

Linear Regression

Prerequisite:

Statistics: mean, median, mode, variance, standard deviation, correlation and covariance.

Exploratory Data Analysis: Data distribution, scatter plot, correlation matrix, heat map.

Objectives:

- Understand what is Linear Regression and motivation behind linear regression
- What is the best fit line and residual of regression
- Least square method to find the best fit line of regression.
- Gradient descent method to find the best fit line of regression

Linear Regression

Linear regression is a way to identify a relationship between two or more variables and use these relationships to predict values for one variable for given value(s) of other variable(s). Linear regression assume the relationship between variables can be modeled through linear equation or an equation of line. The variable which is used in prediction is termed as independent/explanatory/regressor whereas the predicted variable is termed as dependent/target/response variable. Linear regression assumes that independent variables are **related linearly** to response variable.

$$y = c + mx$$

In machine learning and regression literature above equation is used in the form:

$$y = w_0 + w_1x$$

Where w_0 is intercept on y-axis, w_1 is slope of line, these are called parameters or coefficients of regression, x is the explanatory variable and y is the response variable.

Motivational Examples

1. Let's say we have sales data of house prices. For each house we have complete information about the plot size area and the price at which the house was sold. Can we use this information to predict the price of a house for a given plot size area? The problem can be modeled using linear regression with plot size as the explanatory variable and house price as the response variable.

$$HousePrice = w_0 + w_1PlotSize$$

2. Consider a scenario where we have medical data about some patients. The data contains the information of the blood pressure for a patient along with his/her age. Can we use this information to predict the blood pressure level of a patient for some given age? This problem can be solved by using linear regression with age as the explanatory variable and blood pressure as the response variable.

$$\text{BloodPressure} = w_0 + w_1 \text{Age}$$

3. Next consider a problem where we need to predict the price of a used car. The sale price of a used car depends on many attributes, some of them may be mileage (km/litre), model (Maruti, Hundai, Honda, Toyota, Tata), segment (small, medium, luxury). In this scenario, the sale price is the target variable which depends on mileage, model and segment. This problem can be solved using linear regression and can be represented by the following equation.

$$\text{SalePrice}(y) = w_0 + w_1 \text{Mileage} + w_2 \text{Model} + w_3 \text{Segment}$$

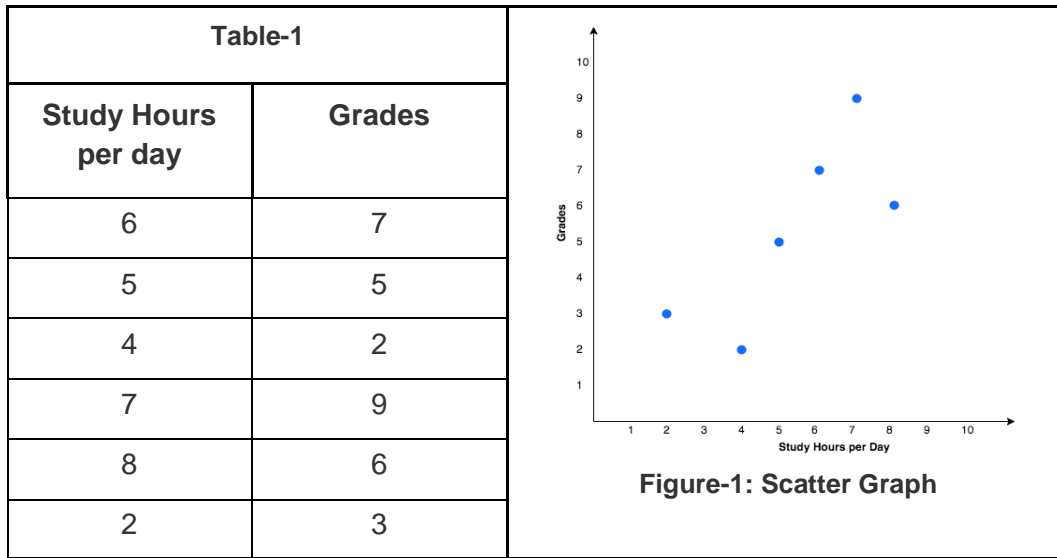
Such equation is called **multiple linear regression** as there are more than one explanatory variables involved in the prediction of target variable. In real world scenarios, we rarely have one explanatory variable, so we use multiple linear regression rather than simple linear regression. However, here we take an example of simple linear regression to understand the fundamentals of regression.

