

```
In [15]: # Expected salary for fresher based on company employers salary
# The scenario is you are a HR officer, you got a candidate with 5 years of experience,
# then what is the best salary you should offer to him?

In [16]: # Importing Packages/Libraries
import numpy as np
from sklearn.linear_model import LinearRegression

In [17]: data = pd.read_csv("E:\Data\SalaryPred\linreq.csv")
data

Out [17]:
   YearsExperience  Salary
0                1.1  39343
1                1.3  46205
2                1.5  37731
3                2.0  43525
4                2.2  38891
5                2.9  56642
6                3.0  60150
7                3.2  54445
8                3.2  64445
9                3.7  57189
10               3.9  63218
11               4.0  55794
12               4.0  56957
13               4.1  57081
14               4.5  61111
15               4.9  67838
16               5.1  66029
17               5.3  83088
18               5.9  81363
19               6.0  93940
20               6.8  91738
21               7.1  98273
22               7.9  101302
23               8.2  113812
24               8.7  109431
25               9.0  105582
26               9.5  116969
27               9.6  112635
28               10.3  122391
29               10.5  121872
```

I. Getting Overview of Dataset

1. getting Datatypes of Columns

```
In [18]: data.dtypes

Out [18]:
YearsExperience    float64
Salary            int64
dtypes: object

In [19]: print(data.info)

Out [19]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 30 entries, 0 to 29
Data columns (total 2 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   YearsExperience  30 non-null      float64
 1   Salary          30 non-null      int64
dtypes: float64(1), int64(1)
memory usage: 612.0 bytes
```

2. getting size of Dataframes(rows x Columns)

```
In [20]: print(data.size)

Out [20]:
(30, 2)
```

3. getting the number of rows and columns in the dataframe

```
In [21]: data.shape

Out [21]:
(30, 2)
```

4. getting the dimentions of the dataframe

```
In [22]: data.ndim

Out [22]:
2
```

5. getting Summary of Dataset

```
In [24]: data.info()

Out [24]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 30 entries, 0 to 29
Data columns (total 2 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   YearsExperience  30 non-null      float64
 1   Salary          30 non-null      int64
dtypes: float64(1), int64(1)
memory usage: 612.0 bytes
```

6. Informaton of Data set

```
In [25]: data.info

Out [25]:
<bound method DataFrame.info of
0                1.1  39343
1                1.3  46205
2                1.5  37731
3                2.0  43525
4                2.2  38891
5                2.9  56642
6                3.0  60150
7                3.2  54445
8                3.2  64445
9                3.7  57189
10               3.9  63218
11               4.0  55794
12               4.0  56957
13               4.1  57081
14               4.5  61111
15               4.9  67838
16               5.1  66029
17               5.3  83088
18               5.9  81363
19               6.0  93940
20               6.8  91738
21               7.1  98273
22               7.9  101302
23               8.2  113812
24               8.7  109431
25               9.0  105582
26               9.5  116969
27               9.6  112635
28               10.3  122391
29               10.5  121872>
```

7.head() Return the first 5 rows of the DataFrame by default.

```
In [26]: data.head()

Out [26]:
   YearsExperience  Salary
0                1.1  39343
1                1.3  46205
2                1.5  37731
3                2.0  43525
4                2.2  38891
```

8. head() Return the first n rows of the DataFrame.

```
In [27]: data.head(10)

Out [27]:
   YearsExperience  Salary
0                1.1  39343
1                1.3  46205
2                1.5  37731
3                2.0  43525
4                2.2  38891
5                2.9  56642
6                3.0  60150
7                3.2  54445
8                3.2  64445
9                3.7  57189
```

9. head() Returning entire size of the DataFrame.

```
In [28]: data.head(n=None)

Out [28]:
   YearsExperience  Salary
0                1.1  39343
1                1.3  46205
2                1.5  37731
3                2.0  43525
4                2.2  38891
5                2.9  56642
6                3.0  60150
7                3.2  54445
8                3.2  64445
9                3.7  57189
10               3.9  63218
11               4.0  55794
12               4.0  56957
13               4.1  57081
14               4.5  61111
15               4.9  67838
16               5.1  66029
17               5.3  83088
18               5.9  81363
19               6.0  93940
20               6.8  91738
21               7.1  98273
22               7.9  101302
23               8.2  113812
24               8.7  109431
25               9.0  105582
26               9.5  116969
27               9.6  112635
28               10.3  122391
29               10.5  121872
```

