

Neural Network Assignment 03 (Tutor Redion)

Priyanka Upadhyay (2581714) - s8prupad@stud.uni-saarland.de

Gopal Bhatrai (7013547) - gobh00001@stud.uni-saarland.de

1.1

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 \ln(x_{i2}) + \beta_3 x_{i3}$$

$$y_i = 10 + 10 x_{i1} + 0.5 \ln(x_{i2}) - 5 x_{i3} \quad (1)$$

(i) x_1 changes from 1 unit

$$y_{i+1} = 10 + 10 x_{i+1} + 0.5 \ln(x_{i2}) - 5 x_{i3}$$

$$\begin{aligned} y_{i+1} - y_i &= [10 + 10 x_{i+1} + 0.5 \ln(x_{i2}) - 5 x_{i3}] \\ &\quad - [10 + 10 x_{i1} + 0.5 \ln(x_{i2}) - 5 x_{i3}] \end{aligned}$$

$$\Rightarrow 10 x_{i+1} - 10 x_{i1}$$

$$[y_{i+1} - y_i] \Rightarrow 10 (x_{i+1} - x_{i1})$$

$$\text{so } x_{i+1} - x_{i1} = 1$$

↓

$$y_{i+1} - y_i = 10 \times 1 = 10$$

so if 1 change in x_1 gives 100% change in y so now 10 unit change in $(y_{i+1} - y_i)$ gives total $10 \times 100\% = 1000\%$ total change in response.

= TRUE

(ii) $x_2 = 1$ unit change

$$\begin{aligned} [y_{i+1} - y_i] &= 0.5 [\ln(x_{i2+1}) - \ln(x_{i2})] \\ &= 0.5 \ln[x_{i2+1} - x_{i2}] \end{aligned}$$

$\text{So } 0.5 \times 100.1 = 50.1$. total change in y

= TRUE

(iv) 100% change in $x_{i2} \rightarrow x_{i2} \rightarrow 2x_{i2}$

$$\begin{aligned}\text{So } y_{i+1} - y_i &= 0.5 \ln [2x_{i2} - x_{i2}] \\ &= 0.5 \ln [x_{i2}]\end{aligned}$$

$\Rightarrow 0.5 \times 100.1 = 50.1$. total change in y

= TRUE

=

(v) Higher debt implies lower future stock values

as y and x_{i3} has inverse relation because of the (-ve) coefficient, we can say that when x_{i3} has the largest value, y would become very small.

= TRUE

(vi) As per the machine learning \rightarrow we know that Bias term gives us an educated guess of the response when we don't have any knowledge about input variables.

So when $IP = 0$

$y = \beta_0 = 10 \rightarrow$ without any input we still have some response value.

1.2 Part b

$$J(\omega) = \text{MSE}_{\text{train}} + \gamma \omega^T \omega$$

ω = weights, γ = Regularization Parameter

Output: $\omega = \left[(x^{\text{Train}})^T x^{\text{Train}} + \gamma I \right]^{-1} x^{\text{Train}}^T y$

writing: $x = x^{\text{Train}}$, $y = y^{\text{Train}}$

ω minimize the distance b/w $J(\omega)$

$$\therefore J(\omega) = \frac{1}{2} \| x\omega - y \|^2 + \gamma \times \frac{1}{2} \| \omega \|^2$$

OR $J(\omega) = \frac{1}{2} (x\omega - y)^T (x\omega - y) + \frac{\gamma}{2} \omega^T \omega$

$$\Rightarrow \frac{1}{2} (x^T \omega^T - y^T) (x\omega - y) + \frac{\gamma}{2} \omega^T \omega$$

$$\Rightarrow \frac{1}{2} \left[x^T x \omega^T \omega - \underbrace{x^T \omega^T y - y^T x \omega + y^T y}_{y^T x = x^T y} \right] + \frac{\gamma}{2} \omega^T \omega$$

$$\Rightarrow \frac{1}{2} \left[x^T x \omega^T \omega - 2 \omega^T x^T y + \underbrace{y^T y}_{x^T x} \right] + \frac{\gamma}{2} \omega^T \omega$$

$$\Rightarrow \omega^T x \frac{1}{2} [2x^T x \omega - 2x^T y] + \frac{\gamma}{2} \omega^T \omega$$

