

```
In [10]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [11]: advertising=pd.read_csv("Advertising.csv")

In [12]: advertising

Out[12]:
   Unnamed: 0    TV    radio  newspaper  sales
0           1  230.1   37.8      69.2   22.1
1           2  44.5   39.3      45.1   10.4
2           3  17.2   45.9      69.3    9.3
3           4 151.5  41.3      58.5   18.5
4           5 190.8  10.8      58.4   12.9
...      ...    ...    ...      ...    ...
195        196   38.2   3.7      13.8    7.6
196        197   94.2   4.9       8.1    9.7
197        198  177.0   9.3       6.4   12.8
198        199  283.6  42.0      66.2   25.5
199        200  232.1   8.6       8.7   13.4

200 rows x 5 columns

In [13]: advertising.shape

Out[13]:
(200, 5)

In [14]: advertising.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column  Non-Null Count  Dtype
---  --
0   Unnamed: 0    200 non-null    int64
1   TV          200 non-null    float64
2   radio       200 non-null    float64
3   newspaper   200 non-null    float64
4   sales       200 non-null    float64
dtypes: float64(4), int64(1)
memory usage: 7.9 KB

In [15]: advertising.describe()

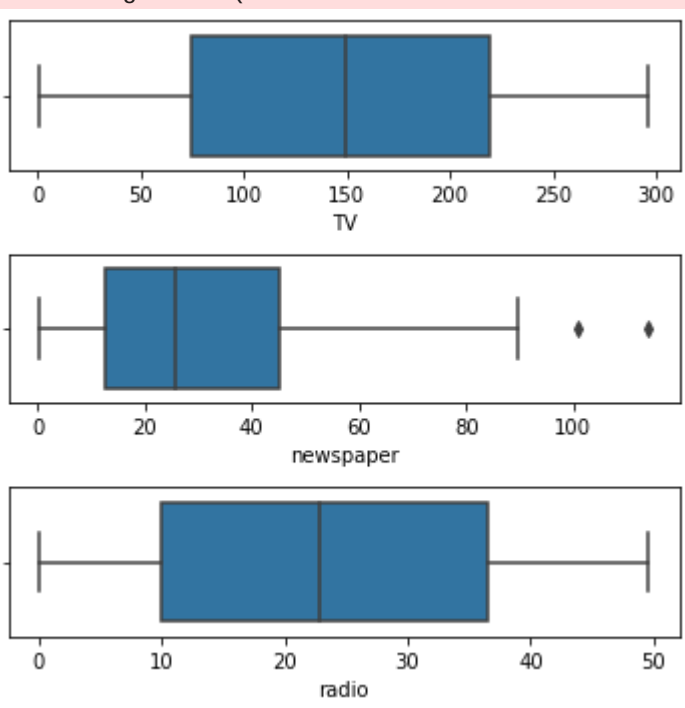
Out[15]:
   Unnamed: 0     TV     radio  newspaper     sales
count  200.000000  200.000000  200.000000  200.000000  200.000000
mean    100.500000  147.042500  23.264000   30.554000  14.022500
std     57.879185   85.854236  14.846809   21.778621  5.217457
min      1.000000    0.700000   0.000000   0.300000  1.600000
25%     50.750000   74.375000   9.975000  12.750000  10.375000
50%    100.500000  140.750000  22.900000  25.750000  12.900000
75%    150.250000  218.025000  36.525000  45.100000  17.400000
max    200.000000  295.400000  49.600000 114.000000  27.000000

In [16]: advertising.isnull().sum()

Out[16]:
Unnamed: 0    0
TV            0
radio         0
newspaper     0
sales         0
dtype: int64

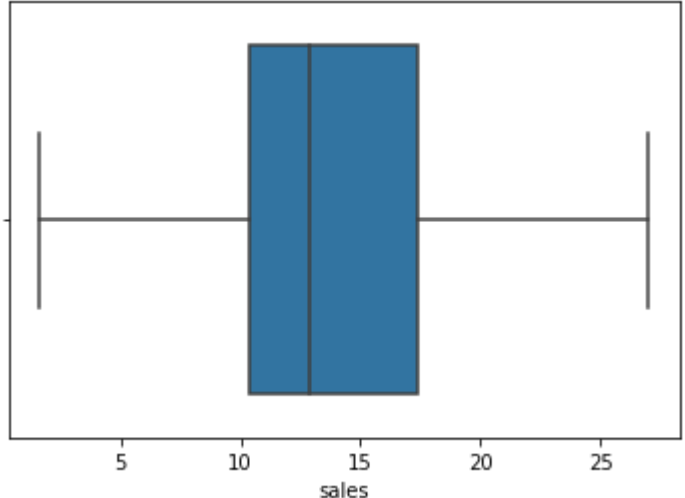
In [19]: # Outlier Analysis
fig, axs = plt.subplots(3, figsize = (5,5))
plt1 = sns.boxplot(advertising['TV'], ax = axs[0])
plt2 = sns.boxplot(advertising['newspaper'], ax = axs[1])
plt3 = sns.boxplot(advertising['radio'], ax = axs[2])
plt.tight_layout()

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keywo
rd will result in an error or misinterpretation.
