

1 What is Machine Learning?

- Machine Learning is a subset of Artificial Intelligence.
- It is focused mainly on the designing of systems, thereby allowing them to learn and make predictions based on some experience which is data in case of Machine.
- How does Machine Learning work?

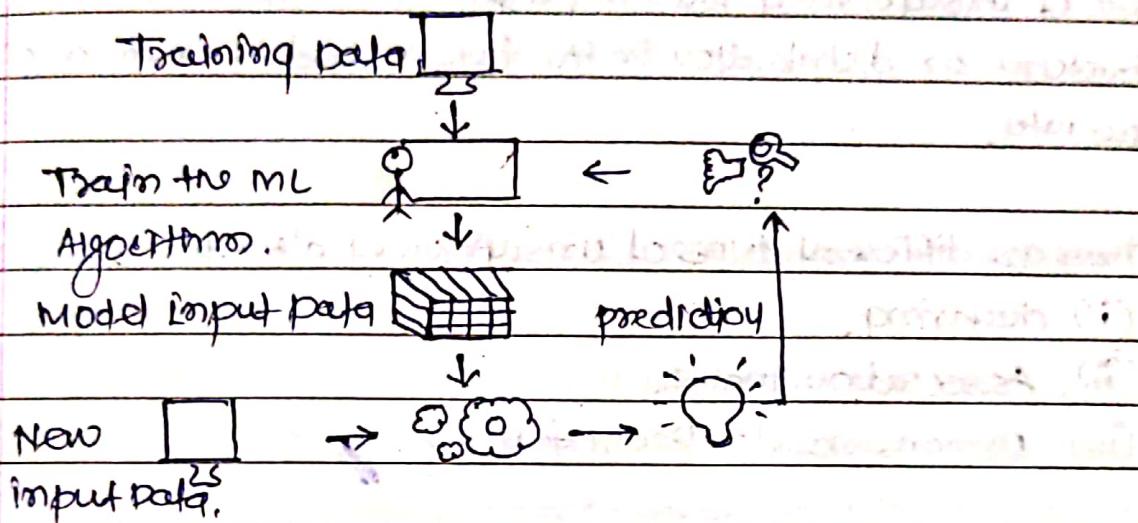


Diagram: How Machine Learning works.

- There are different types of Machine Learning Algorithms. Above are the types of Machine Learning.
 - (i) supervised learning.
 - (ii) Unsupervised learning
 - (iii) Reinforcement Learning.

① supervised learning:

- supervised learning Algorithm is input variable x and output variable y , and use an algorithm to learn mapping function from input to output, that is $y = f(x)$.
- goal is to approximate mapping function so well that whenever new data (input x) good output predicted variable.

(i) Classification

(ii) Regression.

(iii) Unsupervised Learning Algorithms:

- Unsupervised learning algorithm only have input data, no corresponding output.
- goal of unsupervised learning algorithm is model underlying structure or distribution in the data in order to measure the data.
- There are different types of unsupervised algorithm,
 - (i) clustering.
 - (ii). Association Analysis
 - (iii) Dimensionality Reduction.

(iv) Reinforcement Algorithms:

- Reinforcement algorithm soft upon machine automation to determine ideal behaviour, within specific context to rate its performance.
- There are different type of Reinforcement Algorithm.
 - i) Model Free - Q-learning
 - ii) Model Based.

① Linear Regression:

- Linear Regression is a supervised machine learning algorithm where the predicted output is continuous and has constant slope.
- It is used to predict values within a continuous range, (Ex. sales, price), rather than trying to classify them into categories.
- There are two main types:
 - (i) simple linear regression;
 - (ii) Multiple linear regression.

i) Simple Regression:

- Simple linear Regression is used to estimate the relationship between two quantitative variables.
- Use simple linear Regression when you want to know:
 - (i) How strong the relationship between two variables.
 - (ii) The value of the dependent variable at a certain value of independent variable.
- If you have more than one independent variable, use multiple linear Regression.

Assumption of simple linear Regression:

- simple Linear Regression is a parametric test, meaning that it makes a certain assumption about the data.
- Assumptions are:
 - (i) Homogeneity of variance (Homoscedasticity):
 - size of errors in the prediction doesn't change significantly across the values of the independent variable.
 - (ii) Independence of observations:
 - observations in the dataset were collected using statistically valid sampling methods.
 - No hidden relationship among observations.

