

1 CSE 512 Machine Learning - Homework 2

Name: Rohit Bhal

Solar ID: 112073893

NetID email address: rbhal@cs.stonybrook.edu / rohit.bhal@stonybrook.edu

Q1. Parameter Estimation

1.1 MLE

Ans: (i) We have

$$P(X=k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad k \in \{0, 1, 2, \dots\}$$

The log likelihood function of X given λ is:

$$\begin{aligned} \ln(P(X=k|\lambda)) &= \ln\left(\frac{\lambda^k e^{-\lambda}}{k!}\right) \quad (\because \ln(ab) = \ln(a) + \ln(b)) \\ &= \ln \lambda^k + \ln e^{-\lambda} + -\ln(k!) \quad (\ln e = 1) \end{aligned}$$

$$L(X) = k \ln \lambda - \lambda - \ln(k!)$$

(2) MLE for the given λ can be derived as:

$$\frac{\partial L(X)}{\partial \lambda} = 0$$

$$\text{Now, } L(X=k|\lambda) = \sum_{k=1}^N (k \ln \lambda - \lambda - \ln(k!))$$

$$\therefore \frac{\partial (L(X))}{\partial \lambda} = \sum_{k=1}^N \frac{k}{\lambda} - N = 0$$

$$\lambda = \frac{\sum_{k=1}^N k}{N}$$

$$\lambda = \frac{\sum_{k=1}^N k}{N}$$

$$(3) \quad \lambda = \frac{\sum_{i=1}^n x_i}{n}$$

Trip	1	2	3	4	5	6	7
Delay in minutes	4	5	3	5	6	9	10

$$\begin{aligned} \lambda &= \frac{(4+5+3+5+6+9+10)}{7} \\ &= \frac{42}{7} = 6 \end{aligned}$$

1.2 MAP

$$(1) \quad P(\lambda | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}, \quad \lambda > 0$$

$$P(\lambda | x=k) = \frac{P(\lambda) P(x=k|\lambda)}{P(x=k)}$$

The posterior distribution over λ is:

$$\begin{aligned} P(\lambda | x=k) &= \prod_{i=1}^n \frac{P(\lambda) P(x=k_i|\lambda)}{P(x=k_i)} \\ &= \prod_{i=1}^n \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \frac{\left(\lambda^{k_i} e^{-\lambda}\right)}{(k_i!)} \\ &= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{e^{(-\beta\lambda - \lambda \sum_{i=1}^n k_i)}}{\prod_{i=1}^n (k_i!)} \end{aligned}$$

$$P(\lambda | x=k) = \frac{\beta^\lambda}{\Gamma(k)} \frac{e^{-\lambda(\beta+n)}}{\prod_{i=1}^n (k_i!)} \left((\lambda-1) + \sum_{i=1}^n k_i \right)$$

(2) MAP estimate for λ is:

$$\begin{aligned} \frac{\partial \ln(P(\lambda | x=k))}{\partial \lambda} &= 0 \\ &= \frac{\partial \ln\left(\frac{\beta^\lambda}{\Gamma(k)} \frac{e^{-\lambda(\beta+n)}}{\prod_{i=1}^n (k_i!)} \left((\lambda-1) + \sum_{i=1}^n k_i \right)\right)}{\partial \lambda} = 0 \\ &= \frac{\partial \left(\ln \beta^\lambda + \ln e^{-\lambda(\beta+n)} + \ln \lambda^{(\lambda-1) + \sum_{i=1}^n k_i} - \ln \left(\Gamma(k) \prod_{i=1}^n (k_i!) \right) \right)}{\partial \lambda} = 0 \\ &= -(\beta+n) + \frac{(\lambda-1 + \sum_{i=1}^n k_i)}{\lambda} = 0 \\ \lambda(\beta+n) &= (\lambda-1) + \sum_{i=1}^n k_i \\ \therefore \lambda &= \frac{((\lambda-1) + \sum_{i=1}^n k_i)}{(\beta+n)} \end{aligned}$$

