

## Importing Libraries

### # Importing libraries

import Libname as subname

#### Syntax

~~Ex~~  
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt

### # Importing dataset:

① In supervised machine learning } means : In which model the program will know ~~label~~ about the data set or we can say that the data set is given and the programmer knows about that the output }.

→ import pandas as pd

```
dataset = pd.read_csv('file.csv')
X = dataset.iloc[:, :-1].values
Y = dataset.iloc[:, -1].values
```

\* pd.read\_csv : CSV = formatted data set

- clipboard ; - excel ; - feather ; - fast ; - gbo
- h5 ; - html ; - json ; - msgpack ; - parquet
- pickle ; - sas ; - SQL ; - ~~sqlite~~
- SQL - query ; - SQL - table
- Stata ; - tablib

\* Now the whole data set is divided into two part

X and Y :

In which X is the data set which contains the data which is used to train the model. OR in other words the data which X helps to find out result

And the other one Y is the data set which is used to test the ~~data~~ model. It means after training the model then check the result for the ~~data~~.

\* Why do we create X and Y separately?

Because we want to work with Numpy arrays instead of pandas dataframes. Numpy arrays are the most convenient format to work with when doing data preprocessing and building machine learning models. So we create two separate arrays.

\* iloc :-

It locates the column by its index

\* values :-

It return the values of the columns you are taking inside a Numpy arrays

Q Explain indexing in X, Y arrays

: a, a<sub>1</sub>; b, b<sub>1</sub> → denote lines

: c, c<sub>1</sub>; d, d<sub>1</sub> → denote columns

Now a = denote we take all the lines  
b = denote we do not take this, no "b" line  
similar in columns

c = denote we take all the columns  
d = denote we do not take "d" no of columns

\* Program file must store in that folder which contain the dataset

Program's data set stored in the directory in which the program file save

→ we can change the directory on clicking file explorer which is present in the right side of the spyder

## # Finding the missing values

- \* For finding missing value and fixed it by using sklearn library

→ `from sklearn.preprocessing import Imputer  
imputer = Imputer(missing_values='NaN',  
strategy='mean', axis=0)`

`imputer = imputer.fit(X[:, 1:3])  
X[:, 1:3] = imputer.transform(X[:, 1:3])`

help this

b) strategy = 'mean' { default = 'mean' }

by the help of strategy we find the missing value by applying 'mean', 'median', 'median', 'most-frequent'

i) mean:- If replace missing values using the mean along the axis

ii) median: ----- " ----- median

iii) most-frequent: ----- " ----- most-frequent

c) axis = 0 { default = "0" }

If axis = 0 then strategy applied along columns  
If axis = 1 the strategy applied along lines

d) verbose { default = 0 }

\* controls the verbosity of the imputn

e) copy { default = True }

f) What is the difference between fit and transform?

fit:-

It is used to extract some info of the data and on which the object is applied

→ `missing_values = 'NaN'`

The place holder for the missing values. All the blank place replaced by 'nan' by the

transform:-

It is used to apply some transforms

## # Dealing with Categorical data:-

- \* In this section: the machine learning knows only deals with integer values or 0 or 1) not string, character & Name etc so

first convert the character data into integer  
then convert into the con form like 0 and 1

→ from sklearn.preprocessing import LabelEncoder  
labelencoder = X = LabelEncoder()  
① — X[:, 0] = labelencoder — X . fit\_transform (X[:, 0])

After processing, Eq-1: We get some digits like 0, 1, 2 and so on so (However since one is greater than zero and two is greater than 1 the equation in the model will think (according to lecture example Spain has a higher value than Germany and France and Germany has a higher value than France) They are only a 3 categories not a relational order between them)

For solving this the program can create Dummy Variable

→ from sklearn.preprocessing import OneHotEncoder  
one hot encoder = OneHotEncoder (categorical\_features  
one hot encoder = [0])  
X = one hot encoder . fit\_transform (X) . toarray ()

→ LabelEncoder -  $Y = \text{LabelEncoder}()$   
 $Y = \text{LabelEncoder} - Y \cdot \text{fit\_transform}$

Q What do the two 'fit\_transform' methods do?

→ When the 'fit\_transform()' method is called from the 'LabelEncoder()' class, it transforms the categories strings into integers.

Example:

It transforms France, Spain, and Germany into 0, 1, and 2

→ When the 'fit\_transform()' method is called from the 'OneHotEncoder()' class, it creates separate columns for each different label with binary values 0 and 1. Those separate columns are the dummy variables.

\* Categorical data - feature:

in which column  
the fit\_transform is applied

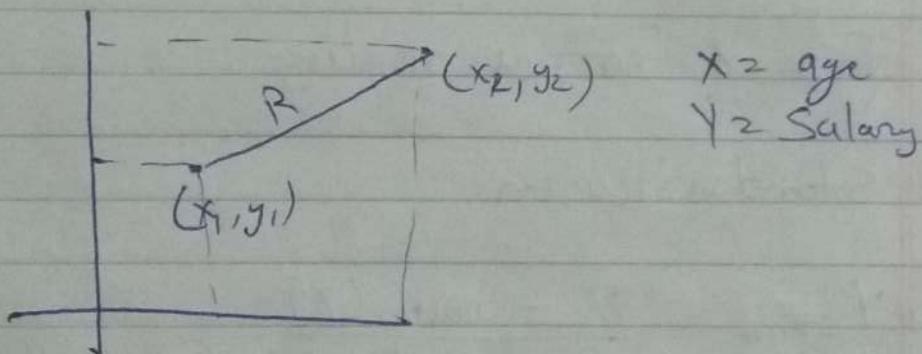
# Splitting the dataset into the training set and test set:

→ `from sklearn.cross_validation import train_test_split  
X_train, X_test, Y_train, Y_test = train_test_split  
(X, y, test_size = 0.25, random_state = 0)`

# Feature scaling :-Q Why <sup>we</sup> use Feature scaling:-

In given data set the age variable (27 to 50) and the Salary variable (48k to 90k) are not on the same scale.

