

Cover page for answers.pdf
CSE512 Fall 2018 - Machine Learning - Homework 2

Your Name: SAIF SULEMAN VAZIR

Solar ID: 112072061

NetID email address: saifsuleman.vazir@stonybrook.edu

Names of people whom you discussed the homework with: .

HOMEWORK 2

Question 1.1:

$X = (X_1, \dots, X_n)$, X_i = delay in minutes

Solⁿ:

i) To find log likelihood function of X given λ

$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmax}} \prod_{k_i=1}^n P(x_i = k_i | \lambda)$$

$$= \prod_{k_i=1}^n \frac{\lambda^{k_i} e^{-\lambda}}{k_i!}$$

\Rightarrow Taking log(natural) on both sides we get:

$$\log \text{ likelihood} = \ln \left[\prod_{k_i=1}^n \frac{\lambda^{k_i} e^{-\lambda}}{k_i!} \right]$$

$$= \ln \left[\frac{\lambda^{k_1}}{k_1!} \cdot \frac{\lambda^{k_2}}{k_2!} \dots \frac{\lambda^{k_n}}{k_n!} e^{-n\lambda} \right]$$

$$= \ln e^{-n\lambda} + \ln \frac{\lambda^{k_1}}{k_1!} + \ln \frac{\lambda^{k_2}}{k_2!} \dots \ln \frac{\lambda^{k_n}}{k_n!}$$

$$\log \text{ likelihood} = -n\lambda + [k_1 \ln \lambda + k_2 \ln \lambda \dots k_n \ln \lambda - \ln k_1 k_2 \dots k_n]$$

$$= -n\lambda + [(k_1 + k_2 + \dots + k_n) \ln \lambda - \ln k_1 k_2 \dots k_n]$$

ii) For MLE:

We differentiate the above eqⁿ and set it to 0.

$$\text{MLE} \Rightarrow \frac{\partial}{\partial \lambda} [-n\lambda + (k_1 + k_2 + \dots + k_n) \ln \lambda - \ln k_1 k_2 \dots k_n] = 0$$

$$\Rightarrow -n + \frac{k_1 + k_2 + \dots + k_n}{\lambda} = 0$$

$$\therefore \hat{\lambda} = \frac{\sum_{i=1}^n k_i}{n}$$

iii) MLE for λ using observed x :

$$\hat{\lambda} = \frac{1}{7} [4 + 5 + 3 + 5 + 6 + 9 + 10]$$

$$\therefore \hat{\lambda} = \frac{1}{7} \times 42 = \underline{\underline{6}}$$

(Question 1.2:

i) For posterior distribution, we know that:

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

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

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

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

$$\propto [\text{const.}] \lambda^{[n \cdot X_{\text{mean}} + \alpha - 1]} e^{-\lambda [n + \beta]}$$

$$\propto \text{Gamma}[n X_{\text{mean}} + \alpha - 1, n + \beta]$$

Taking log we get:

$$\begin{aligned} \log p(\lambda | X) &\propto \log C + \log \lambda^{[n \cdot X_{\text{mean}} + \alpha - 1]} + \log e^{-\lambda [n + \beta]} \\ &= \log C + [n \cdot X_{\text{mean}} + \alpha - 1] \log \lambda - \lambda [n + \beta] \end{aligned}$$

ii) For MAP:

we differentiate above eqⁿ & set it to 0.

$$\therefore \text{MAP} \Rightarrow \frac{\partial}{\partial \lambda} [\log C + (n \cdot X_{\text{mean}} + \alpha - 1) \log \lambda - \lambda [n + \beta]] = 0$$

$$= \frac{n \cdot X_{\text{mean}} + \alpha - 1}{\lambda} - [n + \beta] = 0$$

$$\therefore \text{MAP} \Rightarrow \hat{\lambda}_{\text{map}} = \frac{n \cdot X_{\text{mean}} + \alpha - 1}{n + \beta}$$

Question 1.3:
 i) Let $h = e^{-2\lambda}$, To show that $\hat{h} = e^{-2X}$ is MLE of h