1.3 Estimator Bias

$$(1) \quad \eta = e^{-2\lambda}$$

$$\ln \eta = (-2\lambda) \ln e$$

$$-\frac{1}{2}(\ln \eta) = \lambda$$

$$P(\lambda) = \frac{\lambda^k e^{-\lambda}}{k!} = \frac{(-\frac{1}{2} \ln n)^k e^{-(-\frac{1}{2} \ln n)}}{k!}$$

$$\begin{aligned} \text{Now, } \ln P(\lambda) &= \ln \left(-\frac{1}{2} \ln n \right)^k + \ln e^{\frac{1}{2} \log n} - \ln k! \\ &= k \left(\ln \ln n - \ln 2 \right) + \frac{1}{2} \ln n - \ln k! \end{aligned}$$

$$\begin{aligned} \frac{\partial \ln P(\lambda)}{\partial \lambda} &= 0 \\ &= k \left(\frac{1}{-\ln n} - \frac{1}{n} \right) + \frac{1}{2n} = 0 \end{aligned}$$

$$\frac{-1}{2n} = \frac{k}{\eta^* \ln n}$$

$$\begin{aligned} \ln n &= -2k \\ \eta &= e^{-2k} \quad (k=x) \end{aligned}$$

$$\hat{\eta} = e^{-2k}$$

$\Rightarrow \hat{\eta}$ is the maximum likelihood estimate of λ .

(2) Bias of $\hat{\eta}$ = Expected Value - True Value

$$Ex(\hat{\eta}) = \sum P(\hat{\eta}) \hat{\eta}$$

$$\hat{\eta} = e^{-2x}$$

$$\begin{aligned} Ex(\hat{\eta}) &= \sum_{i=0}^{\infty} e^{-2x_i} \frac{e^{-\lambda} \lambda^{x_i}}{(x_i!)} \\ &= e^{-\lambda} \sum_{i=0}^{\infty} \frac{(e^{-2\lambda})^{x_i}}{(x_i!)} \\ &= e^{-\lambda} \cdot e^{-2\lambda} \\ &= e^{-\lambda(1/e^2 - 1)} = e^{\lambda(1/e^2 - 1)} \\ &= e^{\lambda(\frac{1}{e^2} - 1) - 2\lambda} \end{aligned}$$

$$\begin{aligned} \text{Bias of } \hat{\eta} &= e^{-\lambda(1/e^2 - 1) - 2\lambda} - e^{-\lambda} \\ &= \left(e^{-\lambda(1 - \frac{1}{e^2})} - e^{-2\lambda} \right) \end{aligned}$$

(3) Expectation of an unbiased estimator $f(x)$ is equal to the true mean.

$$E(f(x)) = e^{-2\lambda}$$

$$\sum_{i=0}^{\infty} f(x) \frac{\lambda^{x_i} e^{-\lambda}}{(x_i!)} = e^{-2\lambda}$$

$$\sum_{k=0}^{\infty} f(x) \frac{x^k}{(k!)^{-\lambda}} = e^{-\lambda}$$

Now, we know that

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

$$\sum_{k=0}^{\infty} f(x) \frac{x^k}{(k!)^{-\lambda}} = \sum_{k=0}^{\infty} \frac{(-\lambda)^k}{k!}$$

Now, Both RHS & LHS can be equal iff

$$f(x) = (-1)^x$$

$$\Rightarrow f(x) = (-1)^x.$$

$$\text{for } x=100, f(x)=1$$

$$x=101, f(x)=-1.$$

Now, we can see that the $f(x)$ is a bad estimator of η for the above x values.
The deviation from true value is very large.

2. Ridge Regression & LOOCV

$$\underset{\omega, b}{\text{minimize}} \quad \lambda \|\omega\|^2 + \sum_{i=1}^n (\omega^T x_i + b - y_i)^2$$

Ams (2.1) $\bar{\omega} = [\omega; b] \quad \bar{x} = [x; I_n^T], \quad \bar{I} = [I_k, 0_k; 0_k^T, 0]$

$$C = (\bar{x} \bar{x}^T + \lambda \bar{I}), \quad \bar{d} = \bar{x} y$$

To Prove:

$$\bar{\omega} = C^{-1} \bar{d}$$

$$\begin{aligned} L(\omega) &= \lambda \|\omega\|^2 + \sum_{i=1}^n (\omega^T x_i + b - y_i)^2 \\ &= \lambda \omega^T \omega + \sum_{i=1}^n (\omega^T x_i + b - y_i)^2 \end{aligned}$$

