

kaviyadevi 20106064

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

```
In [16]: #to import dataset
data1=pd.read_csv(r"C:\Users\user\Downloads\6_Salesworkload1 - 6_Salesworkload1.csv")
data1
```

Out[16]:

	MonthYear	Time index	Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLeas
0	10.2016	1.0	United Kingdom	88253.0	London (I)	1.0	Dry	3184.764	0.
1	10.2016	1.0	United Kingdom	88253.0	London (I)	2.0	Frozen	1582.941	0.
2	10.2016	1.0	United Kingdom	88253.0	London (I)	3.0	other	47.205	0.
3	10.2016	1.0	United Kingdom	88253.0	London (I)	4.0	Fish	1623.852	0.
4	10.2016	1.0	United Kingdom	88253.0	London (I)	5.0	Fruits & Vegetables	1759.173	0.
...	...	...	...	...	...	...	...	...	.
7653	6.2017	9.0	Sweden	29650.0	Gothenburg	12.0	Checkout	6322.323	0.
7654	6.2017	9.0	Sweden	29650.0	Gothenburg	16.0	Customer Services	4270.479	0.
7655	6.2017	9.0	Sweden	29650.0	Gothenburg	11.0	Delivery	0	0.
7656	6.2017	9.0	Sweden	29650.0	Gothenburg	17.0	others	2224.929	0.
7657	6.2017	9.0	Sweden	29650.0	Gothenburg	18.0	all	39652.2	0.

7658 rows × 14 columns



```
In [19]: #to display top 5 rows
data=data1.head(200)
data
```

Out[19]:

Country	StoreID	City	Dept_ID	Dept. Name	HoursOwn	HoursLease	Sales units	Turnover	Customers
United Kingdom	88253.0	London (I)	1.0	Dry	3184.764	0.0	398560.0	1226244.0	Na
United Kingdom	88253.0	London (I)	2.0	Frozen	1582.941	0.0	82725.0	387810.0	Na
United Kingdom	88253.0	London (I)	3.0	other	47.205	0.0	438400.0	654657.0	Na
United Kingdom	88253.0	London (I)	4.0	Fish	1623.852	0.0	309425.0	499434.0	Na
United Kingdom	88253.0	London (I)	5.0	Fruits & Vegetables	1759.173	0.0	165515.0	329397.0	Na
...	...	...	...	...	...	...	...	...	...
The Netherlands	95434.0	Den Haag	8.0	Household	2127.372	0.0	58615.0	27960.0	Na
The Netherlands	95434.0	Den Haag	9.0	Hardware	2158.842	0.0	63985.0	554325.0	Na
The Netherlands	95434.0	Den Haag	14.0	Non Food	9887.874	0.0	370250.0	2994267.0	Na
The Netherlands	95434.0	Den Haag	15.0	Admin	5589.072	0.0	55.0	0.0	Na
The Netherlands	95434.0	Den Haag	12.0	Checkout	6781.785	0.0	4510270.0	18356889.0	Na



