

# SUMESH R -20104169

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

```
In [2]: df = pd.read_csv("6_Salesworkload1.csv")
df = df[['MonthYear', 'Time index', 'Country', 'StoreID', 'City', 'Dept_ID',
        'Dept. Name', 'HoursOwn', 'HoursLease', 'Sales units', 'Turnover', 'Area (m2)',
        df
```

```
Out[2]:
```

	MonthYear	Time index	Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLease
<b>0</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	1.0	Dry	3184.764	0.0
<b>1</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	2.0	Frozen	1582.941	0.0
<b>2</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	3.0	other	47.205	0.0
<b>3</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	4.0	Fish	1623.852	0.0
<b>4</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	5.0	Fruits & Vegetables	1759.173	0.0
...	...	...	...	...	...	...	...	...	...
<b>7653</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	12.0	Checkout	6322.323	0.0
<b>7654</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	16.0	Customer Services	4270.479	0.0
<b>7655</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	11.0	Delivery	0	0.0
<b>7656</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	17.0	others	2224.929	0.0
<b>7657</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	18.0	all	39652.2	0.0

7650 rows × 13 columns



```
In [3]: df.head()
```

```
Out[3]:
```

	MonthYear	Time index	Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLease	Sales units
<b>0</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	1.0	Dry	3184.764	0.0	398560.0

	MonthYear	Time index	Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLease	Sales units
1	10.2016	1.0	United Kingdom	88253.0	London (I)	2.0	Frozen	1582.941	0.0	82725.0
2	10.2016	1.0	United Kingdom	88253.0	London (I)	3.0	other	47.205	0.0	438400.0
3	10.2016	1.0	United Kingdom	88253.0	London (I)	4.0	Fish	1623.852	0.0	309425.0
4	10.2016	1.0	United Kingdom	88253.0	London (I)	5.0	Fruits & Vegetables	1759.173	0.0	165515.0

## Data cleaning and pre processing

In [4]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 7650 entries, 0 to 7657
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   MonthYear       7650 non-null   object
1   Time index      7650 non-null   float64
2   Country         7650 non-null   object
3   StoreID         7650 non-null   float64
4   City            7650 non-null   object
5   Dept_ID         7650 non-null   float64
6   Dept. Name      7650 non-null   object
7   HoursOwn        7650 non-null   object
8   HoursLease      7650 non-null   float64
9   Sales units     7650 non-null   float64
10  Turnover        7650 non-null   float64
11  Area (m2)       7650 non-null   object
12  Opening hours   7650 non-null   object
dtypes: float64(6), object(7)
memory usage: 836.7+ KB
```

In [5]:

```
df.describe()
```

Out[5]:

	Time index	StoreID	Dept_ID	HoursLease	Sales units	Turnover
<b>count</b>	7650.000000	7650.000000	7650.000000	7650.000000	7.650000e+03	7.650000e+03
<b>mean</b>	5.000000	61995.220000	9.470588	22.036078	1.076471e+06	3.721393e+06
<b>std</b>	2.582158	29924.581631	5.337429	133.299513	1.728113e+06	6.003380e+06
<b>min</b>	1.000000	12227.000000	1.000000	0.000000	0.000000e+00	0.000000e+00
<b>25%</b>	3.000000	29650.000000	5.000000	0.000000	5.457125e+04	2.726798e+05
<b>50%</b>	5.000000	75400.500000	9.000000	0.000000	2.932300e+05	9.319575e+05
<b>75%</b>	7.000000	87703.000000	14.000000	0.000000	9.175075e+05	3.264432e+06
<b>max</b>	9.000000	98422.000000	18.000000	3984.000000	1.124296e+07	4.271739e+07

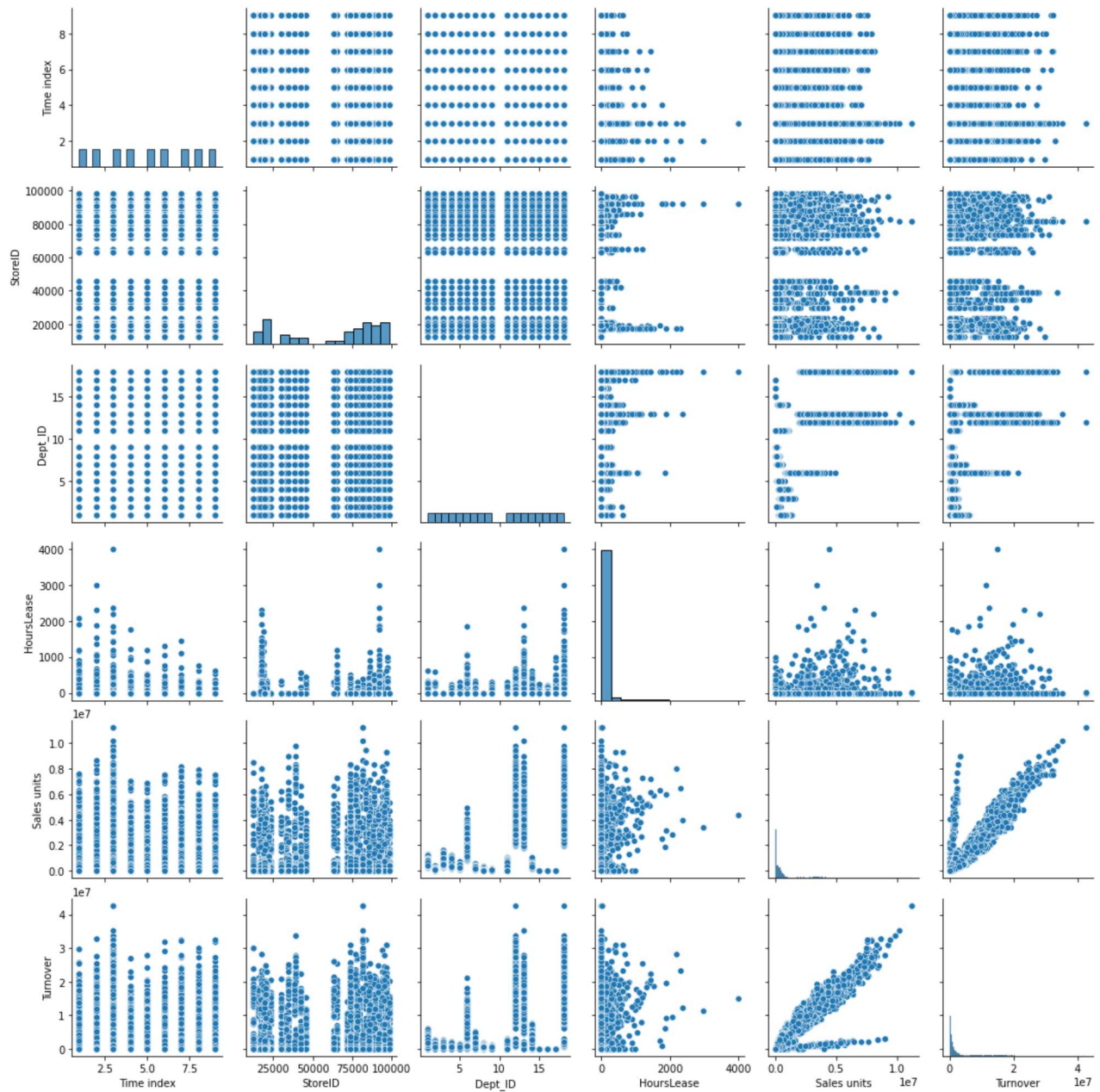
In [6]: `df.columns`

Out[6]: Index(['MonthYear', 'Time index', 'Country', 'StoreID', 'City', 'Dept\_ID', 'Dept. Name', 'HoursOwn', 'HoursLease', 'Sales units', 'Turnover', 'Area (m2)', 'Opening hours'], dtype='object')

## EDA and VISUALIZATION

In [7]: `sns.pairplot(df)`

Out[7]: <seaborn.axisgrid.PairGrid at 0x2a8cf9fedc0>

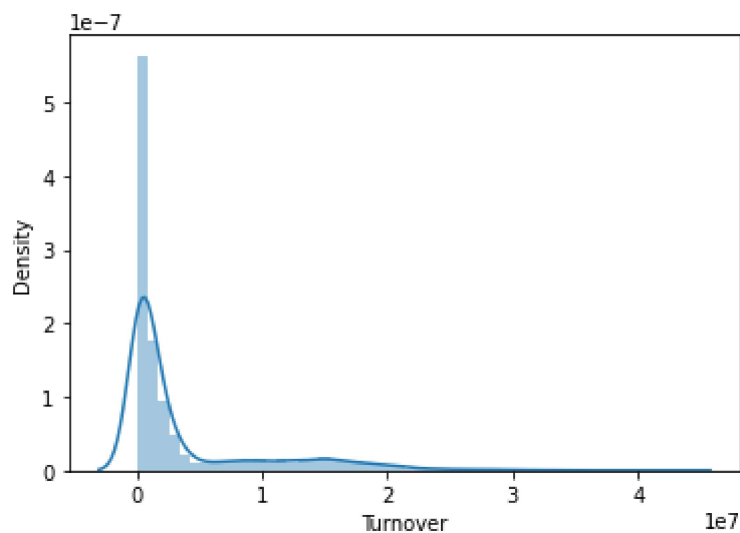