Example: Consider a toy example where we are interested in finding the effect of studying hours per day over grades in examination and predict grades of a student for given study hours. We have sample data about six students for their grades and total study hours per day (Table - 1).

From the given data we get an idea that study hours per day and grades have positive linear relationship. So one can say that if a student spends more hours studying per day he is likely to get good grades in his/her examination. The scatter plot of given data is shown in Figure-1. Scatter plot is a useful tool to judge the strength of relationship between two variables. If scatter plot does not conclude any linear relationship, then fitting a linear model to the data is probably not useful.

A valuable measure to quantify the relationship between two variables is the **correlation coefficient**. The correlation coefficient has range of values between -1 to +1 to indicate the strength of linear relationship where -1 indicates the perfect negative correlation (i.e. higher values of one variable tend to be associated with lower values of the other), +1 indicates the perfect positive correlation (i.e. higher values of one variable tend to be associated with higher values of the other) and 0 shows no correlation between the two variables.

From the scatter plot shown in Figure-1, we get some intuition that there is a positive effect of studying hours per day over grades in exam.



To fit the given data, we can draw multiple lines and choose the best fit line. Let's see how to decide which line is the best fit line. The equation of the linear model is given by:

$$y(\text{Grades}) = w_0 + w_1 X(\text{Study Hours per day})$$

Any line we might come up with will have some fixed intercept w_0 and a slope w_1 . This line may include some data points on it but cannot cover all of them. In our example we have six data points, let us label these points as (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , (x_4, y_4) , (x_5, y_5) and (x_6, y_6) , with values $(6, 7)$, $(5, 5)$, $(4, 2)$, $(7, 9)$, $(8, 7)$ and $(2, 3)$ respectively. For any given point x_i the prediction ($yhat_i$) is given by:

$$yhat_i = w_0 + w_1 x_i$$

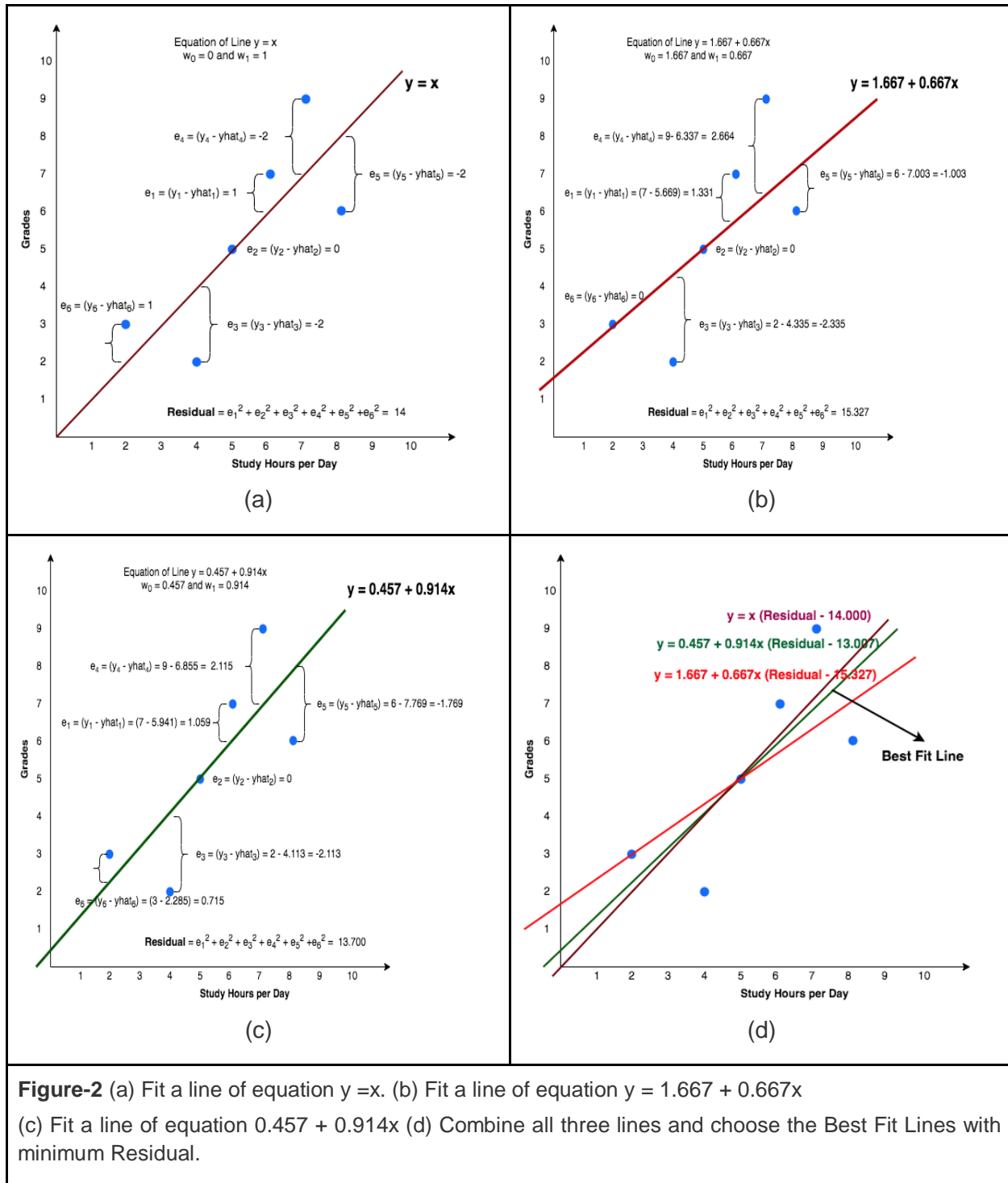
Unless the line passes through (x_i, y_i) the value of $yhat_i$ differs from the observed value of y_i . The difference between the two values is called an **error** or **residual** of regression.

