

CCCC

1.a As it is given in that we have two vectors (y_1, y_2, \dots, y_n) and (x_1, x_2, \dots, x_n) of length n . Now considering that X is a matrix of length $n \times (d+1)$ we cannot make the assumption that the given data set has a non zero mean which means that w_0 cannot be assumed to be zero in this case. So the X matrix shapes up

$$\begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & \dots & x_{2d} \\ 1 & x_{31} & x_{32} & \dots & x_{3d} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$$

As it is clear that the matrix W (weight) should be dimension $(d+1) \times n$. We can say that the weight matrix may look like:

$$\begin{bmatrix} w_0 & w_1 & w_2 & \dots & w_d \\ w_{11} & w_{12} & w_{13} & \dots & w_{1n} \\ w_{21} & w_{22} & w_{23} & \dots & w_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ w_{d1} & w_{d2} & w_{d3} & \dots & w_{dn} \end{bmatrix}$$

Since the model that we are using to fit this data set is a multivariate linear regression model, so as to maintain the property of matrix multiplication of matrix X and W , so the matrix w had to be of the shape mentioned above. Now as the model used is a linear regression we can say that the predicted value of the model can be of the type:

$$\tilde{y} = Xw$$

It is evident that y is a square matrix of dimension $(d+1) \times (d+1)$. As we know that the mean square error is the square of the difference between the expected and the predicted value over all the data points. This can be written as:

$$MSE(w) = \frac{\sum_{i=1}^n (y - Xw)^2}{n}$$

Now this looks similar to the the L2 norm of the term $y - Xw$, if we take the square of the L2 norm of $y - \tilde{y}$ we get:

$$\begin{aligned} & \sum_{i=1}^n ||y - Xw||_2^2 \\ &= \sum_{i=1}^n \sqrt{(y - Xw)^2}^2 = \sum_{i=1}^n (y - Xw)^2 \end{aligned}$$

This result is equivalent to MSE of the function $y - \tilde{y}$.

1.b For getting the optimal values of $MSE(w)$, we can calculate the gradient or the partial derivative w.r.t w . As derived from the above problem we have the $MSE(w)$ as

$$\begin{aligned}MSE(w) &= \frac{1}{n}(y - Xw)^T(y - Xw) \\&= \frac{1}{n}(y^T - X^T w^T)((y - Xw)) \\&= \frac{1}{n}(y^T y - X^T w^T y - y^T Xw - X^T Xw^T w)\end{aligned}$$

Taking gradient of MSE we get:

$$\begin{aligned}\nabla(MSE) &= \frac{\partial MSE}{\partial w} \\&= 2X^T Xw - 2X^T y\end{aligned}$$

For the optimum solution $\nabla(MSE) = 0$. Since it is assumed that $\text{rank}(x) = k(\text{fullrank})$, then $X^T X$ is a positive definite and unique solution of the normal equation is

$$\begin{aligned}X^T X\hat{w} &= X^T y \\ \hat{w} &= X^T y(X^T X)^{-1}\end{aligned}$$

Few other assumptions that we need to consider is as follows :

- (i) X is a non-stochastic matrix
- (ii) X has to be a singular matrix to get its inverse
- (iii) $\lim_{x \rightarrow +\infty} (\frac{X^T X}{n}) = \Delta$ exists and is a non-stochastic and non singular matrix (with finite elements).