

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

```
In [2]: advertising = pd.read_csv("Advertising.csv")
```

```
In [3]: advertising
```

Out[3]:

	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	180.8	10.8	58.4	12.9
...	...	...	...	...	...
195	196	196	38.2	3.7	13.8
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 × 5 columns

```
In [4]: advertising.shape
```

```
Out[4]: (200, 5)
```

```
In [5]: 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 [6]: advertising.describe()
```

Out[6]:

	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	149.750000	22.900000	25.750000	12.900000
75%	150.250000	218.825000	36.525000	45.100000	17.400000
max	200.000000	296.400000	49.600000	114.000000	27.000000

```
In [7]: advertising.isnull().sum()
```

```
Out[7]: Unnamed: 0    0
TV            0
radio         0
newspaper     0
sales         0
dtype: int64
```

```
In [8]: # 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 keyword will result in an error or misinterpretation.

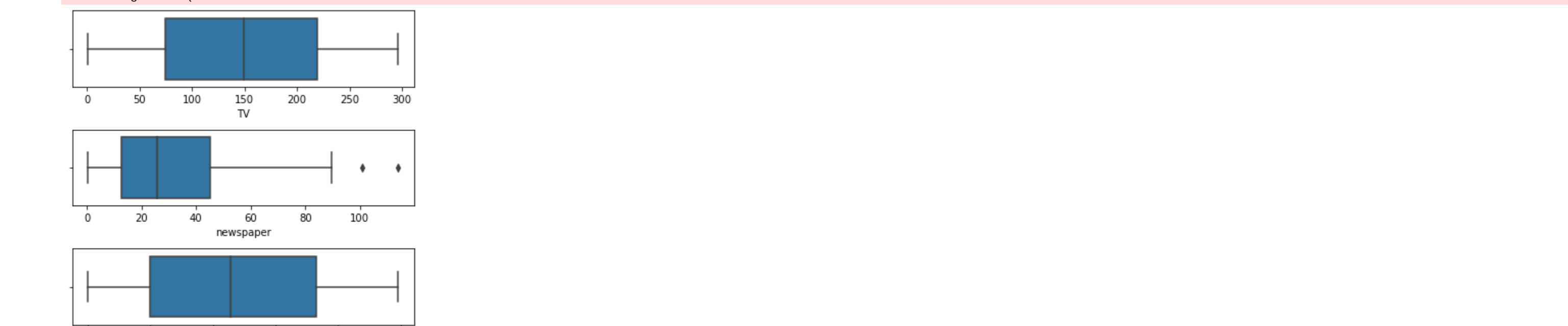
warnings.warn(

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 keyword 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 keyword will result in an error or misinterpretation.

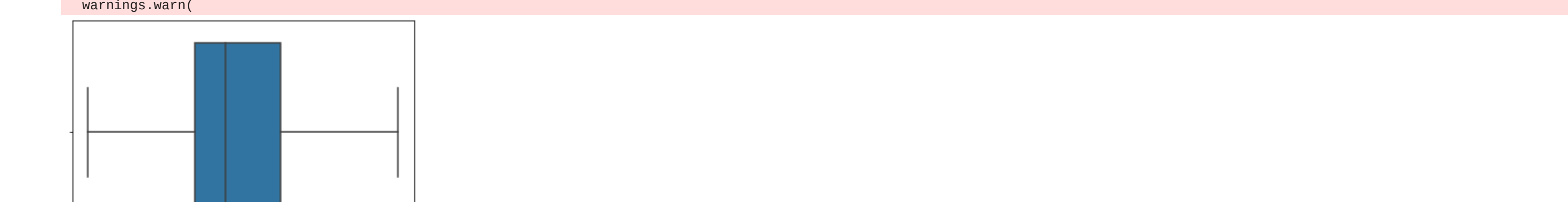
warnings.warn(



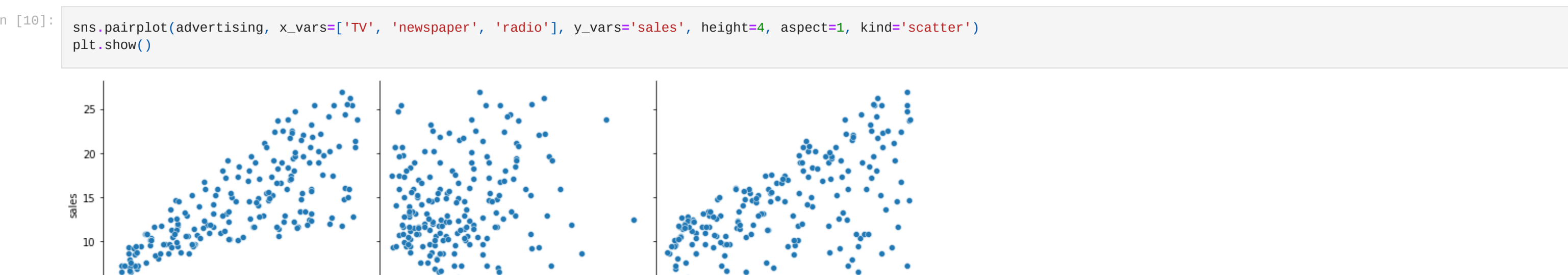
```
In [9]: sns.boxplot(advertising['sales'])
plt.show()
```

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 keyword will result in an error or misinterpretation.

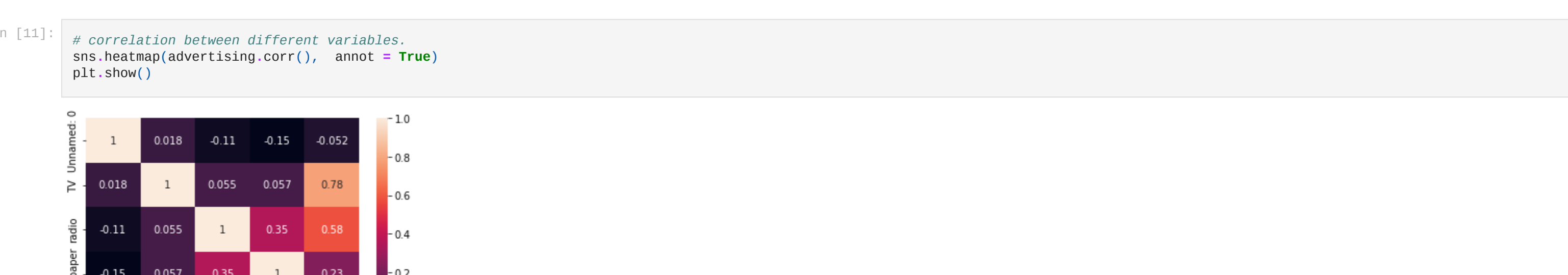
warnings.warn(



