# THE CONCEPT OF LINEAR REGRESSION

## LET'S LEARN MACHINE LEARNING WITH PYTHON

**Distance vs Speed**



© Amit Mishra

- Introduction to Linear Regression
- Example Problem
- Graphical Analysis
- Correlation Analysis
- Applying the Model
- Linear Regression Diagnostics
- How to know which regression model is best fit for the data?
- Predicting Linear Models
- K - Fold Cross Validation

# Definition

Linear regression is statistical tool for modeling the relationship between a scalar **dependent variable** y and one or more **explanatory variables** (or independent variables) denoted by X.

The representation is a linear equation that combines a specific set of input values (x) the solution to which is the predicted output for that set of input values (y). As such, both the input values (x) and the output value are numeric.

# Regression Analysis

Regression analysis is used to:
- Predict the value of a dependent variable based on the value of at least one independent variable.
- Explain the impact of changes in an independent variable on the dependent variable.

**Dependent variable:** the variable we wish to explain
**Independent variable:** The variable used to explain the dependent variable

# Let us take an example:

The goal here is to establish a mathematical equation for distance as a function of speed, so you can use it to predict distance when only the speed of the car is known.

So it is desirable to build a linear regression model with the response variable as distance and the predictor as speed.
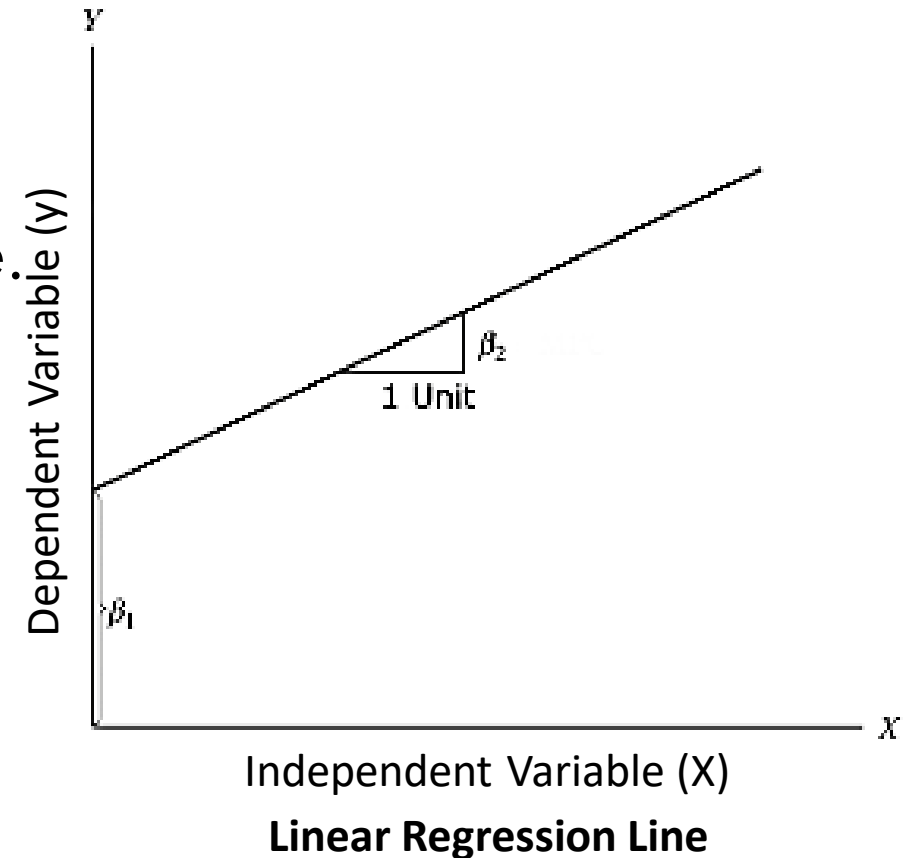
| Speed(X) | Distance(y) |
|----------|-------------|
| 4        | 2           |
| 4        | 10          |
| 7        | 4           |
| 7        | 22          |
| 8        | 16          |
| 9        | 10          |

# Mathematical Equation

The linear regression model fits a linear function to a set of data points. Mathematical equation can be generalized as:

**Y=β1+β2X**

where, β1 is the intercept and β2 is the slope.