## DATA CLEANING AND PREPROCESSING

In [20]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 14 columns):
#   Column          Non-Null Count  Dtype
---  -
0   MonthYear       200 non-null   object
1   Time index      200 non-null   float64
2   Country         200 non-null   object
3   StoreID         200 non-null   float64
4   City            200 non-null   object
5   Dept_ID         200 non-null   float64
6   Dept. Name      200 non-null   object
7   HoursOwn        200 non-null   object
8   HoursLease      200 non-null   float64
9   Sales units     200 non-null   float64
10  Turnover        200 non-null   float64
11  Customer        0 non-null     float64
12  Area (m2)       200 non-null   object
13  Opening hours   200 non-null   object
dtypes: float64(7), object(7)
memory usage: 22.0+ KB
```

In [21]: *#to display summary of statistics*  
data.describe()

Out[21]:

	Time index	StoreID	Dept_ID	HoursLease	Sales units	Turnover	Customer
<b>count</b>	200.0	200.000000	200.000000	200.000000	2.000000e+02	2.000000e+02	0.0
<b>mean</b>	1.0	46739.115000	9.350000	41.23000	9.317313e+05	3.000231e+06	NaN
<b>std</b>	0.0	30654.343517	5.320625	184.09236	1.521370e+06	5.188606e+06	NaN
<b>min</b>	1.0	15552.000000	1.000000	0.00000	0.000000e+00	0.000000e+00	NaN
<b>25%</b>	1.0	18808.000000	5.000000	0.00000	5.200250e+04	2.084858e+05	NaN
<b>50%</b>	1.0	23623.000000	9.000000	0.00000	2.429175e+05	5.771910e+05	NaN
<b>75%</b>	1.0	73949.000000	14.000000	0.00000	9.019388e+05	2.358503e+06	NaN
<b>max</b>	1.0	95434.000000	18.000000	1896.00000	7.476680e+06	2.571973e+07	NaN

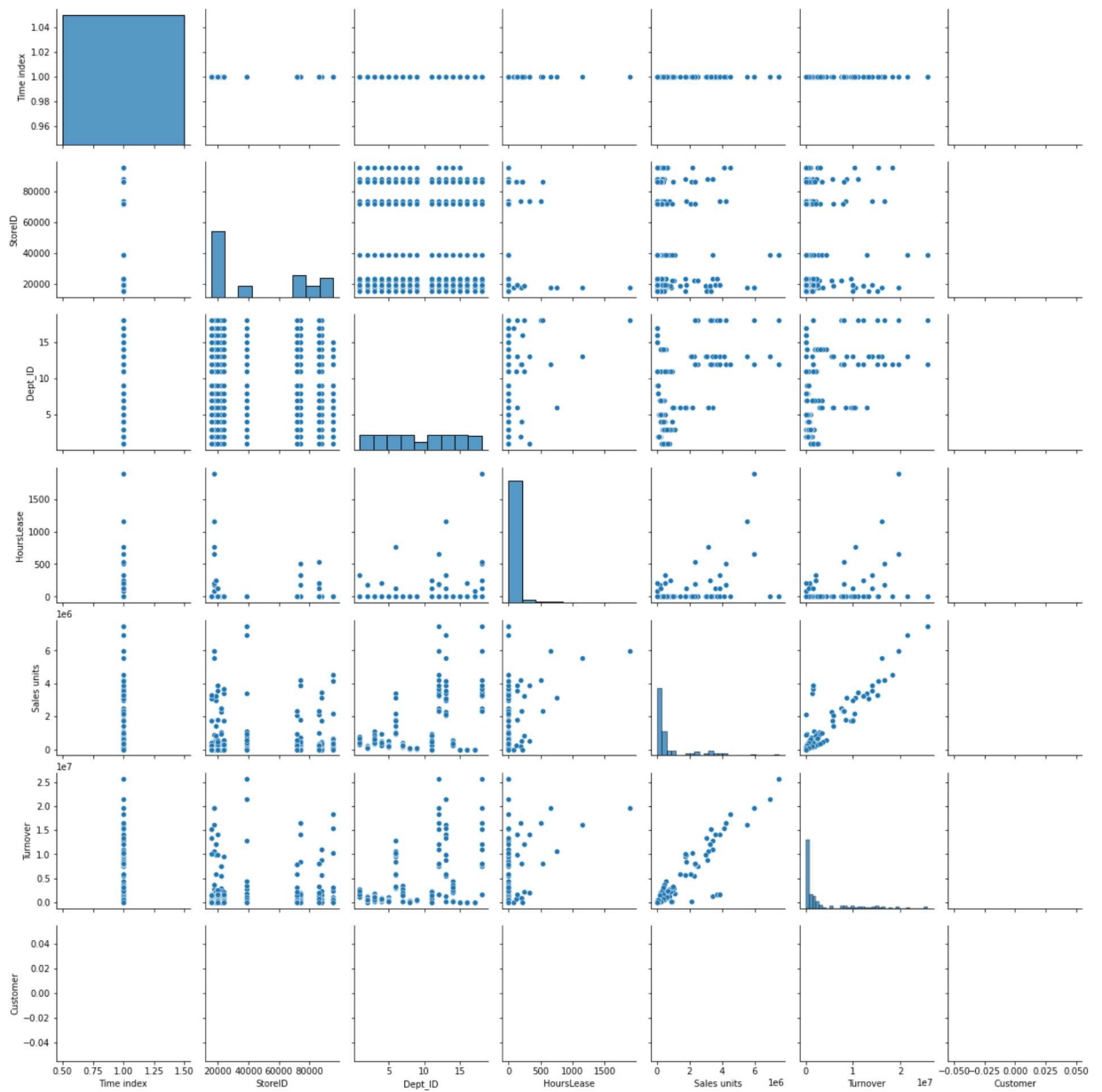
In [22]: *#to display the column heading*  
data.columns

Out[22]: Index(['MonthYear', 'Time index', 'Country', 'StoreID', 'City', 'Dept\_ID',  
                  'Dept. Name', 'HoursOwn', 'HoursLease', 'Sales units', 'Turnover',  
                  'Customer', 'Area (m2)', 'Opening hours'],  
              dtype='object')

## EDA and DATA VISUALIZATION