$$\begin{aligned} \cancel{\frac{\partial L(\omega)}{\partial \omega}} &= (\lambda + \lambda) \omega + 2 \sum_{i=1}^n (\bar{\omega}^T \bar{x}_i + b - y_i) \bar{x}_i = 0 \\ - \cancel{\lambda \omega} &= \cancel{\lambda} \sum_{i=1}^n (\bar{\omega}^T \bar{x}_i + b - y_i) \bar{x}_i^T \\ - \lambda \omega &= \omega^T \bar{x} \end{aligned}$$

$$L(\omega) = \lambda \bar{\omega}^T \bar{\omega} + \sum_{i=1}^n (\bar{\omega}^T \bar{x}_i^T + b - y_i)^2$$

$$\begin{aligned} \cancel{\frac{\partial L(\omega)}{\partial \omega}} &= (\lambda + \lambda) \bar{\omega} + 2 \sum_{i=1}^n (\bar{\omega}^T \bar{x}_i^T + b - y_i) \bar{x}_i = 0 \\ - \cancel{\lambda \bar{\omega}} &= \cancel{\lambda} \sum_{i=1}^n (\bar{\omega}^T \bar{x}_i^T + b - y_i) \bar{x}_i^T \\ - \lambda \bar{\omega} &= \bar{\omega}^T \bar{x}^T \bar{x} - \bar{x} y \\ \bar{\omega} (\bar{x} \bar{x}^T + \lambda \bar{I}) &= \bar{x} y \end{aligned}$$

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

$$\therefore \bar{w} = C^T d$$

$$(2.2) \quad C = (\bar{X}\bar{X}^T + \lambda\bar{I})$$

$$\text{Now, } \bar{X}\bar{X}^T = \sum_{i=1}^n \sum_{j=1}^m x_{ij} x_{ji} \quad \text{for } n \times m \text{ matrix}$$

If we remove x_k row from \bar{X} , we get

$$\begin{aligned} \bar{X}_{-k}\bar{X}_{-k}^T &= \sum_{\substack{i=1 \\ i \neq k}}^n \sum_{\substack{j=1 \\ j \neq k}}^m x_{ij} x_{ji} \\ &= \sum_{i=1}^n \sum_{j=1}^m x_{ij} x_{ji} - \sum_{i=k} \sum_{j=k} x_{ij} x_{ji} \\ &= (\bar{X}\bar{X}^T - x_k x_k^T) \end{aligned}$$

$$C_{(i)} = (\bar{X}_{-i}\bar{X}_{-i}^T + \lambda\bar{I}) \quad (x_i \text{ is the data row removed from } \bar{X})$$

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

$$C_{(i)} = (C - x_i x_i^T)$$

Similarly, $d = \bar{X}y$

$$d_{(i)} = (\bar{X}y - x_i y_i)$$

$$2.3 \quad C_{(i)} = (C - x_i x_i^T)$$

$$C_{(i)}^{-1} = (C - x_i x_i^T)^{-1}$$

Now, we know that :

$$(A + uv^T)^{-1} = A^{-1} - \frac{(A^{-1}uv^TA^{-1})}{(1+v^TA^{-1}u)} \quad (\text{Sherman Morrison formula})$$

We can see that,

$$A = C, \quad u = -x_i, \quad v^T = x_i^T$$

$$\therefore (C - x_i x_i^T)^{-1} = C^{-1} - \left(\frac{C^{-1}(-x_i) x_i^T C^{-1}}{1 + x_i^T C^{-1}(-x_i)} \right)$$

$$C_{(i)}^{-1} = C^{-1} + \left(\frac{(C^{-1} x_i x_i^T C^{-1})}{1 - (x_i^T C^{-1} x_i)} \right)$$

2.4 To Prove :

$$\bar{\omega}_{(i)} = \bar{\omega} + (C^{-1} \bar{x}_i) \frac{(-y_i + \bar{x}_i^T \bar{\omega})}{(1 - \bar{x}_i^T C^{-1} \bar{x}_i)}$$

$$\bar{\omega}_{(i)} = C_{(i)}^{-1} d_{(i)} = \left[C^{-1} + \left(\frac{(C^{-1} \bar{x}_i \bar{x}_i^T C^{-1})}{1 - (\bar{x}_i^T C^{-1} \bar{x}_i)} \right) \right] (d - \bar{x}_i^T y_i)$$

