

~~ML~~ ML Assignment - 1

Q.1 $y = ax + b$ derive formulas for a & b that minimizes the least square error.

Ans Given that the input is one dimensional assume that we have n input points such as

$$x = \{x_1, x_2, x_3, x_4, \dots, x_n\}$$

Now to predict points $Y = \{y_1, y_2, y_3, \dots, y_n\}$ we will use the formula $y = ax + b$.

Lets assume that we have the actual class values as $Z = \{z_1, z_2, z_3, \dots, z_n\}$

The mean square error will be calculated as follows

$$MSE = \frac{1}{n} \sum_{i=1}^n (z_i - y_i)^2$$

substituting $y = ax + b$ we get

$$MSE = \frac{1}{n} \sum_{i=1}^n (z_i - \cancel{(ax_i + b)})^2 \quad -(1)$$

To minimize this error we need to minimize the equation (1)

(P.T.O)

Since n is a constant we ~~can~~ have to minimize the following

$$\sum_{i=1}^n (z_i - (ax_i + b))^2$$

On Expanding the equation we get

$$MSE = \frac{1}{n} \left[(z_1 - (ax_1 + b))^2 + (z_2 - (ax_2 + b))^2 + (z_3 - (ax_3 + b))^2 + (z_4 - (ax_4 + b))^2 \dots \right]$$

$$MSE = \frac{1}{n} \left[(z_1^2 - 2z_1(ax_1 + b) + (ax_1 + b)^2) + (z_2^2 - 2z_2(ax_2 + b) + (ax_2 + b)^2) + (z_3^2 - 2z_3(ax_3 + b) + (ax_3 + b)^2) + (z_4^2 - 2z_4(ax_4 + b) + (ax_4 + b)^2) + \dots \right]$$

On grouping the similar terms together we get

$$MSE = \frac{1}{n} \left[(z_1^2 + z_2^2 + \dots + z_n^2) - 2a(z_1x_1 + \dots + z_nx_n) - 2b(z_1 + z_2 + z_3 + \dots + z_n) + a^2(x_1^2 + x_2^2 + \dots + x_n^2) + 2ab(x_1 + x_2 + x_3 + \dots + x_n) + nb^2 \right]$$

Lets say

$$\hat{z}^2 = z_1^2 + z_2^2 + z_3^2 + \dots z_n^2 / n$$

$$\hat{z} = z_1 + z_2 + \dots z_n / n$$

$$\hat{x}_x = z_1 x_1 + z_2 x_2 + z_3 x_2 \dots z_n x_n / n$$

$$\hat{x} = x_1^2 + x_2^2 \dots x_n^2 / n$$

$$\hat{x} = x_1 + x_2 + x_3 \dots x_n / n.$$

Now we get the MSE equation as follows

~~(2) — MSE = $\hat{z}^2 - 2a\hat{z}_x - 2b\hat{z} + a^2\hat{x}^2 + 2ab\hat{x} + b^2$~~

To minimize this equation we will use partial derivative with respect to a and with respect to b . We will set the partial derivative equal to 0.

$$\therefore \frac{\partial \text{J(MSE)}}{\partial a} = 0 \quad \frac{\partial \text{J(MSE)}}{\partial b} = 0.$$

$$\frac{\partial (\text{MSE})}{\partial (a)} = -2\hat{z}\hat{x} + 2a\hat{x}^2 + 2b\hat{x} - (3)$$

Divide with 2 we get $-z\hat{x} + a\hat{x}^2 + b\hat{x} = 0 - (4)$

$$\frac{\partial (\text{MSE})}{\partial (b)} = -2\hat{z} + 2a\hat{x} + 2b - (5)$$

Divide with 2 we get

$$-z + a\hat{x} + b = 0 - (6)$$

For equation 4 we shift the z term to right & we get

$$a\hat{x}^2 + b\hat{x} = \hat{z}\hat{x}$$

dividing by \hat{x}

$$a\frac{\hat{x}^2}{\hat{x}} + b = \frac{\hat{z}\hat{x}}{\hat{x}} - (7)$$

Similarly for equation (6) we get

$$a\hat{x} + b = \hat{z} - (8)$$

Subtracting (7) from (8)

$$a\left(\hat{x} - \frac{\hat{x}^2}{\hat{x}}\right) = \hat{z} - \frac{\hat{z}\hat{x}}{\hat{x}}$$

$$\therefore a = \frac{\left(\hat{z} - \frac{\hat{z}\hat{x}}{\hat{x}} \right)}{\left(\hat{x} - \frac{\hat{x}^2}{\hat{x}} \right)}$$

