

ASSIGNMENT 2

Aim: Generate a proper 2D dataset of 14 points.
Split the data set into Training Data set and Test Data set.

- i) Perform linear regression analysis with Least Square Method
- ii) Plot the graph for Training MSD and Test MSD and comment on Curve Fitting and Generalization Error
- iii) Verify the Effect of Data set size and Bias - Variance Trade off
- iv) Apply Cross validation and plot the graph for errors
- v) Apply subset selection and plot graph for errors
- vi) Describe your finding

Theory:

When we have single input attribute (or) and we want to use linear regression, this is called simple linear regression.

If we had multiple input attributes, this would be called multiple linear regression. The procedure for linear regression is different and simpler than that for multiple linear regression, so it is a good place to start.

With simple linear regression we want to model our data as follows:

$$y = B_0 + B_1 x$$

This is a line where y is the output variable we want to predict, x is the input variable we know and B_0 and B_1 are coefficients that we need to estimate and move the line around.

Technically, B_0 is called the intercept because it determines where the line intercepts the y -axis. In machine learning we can call this the bias, because it is added to offset all predictions that we make. The B_1 term is called the slope because it defines the slope of the line or how x translates into a y value before we add our bias.

The goal is to find the best estimates for the coefficients to minimise the errors in predicting y from x .

Simple regression is great because rather than having to search for values by trial and error or calculate them analytically using more advanced linear algebra, we can estimate the directly from our data.

We can start off by estimating the value for B_1 as:

$$B_1 = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

where $\text{mean}()$ is the average value for a variable in our dataset. The x_i and y_i refer to the fact that we need to repeat these calculations across all values in our dataset and i refers to i^{th} value of x or y .

$$B_0 = \bar{y} - (B_1 * \bar{x})$$

We can calculate an error for our predictions called the Root Mean Squared Error or RMSE

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

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (\text{Predicted}_i - \text{Actual}_i)^2}{N}}$$

where $\text{sqrt}()$ is the square root function, \hat{y} is the predicted value and y is the actual value, i is the index for a specific instance, n is the number of predictions, because we must calculate the error across all predicted values.

Conclusion: The following objectives with respect to linear regression were implemented in this assignment successfully.