(iii) Normality:

- Data follows Normal Distribution.
- Linear Regression makes one ~~one~~ additional assumption.

(iv) Linearity:

- Relationship between the independent and dependent variable is linear.
- the line of best fit through the data points is a straight line (rather than curve or some sort of grouping factor).
- If data do not meet the Assumptions of Homoscedasticity or Normality, you may able to use a non-parametric test instead.
- If your data violate the assumption of independence of observations (ex. observations are repeated over time), you may able to perform a linear mixed-effects model that accounts addition structure of data.

Perform a simple linear regression:

- the formula for a simple linear regression is
- $$y = B_0 + B_1 X + \epsilon$$
- y - predicted value of dependent variable (y) for given value of independent variable (x).
 - B_0 - intercept, predicted value of y when x is 0.
 - B_1 - Regression coefficient - how much we expect y to change as x increases.
 - x - independent variable (variable we expect influencing y).
 - ϵ - ϵ is the error of the estimate, how much variation there is in our estimate of the regression coefficient.

- Linear Regression finds line of best fit line through data by searching regression coefficient (B_1) that minimizes the total error (e) of model.

① Error calculated in a Linear Regression model?

- Linear Regression most often uses mean-square-error (MSE) to calculate the error of the model. MSE is calculated by.
 - (i) measuring the distance of the observed y values from predicted y -values at each value of x ,
 - (ii) squaring each of these distances,
 - (iii) calculating mean of each of squared distances.
- Linear Regression fits a line to the data by finding the regression coefficient that results in smallest MSE.

② How to find mean?

- Mean or Arithmetic mean of a dataset is the sum of all values divided by total number of values.
- Also called as "Average".

③ What is Regression Model?

- Regression model is a statistical model that estimates the relationship between one dependent variable and one or more independent variables quantitatively using line. (or a plane in the case of two or more independent variables).
- A Regression model can be used when the dependent variable is a quantitative.

(ii) Multiple Linear Regression :-

- Regression models are used to describe relationship between variables by fitting a line to the observed data.
- Regression allows you to estimate how a dependent variable changes as independent variable changes.
- Multiple linear regression is used to estimate the relation between two or more independent variables, and one dependent variable.
- Use multiple linear regression when you want to know,
 - (i) How strong the relationship between two or more independent variable and one dependent variable.
(Ex. How rainfall, temperature, and amount of fertilizer affect crop growth).
 - (ii) The value of dependent variable at a certain value of independent variables.
(Ex. Expected yield of a crop at certain levels of rainfall, temperature and fertilizer addition),

Assumption of multiple linear Regression :-

- Multiple Linear Regression make all assumption as simple linear regression.
- (i) Homogeneity of variance (Homoscedasticity)
- The size of the error in our prediction doesn't change significantly across the values of independent variable.
- (ii) Independence of observations.
 - Observation in the datasets were collected using statistically valid method
 - There is no hidden relationship among these variables

- In multiple linear regression, it is possible that some of independent variable highly correlated with another, so it is important to check these before developing regression model.
- If two variables are too highly correlated ($r^2 > \sim 0.6$), then only one of them should be used in regression model.

(iii) Normality:

- The data follows a Normal Distribution

(iv) Linearity:

- The line of best fit through the data points is a straight line rather than a curve or some sort of grouping factor.

Perform a multiple linear Regression

- Formula for multiple linear regression is,

$$Y = B_0 + B_1 X_1 + B_2 X_2 + \dots + B_n X_n + \epsilon$$

Y = predicted value of dependent variable.

B_0 = y -intercept (value of y when all other parameters are set to 0).

$B_1 X_1$ = Regression coefficient (B_1) of the first independent variable (X_1) [effect that increasing value of independent variable has on the predicted y value]

$\dots =$ same for however many independent variables you are testing.

$B_n X_n$ = Regression coefficient of the last independent variable.