```
In [23]: sns.pairplot(data)
```

```
Out[23]: <seaborn.axisgrid.PairGrid at 0x21129a0b130>
```

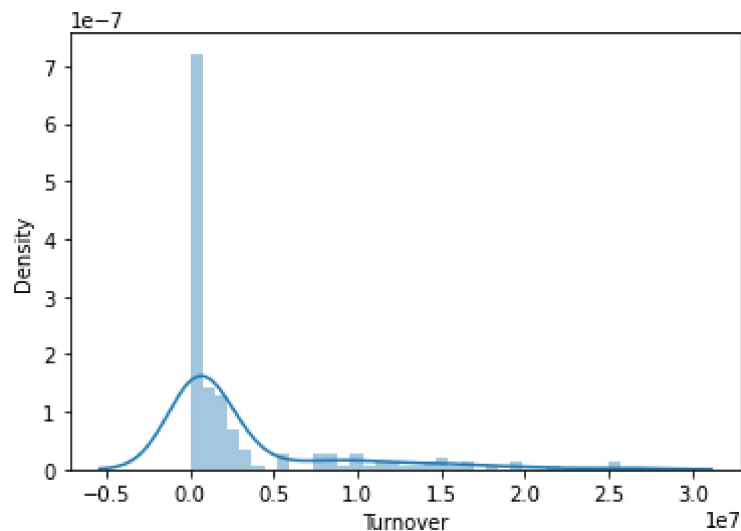


```
In [26]: sns.distplot(data['Turnover'])
```

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)
```