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

$$h = e^{-2\lambda} \quad \ln h = -2\lambda$$

$$\lambda = -\frac{1}{2} \ln h$$

$$\lambda = \ln \frac{1}{\sqrt{h}}$$

$$P(X|h) = \frac{(-\frac{1}{2} \ln h)^x e^{-(-\frac{1}{2} \ln h)}}{x!} \quad (\ln = \text{natural log})$$

$$= \frac{(-\frac{1}{2} \ln h)^x}{x!} h^{1/2}$$

$$\ln P(X|h) = \ln \left[\frac{(-\frac{1}{2} \ln h)^x}{x!} h^{1/2} \right]$$

$$= \frac{1}{2} \ln h - \ln x! + \ln \left[\left(\frac{1}{2} \ln h \right)^x \right]$$

$$= \frac{1}{2} \ln h - \ln x! + x \ln \left(\ln \frac{1}{\sqrt{h}} \right)$$

Differentiating & equating to 0.

$$\text{MLE} \Rightarrow \frac{\partial}{\partial h} \left[\frac{1}{2} \ln h - \ln x! + x \ln \ln \frac{1}{\sqrt{h}} \right] = 0$$

$$\Rightarrow \frac{1}{2h} - 0 + x \left[\frac{1}{\ln h^{-1/2}} \cdot \frac{1}{h^{-1/2}} \cdot -\frac{1}{2} h^{-3/2} \right] = 0$$

$$\frac{1}{2h} = \frac{x}{2} \left[\frac{1}{\ln h^{-1/2}} \cdot h^{1/2} \cdot \frac{1}{h^{3/2}} \right]$$

$$\frac{1}{2h} = \frac{x}{2} \cdot \frac{1}{\ln h^{-1/2}} \cdot \frac{1}{h}$$

$$\ln h^{-1/2} = x \Rightarrow \ln h = -2x$$

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

Hence proved.

ii) Bias = Estimated value - True value

$$\text{bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

$$\text{bias}(\hat{\lambda}) = E(\hat{\lambda}) - \lambda$$

$$= E\left[e^{-2X}\right] - e^{-2\lambda}$$

$$\therefore \text{bias} = \left(\sum_x e^{-2x} \cdot \frac{\lambda^x e^{-\lambda}}{x!} \right) - e^{-2\lambda} \quad [\text{using LOTUS}]$$

$$[E(\hat{\lambda})] =$$

$$\text{bias} = e^{-\lambda} \left[\sum_x \frac{(\lambda/e^2)^x}{x!} \right] - e^{-2\lambda}$$

we know that: $e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$

$$\therefore \text{Bias} = e^{-\lambda} \left[e^{-\lambda/e^2} \right] - e^{-2\lambda}$$

$$= e^{\lambda/e^2 - \lambda} - e^{-2\lambda}$$

$$\text{Bias} = e^{-\lambda} [1 - 1/e^2] - e^{-2\lambda}$$

iii) For unbiased estimate, let our $\hat{\theta}$ be m^X

For unbiased estimate:

Expected value = true value

$$E(m^X) = e^{-\lambda} \quad E(m^X) = e^{-2\lambda}$$

$$\sum_x m^x \frac{\lambda^x e^{-\lambda}}{x!} = e^{-2\lambda} \quad \sum_x m^x \frac{\lambda^x e^{-\lambda}}{x!} = e^{-2\lambda}$$

$$e^{-\lambda} \left[\sum_x \frac{(m\lambda)^x}{x!} \right] = e^{-2\lambda}$$

$$e^{m\lambda} = e^{-\lambda}$$

$$m\lambda = -\lambda$$

$$m = -1$$

$(-1)^X$ is the only unbiased estimate of λ .

Question 2:

$$2.1) \underset{w, b}{\text{minimize}} \quad \lambda \|w\|^2 + \sum_{i=1}^n (w^T x_i + b - y_i)^2$$

$$\text{let } \bar{w} = [w; b] \quad \& \quad \bar{x} = [X; 1_n^T]$$

$$\therefore \text{ we can write above eq}^n \text{ as:}$$

$$\min_{\bar{w}} \quad \lambda [\bar{w}^T \cdot \bar{w} - b^2] + \sum_{i=1}^n (\bar{w}^T x_i - y_i)^2$$

we differentiate the above eqⁿ & set it to 0.