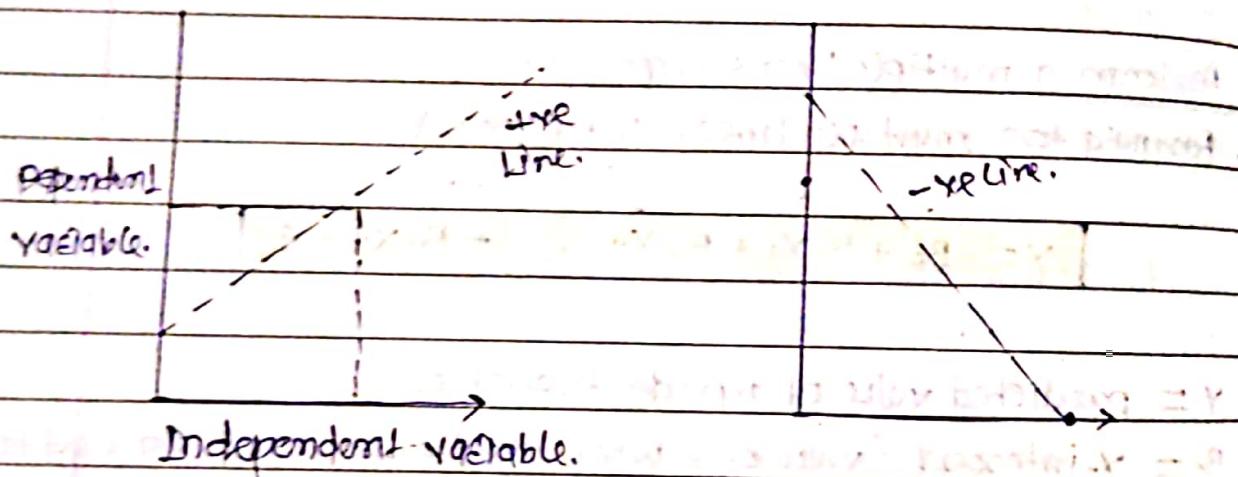
ϵ = model error [how much variability there is our estimate of y ,

- To find best-fit line for each independent variable, multiple linear regression calculates three things,

- (i) Regression coefficient lead to the smallest overall model.
 - (ii) t-statistics of overall model.
 - (iii) Associated p-value (how likely it is that t-statistics have occurred by chance if the Null hypothesis of no relationship between independent and dependent variables was true).
- It then calculates t-statistics and p-value for each regression coefficient in the model.

Algorithm:

- suppose you have dataset independent variable on the x-axis, dependent variable on y axis.

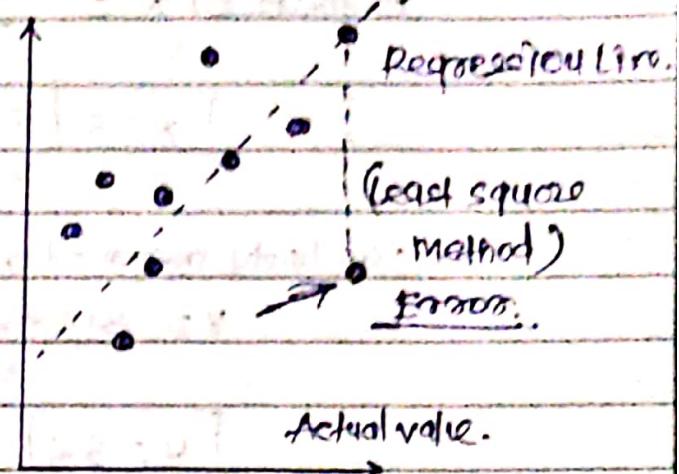
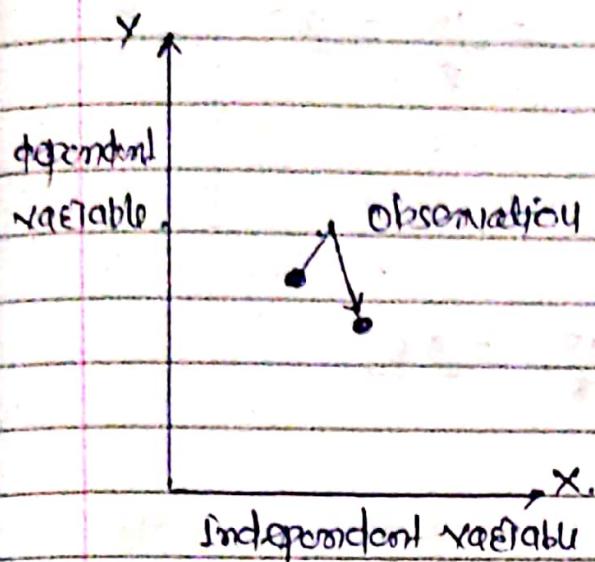


- Independent variable on the x-axis (increasing) and does the dependent variable on the y axis you would get positive linear regression.
- Independent variable on the x-axis (increasing) and other had to dependent variable on the y axis (decreasing) that is Negative linear as a slope of the line is Negative and the particular line that is

$$y = mx + c$$

of linear expression, which shows the relationship between independent variable and dependent variable and this only function is line of linear regression.