The machine learning model and a lot of machine learning models are based on Euclidean method

Euclidean method

$$R = \text{Euclidean distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

The Euclidean distance will be dominated by the salary

Example

$$(48, 75000) ; (27, 48000)$$

$$R = \sqrt{x^2 + y^2}$$

$$y = 31000 ; y^2 = 961000000 \\ x = 21 ; x^2 = 441$$

difference  
the square of salary ( $y^2$ ) dominates the square difference of age

this is the reason we use feature scaling to remove the domination of one variable to another variable.

The feature scaling helps to optimize the accuracy of your model prediction

or The programmer does not apply feature scaling to keep the most interpretation as possible in your model

## → Types of Feature scaling

### ① Standardisation

$$x_{\text{stand}} = \frac{x - \text{mean}(x)}{\text{standard deviation}(x)}$$

### ② Normalisation

$$x_{\text{norm}} = \frac{x - \text{min}(x)}{\text{max}(x) - \text{min}(x)}$$

By applying feature scaling the speed of execution of program increases

→ from sklearn.preprocessing import StandardScaler  
 $\text{sc\_x} = \text{StandardScaler}()$   
 $\text{x\_train} = \text{sc\_x}.fit\_transform(\text{x\_train})$  → complete  
 $\text{x\_test} = \text{sc\_x}.transform(\text{x\_test})$

## Complete Data processing code :-

# Data processing

# Import the libraries

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import pandas as pd
```

# Importing the dataset

```
dataset = pd.read_csv('Data.csv')
```

```
X = dataset.iloc[:, :-1].values
```

```
Y = dataset.iloc[:, 3].values
```

# Taking care of missing data

```
from sklearn.preprocessing import Imputer
```

```
imputer = Imputer(missing_value='NaN', strategy='mean', axis=0)
```

```
imputer = imputer.fit(X[:, 1:3])
```

```
X[:, 1:3] = imputer.transform(X[:, 1:3])
```

# Encoding Categorical data

```
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
```

```
labelencoder_X = LabelEncoder()
```

```
X[:, 0] = labelencoder_X.fit_transform(X[:, 0])
```

```
onehotencoder = OneHotEncoder(categorical_features=[0])
```

```
X = onehotencoder.fit_transform(X).toarray()
```

```
labelencoder_Y = LabelEncoder()
```

```
Y = labelencoder_Y.fit_transform(Y)
```

# splitting the dataset into the training set and test set

```
from sklearn.model_selection import train_test_split  
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=0)
```

## # Feature scaling

```
from sklearn.preprocessing import StandardScaler  
sc_X = StandardScaler()  
X_train = sc_X.fit_transform(X_train)  
X_test = sc_X.transform(X_test)
```

## Part:-2 Regression

### ① Simple linear regression:-

Q What are exactly the coefficients  $b_0$  and  $b_1$  in the SLR equation:

$$\text{Salary} = b_0 + b_1 \text{Experience}$$

Where  $b_0$  and  $b_1$  are the constant coefficient

Salary = dependent variable

Experience = independent variable

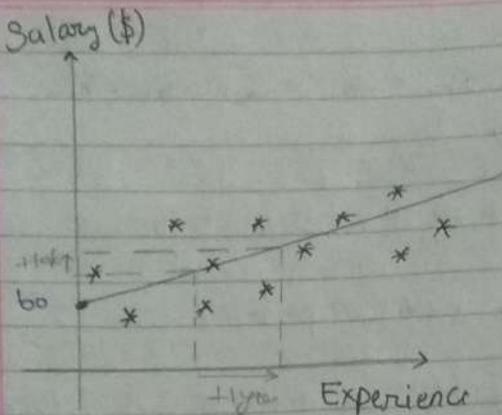
Q Why didn't we apply feature scaling in our SLR model?

It's simply because since  $y$  is a linear combination of the independent variables, the coefficients can adapt their scale to put everything on the same scale.

For example if you have two independent variable  $x_1$  and  $x_2$  and if  $y$  takes values between 0 and 1  $x_1$  takes values between 1 and 10 and  $x_2$  takes values between 10 and 100, then  $b_1$  can be multiplied by 0.1 and  $b_2$  can be multiplied by 0.01 so that  $y$ ,  $b_1 \cdot x_1$  and  $b_2 \cdot x_2$  are all on the same scale.

Q What does 'regressor.fit(X\_train, Y\_train)' do exactly?

The fit method will take the values of  $X_{\text{train}}$  and  $Y_{\text{train}}$  and then will compute the coefficients  $b_0$  and  $b_1$  of the simple linear regression equation ( $y = b_0 + b_1 x$ ) as seen in the.



$$\text{Salary} = b_0 + b_1 * \text{Experience}$$

Complete Simple Linear Regression code

### # Simple Linear Regression

#### # Importing the libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

#### # importing the dataset

```
dataset = pd.read_csv('Salary_Data.csv')
X = dataset.iloc[:, :-1].values
Y = dataset.iloc[:, -1].values
```

```
# splitting the data set into the train-test split
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=1/3, random_state=0)
```

```
# fitting Simple linear regression model
from sklearn.linear_model import LinearRegression
re = LinearRegression()
re.fit(X_train, Y_train) simple linear Model
```

```
# predicting the test set result
Y_pred = re.predict(X_test)
```

# plotting the training set result

```
plt.scatter(X_train, Y_train, color='red')
plt.plot(X_train, re.predict(X_train), color='blue')
plt.title('Salary vs Experience (train set)')
plt.xlabel('Year of Experience')
plt.ylabel('Salary')
plt.show()
```

# plotting the test set result

