution to get the optimal values of  $\lambda$

It can be defined as follows-

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

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

$\because p(x)$  doesn't depend on  $\lambda$ , it can be ignored in the calculations. Substituting the other values, we get

$$\lambda_{MAP} = \arg \max_{\lambda} \left[ \ln p(x|\lambda) + \ln p(\lambda) \right]$$

$$= \arg \max_{\lambda} \left[ \ln \left[ \prod_{i=1}^N \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \right] + \ln \left( \frac{b^a \lambda^{a-1} e^{-b\lambda}}{r[a]} \right) \right]$$

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

ignoring terms without  $\lambda$  again,

$$= \arg \max_{\lambda} \left\{ \ln(\lambda) \underbrace{\sum_{i=1}^N x_i}_{+} - N\lambda + (a-1) \ln \lambda - b\lambda \right\}$$

To get  $\lambda_{MAP}$ ,  $\frac{\partial f}{\partial \lambda} = 0$

i.e,

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

$$\sum_{i=1}^N x_i + a-1 = N+b$$

---

$$\lambda_{MAP}$$

$$\lambda_{MAP} = \frac{\sum_{i=1}^N x_i + a-1}{N+b}$$

d) The Bayes rule states that

$$p(\lambda|x) = \frac{p(x|\lambda) \cdot p(\lambda)}{p(x)}$$

$$p(\lambda|x) \propto p(x|\lambda) \cdot p(\lambda)$$

$$\propto \frac{\lambda^x e^{-\lambda}}{x!} \cdot \frac{b^a \lambda^{a-1} e^{-b\lambda}}{\Gamma[a]}$$

where  $X = \sum_{i=1}^N x_i$

Ignoring the non- $\lambda$  terms, we get

$$p(\lambda|x) \propto [\lambda^x e^{-\lambda} \lambda^{a-1} e^{-b\lambda}]$$

$$\propto \lambda^{(x+a)-1} e^{-\lambda(b+N)}$$

Let  $x+a$  be  $u$  &  $b+N$  be  $v$

$$p(\lambda|x) \propto \lambda^{u-1} e^{-\lambda t}$$

we can multiply constants to this

$$p(\lambda|x) \propto \frac{u^v \lambda^{u-1} e^{-\lambda t}}{\Gamma[u]}$$

This shows that  $p(\lambda|x)$  is proportional to gamma( $u, v$ ).

i.e.,  $p(\lambda|x)$  has a distribution of gamma  $\left( \sum_{i=1}^n x_i + a, b+N \right)$

c) Since we know the distribution is gamma, we can say its mean & variance as follows

$$E[X] = \frac{a}{b} \quad \& \quad \text{Var}[X] = \frac{a}{b^2} \quad \text{for } X \sim \text{Gamma}(a, b)$$

To get its mean

$$E[\lambda] = \frac{\sum_{i=1}^n x_i + a}{N+b}$$

& Variance

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

$$E[\lambda] = \frac{\sum_{i=1}^N x_i + a}{N+b}$$

$$\lambda_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\lambda_{MAP} = \frac{\sum_{i=1}^N x_i + a - 1}{N+b}$$

$$\lambda_{MAP} = \frac{\sum_{i=1}^N x_i + a - \frac{1}{N+b}}{N+b}$$

$$\lambda_{MAP} = E[\lambda] - \frac{1}{N+b}$$

$$\therefore \lambda_{MAP} < E[\lambda]$$

$$\text{If } E[\lambda] \Big|_{a,b=0} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$E[\lambda] \Big|_{a,b=0} = \lambda_{ML}$$

This says that if  $a$  &  $b$  values are 0, we get  $\lambda_{ML}$ .