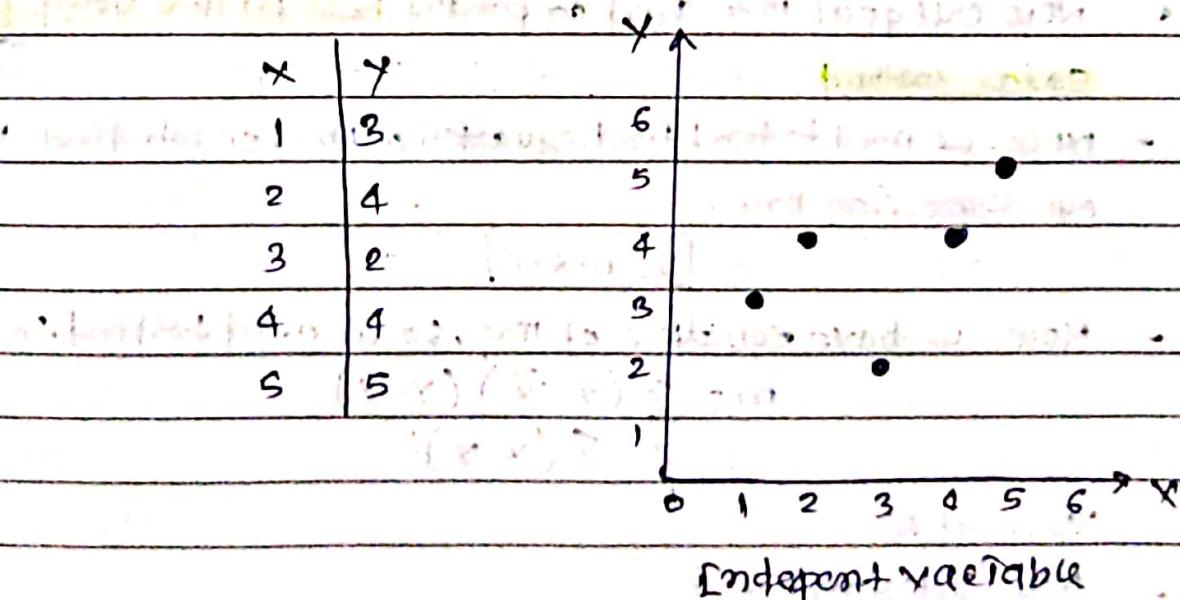
- some data points come to our graph.



- All our data points are plotted now our task is to recall Regression line, or the best fit line.
- Now we get estimated, predicted and Actual value.
- Our main goal is to reduce this error that is difference between estimated value and Actual value, or predicted value.

Mathematical Implementation of:

- consider Data points $x = \{1, 2, 3, 4, 5\}$ on the x axis and $y = \{3, 4, 2, 1, 5\}$ on the y axis,



- Now calculate mean of x and y and plot it on graph,
so mean of $x = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$

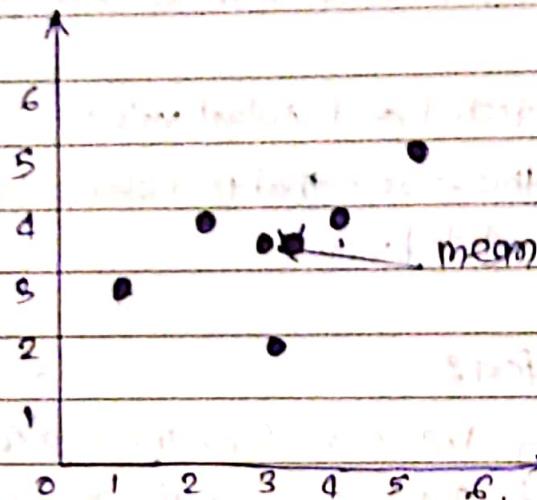
$$x = 3$$

- similarly mean of y ,

$$y = \frac{3+4+2+4+5}{5} = \frac{18}{5} = 3.6$$

$$y = 3.6$$

- Now plot a graph using mean value,



- Now our goal is to find or predict best fit line using Least Square method.