```
plt.scatter(X_test, Y_test, color='red')
plt.plot(X_train, re.predict(X_train), color='blue')
plt.title('Salary vs Experience (test set)')
plt.xlabel('Year of Experience')
plt.ylabel('Salary')
plt.show()
```

① plt.scatter :-

20

this helps to plot points on a graph like (2.2, 39891) or (2.9, 56642)

② plt.plot :-

this helps to draw a linear line which follows this equation  $y = a_0 + a_1 x$

27/01/19

POW

533  
6/17  
1/40

Day 13

Plt·title :-

This is used to provide the title.

Plt·(x\_label, y\_label) :-

helps to label the x and y coordinates respectively.

# Linear Regression function :-  
Parameters :-(1) fit\_intercept :- Default : True (only boolean value)  
(optional)

Whether to calculate the intercept for this model. If set to false, no intercept will be used in calculation (e.g. data is expected to be already centered).

(2) normalize :- Default : false (only boolean value)  
(optional)  
\* Parameter ignored when fit\_intercept is set to false.  
\* If true, the regressors X will be normalized before regression by subtracting the mean and dividing by the L2-norm.(3) copy\_X : Boolean, optional, default = true  
If True, X will be copied; else, it may be overwritten.(4) n\_jobs :- int or None, optional (default = None)  
The number of jobs to use for the computation. This will only provide speedup for n\_targets > 1 and sufficient large problems.2. Multiple linear regression

What are the Multiple linear regression assumption in more details :-

Linearity :- There must be a linear relationship b/w the dependent variable and the independent variable. Scatterplot can show whether there is a linear or curvilinear relationship.Homoscedasticity :- This assumption states that the variance of error terms is similar across the values of the independent variables. A plot of standardised residuals versus predicted values can show whether points are equally distributed across all values of the independent variables.Multivariate Normality :- MLR assumes that the residuals (total difference b/w predicted value and observed value) are approximately normally distributed.

Q

How is the coefficient  $b_0$  related to the dummy variable trap?

Since  $D_2 = 1 - D_1$ , then if you include both  $D_{0,1}$  and  $D_2$  you get:

$$\begin{aligned}y &= b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 D_1 + b_5 D_2 \\&= b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 D_1 + b_5 (1 - D_1) \\&= b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + D_1 (b_4 - b_5) \\&= b_0^* + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4^* D_1,\end{aligned}$$

with  $\begin{cases} b_0^* = b_0 + b_5 \\ b_4^* = b_4 - b_5 \end{cases}$

Q

How can we select the important variable which helps to form best model and locking.

- \* In MLR many independent variables are present in data set; so we have to find the important independent variables which really helps to predict the good result in less the time taken for predicting the result.

If the programmer takes the complete dataset independent variable then the program locks like a geyser and the model takes time to predict the result

Must watch video no. 44  
no. 67

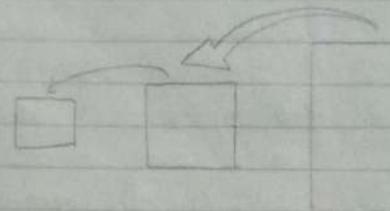
Q

Name the method of building models of MLR

5 Methods of building models

- ① All-in
  - ② Backward Elimination
  - ③ Forward Elimination
  - ④ Bidirectional Elimination
  - ⑤ Score comparison
- Stepwise regression

# Backward Elimination :-



Step

1. Select a significance level to stay in the model (e.g. SL = 0.05)

2. Fit the full model with all possible predictors

Example

```
import statsmodels.formula.api as sm
x = np.append(ann = np.ones((50, 1)).astype(int),
              values = x, axis = 1)
```

$x_{opt} = x[:, [0, 1, 2, 3, 4, 5]]$

$re\_OLS = sm.OLS(endog = Y, exog = x_{opt}).fit()$

$re\_OLS = summary$

- >3 Consider the predictor with the highest P-value  
If  $P > SL$ , go to Step 4, otherwise go to 1.
- 4 Remove the predictor.
- 5 Fit model with out this variable.

# Code:

# Importing the libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

# Importing the data set

```
dataset = pd.read_csv('50_Startups.csv')
X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, -1].values
```

# Importing dealing with categorical data

from sklearn.preprocessing import LabelEncoder, OneHotEncoder

L\_X = LabelEncoder()

X[:, 3] = L\_X.fit\_transform(X[:, 3])

O = OneHotEncoder(categorical\_features=[3])

X = O.fit\_transform(X).toarray()

# Avoiding the dummy variable trap

X = X[:, 1:]

# splitting the dataset into train-test-split

from sklearn.model\_selection import train\_test\_split

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=1/3, random\_state=0)

# fitting the multiple linear regression

from sklearn.linear\_model import LinearRegression

re = LinearRegression()

re.fit(X\_train, y\_train)

# Predict the value:  
 $y_{pred} = \text{re.predict}(x\_test)$

# Building the optimal model using backward Elimination

import statsmodels.formula.api as sm  
 $X = np.append(\text{array}([1]), \text{array}(X), axis=1)$   
 $\text{value} = X, axis=1$

①  $X\_opt = X[:, [0, 1, 2, 3, 4, 5]]$   
 $re\_OLS = sm.OLS(endog=y, exog=X\_opt).fit()$   
 $re\_OLS.summary()$

②  $X\_opt = X[:, [0, 1, 2, 4, 5]]$   
 $re\_OLS = sm.OLS(endog=y, exog=X\_opt).fit()$   
 $re\_OLS.summary()$

③  $X\_opt = X[:, [0, 1, 4, 5]]$   
 $re\_OLS = sm.OLS(endog=y, exog=X\_opt).fit()$   
 $re\_OLS.summary()$

④  $X\_opt = X[:, [0, 4, 5]]$   
 $re\_OLS = sm.OLS(endog=y, exog=X\_opt).fit()$   
 $re\_OLS.summary()$

⑤  $X\_opt = X[:, [4, 5]]$   
 $re\_OLS = sm.OLS(endog=y, exog=X\_opt).fit()$   
 $re\_OLS.summary()$

\* Only 4 and 5 attributes have p-value is less than 5%

\* toarray() :- Return a dense ndarray representation of this matrix

## # Polynomial Regression:

① Is Polynomial Regression a linear or non-linear model?

That depends on what you are referring to. Polynomial Regression is linear on the coefficients since we have only power of the coefficients (all the coeff have the power 1:  $b_0, b_1, \dots, b_n$ ). However, polynomial regression is a non-linear function of the input  $X$ , since we have the inputs raised to several powers:  $x, x^2, \dots$ . That is how we can also see the polynomial regression as a non-linear model.

② Why didn't we apply feature scaling in our model?

It's simply because, since  $y$  is a linear combination of  $X$  and  $X^2$ , the coeff. can adapt their scale to fit everything on the same scale. For example if  $y$  takes values between 0 and 1,  $X$  takes values between 1 and 10 and  $X^2$  takes values between 1 and 100, then  $b_1$  can be multiplied by 0.01 and  $b_2$  can be multiplied by 0.01 so that  $y, b_0x_1$  and  $b_2x_2$  are all on the same scale.

3. How do we find the best degree?

Test several degrees and you see which one gives you the best fit.

4. What's the reason behind to create a second linear regression 'lin\_rgy\_2'?

Because 'lin\_rgy' is already fitted to  $X$  and  $y$  and now we want to fit a new linear model to  $X_{-poly}$  and  $y$ . So we have to create a new regressor object.

Why we create linear-regressor model in polynomial linear regression?

This is because we can create / compare the outcomes which comes from both different model.

Why we create  $X_{-poly}$  variable?

Because we ~~transform~~ transform the variable ' $X$ ' into the  $X^n$  where  $n$  is the degree, which is feeded by the programme according to the best fitting.

After applying the fit-transform method  
index 0 is a dummy variable, index 1 is the  $X$  variable, and index 2 is the  $X^n$  variable

\* Feature scaling apply only when the algorithm uses must based on Euclidean distance.

Raw Code:

# Polynomial regression

# importing the libraries

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

# importing the dataset