Using this equation/value of a to find value of b we get

$$b = \hat{z} - \frac{\left(\hat{z} - \frac{\hat{z}\hat{x}}{\hat{x}} \right)}{\left(\hat{x} - \frac{\hat{x}^2}{\hat{x}} \right)}$$

$$\therefore a = \frac{\hat{z}\hat{x} - \hat{z}\hat{x}}{(\hat{x})^2 - \hat{x}^2} \quad - \quad (9)$$

$$b = \hat{z} - \hat{x} \left(\frac{\hat{z}\hat{x} - \hat{z}\hat{x}}{(\hat{x})^2 - \hat{x}^2} \right) \quad - \quad (10)$$

These values of a and b will minimize the \leq mean square error.

Q.2 Regularized Linear Regression

(a) What will happen if we overestimate the value of λ for LASSO? What would happen to the number of non-zero elements of ω .

Ans. The sparsity in weights will increase if λ is overestimated in LASSO. The formula in LASSO is as follows

$$\lambda \|\beta\|_1 - (1)$$

Now in equation since λ is multiplied by $\|\beta\|_1$ and since the algorithm tries to decrease or minimize the error, it will lead to many β values ~~not~~ zero and this will lead to omission of many non-zero elements. Although this helps in decreasing the model complexity, it will lead to underfitting.

(b) What if underestimate the value of λ for RIDGE. What happens to value of ω ?

Answer) The formula for in RIDGE is given as

$$\lambda \|\beta\|_2^2$$

which is

$$\lambda (\beta_1^2 + \beta_2^2 + \beta_3^2 + \dots + \beta_n^2)$$

(2b) (continued)

Now if we under-estimate the value of λ then the term λ will be negligible and it will give the same equation as linear regression OLS. Since OLS reduces the bias, it will be beneficial but it will lead to overfitting.

(c) Calculate the partial derivative of penalty terms for LASSO and RIDGE.

Ans) LASSO:

$$J(\omega) = \frac{1}{N} \sum_{i=1}^N (y_i - (\omega_0 + \omega^T x_i))^2 + \lambda \|\omega\|_1$$

$$\text{Penalty term} = \lambda \|\omega\|_1 = \lambda \sum_{i=1}^d |\omega_i|$$

$$\frac{\partial (\text{Penalty})}{\partial \omega} = \begin{cases} \lambda & \text{if } \omega_i > 0 \\ -\lambda & \text{if } \omega_i < 0 \end{cases}$$

Ridge:-

$$\text{Penalty} = \lambda \|\omega\|_2^2$$

$$\frac{\partial (\text{Penalty})}{\partial \omega} = 2\lambda \omega_i$$

(d) For what values of w_i will their behaviors differ.

Ans The Lasso coefficients will be reduced by constant factor. Their values will be lower than ridge coefficients and will become zero faster.

(ii) The formula for penalty in ridge is $\lambda \|w\|_2^2$ which means $\lambda \cdot |w_i|^2$

The formula for penalty term in Lasso is $\lambda \|w\|_1$, which means $\lambda \cdot |w_i|$

(iii) From the above formulas it can be seen that with negative w_i ($-w_i$) the ridge factor will square and the effect of negative sign won't be there but in Lasso the negative term will affect the penalty terms value and in the term affect the values of β .

(iv) In terms of partial derivative of the w_i value won't affect LASSO as the partial derivative ~~does not have~~ is λ or $-\lambda$

(P.T.O)

(iv) For negative or positive values of w_i the partial derivative terms in lasso will give λ or $-\lambda$ but in ridge it will reduce or increase the penalty term by 2λ . For value

(v) For $w_i = 0$ the penalty term will be 0 for both, the ridge coefficient will be reduced factor of linear regression.

and an additional note

Maths

maths

for a given dataset and a given condition
the Lasso regression finds the best fit
by adding a constraint that the sum of the absolute values of the coefficients must be less than or equal to a constant. This constraint tends to shrink the coefficients towards zero, which can help in feature selection and can also reduce the variance of the estimates. In contrast, Ridge regression adds a constraint that the sum of the squares of the coefficients must be less than or equal to a constant. This constraint does not shrink the coefficients towards zero, but it increases the variance of the estimates. Both methods are useful in different situations depending on the specific requirements of the analysis.