- Now, we need to first find equation of line, so lets find the equation of regression line.

$$y = mx + c$$

- Now, we have equation of line, so we need to find line,

$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

value of m .

- Now calculate $x - \bar{x}$

x	y	(x - \bar{x})	(y - \bar{y})	(x - \bar{x})^2
1	3	(1 - 3) = -2	(3 - 3.6) = -0.6	4
2	4	(2 - 3) = -1	(4 - 3.6) = 0.4	1
3	2	(3 - 3) = 0	(2 - 3.6) = -1.6	0
4	4	(4 - 3) = 1	(4 - 3.6) = 0.4	1
5	5	(5 - 3) = 2	(5 - 3.6) = 1.4	4

mean 3 3.6

$$\sum = 10 \quad \sum = 18$$

- Now we need product of $(x - \bar{x})$ and $(y - \bar{y})$

$$(x - \bar{x}) \cdot (y - \bar{y})$$

$$(-2) \cdot (-0.6) = 1.2$$

$$(-1) \cdot (0.4) = -0.4$$

$$(0) \cdot (-1.6) = 0$$

$$(1) \cdot (0.4) = 0.4$$

$$(2) \cdot (1.4) = 2.8$$

- Now calculate sum of last two columns.

$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})} = \frac{4}{10} = 0.4$$

- Now, we have x as 3.6 mean of $m = 0.4$.

Now $y = m x + c$

$$3.6$$

Now, $m = 0.4$

$$3.6 = 0.4 * 3 + c$$

$$c = 2.4$$

$$3.6 = 1.2 + c$$

$$y = 0.4x + 2.4$$

$$3.6 - 1.2$$

so this is the Regression line.

$$c = 2.4$$

Now, $X = \{1, 2, 3, 4, 5\}$.

$$m = 0.4$$

$$c = 2.4.$$

Now, calculate predicted value for y .

$$Y = mx + c$$

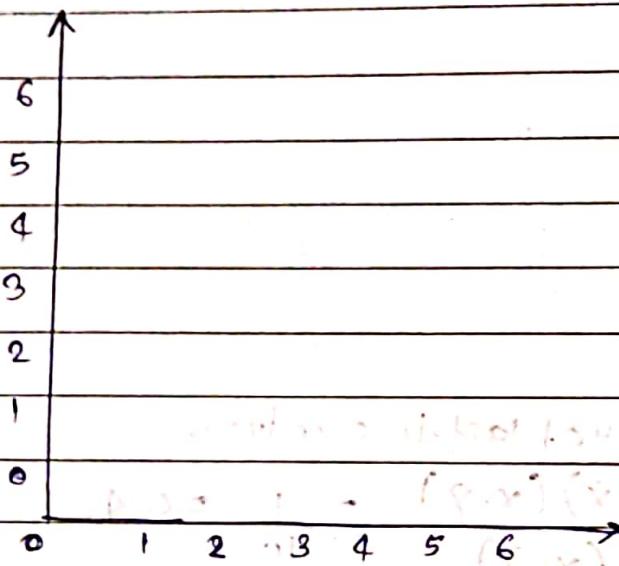
$$Y = 0.4 * 1 + 2.4 = 2.8$$

$$Y = 0.4 * 2 + 2.4 = 3.2$$

$$Y = 0.4 * 3 + 2.4 = 3.6$$

$$Y = 0.4 * 4 + 2.4 = 4.0$$

$$Y = 0.4 * 5 + 2.4 = 4.4$$



- Now calculate distance between actual value and predicted value and reduce distance. (Least square method).
- so it will perform N number of iteration for different values. It will calculate equation of line.

$$[Y = mx + c]$$

- so value of m changes the line is changing. If it changes from one and it will perform n number of iteration, after iteration it will calculate the predicted value.