```
data = pd.read_csv('Position_Salaries.csv')
```

```
X = data.iloc[:, 1:2].values
```

```
y = data.iloc[:, 2].values
```

# Fitting the linear regression to the dataset

```
from sklearn.linear_model import LinearRegression
```

```
lin_regr = LinearRegression()
```

```
lin_regr.fit(X, y)
```

# Fitting the Polynomial Regression to the dataset

```
from sklearn.preprocessing import PolynomialFeatures
```

```
poly_reg = PolynomialFeatures(degree = 2)
```

```
poly_reg.fit(X_poly, y)
```

```
X_poly = poly_reg.fit_transform(X)
```

```
lin_regr_2 = LinearRegression()
```

```
lin_regr_2.fit(X_poly, y)
```

# Visualising the linear Regression result

```
plt.scatter(X, y, color = 'red')
```

```
plt.plot(X, lin_regr.predict(X), color = 'blue')
```

```
plt.title('')
```

```
plt.xlabel('')
```

```
plt.ylabel('')
```

plt.show()

# Visualising the polynomial regression result

```
plt.scatter(X, y, color = 'red')
```

```
plt.plot(X, lin_regr_2.predict(poly_reg.fit_transform(X)), color = 'blue')
```

```
plt.title('')
```

```
plt.xlabel('')
```

```
plt.ylabel('')
```

```
plt.show()
```

# Visualising the polynomial regression result (for high quality)

```
X_grid = np.arange(min(X), max(X), 0.1)
```

```
X_grid = X_grid.reshape(len(X_grid), 1)
```

```
plt.scatter(X, y, color = 'red')
```

```
plt.plot(X_grid, lin_regr_2.predict(poly_reg.fit_transform(X_grid)), color = 'blue')
```

```
plt.title('')
```

```
plt.xlabel('')
```

```
plt.ylabel('')
```

```
plt.show()
```

# predicting the result by linear regression

```
lin_regr.predict(6.5)
```

# predicting the result by Polynomial regression

```
lin_regr_2.predict(poly_reg.fit_transform(6.5))
```

## # Polynomial Feature function:

① Generates X-poly matrix

Parameters:

② degree: integer Default = 2  
The degree of the PF

③ interaction\_only: boolean Default = False  
If true, only interaction features are produced.  
Features are produced: features that are product  
of at most degree distinct input features

## # SVR : Support vector regression

1. When should i use SVR?

You should use SVR if a linear model like linear regression doesn't fit very well your data. This would mean you are dealing with a non-linear problem.

Why do we need to 'sc\_y inverse\_transform'?

We need the inverse-transform method to go back to the original scale. Indeed we applied feature scaling so we get this scale around '0' and if we make a prediction without inverting the scale we will get the scaled predicted salary. And of course we want the real salary, not the scaled one, so we have to use inverse transform method.

- \* Support vector machines (SVM) support linear and non-linear regression that we can refer to as SVR
- \* SVR performs linear regression in a higher dimensional space
- \* In SVR each point in train set represent its own dimension
- \* In SVR we can't apply limit in the margins of  $b$  which is  $\frac{y_i - y_j}{2}$  between the original support vector regression line and the line pass from the element of two different cluster's element. The element is the nearest element of both the cluster. We can limit the margin by Epsilon

## Building SVR Model

- Step 1 (optional) Collect a training set  $T = \{(\vec{x}, \vec{y})\}$
- ② Choose a kernel and its parameters as well as any regularization method.
  - ③ Form the correlation matrix,  $K$
  - ④ Train your machine, exactly or approximately, to get contraction coefficients  $\{q_i\}_{i=1}^n$
  - ⑤ Use these coefficients, create your estimator  

$$f(\vec{x}, \vec{q}, \vec{x}') = \vec{y}$$
- \* Kernel are
- ① Gaussian
  - ② Regularization
  - ③ Noise

Raw code:-

```
# SVR
# importing the libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# importing the dataset
dataset = pd.read_csv('Position_Salary.csv')
X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, -1].values
```

```
# Feature Scaling:
from sklearn.preprocessing import StandardScaler
sc_x = StandardScaler()
```

```
sc_y = StandardScaler()
X = sc_x.fit_transform(X)
y = sc_y.fit_transform(y)
```

# Fitting SVR model.