```
In [8]: sns.distplot(df["Turnover"])
```

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)

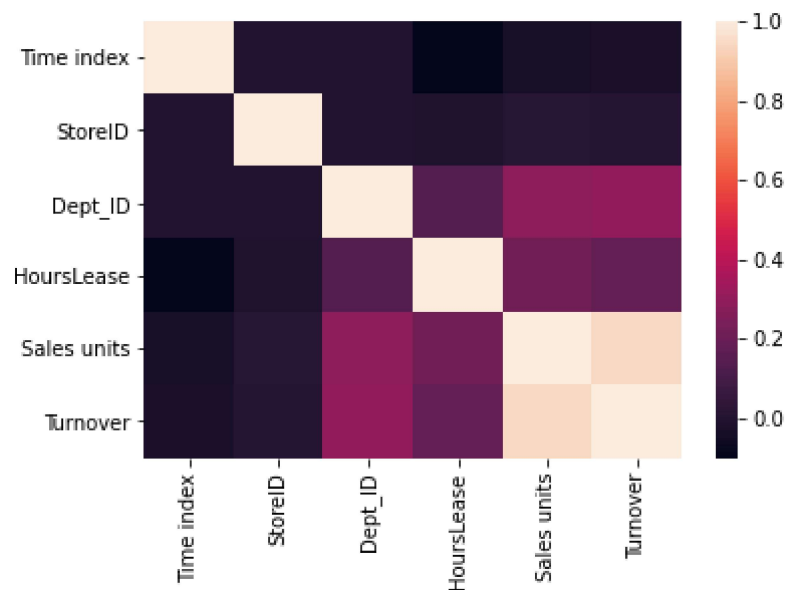
```
Out[8]: <AxesSubplot:xlabel='Turnover', ylabel='Density'>
```



```
In [9]: df1 = df[['MonthYear', 'Time index', 'Country', 'StoreID', 'City', 'Dept_ID',
                'Dept. Name', 'HoursOwn', 'HoursLease', 'Sales units', 'Turnover',
                'Area (m2)', 'Opening hours']]
```