$$\frac{\partial}{\partial \bar{w}} [\lambda (\bar{w}^T \cdot \bar{w} - b^2) + \sum_{i=1}^n (x_i^T \bar{w} - y_i)^2] = 0$$

[transpose of scalar is equal to the original & $(\bar{w}^T x_i - y_i)$ is a scalar]

$$\therefore \lambda [2 \bar{w}] + \frac{\partial}{\partial \bar{w}} [\| \bar{X}^T \cdot \bar{w} - y \|^2] = 0$$

$$\therefore 2\lambda \bar{w} + \frac{\partial}{\partial \bar{w}} [(\bar{X}^T \bar{w} - y)^T (\bar{X}^T \bar{w} - y)] = 0$$

$$\therefore 2\lambda \bar{w} + \frac{\partial}{\partial \bar{w}} [(\bar{w}^T \bar{X} - y^T)(\bar{X}^T \bar{w} - y)] = 0$$

$$2\lambda \bar{w} + \frac{\partial}{\partial \bar{w}} [\bar{w}^T \bar{X} \bar{X}^T \bar{w} - \bar{w}^T \bar{X} y - y^T \bar{X}^T \bar{w} + y^T y] = 0$$

[scalar, hence take transpose]

$$2\lambda \bar{w} + \frac{\partial}{\partial \bar{w}} [\bar{w}^T \bar{X} \bar{X}^T \bar{w} - 2y^T \bar{X}^T \bar{w} + y^T y] = 0$$

$$2\lambda \bar{w} + [\bar{X} \bar{X}^T + (\bar{X} \bar{X}^T)^T] \cdot \bar{w} - 2\bar{X} y = 0$$

$$\cancel{2\lambda \bar{w}} + \cancel{2\bar{X} \bar{X}^T} \cdot \bar{w} = \cancel{2\bar{X} y}$$

$$[\bar{X} \bar{X}^T + \lambda \bar{I}] \bar{w} = \bar{X} y$$

$$\therefore \bar{w} = \underbrace{[\bar{X} \bar{X}^T + \lambda \bar{I}]^{-1}}_C \cdot \underbrace{\bar{X} y}_d$$

$$\bar{w} = C^{-1} \cdot d$$

Hence proved

$$C = \underbrace{\bar{X}}_{(k+1) \times n} \underbrace{\bar{X}^T}_{n \times (k+1)} + \lambda \underbrace{\bar{I}}_{(k+1) \times (k+1)}$$

∴ initial dimensions of $C \Rightarrow (k+1) \times (k+1)$

even when we remove one column (say x_i)

$$\begin{aligned} C &= \bar{X}_i \bar{X}_i^T + \lambda \bar{I}_i \\ &= \underbrace{(k+1) \times (n-1)}_{(k+1) \times (k+1)} \underbrace{(n-1) \times (k+1)}_{(k+1) \times (k+1)} \Rightarrow (k+1) \times (k+1) \\ &\Rightarrow (k+1) \times (k+1) \end{aligned}$$

dimensions remain same

When we see the matrix multiplications:

$$\begin{bmatrix} x_1 & \dots & x_i & \dots & x_n \\ x_{12} & & x_{i2} & & x_{n2} \\ \vdots & & \vdots & & \vdots \\ x_{1k+1} & & x_{ik+1} & & x_{nk+1} \end{bmatrix} \cdot \begin{bmatrix} x_1 & \dots & x_{1k+1} \\ \vdots & & \vdots \\ x_{1n} & \dots & x_{1k+1} \end{bmatrix}$$

or terms of column vectors:

$$\begin{bmatrix} x_1 & x_2 & \dots & x_i & \dots & x_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{bmatrix}$$

∴ When we remove x_i , only terms missing from multiplication is $x_i \cdot x_i^T$

$$C(i) = C - \bar{x}_i \bar{x}_i^T$$

similarly for d:

$$\begin{aligned} d &= \bar{X}^T y \\ &= \begin{bmatrix} x_1 & x_2 & \dots & x_i & \dots & x_n \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} \end{aligned}$$

when we remove x_i , only term missing is $x_i \cdot y_i$

$$d(i) = d - \bar{x}_i y_i$$

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

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

$$\therefore C_{(i)} = (\bar{X}\bar{X}^T - \bar{x}_i\bar{x}_i^T + \lambda \bar{I})^{-1}$$

$$\therefore C_{(i)} = [C - \bar{x}_i\bar{x}_i^T]^{-1}$$

$$\therefore C_{(i)}^{-1} = [C - \bar{x}_i\bar{x}_i^T]^{-1}$$

By Sherman morrison formula:

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}u \cdot v^T A^{-1}}{1 + v^T A^{-1}u}$$

Substitute $A = C$ $u = (-\bar{x}_i)$ $v^T = \bar{x}_i^T$

$$\therefore C_{(i)}^{-1} = C^{-1} - \frac{C^{-1}(-\bar{x}_i)(\bar{x}_i^T)C^{-1}}{1 + (\bar{x}_i^T)C^{-1}(-\bar{x}_i)}$$

$$\therefore C_{(i)}^{-1} = C^{-1} + \frac{C^{-1}\bar{x}_i\bar{x}_i^T C^{-1}}{1 - \bar{x}_i^T C^{-1}\bar{x}_i}$$

iv) Now, $\bar{w} = C^{-1}d$

$$\therefore \bar{w}_{(i)} = C_{(i)}^{-1}d_{(i)}$$

$$= \left[C^{-1} + \frac{C^{-1}\bar{x}_i\bar{x}_i^T C^{-1}}{1 - \bar{x}_i^T C^{-1}\bar{x}_i} \right] [d - \bar{x}_i y_i]$$

$$\therefore \bar{w}_{(i)} = \underbrace{C^{-1}d}_{\bar{w}} + \frac{(C^{-1}\bar{x}_i\bar{x}_i^T C^{-1})d}{1 - \bar{x}_i^T C^{-1}\bar{x}_i} - (C^{-1}) \cdot \bar{x}_i y_i$$

$$= \frac{(C^{-1}\bar{x}_i\bar{x}_i^T C^{-1}) \cdot \bar{x}_i y_i}{1 - \bar{x}_i^T C^{-1}\bar{x}_i} \rightarrow \bar{w}$$

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

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

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

[y_i can be placed anywhere as it is scalar]

$$\therefore \bar{w}(x_i) = \bar{w} + (C^{-1} \bar{x}_i) \left[\frac{-y_i^0 + \bar{x}_i^0^T \bar{w}}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} \right]$$

v) $(\bar{w}^T \bar{x}_i - y_i) \Rightarrow$ To calculate this we take $\frac{1}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}$ as λ because it is a scalar value

$$\therefore \bar{w}^T \bar{x}_i - y_i = (\bar{w} + (C^{-1} \bar{x}_i) \lambda)^T \bar{x}_i - y_i$$

$$= (\bar{w}^T + \lambda^T (C^{-1} \bar{x}_i)^T) \bar{x}_i - y_i$$

$$= \bar{w}^T \bar{x}_i - y_i + \lambda^T (\bar{x}_i^T C^{-1}) \bar{x}_i$$

$$\Rightarrow \bar{w}^T \bar{x}_i - y_i + \left[\frac{-y_i^0 + \bar{x}_i^0^T \bar{w}}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} \right]^T (\bar{x}_i^T C^{-1} \bar{x}_i)$$

$$= \bar{w}^T \bar{x}_i - y_i + \frac{-y_i^0 + \bar{w}^T \bar{x}_i}{1 - \bar{x}_i^T C^{-1} \bar{x}_i} (\bar{x}_i^T C^{-1} \bar{x}_i)$$

$$= \left[\bar{w}^T \bar{x}_i - (\bar{w}^T \bar{x}_i \bar{x}_i^T C^{-1} \bar{x}_i) - y_i^0 + y_i^0 \bar{x}_i^T C^{-1} \bar{x}_i + \bar{w}^T \bar{x}_i \bar{x}_i^T C^{-1} \bar{x}_i - y_i^0 \bar{x}_i^T C^{-1} \bar{x}_i \right] \frac{1}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}$$

$$= \frac{\bar{w}^T \bar{x}_i - y_i^0}{(1 - \bar{x}_i^T C^{-1} \bar{x}_i)} \rightarrow \left[\text{Taking transpose as it is scalar} \right. \\ \left. \& \text{ } T(\text{scalar}) = \text{original} \right]$$