So ~~$\partial J(\omega) / \partial \omega$~~ for closed form put $\frac{\partial J(\omega)}{\partial \omega} = 0$
 ~~$\partial J(\omega) / \partial \omega$~~ (Resulting matrix is full rank)

$$\frac{\omega^T}{2} [2x^T x \omega - 2x^T y] + \frac{\omega^T}{2} \gamma \omega = 0$$

$$\Rightarrow x^T x \omega - 2x^T y + \gamma \omega = 0$$

$$\Rightarrow \omega (x^T x + \gamma) - x^T y = 0$$

$$w = \frac{x^T y}{x^T x + \gamma}$$

$$w = [x^T x + \gamma I]^{-1} x^T y$$

putting back again $x = x^{\text{train}}$, $y = y^{\text{train}}$

$$\text{weights } w = [x^{\text{train}}{}^T x^{\text{train}} + \gamma I]^{-1} x^{\text{train}}{}^T y^{\text{train}}$$

Answer

[Reference:-]

Ridge Regression, Wessel N. Van → Page 8

Problem - 4 Logistic Regression

- 1 → Logistic Regression performs when our Response variable is categorical, predictors variable can be anything continuous or categorical. Since Linear Regression work on a very strong assumptions which is that the outcome variable is continuous & using Linear Regression such case violates the strong assumption & therefore can not be used.
- 2 → Logistic Regression is supervised learning method.
- 3 → We can not use because then cost function would become non-convex function which is not our desired goal.

5 → for a simple classifier, the outcome Y has 2 category 0 or 1 / yes or No but nothing in the middle of these values. But this is not true of all various input values X . So in the Logistic Regression the output can be seen as probability since it matches to training distribution.

4 → Logistic Regression is trained by feeding on input data and a binary class to which this data belongs

Problem-5 EigenDecomposition:-

$$M = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix}$$

→ Symmetric $x = x^T$ and M is a symmetric matrix. It implies that matrix is orthogonal so M^{-1} could be replaced by ~~M^T~~ M^T which is much easier than calculating M^{-1}

→ Yes, M is Singular. It implies that matrix is diagonalizable

→ As we have also solved it in the assignment 02!.

$$MV = \lambda V$$

$$[1-\lambda] \cdot v = 0$$

$$\begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix} = 0$$

$$\Rightarrow [1-\lambda] = 0$$

$$= \begin{bmatrix} 1-\lambda & -1 & 0 \\ -1 & 2-\lambda & -1 \\ 0 & -1 & 1-\lambda \end{bmatrix} = 0$$

$$\Rightarrow (1-\lambda)[(2-\lambda)(1-\lambda) - 1] + (-1)[(1-\lambda)]$$

+ 0

$$\Rightarrow (1-\lambda)(2-\lambda)(1-\lambda) - (1-\lambda) - (1+\lambda) = 0$$

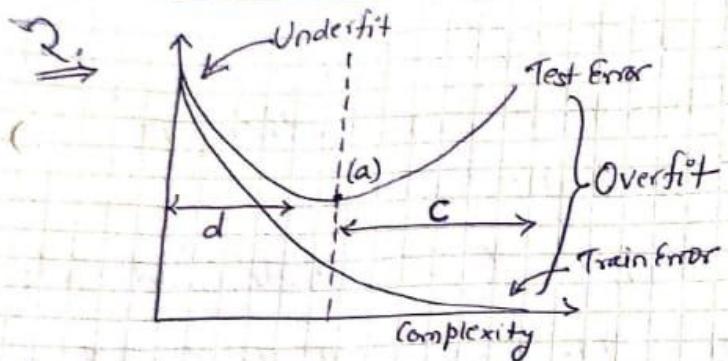
$$\Rightarrow \lambda = 3, 1, 0$$

$$\lambda_1 = 3, \quad v_1 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$$

$$\lambda_2 = 1 = v_2 = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

$$\lambda_3 = 0 = v_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Ans



* Point (a) represents the ideal point where we should operate ML algorithm. Because after that test error starts to go up and train error down \rightarrow Overfitting.

- After the Point (a), the range shown by C is Overfitting. the model became highly complex, we should reduce down its complexity.

Like \Rightarrow Minimize the degree of Polynomial.

Manipulate Regularization Const. (i.e. increase the value of it)

- The point where both train and test error is very high the model is Underfitting, because model is not complex enough to grab the pattern of underlying dataset.