Dependent Variable (y)

$\beta_2$

1 Unit

$\beta_1$

Independent Variable (X)

**Linear Regression Line**

**The population regression model representation**

Population
y intercept

Population
Slope
Coefficient

Independent
Variable

Random
Error term,
or residual

Dependent
Variable

$$y = \beta_0 + \beta_1 x + \varepsilon$$

Linear component

Random Error
component

## Mathematical Equation

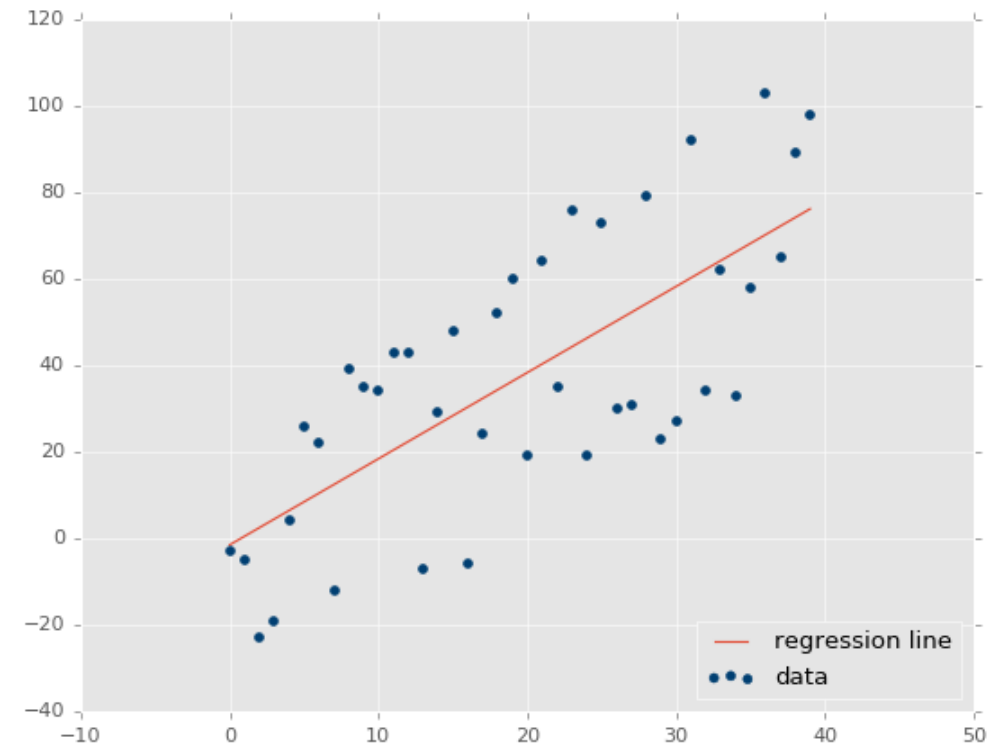The linear regression model fits a linear function to a set of data points. The form of the function for multiple linear regression is:

$$Y = β0 + β1*X1 + β2*X2 + … + βn*Xn$$

Where Y is the target variable, and X1, X2, ... Xn are the predictor variables and β1, β2, … βn are the coefficients that multiply the predictor variables. β0 is constant.
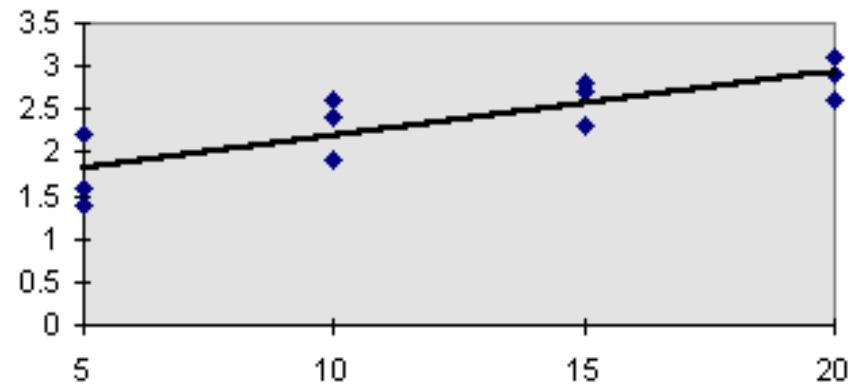
**The goal is to take continuous data, find the equation that best fits the data**, and be able forecast out a specific value. With simple linear regression, you are just simply doing this by creating a best fit line:
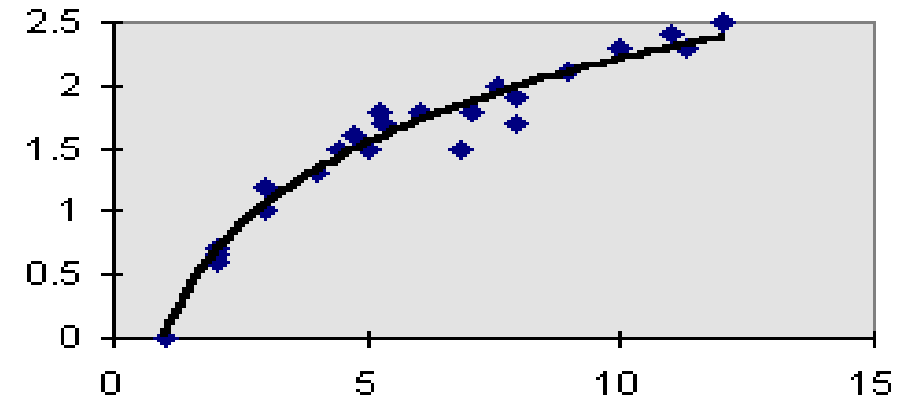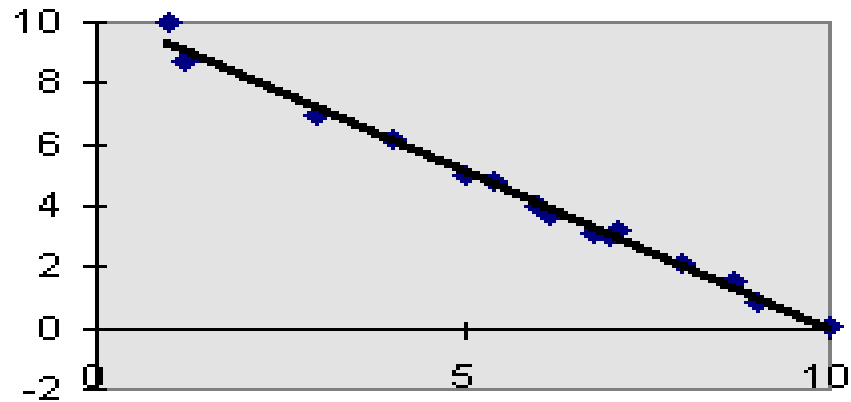
**Positive Linear Relationship**



