# Assignment - 04 - Simple Linear Regression - 2

```
In [1]:  1  # import libraries
         2  import pandas as pd
         3  import numpy as np
         4  import seaborn as sns
         5  import statsmodels.formula.api as smf
```

```
In [2]:  1  # import dataset
         2  dataset=pd.read_csv('Downloads\\Salary_Data.csv')
         3  dataset
```

Out[2]:

| | YearsExperience | Salary |
|---|---|---|
| 0 | 1.1 | 39343.0 |
| 1 | 1.3 | 46205.0 |
| 2 | 1.5 | 37731.0 |
| 3 | 2.0 | 43525.0 |
| 4 | 2.2 | 39891.0 |
| 5 | 2.9 | 56642.0 |
| 6 | 3.0 | 60150.0 |
| 7 | 3.2 | 54445.0 |
| 8 | 3.2 | 64445.0 |
| 9 | 3.7 | 57189.0 |
| 10 | 3.9 | 63218.0 |
| 11 | 4.0 | 55794.0 |
| 12 | 4.0 | 56957.0 |
| 13 | 4.1 | 57081.0 |
| 14 | 4.5 | 61111.0 |
| 15 | 4.9 | 67938.0 |
| 16 | 5.1 | 66029.0 |
| 17 | 5.3 | 83088.0 |
| 18 | 5.9 | 81363.0 |
| 19 | 6.0 | 93940.0 |
| 20 | 6.8 | 91738.0 |
| 21 | 7.1 | 98273.0 |
| 22 | 7.9 | 101302.0 |
| 23 | 8.2 | 113812.0 |
| 24 | 8.7 | 109431.0 |
| 25 | 9.0 | 105582.0 |
| 26 | 9.5 | 116969.0 |
| 27 | 9.6 | 112635.0 |
| 28 | 10.3 | 122391.0 |
| 29 | 10.5 | 121872.0 |

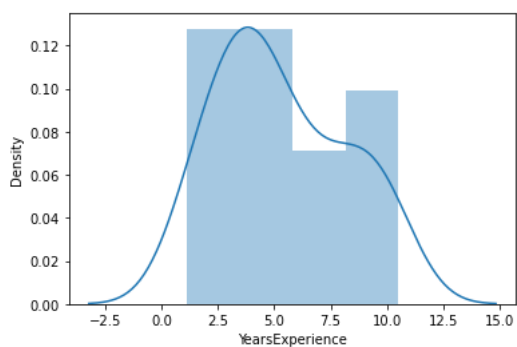## EDA and Data Visualization

```
In [3]:  1  dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 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     float64
dtypes: float64(2)
memory usage: 608.0 bytes
```

In [7]:
```python
1  import warnings
2  warnings.filterwarnings('ignore')
3  sns.distplot(dataset["YearsExperience"])
```
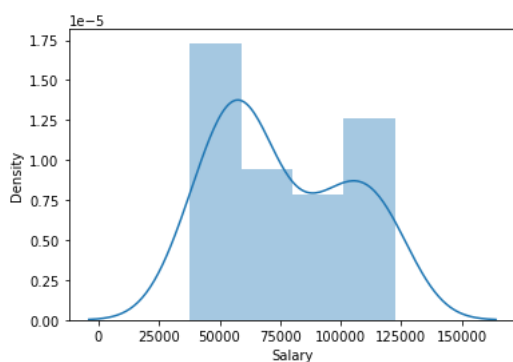
Out[7]: <AxesSubplot:xlabel='YearsExperience', ylabel='Density'>



In [8]:
```python
1  sns.distplot(dataset["Salary"])
```
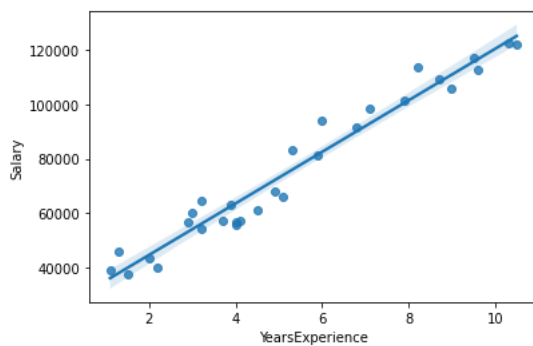
Out[8]: <AxesSubplot:xlabel='Salary', ylabel='Density'>



## Correlation Analysis

In [10]:
```python
1  dataset.corr()
```

Out[10]:

|                 | YearsExperience | Salary   |
|-----------------|-----------------|----------|
| YearsExperience | 1.000000        | 0.978242 |
| Salary          | 0.978242        | 1.000000 |

In [11]:
```python
1  sns.regplot(x=dataset['YearsExperience'],y=dataset['Salary'])
```

Out[11]: <AxesSubplot:xlabel='YearsExperience', ylabel='Salary'>



## Model Building

```
In [12]:    1  model=smf.ols("Salary~YearsExperience",data=dataset).fit()
```

## Model Testing

```
In [13]:    1  # Finding Coefficient Parameters
            2  model.params
```

```
Out[13]:  Intercept          25792.200199
          YearsExperience     9449.962321
          dtype: float64
```

```
In [15]:    1  # Finding Pvalues and tvalues
            2  model.pvalues , model.tvalues
```

```
Out[15]:  (Intercept          5.511950e-12
           YearsExperience    1.143068e-20
           dtype: float64,
           Intercept          11.346940
           YearsExperience    24.950094
           dtype: float64)
```

```
In [16]:    1  # Finding Rsquared values
            2  model.rsquared , model.rsquared_adj
```

```
Out[16]:  (0.9569566641435086, 0.9554194021486339)
```

## Model Prediction

```
In [19]:    1  # Manual prediction for say 3 years
            2  Salary = 5792.200199 + 9449.962321
            3  Salary
```

```
Out[19]:  15242.162520000002
```

```
In [20]:    1  # Automatic prediction for say 3 and 5 years
```

```
In [21]:    1  new_data=pd.Series([3,5])
            2  new_data
```

```
Out[21]:  0    3
          1    5
          dtype: int64
```

```
In [26]:    1  data_pred=pd.DataFrame(new_data,columns=["YearsExperience"])
            2  data_pred
```

Out[26]:

|   | YearsExperience |
|---|---|
| 0 | 3 |
| 1 | 5 |

```
In [27]:    1  model.predict(data_pred)
```

```
Out[27]:  0    54142.087163
          1    73042.011806
          dtype: float64
```

```
In [ ]:     1
```