```
from sklearn.svm import SVR
r = SVR(kernel='rbf')
r.fit(X, y)
```

# predicting a new result

```
y_pred = r.predict(6.5)
```

```
y_pred = sc_y.inverse_transform(y_pred)
```

# visualising the graph

## # Decision Tree Regression

Q

How does the algo split the data points?

It uses reduction of standard deviation of the predictions. In other words, the standard deviation is decreased right after a split. Hence building a decision tree is all about finding the attribute that returns the highest standard deviation reduction.

Q

What is the Information Gain and how does it work in Decision tree?

The information gain in Decision tree is exactly the standard Deviation Reduction we are looking to reach. We calculate by how much the standard deviation decreases after each split. Because the more the standard deviation is decreased after a split; the more homogeneous the child nodes will be.

Q

What is the Entropy and how does it work in Decision trees?

The Entropy measures the disorder in a set, here in a part resulting from a split. So the more homogeneous is your data in a part, the lower will be the entropy. The more you have splits, the more you have chance to find parts in which your data is homogeneous, and therefore the lower will be the entropy (close to 0) in these parts. However you might still find some nodes where the data is not homogeneous, and therefore the entropy would not be that small.

- \* Decision tree is not much sense in 1D.
- \* In Decision tree & data can be divide by applying the splits

## # Code

# importing the libraries

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

# importing the dataset

data = pd.read\_csv('Position\_Salaries.csv')

X = data.iloc[:, 1:-1].values

y = data.iloc[:, -1].values

# Implementing the Decision tree Regression alg.

from sklearn.tree import DecisionTreeRegressor

regressor = DecisionTreeRegressor(random\_state=0)

regressor.fit(X, y)

# predicting the result

y\_pred = regressor.predict(6.5)

~~#~~ # Visualizing the Regression result - (1)

plt.scatter(X, y, color='red')

plt.plot(X, regressor.predict(X), color='blue')

plt.title(' )

plt.xlabel(' )

plt.ylabel(' )

plt.show()

# Visualizing the result (for higher resolution and

1\_grid = np.arange(min(X), max(X), 0.1)

```
x-grid = x-grid.reshape((len(x-grid), 1))
plt.scatter(x, y, color = 'red')
plt.plot(x-grid, r.predict(x-grid), color = 'blue')
plt.title(' ')
plt.xlabel(' ')
plt.ylabel(' ')
plt.show()
```

- \* In x-grid ; 0.1. Explain the resolution . Resolution can be increases by increasing the 0.1 value like 0.01.

## # Random Forest Regression

Q

What is the advantage and drawback of Random Forests compared to Decision Trees?

Advantage :- Random Forests can give you a better prediction power than Decision Trees.

Drawback - Decision tree will give you more interpretability than Random forest, because you can plot the graph of a Decision tree to see the different splits made to the prediction, as seen in the Intuition Lecture. That's something you can't do with Random forest.

Q

When to use Random Forest and when to use the other models?

- \* The best answer to that question is: try them all
- \* First you need to figure out whether your problem is linear or non-linear.
  - a) If the problem is linear, you should go for Simple Linear Regression if you only have one feature and Multiple Linear Regression if you have several features.
  - b) If the problem is non linear, you should go for Polynomial Regression, SVR, Decision tree, and Random Forest.

The method consists of using a very relevant technique that evaluates your model's performance, called k-Fold Cross Validation, and then picking the model that shows the best result. Feel free to

d

How do I know how many trees I should use?  
First, I would recommend to choose the number of trees by experimenting. It usually takes less time than we think to figure out a best value by tweaking and tuning your model manually. That's actually what we do in general when we build a Machine Learning model: we do it in several shots, by experimenting several values of hyperparameters like the no. of trees.

k-Fold Cross validation and Grid search which are powerful techniques that you can use to find the optimal value of a hyperparameter, like here the no. of trees.

## # Code:-

```
# importing the libraries
```

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
# importing the dataset
```

```
data = pd.read_csv('Position_Salaries.csv')
```

```
X = data.iloc[:, 1:2].values
```

```
y = data.iloc[:, 2].values
```

```
# implementing the Random forest model.
```

```
from sklearn.ensemble import RandomForestRegressor
```

```
r = RandomForestRegressor(n_estimators=100,  
random_state=0)
```

```
r.fit(X, y)
```

```
# predicting the results
```

```
y_pred = r.predict(6.5)
```

```
# visualising the result (high resolution)
```

```
X_grid = np.arange(min(X), max(X), 0.1)
```

```
X_grid = X_grid.reshape((len(X_grid), 1))
```

```
plt.scatter(X, y, color='red')
```

```
plt.plot(X_grid, r.predict(X_grid), color='blue')
```

```
plt.title('')
```

```
plt.xlabel('')
```

```
plt.ylabel('')
```

```
plt.show()
```

- \* n-estimators = define the no. of trees
- \* no. of trees ↑ the prediction & result will more accurate

## # R-squared:

[www.endeavocareers.com](http://www.endeavocareers.com)

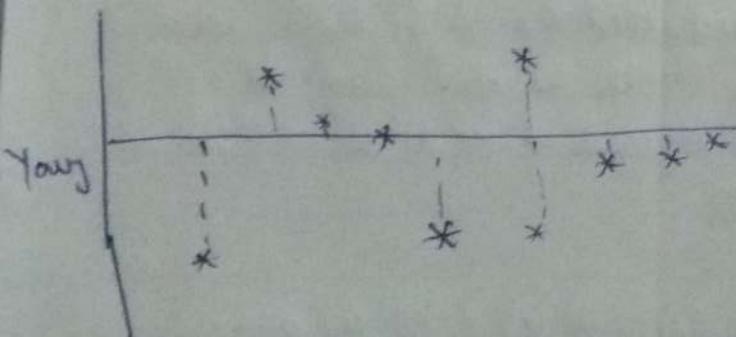
R-squared Explained that how good is your line compared to the you know average line in simple linear regression

$$SS_{\text{res}} = \sum (y_i - \hat{y}_i)^2 = SS_{\text{residual}}$$

$$SS_{\text{tot}} = \sum (y_i - \bar{y}_{\text{avg}})^2 = SS_{\text{total}}$$

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

- \*  $R^2$  can be (-ve) that means direction is wrong and model is broken
- \*  $R^2 = 1$  best case, but not happen • Ideal case



- \*  $R^2$  is closer to 1 it is good for model

"If you can dream it, you can achieve it."