Q.3a Write down the equation for weight update step.

$$\omega = \omega + \alpha \frac{\partial l(\omega)}{\partial \omega}$$

'+' indicates maximizing Likelihood.

$$\frac{\partial l(\omega)}{\partial \omega_j} = \left(\frac{y}{\text{sigm}(\omega^T x)} - \frac{(1-y)}{1-\text{sigm}(\omega^T x)} \right) \frac{\partial \text{sigm}(\omega^T x)}{\partial \omega_j}$$

$$= \cancel{y} + \cancel{(1-y)}$$

$$= \left(\frac{y}{\text{sigm}(\omega^T x)} - \frac{(1-y)}{1-\text{sigm}(\omega^T x)} \right) g(\omega^T x)$$

$$\times (1 - \text{sigm}(\omega^T x)) \frac{\partial (\omega^T x)}{\partial \omega}$$

$$= (y(1 - h_{\omega}(x)) - (1-y)(h_{\omega}(x))) x_j$$

$$= (y - h_{\omega}(x)) x_j$$

\therefore Batch gradient step

$$\omega_j = \omega_j + \alpha \sum_{i=1}^M (y_i - \text{sigm}(\omega^T x)) x_j$$

(3b) Does it always converge? what is a suitable stopping condition?

Ans) (i) A good stopping condition can be to check if the difference between the updated weights and previous weights is negligible

$$\text{i.e } (\omega' - \omega) < \epsilon \quad 0 < \epsilon < 1$$

(ii) Another good stopping condition would be a fixed number of iteration. Study shows that ~~iteration~~ between 5000 to 10000 iteration is a good stopping value.

(iii) If the learning rate α is not set properly, it will never converge.

(iv) Line search algorithm is good solution for finding a good α rate which will give a better solution.

Q.5a. Prove that for any arbitrary tree with unequal branching ratios, there exists a binary tree that implements same classification functionality.

Ans) In this proof we will only be interested in trees with nodes who have children or subnodes greater than 2.

(i) We start by from the root node and traverse the tree to find node with children greater than 2

(ii) If we don't find such a node the tree is a binary tree.

(iii) If we find a tree with children greater than 2, say n children we will ~~skip~~ build the tree as follows

(a) Keep the old node as it is

(b) Keep the leftmost child of the old node as ~~it~~ it is.

(c) Create a new right node and assign ~~the~~ $n-1$ children of the old node as children of the new right node.

(d) Do this recursively till all nodes ~~are~~ have children less than or equal to two.

(iv) Thus we have converted to functionally equivalent binary tree.

5(b) Consider a tree with just 2 levels, a root node connected to B leaf nodes ($B >= 2$)
what are upper & lower limits on functionally equivalent binary tree as function of B

Ans The maximum number of leaf nodes in a binary tree of depth d is equivalent to 2^d .

\therefore Since the given tree which has B leaf nodes should be less than or equal to 2^d . i.e. $B \leq 2^d$

Taking log we get

$$\log_2 B \leq d$$

To this we add 1 for the root node

$$d = \lceil \log_2 B + 1 \rceil - 1 \quad (\text{Minimum Value})$$

The upper limit for binary tree with depth d will be

$$d \leq B$$

(5c) For binary tree with B terminal nodes, the number of nodes will be $2B - 1$.

Since $B \geq 2$ the minimum number of nodes will be 3.

Now for a tree with $B+1$ ~~nodes~~ terminal nodes ~~the~~ with depth d . To convert this to binary tree we need to have $2B - 1$ nodes. To convert we take $n-1$ right children, add a new node and make the $n-1$ ~~children~~ children of the new node so for example if we have 2 level tree with 1 extra node we get binary tree as

$$2(B+1) - 2(B-1) + 2 = 0$$

For $B+1$ the maximum number of nodes should be $2(B+1) - 1$.

\therefore It will be $2B+1$.