## Problem 2 - (20 points)

a) Given -

Data -  $(x_i, y_i), i = 1, 2, \dots, n$ ,  $x \in \mathbb{R}^d$  &  $y \in \mathbb{R}$

Model -  $y_i \stackrel{iid}{\sim} N(x_i^T \omega, \sigma^2)$

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

Mean of  $\omega_{RR}$  ( $E[\omega_{RR}]$ ) -

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

Since in this given model, we assume that we know the  $X$  matrix &  $y$  matrix is the one that is random. we get

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

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

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

$$E[\omega_{RR}] = (\lambda(x^T x)^{-1} + I)^{-1} \omega$$

$$E[\omega_{RR}] = Z \omega$$

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

Variance of  $\omega_{RR}$  ( $\text{Var}(\omega_{RR})$ )

$$\text{Var}(\omega_{RR}) = E[(\omega_{RR} - E[\omega_{RR}]) (\omega_{RR} - E[\omega_{RR}])^T]$$

$$\text{Var}(\omega_{RR}) = E[\omega_{RR} \omega_{RR}^T] - E[\omega_{RR}] E[\omega_{RR}]^T \longrightarrow ①$$

$$\text{w.r.t } \omega_{RR} = (\lambda I + x^T x)^{-1} x^T y = ((x^T x)^{-1} \lambda + I)^{-1} (x^T x)^{-1} x^T y$$

$$\text{& } E[yy^T] = \sigma^2 I + X \omega \omega^T X^T \quad \omega_{RR} = \underline{\underline{Z(x^T x)^{-1} x^T y}}$$

Applying these results to ①

$$\begin{aligned} \text{Var}(\omega_{RR}) &= Z(x^T x)^{-1} x^T E[yy^T] x (x^T x)^{-1} Z^T - Z \omega \omega^T Z^T \\ &= Z(x^T x)^{-1} x^T [I^2 + X \omega \omega^T X^T] x (x^T x)^{-1} Z^T - Z \omega \omega^T Z^T \end{aligned}$$

$$= \sigma^2 \left[ z \left( x^T x \right)^{-1} x^T \cancel{x} \left( x^T \cancel{x} \right)^{-1} z^T \right] + z \left( x^T \cancel{x} \right)^{-1} x^T \cancel{x} w w^T \cancel{x} x \left( x^T \cancel{x} \right)^{-1} z^T - z w w^T z^T$$

$$\text{Var}(x) = \sigma^2 z \left( x^T x \right)^{-1} z^T$$

b)  $\omega_{RR} = (\lambda I + x^T x)^{-1} x^T y \quad , \quad \omega_{LS} = (x^T x)^{-1} x^T y$

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

$$= (\lambda I + x^T x)^{-1} x^T x \omega_{LS}$$

$$= \left[ (x^T x) \left( \lambda (x^T x)^{-1} + I \right) \right]^{-1} x^T x \omega_{LS}$$

$$= \left[ \lambda (x^T x)^{-1} + I \right]^{-1} (x^T x)^{-1} x^T x \omega_{LS}$$

$$\omega_{RR} = \left[ \lambda (x^T x)^{-1} + I \right]^{-1} \omega_{LS}$$

$$\text{w.r.t } X = USV^T$$

$$X^T X = V S U^T U S V^T$$

$$\underline{X^T X = V S^2 V^T}, (X^T X)^{-1} = V S^{-2} V^T$$

Substituting that in the above equation, we get

$$w_{RR} = \left[ \lambda V S^{-2} V^T + I \right]^{-1} w_{LS}$$

$$w_{RR} = V \left[ \lambda S^{-2} + I \right]^{-1} V^T w_{LS}$$

It can be re-written as

$$\underline{\underline{w_{RR} = V M V^T}} w_{LS} \quad \text{, where } M = \frac{s_{ii}^2}{\lambda + s_{ii}^2}$$

### Problem - 3 (Coding) - (30 points)

a)

- (a) For  $\lambda = 0, 1, 2, 3, \dots, 5000$ , solve for  $w_{RR}$ . (Notice that when  $\lambda = 0$ ,  $w_{RR} = w_{LS}$ .) In one figure, plot the 7 values in  $w_{RR}$  as a function of  $d(\lambda)$ . 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  $x$ .<sup>2</sup>