$$e_i = y_i - yhat_i$$

The **best line** is the line which minimizes the **sum of the squared error**:

$$\sum e_i^2 = (y_1 - yhat_1)^2 + (y_2 - yhat_2)^2 + \dots \dots \dots + (y_n - yhat_n)^2$$

Following graphs illustrate the process to find the best line of regression.



Methods to Find Best Fit Line - We can use two different methods to find the best fit line of regression:

1. Ordinary Least Squares (OLS).
2. Gradient Descent.

Least Square:

Let the equation of regression line of y on x be:

$$y = w_0 + w_1x$$

According to the least square principle, the equations to estimate the values of w_0 and w_1 are:

$$\sum_{i=1}^n y_i = nw_0 + w_1 \sum_{i=1}^n x_i \dots \dots (1)$$

Multiply the equation (1) with x_i on both sides,

$$\sum_{i=1}^n x_i y_i = w_0 \sum_{i=1}^n x_i + w_1 \sum_{i=1}^n x_i^2 \dots \dots (2)$$

Dividing equation (1) by n we get,

$$\frac{1}{n} \sum_{i=1}^n y_i = w_0 + \frac{w_1}{n} \sum_{i=1}^n x_i$$

$$\bar{y} = w_0 + w_1 \bar{x} \dots \dots (3)$$

Thus we can say the line of regression will always pass through the points (\bar{x}, \bar{y}) .

Now we need to estimate the values for w_0 and w_1 ,

We know,

$$cov(x, y) = \frac{1}{n} \sum_i x_i y_i - \bar{x} \bar{y} \Rightarrow \frac{1}{n} \sum_i x_i y_i = cov(x, y) + \bar{x} \bar{y} \dots \dots (4)$$

Also,

$$var(x) = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \Rightarrow \frac{1}{n} \sum_{i=1}^n x_i^2 = var(x) + \bar{x}^2 \dots \dots (5)$$

Dividing equation (2) by n and using equation (4) and (5),

$$cov(x, y) + \bar{x} \bar{y} = w_0 \bar{x} + w_1 (var(x) + \bar{x}^2) \dots \dots (6)$$

By solving equation (3) and (6) we get,

$$w_1 = \frac{cov(x, y)}{var(x)} \dots \dots (7)$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

The straight line defined by $y = w_0 + w_1x$ which has the least sum of squared errors for variations in w_0 and w_1 , is called the line of regression of y on x . Let us try these equations to estimate best fit line on our data given in Table-1.

