

# Problem Statement

A real estate agent want help to predict the house price for regions in USA.He gave us the dataset to work on to use linear regression model.Create a model that helps him to estimate of what the house would sell for

## Import libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: # To import dataset
df=pd.read_csv('placement csv')
df
```

Out[2]:

	cgpa	placement_exam_marks	placed
0	7.19	26.0	1
1	7.46	38.0	1
2	7.54	40.0	1
3	6.42	8.0	1
4	7.23	17.0	0
...	...	...	...
995	8.87	44.0	1
996	9.12	65.0	1
997	4.89	34.0	0
998	8.62	46.0	1
999	4.90	10.0	1

1000 rows × 3 columns

```
In [3]: # To display top 10 rows
df.head(10)
```

Out[3]:

	cgpa	placement_exam_marks	placed
0	7.19	26.0	1
1	7.46	38.0	1
2	7.54	40.0	1
3	6.42	8.0	1
4	7.23	17.0	0
5	7.30	23.0	1
6	6.69	11.0	0
7	7.12	39.0	1
8	6.45	38.0	0
9	7.75	94.0	1

## Data Cleaning and Pre-Processing

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 3 columns):
#   Column                Non-Null Count  Dtype
---  -
0   cgpa                  1000 non-null  float64
1   placement_exam_marks 1000 non-null  float64
2   placed                1000 non-null  int64
dtypes: float64(2), int64(1)
memory usage: 23.6 KB
```

```
In [5]: df.describe()
```

Out[5]:

	cgpa	placement_exam_marks	placed
<b>count</b>	1000.000000	1000.000000	1000.000000
<b>mean</b>	6.961240	32.225000	0.489000
<b>std</b>	0.615898	19.130822	0.500129
<b>min</b>	4.890000	0.000000	0.000000
<b>25%</b>	6.550000	17.000000	0.000000
<b>50%</b>	6.960000	28.000000	0.000000
<b>75%</b>	7.370000	44.000000	1.000000
<b>max</b>	9.120000	100.000000	1.000000

```
In [6]: df.columns
```

```
Out[6]: Index(['cgpa', 'placement_exam_marks', 'placed'], dtype='object')
```

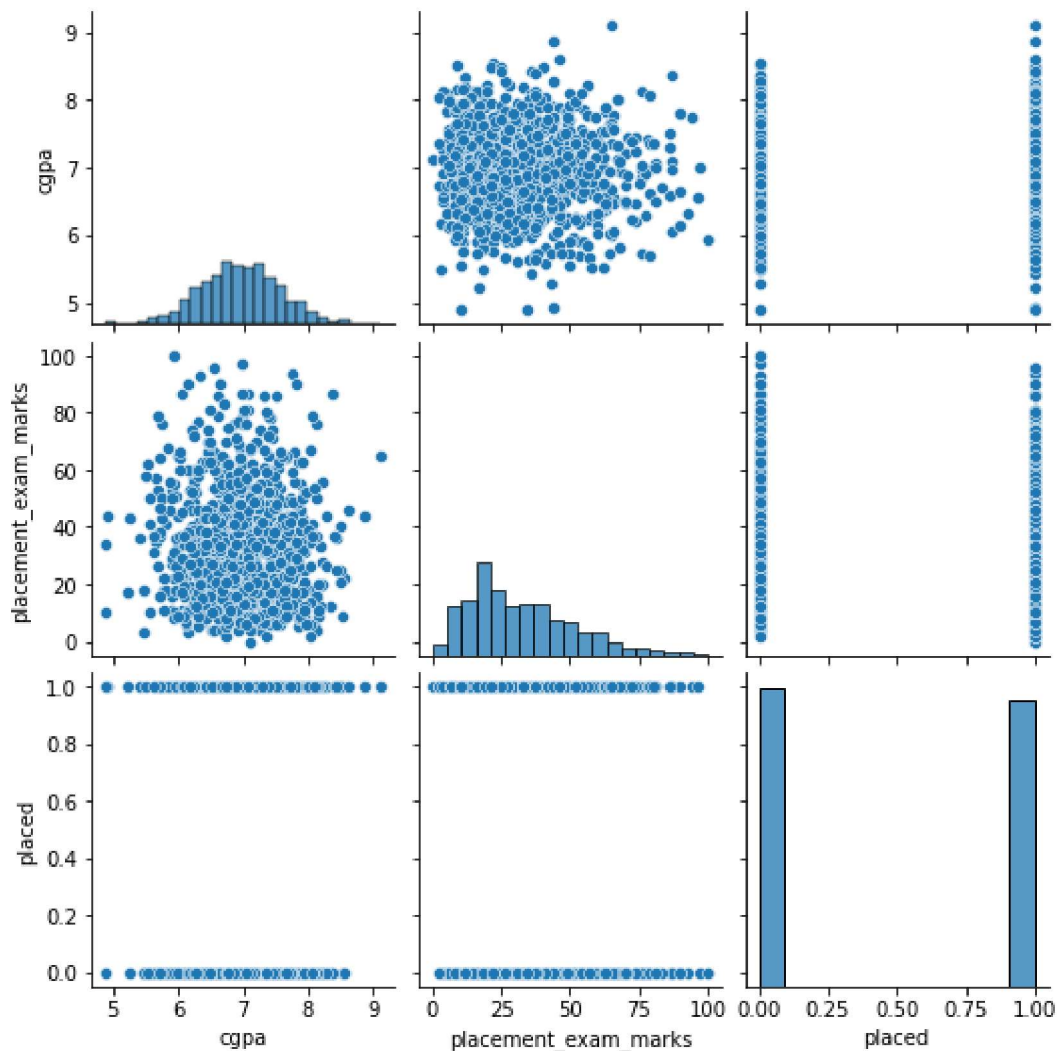
```
In [7]: a = df.dropna(axis='columns')  
a.columns
```

```
Out[7]: Index(['cgpa', 'placement_exam_marks', 'placed'], dtype='object')
```