warnings.warn(
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keywo
rd will result in an error or misinterpretation.
warnings.warn(
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keywo
rd will result in an error or misinterpretation.
rd will result in an error or misinterpretation.
warnings.warn(

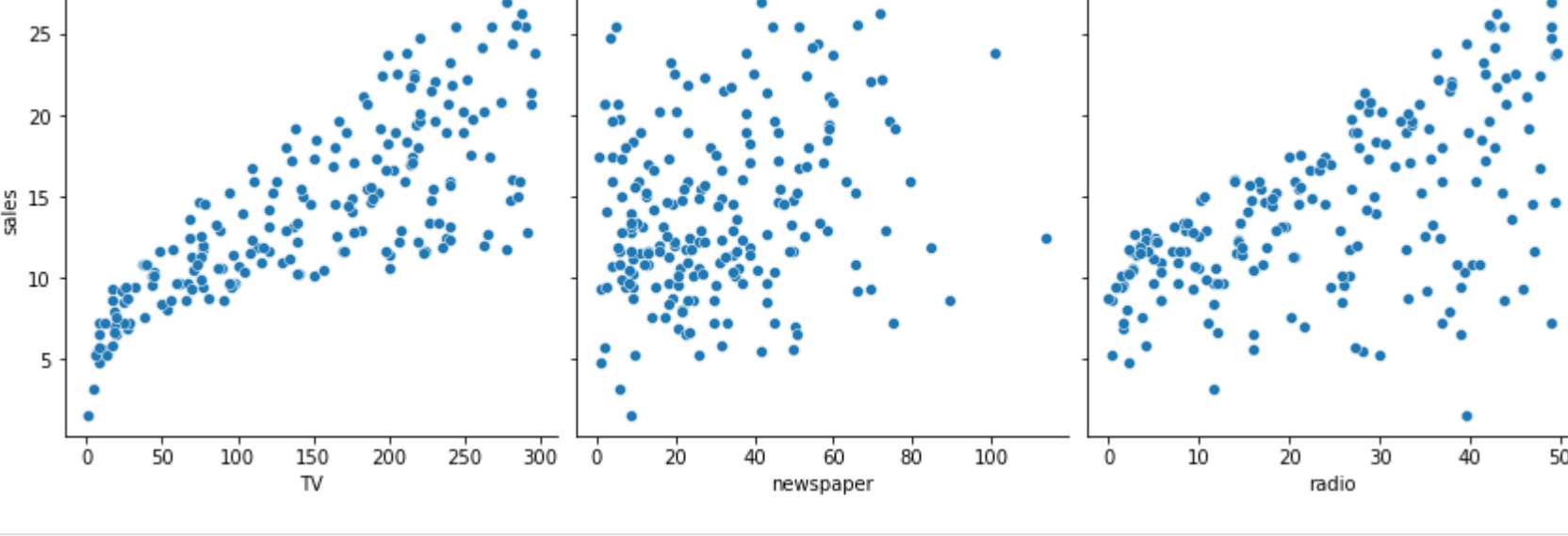


In [20]: sns.boxplot(advertising['sales'])
plt.show()

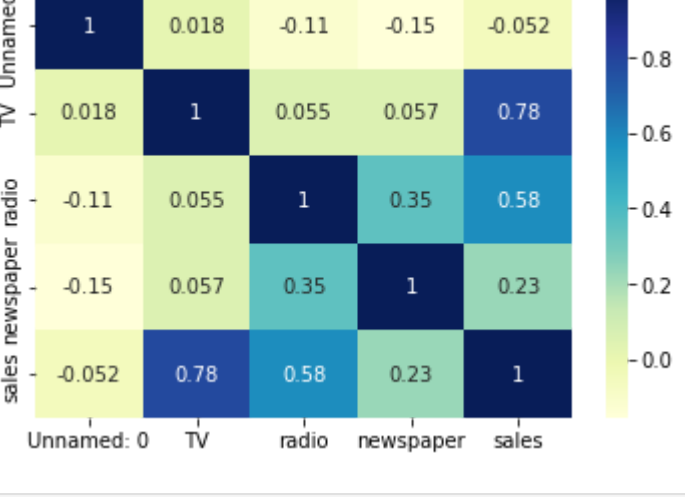
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keywo
rd will result in an error or misinterpretation.
warnings.warn(



In [23]: sns.pairplot(advertising, x_vars=['TV', 'newspaper', 'radio'], y_vars='sales', height=4, aspect=1, kind='scatter')
plt.show()



In [24]: # Let's see the correlation between different variables.
sns.heatmap(advertising.corr(), cmap="YlGnBu", annot = True)
plt.show()



In [25]: x = advertising['TV']
y = advertising['sales']

In [26]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.7, test_size = 0.3, random_state = 100)

In [27]: X_train.head()

Out[27]:
74    213.4
3      151.5
185    285.0
26    142.9
90     134.3
Name: TV, dtype: float64

In [28]: y_train.head()

Out[28]:
74     17.0
3      18.5
185    22.6
26     15.0
90     11.2
Name: sales, dtype: float64

In [29]: import statsmodels.api as sm

In [30]: X_train_sm = sm.add_constant(X_train)

C:\ProgramData\Anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only
x = pd.concat(X[:order], 1)

In [31]: lr = sm.OLS(y_train, X_train_sm).fit()

In [32]: lr.params

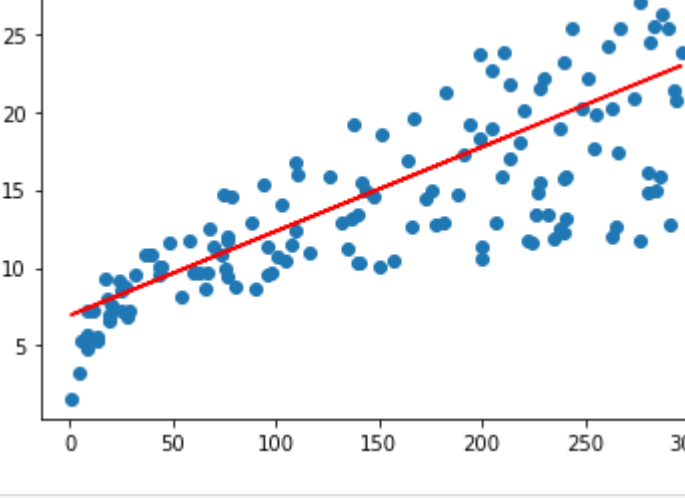
Out[32]:
const    6.989666
TV       0.046497
dtype: float64

In [33]: print(lr.summary())

OLS Regression Results
=====
Dep. Variable:  sales  R-squared:  0.613
Model:  OLS  Adj. R-squared:  0.611
Method:  Least Squares  F-statistic:  219.0
Date:  Sun, 19 Jun 2022  Prob (F-statistic):  2.84e-38
Time:  23:19:06  Log-Likelihood:  -370.62
No. Observations:  140  AIC:  745.2
DF Residuals:  138  BIC:  745.1
DF Model:  1
Covariance Type:  nonrobust
=====
coef    std err    t    P>|t|    [0.025    0.975]
-----
const    6.9897    0.548    12.762    0.000    5.907    8.073
TV       0.0465    0.003    14.798    0.000    0.040    0.053
=====
Omnibus:  0.995  Durbin-Watson:  1.963
Prob(Omnibus):  0.608  Jarque-Bera (JB):  0.970
Skew:  -0.088  Prob(JB):  0.616
Kurtosis:  2.593  Cond. No.  320.
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

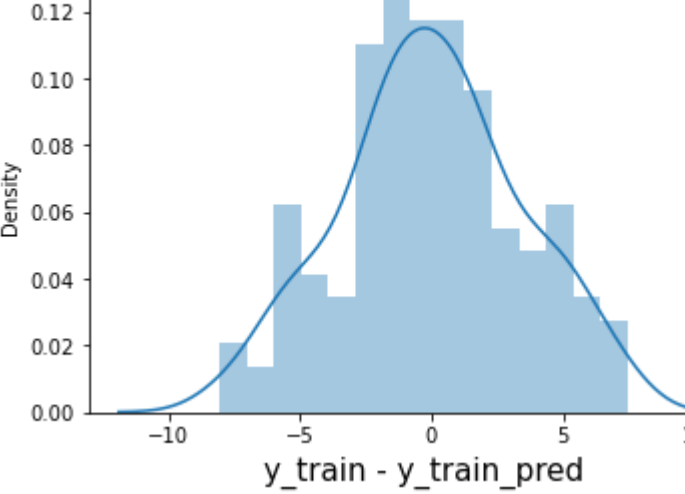
In [34]: plt.scatter(X_train, y_train)
plt.plot(X_train, 6.940 + 0.054*X_train, 'r')
plt.show()



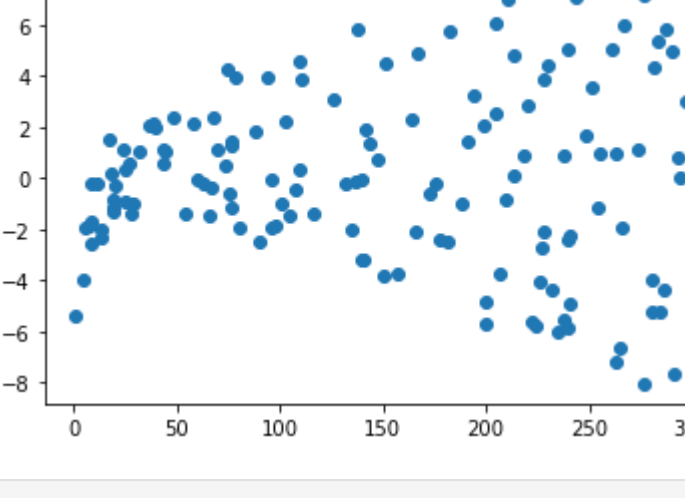
In [35]: y_train_pred = lr.predict(X_train_sm)
res = (y_train - y_train_pred)

In [36]: fig = plt.figure()
sns.distplot(res, bins = 15)
fig.suptitle('Error Terms', fontsize = 15) # Plot heading
plt.xlabel('y_train - y_train_pred', fontsize = 15) # X-label
plt.show()

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: 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 flexibi
lity) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



In [37]: plt.scatter(X_train,res)
plt.show()



In [38]: # Add a constant to X_test
X_test_sm = sm.add_constant(X_test)

C:\ProgramData\Anaconda3\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only
x = pd.concat(X[:order], 1)

In [39]: # Predict the y values corresponding to X_test_sm
y_pred = lr.predict(X_test_sm)

In [40]: y_pred.head()

Out[40]:
126    7.352345
184    18.865337
99     13.276109
92     17.112141
111    18.228077
dtype: float64

In [41]: from sklearn.metrics import mean_squared_error
from sklearn.metrics import r2_score

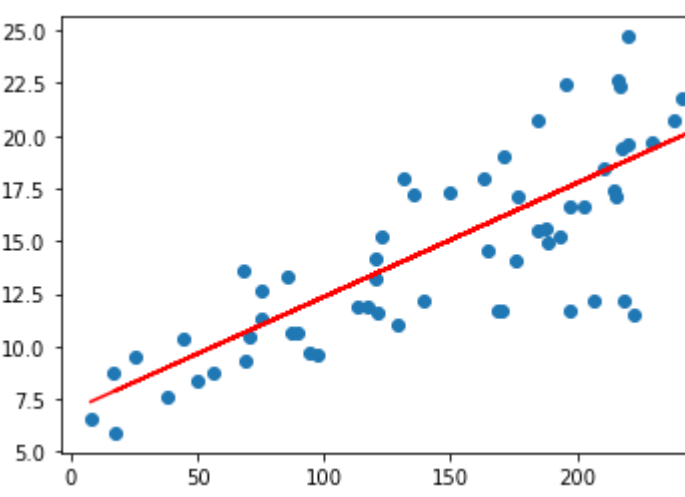
In [42]: np.sqrt(mean_squared_error(y_test, y_pred))

Out[42]:
2.8241456288327016

In [43]: # Checking the R-squared on the test set
r_squared = r2_score(y_test, y_pred)
r_squared

Out[43]:
0.59429872677833

In [44]: plt.scatter(X_test, y_test)
plt.plot(X_test, 6.940 + 0.054 * X_test, 'r')
plt.show()


```