## Hierarchical clustering

POORU

- \* Two types of Hierarchical clustering,

- ① Agglomerative } bottom up approach?
- ② Divisive } Top to bottom approach?

### ① Agglomerative Hierarchical clustering:-

Step<sup>1</sup>

① Make each data point a single-point cluster  $\Rightarrow$  That forms N clusters

② Takes the two closest data points and make them one cluster  $\rightarrow$  That forms  $N-1$  clusters

③ Takes the two closest clusters and make them one cluster  $\rightarrow$  That forms  $N-2$  clusters

④ Repeat step 3 until there is only one cluster

Fin

- \* At the end of all these steps only a large cluster will form.

- \* In each repetition only single combine cluster formed and at the end of this  $N-n$  {when 'n' is the no. of steps repeated}.

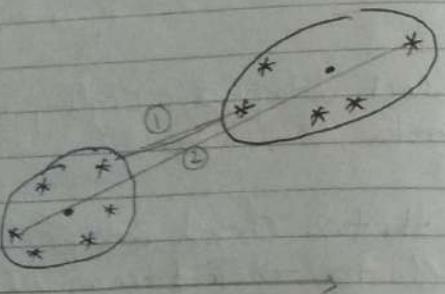
**POW**  
formation  
of cluster  
b/w  
two  
different  
clusters

\* Ways to calculate distance b/w two clusters

- ① closet points
- ② furthest Points
- ③ Average Distance

{sum of all the distance of points lied in two different clusters}

④ Distance B/W centroids



\* Distance b/w two cluster is important

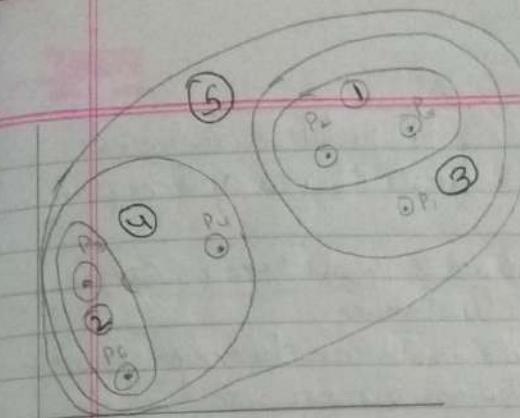
# Dendograms:-

\* Dendogram is kind of like memory of the AC alg.

\* In Dendograms

a) Horizontal line  $\Rightarrow$  connecting one <sup>represent</sup> a one cluster

b) vertical line  $\Rightarrow$  difference of property of two different connecting two cluster



Cluster diagram

Properties  
Difference

Dendrogram

\* Dendrogram can be formed by algo automatically

\* Somm Imp. Questions:-

① What is the point of Hierarchical clustering if it always leads to one cluster per observation point.  
The main point of HC is to make the dendrogram, because you need to start with one single cluster then work your way down to see the different combination of clusters until having a number of clusters equal to the number of observations. And it's the dendrogram itself that allows to find the best clustering configuration.

When you are comparing the distance b/w two clusters or a cluster and a point, how exactly is it measured?  
Are you taking the centroid in the cluster and measuring the distance?

Exactly, the metric is the Euclidean distance b/w the centroid of the first cluster and the point.

- Ques 3. Do we also need to perform feature scaling of Hierarchical clustering?  
Ans. Yes because the equation of the clustering problem involve the Euclidean distance. Anytime the model equations involve the Euclidean distance, you should apply feature scaling.

- Imp. Should we use the dendrogram or the elbow method to find that optimal number of clusters?  
\* You should use both, just to double check the optimal solution.  
\* Elbow order :- elbow method > dendrogram  
\* If you really only have time for one.  
\* Dendrogram is not always the easiest way to find the optimal number of clusters. But with the elbow method it is very easy, since the elbow is most of the time very obvious to spot.

### # Dendrogram:-

```
import scipy.cluster.hierarchy as sch
d = sch.dendrogram(linkage(X, method='ward'))
plt.title('Dendrogram')
plt.xlabel('Customer')
plt.ylabel('Distance')
plt.show()
```

# Ward method

- \* method = 'Ward' :-
- ① For method 'single' an optimized algorithm based on minimum spanning tree is implemented. It has time complexity  $O(n^2)$ .
- ② Other method like 'complete', 'average', 'weighted' or 'ward' an algorithm called nearest-neighbours chain is implemented. It also has time complexity  $O(n^2)$ .
- ③ Remaining methods a naive algorithm is implemented with  $O(n^3)$  time complexity.
- ④ Method 'centroid', 'median', and 'ward' are correctly defined only if Euclidean pairwise metric is used.

If y is passed as precomputed pairwise distance, then it is used. If y is passed as precomputed pairwise distance, then it is a user responsibility to assure that these distance are in fact Euclidean, otherwise the produced result will be incorrect.

How the clusters are formed through the ward method of the agglomerative clustering (?) python class? HC is a general family of clustering algo that build nested clusters by merging or splitting them successively. This hierarchy of clusters is present represented as a tree. The root of the tree is the unique clusters.

that gather all the samples, the leaves being the clusters with only one sample. The agglomerative clustering ( ) class performs HAC using bottom up approach.

The ward method minimizes the sum of squared differences within all clusters. It is a variance-minimizing approach and in this sense is similar to the k-means ~~to the~~ objective function but tackled with an AHC approach.