$$\begin{aligned}
&= C^{-1}d - C^{-1}\bar{x}_i y_i + \frac{(C^{-1}\bar{x}_i \bar{x}_i^T C^{-1})d}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)} - \frac{(C^{-1}\bar{x}_i \bar{x}_i^T C^{-1})(\bar{x}_i^T y_i)}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)} \\
&= \bar{\omega} + C^{-1}\bar{x}_i \left(-y_i + \frac{(C^{-1}\bar{x}_i \bar{x}_i^T C^{-1})d}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)} - \frac{(\bar{x}_i^T C^{-1}\bar{x}_i y_i)}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)} \right) \\
&= \bar{\omega} + C^{-1}\bar{x}_i \underbrace{\left(-y_i + y_i \cancel{\bar{x}_i^T C^{-1}\bar{x}_i} + \bar{x}_i^T C^{-1}d - \cancel{\bar{x}_i^T C^{-1}\bar{x}_i y_i} \right)}_{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)}
\end{aligned}$$

$$\bar{\omega}_{(i)} = \bar{\omega} + C^{-1}\bar{x}_i \frac{(-y_i + \bar{x}_i^T \bar{\omega})}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)}$$

Indence proved.

2.5 To Prove

$$\omega_{(i)}^T \bar{x}_i - y_i = \frac{(\bar{\omega}^T \bar{x}_i - y_i)}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)}$$

$$\begin{aligned}
\omega_{(i)}^T \bar{x}_i - y_i &= \left[\bar{\omega} + C^{-1}\bar{x}_i \frac{(-y_i + \bar{x}_i^T \bar{\omega})}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)} \right] \bar{x}_i - y_i \\
&= \left[\frac{\bar{\omega} \bar{x}_i - \bar{\omega} \bar{x}_i^T C^{-1}\bar{x}_i \bar{x}_i + (-C^{-1}\bar{x}_i y_i) \bar{x}_i + C^{-1}\bar{x}_i \bar{x}_i^T \bar{\omega} \bar{x}_i}{(1 - \bar{x}_i^T C^{-1}\bar{x}_i)} \right] - y_i
\end{aligned}$$

$$\begin{aligned}
 &= \frac{(\bar{\omega}^T \bar{x}_i - \bar{\omega}^T C^T \bar{x}_i \bar{x}_i - C^T \bar{x}_i y_i + C^T \bar{x}_i \bar{x}_i^T \bar{w} \bar{x}_i - y_i + \bar{x}_i^T C^T \bar{x}_i y_i)}{(1 - \bar{x}_i^T C^T \bar{x}_i)} \\
 &= \frac{(\bar{\omega}^T \bar{x}_i - y_i)}{(1 - \bar{x}_i^T C^T \bar{x}_i)}
 \end{aligned}$$

Hence, proved

2.6

$$LOOCV : \sum_{i=1}^n (\bar{\omega}_{(i)}^T \bar{x}_i - y_i)^2$$

$$(\bar{\omega}_{(i)}^T \bar{x}_i - y_i) = \frac{(\bar{\omega}^T \bar{x}_i - y_i)}{(1 - \bar{x}_i^T C^T \bar{x}_i)}$$

C is a matrix of (k, k)

\therefore To calculate inverse of C , we need $O(k^3)$ time.

Also, to calculate $\bar{\omega}^T \bar{x}_i$, we need $O(k)$ time.

To calculate $\bar{x}_i^T C^T \bar{x}_i$, we need $O(k^2)$ time, but this is calculated n times.

Maximum Time Complexity = $\text{Max}(O(k^3), O(nk^2))$

If we use the naive way of calculating LOOCV error, we will be calculating C^{-1} for n time.

$$\text{Time Complexity} = O(nk^3).$$

3. Programming

3.2

The $\lambda = 0.7758246$ achieves the best performance.

The objective value for this λ is:

$$21217.73885524$$

The regularization term ($\lambda \|w\|_2$) is:

$$10127.32674793$$

The RMSE for the above λ is using LOOCV

$$\approx 1.959509042$$

The sum of square of errors is:

$$11090.4121$$

The root mean square error for the data:

$$1.053110$$

3. The top 10 most important features & the top 10 least important features are very hard to define.

If I have to choose, I would choose the features according to the significance it has to the obj. function.