To estimate w_0 and w_1 , we need to find covariance between x and y , variance of x and mean of x and y variables. For given data we get,

$$\bar{x} = \frac{6 + 5 + 4 + 7 + 8 + 2}{6} = 5.333$$

$$\bar{y} = \frac{7 + 5 + 2 + 9 + 6 + 3}{6} = 5.333$$

$$cov(x, y) = 4.267$$

$$var(x) = 4.667$$

when we substitute these values in equation (7) and (8) we get,

$$w_0 = 0.4571 \text{ and } w_1 = 0.9143$$

which are exactly the same as shown in Figure-2(c) for the line $y = 0.457 + 0.914x$, which gives the minimum residual among all the lines.

Performance metric for least square regression:

Performance metrics are the tools to quantify and compare the efficiency of any machine learning model. Least square regression uses R^2 (R-squared) and R_{adj}^2 (Adjusted R-squared) metrics to measure the performance of a regression model. Both of these metrics denotes the power of selected independent variable(s) to explain the variation in response variable. The equations of R^2 and R_{adj}^2 are given by:

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

The numerator term gives the average of squares of residuals and denominator gives the variance in the y (response) variable. A small value for R^2 or higher mean residual error denotes poor model.

$$R_{adj}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - k - 1}$$

Where n is the total number of observations in data and k is the number of explanatory variables. Adjusted R-square is a slight improvement over R-squared by adding an additional term to it. The problem with R^2 is that, it increases with increase in number of explanatory variables in the model irrespective of whether the added variable(s) are significantly contributing in prediction or not. On the contrary, the value of R_{adj}^2 only increases if useful variables are added to the model otherwise it might decrease with the addition of insignificant variables. The relation between R^2 and R_{adj}^2 is:

$$R_{adj}^2 \leq R^2$$

Gradient Descent:

Let the equation of regression line of y on x be:

$$y = w_0 + w_1x$$

This straight line tries to approximate the relationship between x and y for given set of data. By varying the values of w_0 and w_1 we can find the best fit line. By the above discussion, we know that the best fit line is one which minimizes the total error in prediction. Gradient descent method defines a cost function (also called loss function) with parameters w_0 and w_1 and uses a systematic approach to optimize the values of parameters to get the minimum value of the cost function. Let's dive into mathematics of the algorithm.

The model be defined as:

$$y = w_0 + w_1x$$

Now define the cost function of gradient descent as Mean Squared Error (MSE):

$$\text{cost}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \text{yhat}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1x_i)^2$$

The cost function includes two parameters w_0 and w_1 , which controls the value of cost function. As we know that the derivative gives us the rate of change in one variable with respect to others, so we can use partial derivatives to find the impact of individual parameter over the cost function.

The principle of gradient descent is that we always make progress in the direction where the partial derivatives of w_0 and w_1 are steepest. If the derivatives of parameters become zero or very less, point the situation of either maxima or minima on the surface of cost function. The process of gradient descent is started with a random initialization of w_0 and w_1 . Every iteration of gradient descent improves in the direction of optimal values for w_0 and w_1 parameters which will have minimum cost function value. Following figure illustrates the process of optimization.

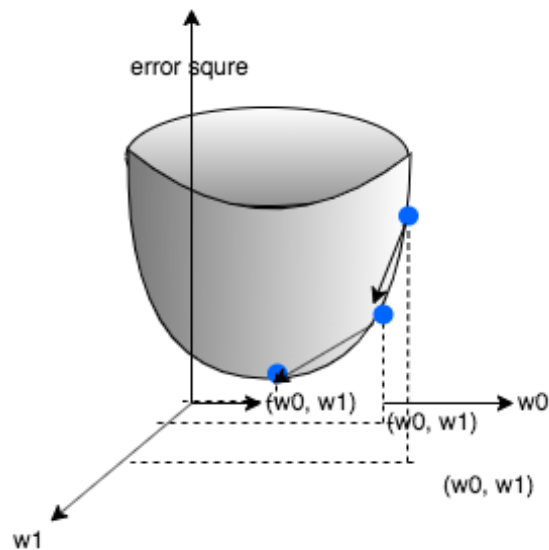


Figure 3: Gradient Descent Iteration

Gradient descent works in following steps:

1. Random initialization of parameters.
2. Calculate the partial derivatives of the cost function with respect to each parameter, which are called gradients.
3. Update the parameters in the opposite direction of gradients.
4. Repeat step 3 and 4 till maximum iteration reached or minimum cost function value achieved.