**Relationship NOT Linear**



**Negative Linear Relationship**



**No Relationship**



© Amit Mishra

The Least square regression line is the unique line such that the sum of squared vertical error (y) distances between the data points and the line is the smallest possible.



**Linear Regression Line**

Hypothesis $\qquad h_\theta(x) = \theta_0 + \theta_1 x$

Parameters $\qquad \theta_0, \theta_1$

Cost Function $\qquad J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Goal $\qquad \underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

- Error values ($\varepsilon$) are statistically independent
- Error values are normally distributed for any given value of x
- The probability distribution of the errors is normal
- The probability distribution of the errors has constant variance
- The underlying relationship between the x variable and the y variable is linear

Correlation analysis studies the strength of relationship between two continuous variables. It involves computing the correlation coefficient between the two variables.

Correlation is a statistical measure that shows the degree of linear dependence between two variables.

Correlation can take values between -1 to +1. If one variables consistently increases with increasing value of the other, then they have a strong positive correlation (value close to +1).
If one consistently decreases when the other increase, they have a strong negative correlation (value close to -1).

Correlation is only an aid to understand the relationship.
Simply use the cor() function with the two numeric variables as arguments.

**Residual Sum of Squares**

$RSS = \Sigma(Y_i - Y_{fitted})^2$

**Explained Sum of Squares**

$ESS = \Sigma(Y_{fitted} - Y_{mean})^2$

Intercept ($\beta_1$)

Y

Actual $Y_i$

$Y_{fitted}$

Residual

**Total Sum of Squares**

$TSS = \Sigma(Y_i - Y_{mean})^2$

$Y_{mean}$

X

© Amit Mishra

Total information in a variable is the amount of variation it contains.

R Sq.=1 − RSS/TSS

Where, RSS is the Residual Sum of Squares given by

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

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

**data**

an optional data frame, list or environment (or object coercible by as.data.frame to a data frame) containing the variables in the model. If not found in data, the variables are taken from environment(formula), typically the environment from which lm is called.

# Regression Model Summary

| Statistic | CRITERION |
|---|---|
| R-Squared | Higher the better |
| Adj. R-Squared | Higher the better |
| Std. Error | Closer to zero the better |
| MAPE (Mean absolute percentage error) | Lower the better |
| MSE (Mean squared error) | Lower the better |
| | |

The principle of least squares is one of the methods for finding a curve fitting a given data. Say (x1, y1), (x2, y2), …. (xn, yn) be n observations from an experiment.

We are interested in finding a curve.

$$y = f(x)$$

Closely fitting the given data of size 'n'. Now at x=x1 while the observed value of  y is  y1, the expected value of y from the curve (above) is f(x1). Let us define the residual by ..

$$e1 = y1 - f(x1)$$

Likewise, the residuals at all other points are given by ..

$$e2 = y2 - f(x2)$$
$$........$$
$$en = yn - f(xn)$$

Some of the residuals may be positive and some may be negative. We would like to find the curve fitting the given data such that the residual at any xi is as small as possible. Now since some of the residuals are positive and others are negative and as we would like to give equal importance to all the residuals it is desirable to consider sum of the squares of these residuals, say and thereby find the curve that minimizes E.

Thus, we consider …

$$E = \sum_{i=1}^{n} \left( e_i^2 \right)$$

# HAPPY LEARNING