\* The only purpose of making first the dendrogram is to have an idea of the optimal number of clusters. Then when we apply AHC, we can input the optimal number of clusters that we found thanks to the dendrogram.

\* Affinity:- It refers to the distance; which is the euclidean distance.

### # Elbow method:-

from sklearn cluster import KMeans  
wcss = []

for i in range(1, 11):

kmeans = KMeans(n\_clusters=i, init='  
k-means++', random\_state=42)

kmeans.fit(x)

wcss.append(kmeans.inertia\_)

plt.plot(range(1, 11), wcss)  
plt.title('The Elbow method')  
plt.xlabel('No. of clusters')  
plt.ylabel('WCSS')  
plt.show()

### # K-Means:

Where can we apply clustering alg. in real life?

- ① Market Segmentation
- ② Medicine with for example tumor detection
- ③ Fraud detection
- ④ to simply identify some clusters of your ~~go~~ customer in your company or business

Q How does the perpendicular link trick work when k = 3

\* This trick is mainly used only in 2D or 3D space.

\* In high dimensional data we use Euclidean distance to perform the clustering

\* In very high dimensional data, we could do a trick, which is an "engulfing sphere" i.e. starting at every centroid, start growing a spherical space outward radially. Until the sphere intersects everything the sphere engulfs belongs to the same clusters.

Q) How does the elbow method work when more than 2 features are involved?

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

Q) Name the another method / way to the Elbow method to find the ultimate no. of clusters?

Explaining in detail & can be done through Parameter tuning with teacher grid search.

## # Hierarchical clustering:-

# importing the libraries

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

```
data = pd.read_csv('Mall_Customers.csv')
X = data.iloc[:, [3, 4]].values
```

```
import scipy.cluster.hierarchy as sch
```

```
d = sch.dendrogram(X, linkage(X, method='ward'))
```

```
plt.title('Dendrogram')
```

```
plt.xlabel('Customers')
```

```
plt.ylabel('Distance')
```

```
from sklearn.cluster import AgglomerativeClustering
```

```
hc = AgglomerativeClustering(n_clusters=5, affinity
```

```
= 'euclidean', linkage='ward')
```

```
y_hc = hc.fit_predict(X)
```

```
plt.scatter(X[y_hc == 0, 0], X[y_hc == 0, 1], s=100,
```

```
c='red', label='Cluster 1')
```

```
plt.scatter(X[y_hc == 1, 0], X[y_hc == 1, 1], s=100,
```

```
c='blue', label='Cluster 2')
```

~~plot~~ - title (Creates a custom title)  
~~plot~~ - subtitle (Creates a custom subtitle)  
~~plot~~ - ylabel (Creates a custom y-axis label)  
~~plot~~ - legend ()  
~~plot~~ - show ()

K

12

P

www.ankuravashisth.com /  
KK /

# Par

A priori

What are  
support

Given two  
relations

blw s

Step 1 find support

Step 2 find confidence

Step 3 find lift

\* Why we create "transaction" name list.

## # Part 5 : Association

# A priori :

Q What are three essential relation ship b/w the support, confidence and lift?

Given two movies  $M_1$  and  $M_2$  here are three essential relations to remember :-

Relation b/w support and the confidence

$$\text{confidence}(M_1 \rightarrow M_2) = \frac{\text{support}(M_1, M_2)}{\text{support}(M_1)}$$

b/w lift and the support:-

$$\text{lift}(M_1 \rightarrow M_2) = \frac{\text{support}(M_1, M_2)}{\text{support}(M_1) + \text{support}(M_2)}$$

b/w lift and the confidence:

$$\text{lift}(M_1 \rightarrow M_2) = \frac{\text{confidence}(M_1 \rightarrow M_2)}{\text{support}(M_2)}$$

# Code:

# importing the libraries

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

# importing the dataset

```
data = pd.read_csv("Market_Basket_Optimisation  
.csv", header = None)
```

transactions = []

```
for i in range(0, 750):
```

```
    transactions.append(data.values[i, :])
```

# Training the data

```
from apyori import apriori
```

```
rules = apriori(transactions, min_support=0.03,  
confidence=0.2, min_lift=3; min_lg
```

```
results = list(rules)
```

# The Multi-Armed Bandit Problem

[www.endeavorcareers.com](http://www.endeavorcareers.com)

- We have  $d$ -arms. For Example, arms are ads that we display to users each time they connect to a web page.
- Each time a user connects to this web page that makes a round.
- At each round  $n$ , we choose one ad to display to the user.
- At each round  $n$ , ad  $i$  gives reward  $r_i(n) \in \{0, 1\}$ :
  - $r_i(n) = 1$  if the user clicked on the ad  $i$ ,
  - 0 if the user didn't
- Our goal is to maximize the total reward we get over many rounds.

## ① Upper confidence Bound:

Q Why does a single round can have multiple 1s for different ads?

Each round corresponds to a user connecting to a web page, and the user can only see one ad when he/she connects to the web page, the ad that is being shown to him/her. If he/she clicks on the ad, the reward is 1, otherwise it's 0. And there can be multiple ones because there are several ads that the user would love to click on. But only a crystal ball would tell us which ads the user would click on. And the dataset is exactly that crystal ball.

Q In the first Intuition lecture, could you please explain why  $D_5$  is the best distribution? Why is it not  $D_3$ ?

In this situation, 0 is loss and 1 is gain, or win. The  $D_5$  is the best because it is skewed so we will have average outcome close to 1, meaning that there we have more wins, or gain. And actually all casino machine nowadays are carefully programmed to have distribution like  $D_1$  or  $D_3$ . But it is a good concrete example.

Q Does the UCB strategy really come into effect during the first rounds?

Not at the beginning, first we must have some first insights of the users response to the ads. That is just the

the beginning of the strategy. In real world for the first ten users connecting to the webpage, you would show ad<sub>1</sub> to the first user, ad<sub>2</sub> to the second user, ad<sub>3</sub> to the third user, ..., ad<sub>10</sub> to the 10<sup>th</sup> user. Then the algo starts.