## EDA and Visualization

```
In [8]: sns.pairplot(a)
```

```
Out[8]: <seaborn.axisgrid.PairGrid at 0x2b9a521daf0>
```

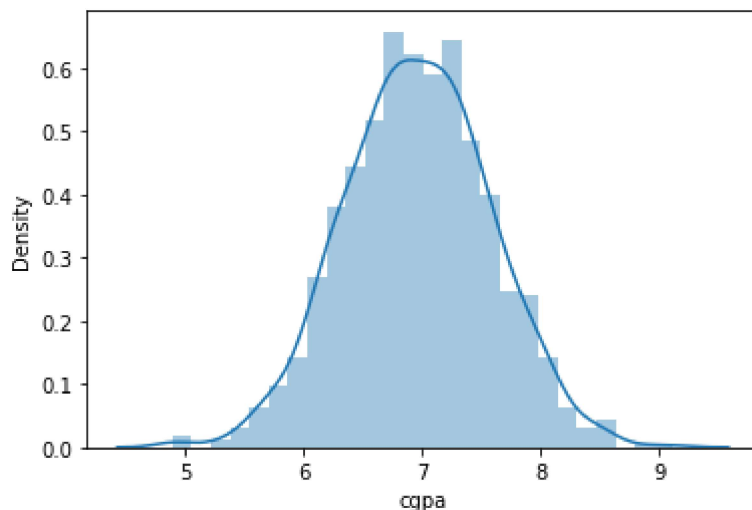


```
In [11]: sns.distplot(a['cgpa'])
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

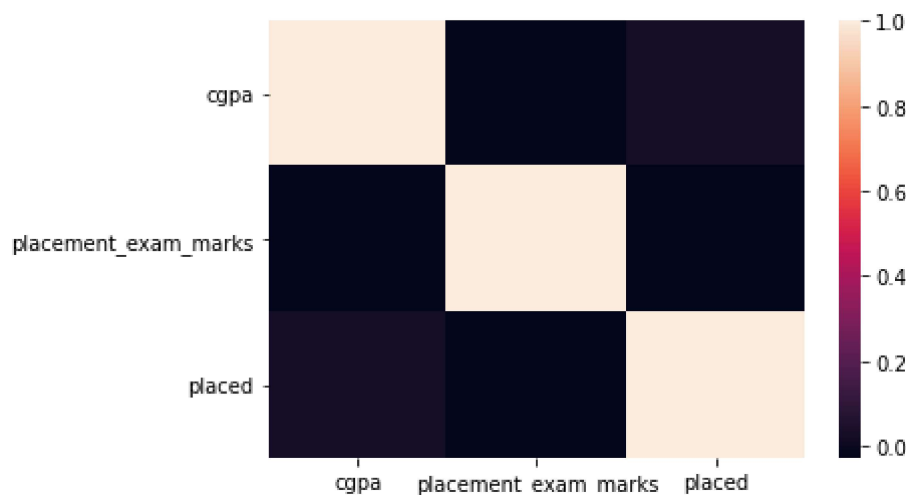
```
Out[11]: <AxesSubplot:xlabel='cgpa', ylabel='Density'>
```



```
In [12]: a1=a[['cgpa', 'placement_exam_marks', 'placed']]
```

```
In [13]: sns.heatmap(a1.corr())
```

```
Out[13]: <AxesSubplot:>
```



## To Train the Model - Model Building

We are going to train Linear Regression model; We need to split out data into two variables x and y where x is independent variable (input) and y is dependent on x (output). We could ignore address column as it is not required for our model.

```
In [14]: x=a1[['placement_exam_marks', 'placed']]
y=a1['cgpa']
```

## To split my dataset into training and test data

```
In [15]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3)
```

```
In [16]: from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(x_train,y_train)
```

Out[16]: LinearRegression()

```
In [17]: print(lr.intercept_)
```

6.974174068263444

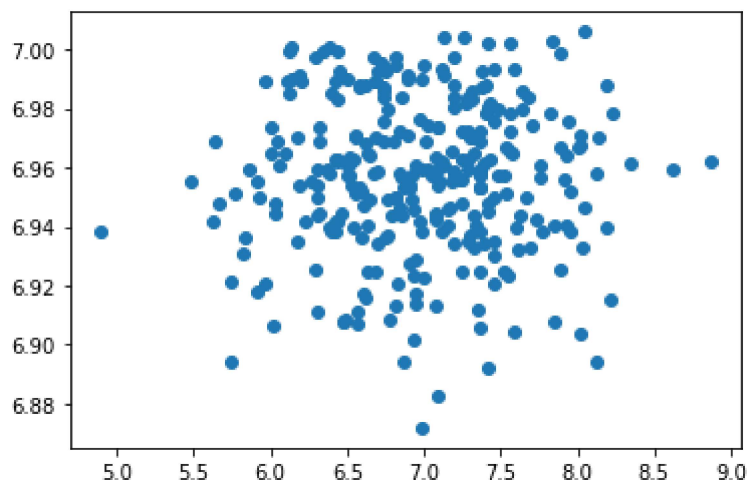
```
In [18]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[18]:

	Co-efficient
placement_exam_marks	-0.001056
placed	0.034016

```
In [19]: prediction=lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[19]: <matplotlib.collections.PathCollection at 0x2b9a7bc6a30>



In [20]: `print(lr.score(x_test,y_test))`

-0.00033406734904151136

In [ ]: