Exp no: 3   sales prediction and customer segmentation codes
            data regression and demans clustering

Aim:
To predict sales based on advertising channels
for advertising budget values (TV) radio and
newspaper using linear regression additionally
to identify distinct patterns in customer data
through k-means clustering

Algorithm:

Step 1: load the advertising.csv containing TV,
radio, news paper, sales and data.

Step 2: input data: select features (TV, radio, newspaper)
and data (sales), then split into inputs,
split the ~~model~~.

Step 3: Train and evaluate linear regression, fit the
model on training data, predict sales on
test data and calculate mean square error

Step 4: visualize result. plot actual vs predicted
sales and display which the contribute model
segments and inputs.

Program

import pandas as pd
import matplotlib as plt
import sklearn as sns.
from sklearn.model.selection import train_test
from sklearn.nlbn. import split mean_square_mean_split
from sklearn.clube import k-mean.
from sklearn. preprocessing input StandardScale

output

17.5

15.5

13.5

11.5

7.5

5   10   15   20   25

Linear Regression   MSE : 4. 52 25525

```python
df = pd.read_csv (advertising_csv)
print (df.head()) print(df.dub())
x = df [['TV', 'radio', 'newspaper']]
y = df ['sales']
x_train, x_test, y_train, y_test = train_test = split (x,y,
                                                  test, size = 0.2,
                                                  random_state = )

model = Linear regression ()
model.fit (x_train, y_train)
y_pred = model.predict (x_test)
mse = mean.predict (y_test)
print (" Linear Regression MSE:", mse)
plt.figure (figsize =(8-5))
sns.scatterplot (x=y, test, y= y_pred)
plt.ylabel (' Actual sales")
plt.title (" Linear regression Actual vs predicted sales)
plt.show()
scaler = standard scaler ()
plt.figure (fig size =(8,16)]
sns.relat.plot (data= df, x='TV', y='sales', row='test',
                pd Col = ("star")
.plt.show()
```

result:-

The linear regression model what is mean square (msels) of Aug shows its is good for better calculated as predicted value mean, illustrating predicting sales value value (selects using on pattern 3 data in vita slows patterns, advertising display and news.