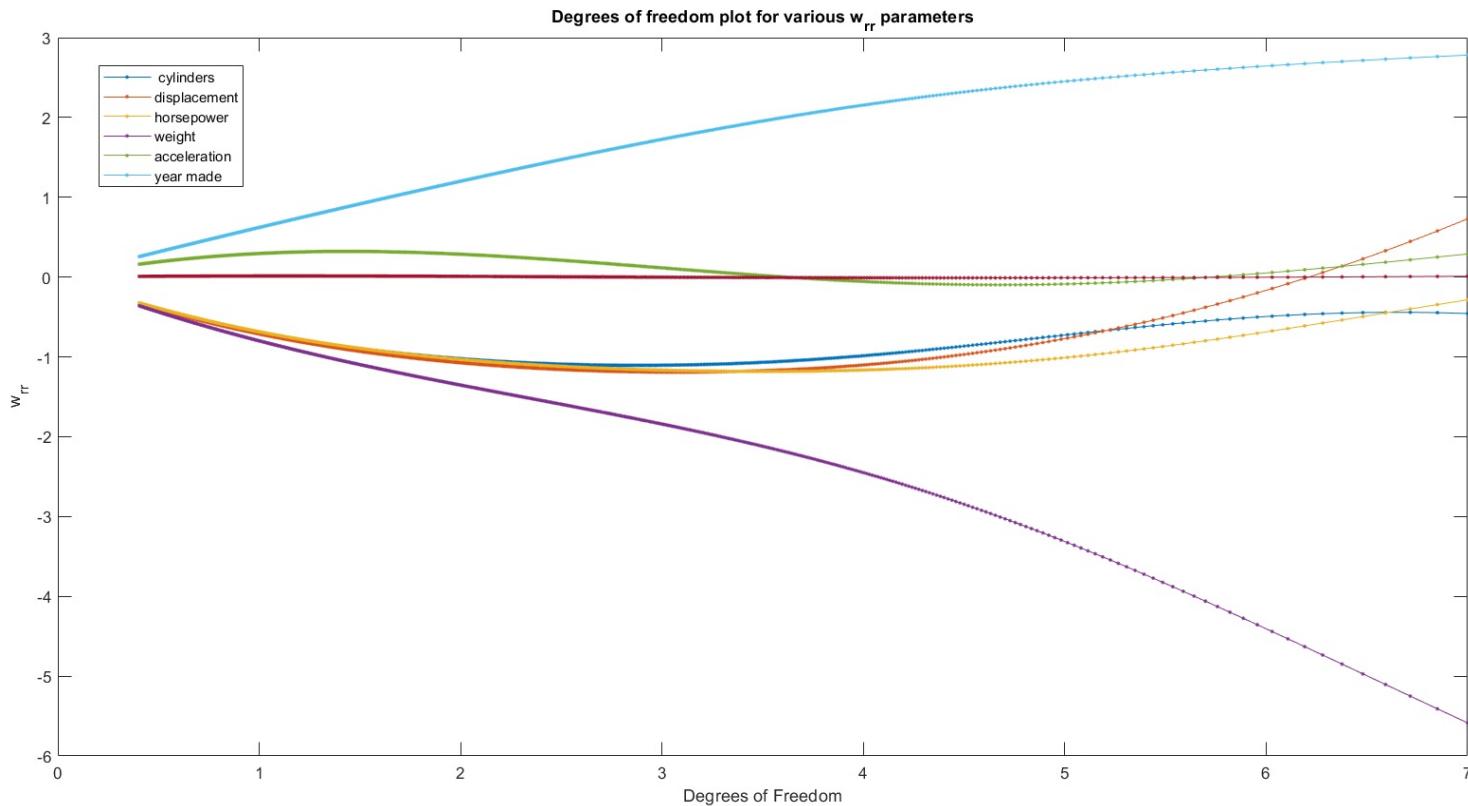


Figure-1

b)

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

We clearly see that 'weight' & 'year-made' parameters have deviated the most, hence are key features for the estimation of miles per gallon of a car.

At  $\lambda=0$ , i.e., the Least squares method the above 2 parameters have the highest absolute

value weights which suggests it's importance. Also, as  $\lambda$  value increases, it's degree of freedom decreases rapidly, thus furthering it's importance.

The information above says that as the 'year made' parameter increases, the miles per gallon increases, which says us that newer cars have higher & better mileage than older cars. The 'weight' parameter of car says that as the weight of the car increases, the miles per gallon decreases & is inversely proportional. This is true for all values of  $\lambda$ .

c)

- (c) For  $\lambda = 0, \dots, 50$ , predict all 42 test cases. Plot the root mean squared error (RMSE)<sup>3</sup> on the test set as a function of  $\lambda$ —not as a function of  $df(\lambda)$ . What does this figure tell you when choosing  $\lambda$  for this problem (and when choosing between ridge regression and least squares)?