10.tail() Return the last 5 rows of the DataFrame by default.

```
In [29]: data.tail()

Out [29]:
   YearsExperience  Salary
26               9.5  116969
27               9.6  112635
28               10.3  122391
29               10.5  121872
```

12. tail() Return the last 10 rows of the DataFrame.

```
In [31]: data.tail(10)

Out [31]:
   YearsExperience  Salary
20               6.8  91738
21               7.1  98273
22               7.9  101302
23               8.2  113812
24               8.7  109431
25               9.0  105582
26               9.5  116969
27               9.6  112635
28               10.3  122391
29               10.5  121872
```

13. tail() Returning entire rows of the DataFrame.

```
In [32]: data.tail(n=None)

Out [32]:
   YearsExperience  Salary
0                1.1  39343
1                1.3  46205
2                1.5  37731
3                2.0  43525
4                2.2  38891
5                2.9  56642
6                3.0  60150
7                3.2  54445
8                3.2  64445
9                3.7  57189
10               3.9  63218
11               4.0  55794
12               4.0  56957
13               4.1  57081
14               4.5  61111
15               4.9  67838
16               5.1  66029
17               5.3  83088
18               5.9  81363
19               6.0  93940
20               6.8  91738
21               7.1  98273
22               7.9  101302
23               8.2  113812
24               8.7  109431
25               9.0  105582
26               9.5  116969
27               9.6  112635
28               10.3  122391
29               10.5  121872
```

II. Data Pre-processing

1. Checking NULL Values and count

```
In [33]: print(data.isna())

Out [33]:
YearsExperience  Salary
0             False      False
1             False      False
2             False      False
3             False      False
4             False      False
5             False      False
6             False      False
7             False      False
8             False      False
9             False      False
10            False      False
11            False      False
12            False      False
13            False      False
14            False      False
15            False      False
16            False      False
17            False      False
18            False      False
19            False      False
20            False      False
21            False      False
22            False      False
23            False      False
24            False      False
25            False      False
26            False      False
27            False      False
28            False      False
29            False      False

In [34]: null_val_count_data.isna().sum()

Out [34]:
YearsExperience    0
Salary            0
dtypes: int64 2
```

2. Categorical to Numerical Values

```
In [35]: print(data.isna().astype(int))

Out [35]:
YearsExperience  Salary
0                0      0
1                0      0
2                0      0
3                0      0
4                0      0
5                0      0
6                0      0
7                0      0
8                0      0
9                0      0
10               0      0
11               0      0
12               0      0
13               0      0
14               0      0
15               0      0
16               0      0
17               0      0
18               0      0
19               0      0
20               0      0
21               0      0
22               0      0
23               0      0
24               0      0
25               0      0
26               0      0
27               0      0
28               0      0
29               0      0
```

3. Reading Particular Column data

```
In [36]: data["YearsExperience"]

Out [36]:
0      1.1
1      1.3
2      1.5
3      2.0
4      2.2
5      2.9
6      3.0
7      3.2
8      3.2
9      3.7
10     3.9
11     4.0
12     4.0
13     4.1
14     4.5
15     4.9
16     5.1
17     5.3
18     5.9
19     6.0
20     6.8
21     7.1
22     7.9
23     8.2
24     8.7
25     9.0
26     9.5
27     9.6
28    10.3
29    10.5
Name: YearsExperience, dtype: float64

In [37]: data["YearsExperience"].astype(int)

Out [37]:
0      1
1      1
2      1
3      2
4      2
5      2
6      3
7      3
8      3
9      3
10     3
11     4
12     4
13     4
14     4
15     4
16     5
17     5
18     5
19     6
20     6
21     7
22     7
23     8
24     8
25     9
26     9
27     9
28    10
29    10
Name: YearsExperience, dtype: int32
```