Partial derivatives:

We have,

$$error = cost(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)^2$$

Partial derivative w. r. t. w_0 and w_1 :

$$\frac{\partial cost(w_0, w_1)}{\partial w_0} = \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)(-2) = \frac{-2}{n} \sum_{i=1}^n error_i$$

$$\frac{\partial cost(w_0, w_1)}{\partial w_1} = \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)(-2x_i) = \frac{-2}{n} \sum_{i=1}^n error_i * x_i$$

Parameter updates:

$$w_0 = w_0 - \text{lr} \frac{\partial cost(w_0, w_1)}{\partial w_0}$$

$$w_1 = w_1 - \text{lr} \frac{\partial cost(w_0, w_1)}{\partial w_1}$$

lr is the learning rate which controls the step size of parameter update.

Let's run it on our example:

X:	6	5	4	7	8	2
y:	7	5	2	9	6	3

Let's initialize both coefficient w_0 and w_1 with 0.0,

$$w_0 = 0.0$$

$$w_1 = 0.0$$

Iteration #1:

$$\hat{y}_i = 0.0 + 0.0 x_i$$

Calculate gradients:

$$\frac{\partial cost(w_0, w_1)}{\partial w_0} = \frac{-2}{6} (7 + 5 + 2 + 9 + 6 + 3) = -10.6667$$

$$\frac{\partial cost(w_0, w_1)}{\partial w_1} = \frac{-2}{6} (7 * 6 + 5 * 5 + 2 * 4 + 9 * 7 + 6 * 8 + 3 * 2) = -64$$

Update parameters: (lr = 0.01)

$$w_0 = w_0 - \text{lr} \frac{\partial cost(w_0, w_1)}{\partial w_0} = 0.0 - 0.01 (-10.6667) = 0.1066$$

$$w_1 = w_1 - \text{lr} \frac{\partial cost(w_0, w_1)}{\partial w_1} = 0.0 - 0.01 (-64) = 0.64$$

Iteration #2:

$$\hat{y}_i = 0.1067 + 0.64 x_i$$

Calculate gradients:

$$\frac{\partial \text{cost}(w_0, w_1)}{\partial w_0} = \frac{-2}{6} (3.0533 + 1.6933 - 0.6667 + 4.4133 + 0.7733 + 1.6133) = -3.6266$$

$$\begin{aligned} \frac{\partial \text{cost}(w_0, w_1)}{\partial w_1} &= \frac{-2}{6} (3.0533 * 6 + 1.6933 * 5 - 0.6667 * 4 + 4.4133 * 7 + 0.7733 * 8 + 1.6133 * 2) \\ &= -21.475 \end{aligned}$$

Update parameters: (lrate = 0.01)

$$w_0 = w_0 - \text{lrate} \frac{\partial \text{cost}(w_0, w_1)}{\partial w_0} = 0.1067 - 0.01 (-3.6266) = 0.14296$$

$$w_1 = w_1 - \text{lrate} \frac{\partial \text{cost}(w_0, w_1)}{\partial w_1} = 0.64 - 0.01 (-21.475) = 0.8547$$

Similarly, all the iterations for gradient descent are performed till the minimum value of cost function of error is achieved or some finite iterations are reached.

Multiple Linear Regression:

Till now we have discussed the case of simple linear regression i.e. only one explanatory variable. But in real scenarios, the target variable depends on multiple explanatory variables which needs to be catered during the development of linear regression model. The model is expressed as:

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 \dots \dots \dots + w_n x_n$$

Where $x_1, x_2, x_3, \dots, x_n$ are explanatory variables and y is the target variable.

Evaluation of Linear regression model:

Evaluation helps to judge the performance of any machine learning model that would provide best results to our test data. Fundamentally three types of evaluation metrics are used to evaluate linear regression model.

- R2 measure (discussed with least square method)
- Mean Absolute Error (MAE)
- Root Mean Square Error (RMSE)

Mean Absolute Error (MAE):

Mean absolute error is the average of the differences between actual and predicted values of the target variable.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Root Mean Square Error(RMSE):

Root mean square error is the square root of average of the squared differences between actual and predicted values of the target variable.

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

Pros and cons of Linear Regression:

Pros:

- Linear regression models are very simple and easy to implement.
- It is computationally very efficient.
- Output's coefficients are easy to interpret.

Cons:

- Linear regression models are largely affected by the presence of outlier in the training data.
- It assumes linear relationship between target and explanatory variables which is sometimes is not true.
- It assumes independence between attributes which is not true for real word scenarios.