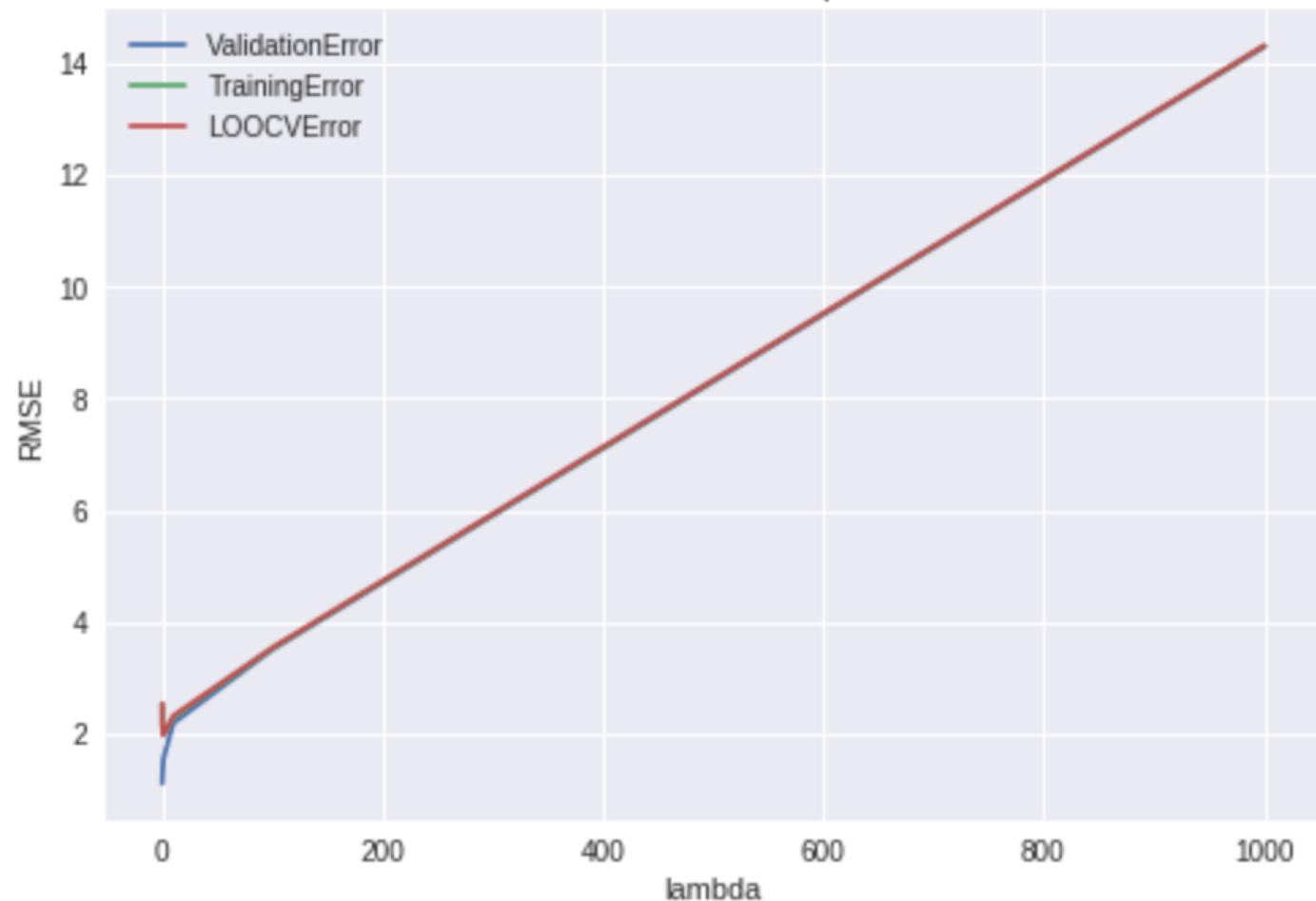
The top 10 most important features :

- | | |
|---------------|---------------|
| 1. Age Years | 5. Area |
| 2. Wine Price | 7. Region |
| 3. Rare | 8. Versatile |
| 4. Delightful | 9. Affordable |
| 5. Pleasing | 10. Smells |

Top 10 least impt. features:

- | | |
|--------------------|------------------|
| 1. tar | 6. Sonoma County |
| 2. Ava | 7. Tiny |
| 3. Pineapple Pears | 8. Oil |
| 4. Grape | 9. Home |
| 5. Leaner | 10. Rocks |

Lambda Vs Root Mean Square Error



Lambda Vs Root Mean Square Error