```
In [10]: sns.heatmap(df1.corr())
```

Out[10]: <AxesSubplot:>



```
In [11]: df1.fillna(1)
```

```
Out[11]:
```

	MonthYear	Time index	Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLease
0	10.2016	1.0	United Kingdom	88253.0	London (I)	1.0	Dry	3184.764	0.0
1	10.2016	1.0	United Kingdom	88253.0	London (I)	2.0	Frozen	1582.941	0.0

	MonthYear	Time index	Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLease	
<b>2</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	3.0	other	47.205	0.0	
<b>3</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	4.0	Fish	1623.852	0.0	
<b>4</b>	10.2016	1.0	United Kingdom	88253.0	London (I)	5.0	Fruits & Vegetables	1759.173	0.0	
...	...	...	...	...	...	...	...	...	...	
<b>7653</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	12.0	Checkout	6322.323	0.0	3
<b>7654</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	16.0	Customer Services	4270.479	0.0	
<b>7655</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	11.0	Delivery	0	0.0	
<b>7656</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	17.0	others	2224.929	0.0	
<b>7657</b>	06.2017	9.0	Sweden	29650.0	Gothenburg	18.0	all	39652.2	0.0	3

7650 rows × 13 columns

In [12]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 7650 entries, 0 to 7657
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   MonthYear       7650 non-null   object
1   Time index      7650 non-null   float64
2   Country         7650 non-null   object
3   StoreID         7650 non-null   float64
4   City            7650 non-null   object
5   Dept_ID         7650 non-null   float64
6   Dept. Name      7650 non-null   object
7   HoursOwn        7650 non-null   object
8   HoursLease      7650 non-null   float64
9   Sales units     7650 non-null   float64
10  Turnover        7650 non-null   float64
11  Area (m2)       7650 non-null   object
12  Opening hours   7650 non-null   object
dtypes: float64(6), object(7)
memory usage: 836.7+ KB
```

In [13]:

```
x = df1[['Time index', 'StoreID', 'Dept_ID',
         'HoursLease', 'Sales units']]
y = df1['Turnover']
```

## split the data into training and test data

In [14]:

```
x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.3)
```

```
In [15]: lr = LinearRegression()  
lr.fit(x_train, y_train)
```

```
Out[15]: LinearRegression()
```

```
In [16]: lr.intercept_
```

```
Out[16]: -232128.03611873696
```

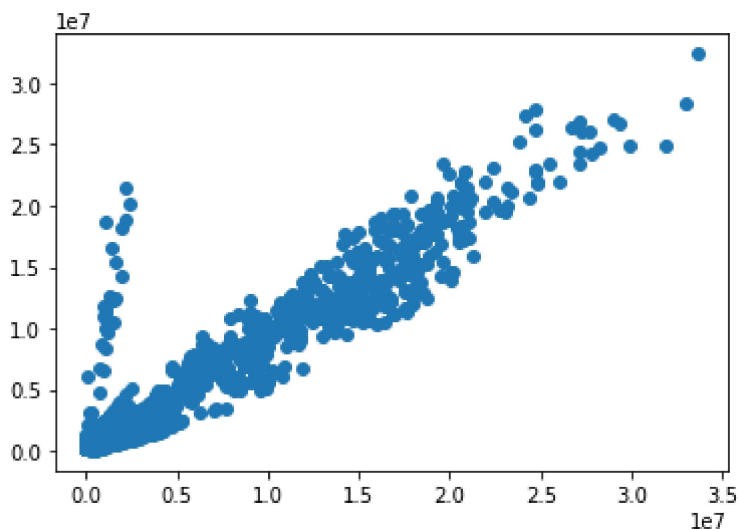
```
In [17]: coeff = pd.DataFrame(lr.coef_, x.columns, columns = ['Co-efficient'])  
coeff
```

```
Out[17]:
```

	Co-efficient
Time index	25022.645451
StoreID	-0.388014
Dept_ID	36768.421426
HoursLease	-525.132851
Sales units	3.246580

```
In [18]: prediction = lr.predict(x_test)  
plt.scatter(y_test, prediction)
```

```
Out[18]: <matplotlib.collections.PathCollection at 0x2a8d2b87880>
```



```
In [19]: lr.score(x_test,y_test)
```

```
Out[19]: 0.912993028532293
```