$$\therefore \bar{w}^T (x_i) \bar{x}_i - y_i = \frac{\bar{w}^T \bar{x}_i - y_i^0}{1 - \bar{x}_i^T C^{-1} \bar{x}_i}$$

vi) LOOCV error $\Rightarrow \sum_{i=1}^n (\bar{w}_{(i)}^T \bar{x}_i - y_i)^2$

According to the formula in section 2.5

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

Here, we calculate C^{-1} only once, which has a complexity of $O(k^3)$

Multiplications require $O(k^2)$ time. [Additions & subtraction are linear]
We will perform these multiplications n times.

\therefore Total complexity $\Rightarrow O(k^3 + n(k^2))$ [kH & k for big k]
 ~~$O(k^3)$~~

For usual method:

we will calculate $(c_i)^{-1}$ everytime
and use $w(c_i) = (c_i)^{-1} \cdot d(i)$ &

This requires $O(n \times k^3)$ for all inverse calculation, and requires $O(n \times k^2)$ for multiplications as well.

\therefore Total complexity would be $O(n(k^3 + k^2)) \approx O(nk^3)$
which is higher than the previous one.

3.2.1 The values of rmse on training, validation and loocv errors for different Lamda are:

1. Lambda 0.01-
 - 1.1.RMSE train-1.121
 - 1.2. RMSE validation- 2.579
 - 1.3.RMS loocv error-2.580
2. Lambda 0.1-
 - 2.1.RMSE train- 1.224
 - 2.2. RMSE validation- 2.157
 - 2.3.RMS loocv error- 2.182
3. Lambda 1 –
 - 3.1.RMSE train- 1.578
 - 3.2. RMSE validation- 1.997
 - 3.3.RMS loocv error- 2.009
4. Lambda 10-
 - 4.1.RMSE train- 2.19
 - 4.2. RMSE validation- 2.348
 - 4.3.RMS loocv error- 2.32
5. Lambda 100 –
 - 5.1.RMSE train- 2.971
 - 5.2. RMSE validation- 3.017
 - 5.3.RMS loocv error- 2.997
6. Lambda 1000-
 - 6.1. RMSE train-3.332
 - 6.2. RMSE validation- 3.345
 - 6.3.RMS loocv error- 3.335

3.2.2 Value of Lambda=1 achieves the best LOOCV performance. The values are:

Objective function: 17200.94

RMSE training value: 1.578

Regularization Term: 4749.95

3.2.3 The most important features are :

1. infused 2. Pineapple orange 3.red 4. Flavors nice 5. Sweet black
6. little heavy 7. New French 8. Future 9. Currant cola 10.cocktail

The least important features are:

- 1.offers 2. Light body 3. Highlights 4. Franc petit verdot 5.framed 6.tannins frame
7. tannins finish 8. Sour 9. Black cherry 10. Oakville

The screenshot contains the weights of these features. We can see that the absolute weights of important features are high whereas those of unimportant features are very low.

```
C:\Users\Saif\Desktop\ms\sem 1\ML\HW2\pandas>python hw2_pandas.py
1
2.00947406089
For best lambda:
('Cost of objective function ', array([[ 17200.94056872]]))
('RMSE training value ', 1.5780360753177243)
('Regularization term ', array([[ 4749.9512937]]))
('max weights are ', [6.998953378204021, 5.663260515333284, 5.636523454586268, 5.280346413180723, 5.194727564872672, 5.128505094493505,
5.078891240481596, 4.848258144674105, 4.786819319038614, 4.73595309040968])
('max weight indices are ', [184, 754, 773, 2368, 1272, 642, 186, 1924, 2835, 2068])
('min weights are ', [0.00014231005513920536, 0.0004945430905536341, 0.0007847465413199362, 0.0011453256827849145, 0.00178610341964713
37, 0.0018535941321715654, 0.0026979242565943196, 0.0037053503599508986, 0.0037468668994993237, 0.004492472102171519])
('min weight indices are ', [1045, 2310, 95, 2459, 1755, 1477, 404, 2388, 1085, 1730])
```