5. Checking Particular Column data NULL Values

```
In [38]: data["YearsExperience"].isna().astype(int)

Out [38]:
0      0
1      0
2      0
3      0
4      0
5      0
6      0
7      0
8      0
9      0
10     0
11     0
12     0
13     0
14     0
15     0
16     0
17     0
18     0
19     0
20     0
21     0
22     0
23     0
24     0
25     0
26     0
27     0
28     0
29     0
Name: YearsExperience, dtype: int32

In [39]: data["YearsExperience"].isna().astype(int).sum()

Out [39]:
0
```

6. Checking Particular Column data

```
In [40]: data["Salary"]

Out [40]:
0      39343
1      46205
2      37731
3      43525
4      38891
5      56642
6      60150
7      54445
8      64445
9      57189
10     63218
11     55794
12     56957
13     57081
14     61111
15     67838
16     66029
17     83088
18     81363
19     93940
20     91738
21     98273
22     101302
23     113812
24     109431
25     105582
26     116969
27     112635
28     122391
29     121872
Name: Salary, dtype: int64

In [41]: # Fit a copy of dataset include last column
# i.e first column of dataset
X = data.iloc[:, :-1].values

In [42]: array([[ 1.1],
 [ 1.3],
 [ 1.5],
 [ 2.0],
 [ 2.2],
 [ 2.9],
 [ 3.0],
 [ 3.2],
 [ 3.2],
 [ 3.7],
 [ 3.9],
 [ 4.0],
 [ 4.0],
 [ 4.1],
 [ 4.5],
 [ 4.9],
 [ 5.1],
 [ 5.3],
 [ 5.9],
 [ 6.0],
 [ 6.8],
 [ 7.1],
 [ 7.9],
 [ 8.2],
 [ 8.7],
 [ 9.0],
 [ 9.5],
 [ 9.6],
 [10.3],
 [10.5]])

In [43]: # Fit array of dataset in column 1st
Y = data.iloc[:, 1].values
Y

Out [43]:
array([ 39343,  46205,  37731,  43525,  38891,  56642,  60150,  54445,
        64445,  57189,  63218,  55794,  56957,  57081,  61111,  67838,
        66029,  83088,  81363,  93940,  91738,  98273, 101302, 113812,
        109431, 105582, 116969, 112635, 122391, 121872], dtype=int64)
```

```
In [44]: # # # Fitting Data
# 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=1/3, random_state=0)

In [45]: # # # Fitting Simple Linear Regression to the Training set
# Predicting the result of 5 Years Experience
# Fitting Simple Linear Regression to the Training set
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)

Out [45]:
LinearRegression

In [46]: # Predicting the result of 5 Years Experience
from sklearn.linear_model import LinearRegression
y_pred = regressor.predict(X_test)

In [47]: array([ 40503.28550771, 122979.29940819,  60124.38623805,  62285.38777221,
        135022.48453269, 108125.8914952 , 115527.23969801,  66199.96220102,
        18349.49123286, 100649.2375447 ])

In [48]: # Predicting the result of 5 Years Experience
y_pred = regressor.predict([5])

Out [48]:
array([73545.90445641])

In [49]: # # # Intercept and slope of a simple linear regression model in Python using scikit-learn
# Print the intercept and slope
from sklearn.linear_model import LinearRegression
print("Intercept:", regressor.intercept_)
print("Slope:", regressor.coef_)

print(regressor.get_params())

Intercept: 26816.10224403119
Slope: 1945.76243122

In [50]: # Copy X = True, fit_intercept = True, n_jobs = None, positive = False

In [51]: # # # Visualizing the Training set results
# Visualizing the Train set results
plt.scatter(X_train, y_train, color='red')
plt.plot(X_train, regressor.predict(X_train), color='blue')
plt.xlabel("Years of Experience")
plt.ylabel("Salary")
plt.show()

In [52]:
```



```
In [53]: # # # Visualizing the Test set results
# Visualizing the Test set results
plt.scatter(X_test, y_test, color='red')
plt.plot(X_test, regressor.predict(X_test), color='blue')
plt.xlabel("Years of Experience")
plt.ylabel("Salary")
plt.show()

In [54]:
```