- According to Libre and compodo distance of Actual values to the predicted value, and value of M for which the distance between Actual and predicted value is minimum will be selected as best fit line.
- Now we have calculated the best fit line, now its time to check goodness of fit or check how good our model performing.
- In order to do that we have a method called **R-squared method**.
 [R squared method use for check goodness of fit in Regression line]

What is R squared ?

- R-squared value is a statistical measure of how close data are to the fitted regression line.
- It also known as coefficient of determination or coefficient of multiple determination.
- In general it is considered as that high-square value model is a good model but you can also have a lower squared value for a good model as well as higher squared value for a model that does not fit at all.

calculation of R^2 :

- $\text{Distance} = \text{Actual} - \text{mean}$ or $\text{Distance} = \text{predicted} - \text{mean}$.
 this is nothing but R^2 .

$$R^2 = \frac{\sum (Y_p - \bar{Y})^2}{\sum (Y - \bar{Y})^2}$$

Where,

y is Actual value.

Y_p is predicted value.

\bar{Y} is mean value of y (S.E.).

x	y	$y - \bar{y}$	$(y - \bar{y})^2$	y_p	$(y_p - y)$	$(y_p - y)^2$
1	3	-0.6	0.36	2.8	-0.8	0.64
2	4	-0.4	0.16	3.2	-0.4	0.16
3	2	-1.6	2.56	3.6	0	0
4	4	0.4	0.16	4.0	-0.4	0.16
5	5	1.4	1.96	4.4	0.8	0.64
	3.6		$\Sigma(5.2)$		$\Sigma(1.6)$	$\Sigma 1.6$

$$R^2 = \frac{\Sigma(y_p - \bar{y})^2}{\Sigma(y - \bar{y})^2} = \frac{\Sigma 1.6}{\Sigma 5.2} = 0.3$$

- this is not a good fit. It suggests that the data points are far away from Regression line.
- When increase value of R^2 to 0.7 the actual value would lie closer the Regression line, $R^2 \approx 0.7$
- When it comes to $R^2 \approx 0.9$ it comes to more close, and value approximately equals to 1 , then actual values lies on Regression line itself.

For Example:

- In this case if you get very low value of R squared suppose $R^2=0.2$ so in that case, actual value are far away from Regression line
- there are outliers present in the data.

RSS (Residuals).

- Residual for a point in the data is the difference between the actual value and the value predicted by our linear regression model.

$$\text{Residual} = \text{actual} - \text{predicted} = y - \hat{y}$$

- Residual plot tell us whether regression model is right fit or not (data).
- It is actually an assumption of ~~data~~ Regression model that there is no trend in Residual plot.
- Using residual value, we can determine sum of square of residuals, also known as residual sum of square or RSS.

$$RSS = \sum (y_i - \hat{y}_i)^2$$

- Lower value of RSS, the better is model prediction. or
- Regression line fits a line of best fit if it minimizes the RSS value.
- The value depends on a scale of the target variable.

Total sum of squares :

- TSS or total sum of squares gives the total variation in Y.
- similar to variance of Y,
- variance is a average of sum of difference between actual values and data points.

calculated R-squared :

$$R\text{-squared} = \frac{(TSS - RSS)}{TSS}$$

$$= \frac{\text{Explained variation}}{\text{Total variation}}$$

$$= \frac{1 - \text{Unexplained variation}}{\text{Total variation}}$$

Adjusted R-squared :

- Adjusted R-squared takes into account the number of independent variable used for predicting the target variable.
- so, determine whether adding new variable to the model actually increases the model fit.

$$\text{Adjusted } R^2 = \left\{ 1 - \frac{[(1-R^2)(n-1)]}{(n-k-1)} \right\}$$

where,

$n \rightarrow$ number of data points in our dataset.

$k \rightarrow$ number of independent variable.

$R \rightarrow$ R-squared value determined by model.

- IF R-squared does not increase significantly on the addition of new independent variable, then value of Adjusted R-squared will actually decrease.

$$\text{Adjusted } R^2 = \left\{ 1 - \frac{[(1-R^2)(n-1)]}{(n-k-1)} \right\}$$

- If on adding new independent variable we see significant increase R-squared value, then Adjusted R-squared value will increase.

$$\text{Adjusted } R^2 = \left\{ 1 - \frac{[(1-R^2)(n-1)]}{(n-k-1)} \right\}$$