For the plot below we see that as the  $\lambda$  values increases the RMSE values increases as well. This shows us that regularization by ridge regression isn't exactly helping the case & the solution obtained by least square, i.e., at  $\lambda=0$ , is the best solution in this particular case.

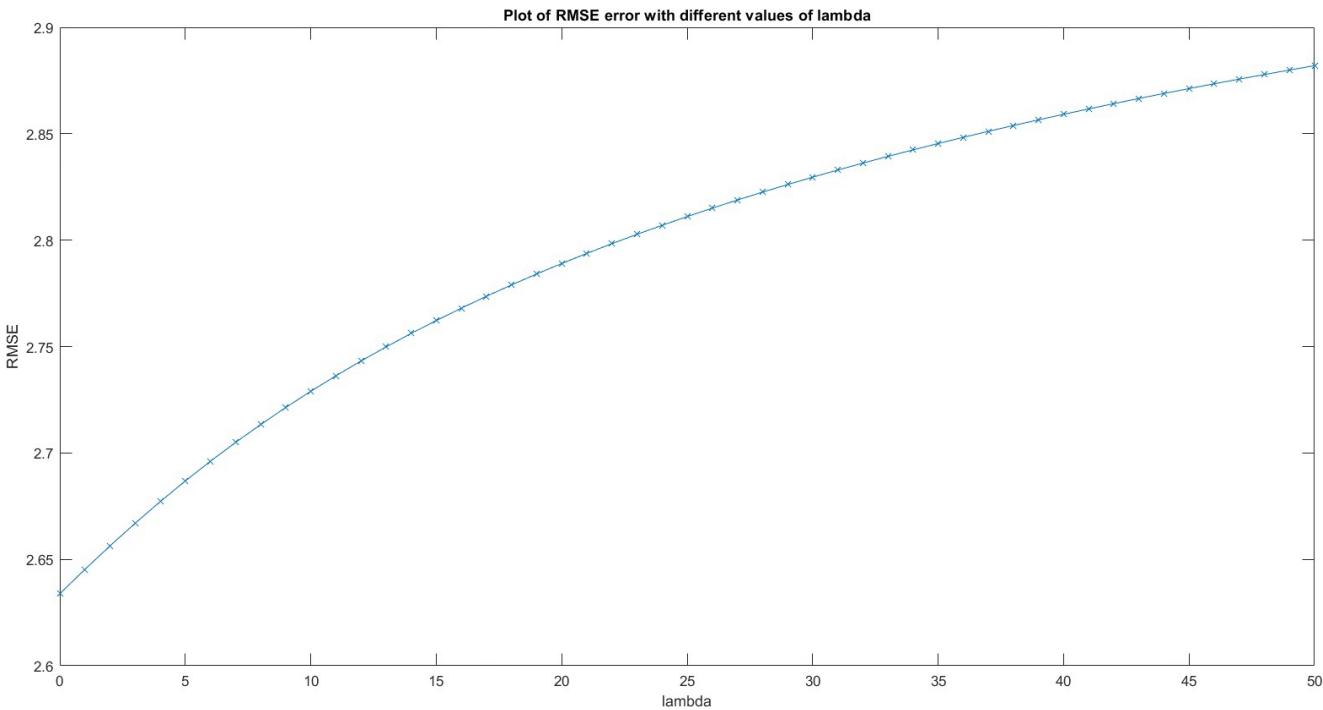


Figure - 2

d)

- (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?

Based on the plot, we see that for  $p=3$  we get the lowest root mean square error. At  $\lambda=51$  &  $p=3$  we get the least error of 2.10.

For  $p=1$ , as  $\lambda$  increases, the RMSE value also increases thus suggesting us that ridge

Regression isn't the way to proceed since the error is getting larger compared to least square method. Hence, the best value for  $\lambda$  is 0.

For  $p=2$ , we see that as  $\lambda$  increases, there's a decrease in the error till  $\lambda=49$ , where we get the least error 2.012, before it starts to increase again. This suggests us that if we are to use  $p=2$  for model estimation, then the corresponding  $\lambda$  value for the best results is  $\lambda=49$ .

Hence, we see that introduction of higher order polynomial terms would require us to regularize the result since  $\lambda \neq 0$ .