(b) Ridge Regression $\Rightarrow \|\theta^T X - b\|_2^2 + \lambda \|\theta\|_2^2$

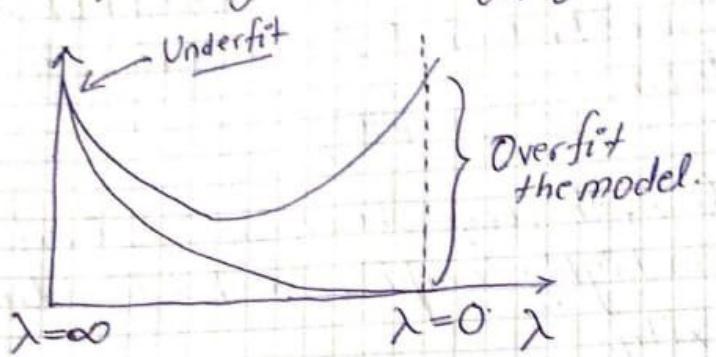
Here $\theta \in \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}, X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \in \mathbb{R}^{P+1}$

If I increase $\lambda \Rightarrow$

Our Job is to minimize the loss. f^n above, If I increase λ , then θ has to be small, then and only then the loss f^n will go down.

If $\lambda = \text{Very high} \rightarrow \theta \approx 0$, Model Underfits

$\lambda = \text{Very Small} \rightarrow \text{No restriction}$, Model Overfits



Q.3 here $K = \# \text{Instances}$. i.e. Training on $n-1$ instances and predicting on 1 instance.

So we trained on $n-1$ instances and then we calculated the score on one remaining instance we got MSE. Here I assume that error 35 we got was on Validation points. Not the average of $[RMSE_1, RMSE_2, \dots, RMSE_n]$. So when I reshuffle and again do LOOCV test then two Case might happen \Rightarrow

1. We again Selected the same Validation point. In that case I'll get 35 again.

2. In the Case when I selected another

Point as the Validation point. :

Since we used $N-1$ instances for training, and we assume it captures almost all of the variance of the data. The MSE that we'll get on the Validation point will be more likely close to 35. (Assuming it's not an Outlier.)

$$3. \quad L = \underset{\lambda}{\operatorname{argmax}} P(x)$$

$$L = \log \underset{\lambda}{\operatorname{argmax}} \frac{e^{-\lambda} \lambda^n}{n!}$$

$$L = \log \prod_{i=1}^n \frac{e^{-\lambda} \cdot \lambda^{x_i}}{x_i!} \quad \text{Because R.V are i.i.d.}$$

$$L = \sum_{i=1}^n \log(e^{-\lambda} \cdot \lambda^{x_i}) - \log(x_i!)$$

$$L = \sum_{i=1}^n \log(e^{-\lambda}) + \log(\lambda^{x_i}) - \log(x_i!)$$

$$L = \sum_{i=1}^n -\lambda + x_i \cdot \log \lambda - \log(x_i!)$$

$$\frac{\partial L}{\partial \lambda} = \sum_{i=1}^n \left[-1 + \frac{x_i}{\lambda} \right] = 0$$

$$-\eta + \frac{1}{\lambda} \cdot \sum_{i=1}^n x_i = 0$$

$$\frac{1}{\lambda} \sum_{i=1}^n x_i = \eta \quad \text{Hence,}$$

$$\left[\lambda = \frac{1}{\eta} \cdot \sum_{i=1}^n x_i \right]$$

(b) Since we now assumed $x_i = 0$ or $x_i > 0$.

$$\text{So by Part (a)} \Rightarrow \lambda = \frac{1}{n} \sum_{i=1}^n x_i \quad \because x_i > 0 \quad \text{or} \\ x_i = 0$$

If I want to Interpret λ in terms
of y_i 's, I can \Rightarrow

$$\lambda = \frac{1}{n} \sum_{i=1}^n y_i \quad ; \quad y_i = 1 \rightarrow x_i > 0 \\ y_i = 0 \rightarrow x_i = 0$$

So Say All the $x_i > 0$ then

$$[\lambda = 1]$$

$$\text{if all } x_i = 0 \forall i : \boxed{\lambda = 0}$$

Because while doing Maximum Likelihood Estimation
we took all the random variables into account.
and now by getting the final expression, we can
find value of λ based on cond' on Random Variables.