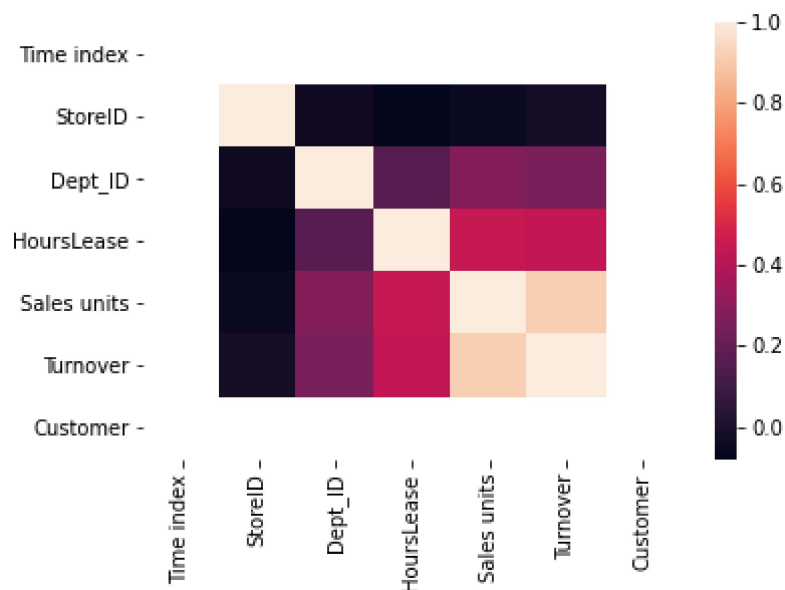
```
Out[26]: <AxesSubplot:xlabel='Turnover', ylabel='Density'>
```



```
In [28]: df=data[['MonthYear', 'Time index', 'Country', 'StoreID', 'City', 'Dept_ID',  
                'Dept. Name', 'HoursOwn', 'HoursLease', 'Sales units', 'Turnover',  
                'Customer', 'Area (m2)', 'Opening hours']]
```

```
In [29]: sns.heatmap(df.corr())
```

```
Out[29]: <AxesSubplot:>
```



## TRAINING MODEL

```
In [39]: x=df[['MonthYear', 'Time index', 'Dept_ID', 'HoursOwn', 'HoursLease', 'Sales units']]
         y=df[['Turnover']]
```

```
In [40]: #to split my dataset into training and test
         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 [41]: from sklearn.linear_model import LinearRegression

         lr=LinearRegression()
         lr.fit(x_train,y_train)
```

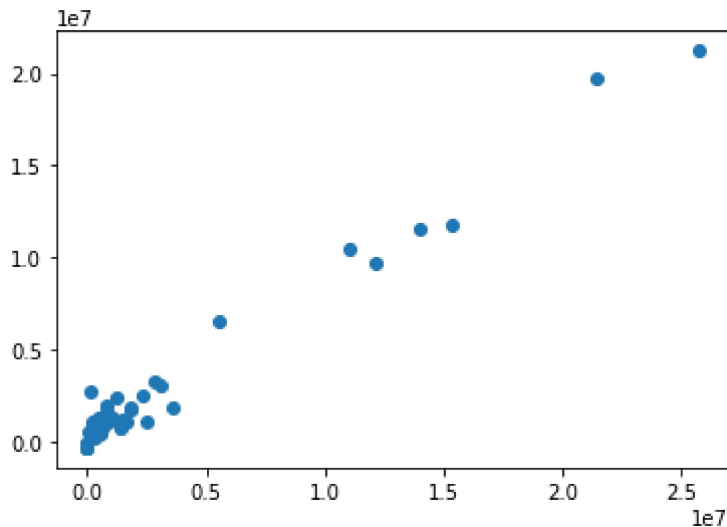
```
Out[41]: LinearRegression()
```

```
In [42]: #to find intercept  
print(lr.intercept_)
```

```
[499205.29783988]
```

```
In [45]: prediction = lr.predict(x_test)  
plt.scatter(y_test,prediction)
```

```
Out[45]: <matplotlib.collections.PathCollection at 0x2112c8c63d0>
```



```
In [46]: print(lr.score(x_test,y_test))
```

```
0.9541744830746266
```

## RIDGE AND LASSO REGRESSION

```
In [47]: from sklearn.linear_model import Ridge,Lasso
```

```
In [48]: rr=Ridge(alpha=10)  
rr.fit(x_train,y_train)
```

```
Out[48]: Ridge(alpha=10)
```

```
In [49]: rr.score(x_test,y_test)
```

```
Out[49]: 0.9542049571246298
```

```
In [50]: la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

```
Out[50]: Lasso(alpha=10)
```

```
In [37]: la.score(x_test,y_test)
```

```
Out[37]: 0.953355166182402
```