The main purpose of this algo is not only to maximize the revenue but also to minimize the cost.

## \* Upper Confidence Bound algo:-

Step:-

1. At each round  $n$ , we consider two numbers for each ad  $i$ :
  - $N_i(n)$  - the number of times the ad  $i$  was selected up to round  $n$ ,
  - $R_i(n)$  - the sum of rewards of the ad  $i$  upto round  $n$ .
2. From these two numbers we compute
  - the average reward of ad  $i$  upto round  $n$

$$\bar{r}_i(n) = \frac{R_i(n)}{N_i(n)}$$

- the confidence interval  $[\bar{r}_i(n) - \Delta_i(n), \bar{r}_i(n) + \Delta_i(n)]$  at round  $n$  with

$$\Delta_i(n) = \sqrt{\frac{3 \log(n)}{2 N_i(n)}}$$

3. We select the ad  $i$  that has the maximum UCB  
 $\bar{r}_i(n) + \Delta_i(n)$

## # Random Selection

### # importing the Libraries

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

### # importing the dataset

```
data = pd.read_csv('Ads_CTR_Optimisation.csv')
```

### # Implementing Random Selection

```
import random
```

```
N = 10000
```

```
d = 10
```

```
ads_selected = []
```

```
total_reward = 0
```

```
for n in range(0, N):
```

```
    ad = random.randrange(d)
```

```
    ads_selected.append(ad)
```

```
    reward = data.values[n, ad]
```

```
    total_reward = total_reward + reward
```

### # Visualizing the result - Histogram

```
plt.hist(ads_selected)
```

```
plt.title('Histogram of ads selection')
```

```
plt.xlabel('Ads')
```

```
plt.ylabel('Number of times each ad was selected')
```

```
plt.show()
```

\* Result = total\_reward = 1193

# UCB algo ↗ there is no package | Libraries to direct implement UCB)

# importing the libraries

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

# importing the dataset

data = pd.read\_csv('AdsCTROptimisation.csv')

# Implementing UCB

import math

N = 10000

d = 10

ad\_selected = []

numbers\_of\_selections = [0] \* d

sums\_of\_rewards = [0] \* d

for n in range(0, N):

ad = 0

max\_upper\_bound = 0

for i in range(0, d):

if (numbers\_of\_selections[i] > 0):

average\_reward = sums\_of\_rewards[i] / numbers\_of\_selections[i]

delta\_i = math.sqrt(3/2 \* math.log(n+1) / numbers\_of\_selections[i])

upper\_bound = average\_reward + delta\_i

else:

upper\_bound = 1e400

if upper\_bound > max\_upper\_bound:

max\_upper\_bound = upper\_bound.

ads - selected : append (ad)

numbers - f - selections [ad] = numbers - f - selection  
+ 1

reward = data - values [n, ad]

sums - f - rewards [ad] = sums - f - rewards [ad]  
+ reward

total - reward = total - reward + reward

# Visualising the result by histogram

plt - hist (ads - selected )

plt - title ( ' ' )

plt - xlabel ( ' ' )

plt - ylabel ( ' ' )

plt - show()

\* reward / result - total - reward = 2178

\* Thompson Sampling is better than UCB.

## 2) Thompson Sampling:-

- \* Thompson Sampling algo is probabilistic or the UCB is deterministic.

Q Why is the yellow mark the best choice, and not the green mark?

The yellow mark is the best choice because it is the furthest from the origin on the x-axis which therefore means that it has the highest estimated return.

Q I don't understand how Thompson Sampling can accept delayed feedback. Please Explain

When doing Thompson Sampling, we can still perform updates in our algo. While we are waiting for the results of an experiment in the real world. This would not hinder our algo from working. This is why it can accept delayed feedback.

Q Comparison b/w UCB & Thompson algo.

UCB

- Deterministic

- Requires update at every round.

Thompson Sampling

Probabilistic

can accommodate delayed feedback.

- Better empirical evidence

## Thompson sampling algo:-

Step 1 At each round  $n$ , we consider two numbers for each ad  $i$ :

- $N_i^1(n)$  - the number of time the ad  $i$  got reward up to round  $n$
- $N_i^0(n)$  - the number of time the ad  $i$  got reward 0 up to round  $n$

Step 2 For each ad  $i$ , we take a random draw from the distribution below

$$\theta_i(n) = \beta(N_i^1(n) + 1, N_i^0(n) + 1)$$

Step 3 We select the ad that has the highest  $\theta_i(n)$

## # Thompson sampling algo:

```
# Importing the libraries
```

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
# importing the dataset
```

```
data = pd.read_csv('Ads_CTR_Optimization.csv')
```

```
# Implementing Thompson sampling algo
```

```
import random
```

```
N = 10000
```

```
d = 10
```

```
ads_selected = []
```

```
numbers_of_rewards_1 = [0] * d
```

```
numbers_of_rewards_0 = [0] * d
```

```
total_reward = 0
```

```
for n in range(0, N):
```

```
    ad = 0
```

```
    max_random = 0
```

```
    for i in range(0, d):
```

```
        random_beta = random.betavariate(numbers_of_rewards_1[i] + 1, numbers_of_rewards_0[i] + 1)
```

```
        if random_beta > max_random:
```

```
            max_random = random_beta
```

```
            ad = i
```

```
    ads_selected.append(ad)
```

```
    reward = data.values[n, ad]
```

```
    if reward == 1:
```

numbers\_f\_rewards\_1[ad] = numbers\_f\_rewards  
[ad] + 1

else :

numbers\_f\_rewards\_0[ad] = numbers\_f\_rewards  
- 0[ad] + 1

total\_reward = total\_reward + reward

# visualising the result - Histogram

plt.hist(ads\_selected)

plt.title('')

plt.xlabel('')

plt.ylabel('')

plt.show()

\* result = total\_reward = 2613