```
In [10]: sns.pairplot(advertising, x_vars=['TV', 'newspaper', 'radio'], y_vars='sales', height=4, aspect=1, kind='scatter')
plt.show()
```



```
In [11]: # correlation between different variables.
sns.heatmap(advertising.corr(), annot = True)
plt.show()
```



```
In [12]: X = advertising['TV']
y = advertising['sales']
```

```
In [13]: 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 [14]: X_train.head()
```

```
Out[14]: 74    213.4
3       151.5
185     285.0
26      142.9
90      134.3
Name: TV, dtype: float64
```

```
In [15]: y_train.head()
```

```
Out[15]: 74      17.0
3       16.5
185     22.6
26      15.0
90      11.2
Name: sales, dtype: float64
```

```
In [16]: import statsmodels.api as sm
```

```
In [17]: 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 [18]: lr = sm.OLS(y_train, X_train_sm).fit()
```

```
In [19]: lr.params
```

```
Out[19]: const    6.989666
TV         0.046497
dtype: float64
```

```
In [20]: 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:                   Mon, 20 Jun 2022      Prob (F-statistic):      2.84e-30
Time:                   21:32:35              Log-Likelihood:       -379.62
No. Observations:       140                  AIC:              745.2
Df Residuals:           138                  BIC:              751.1
Df Model:                1
Covariance Type:         nonrobust
=====
               coef      std err          t      Pr>|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:       1.983
Prob(Omnibus):           0.608    Jarque-Bera (JB):      0.979
Skew:                    0.698    Prob(JB):          0.616
Kurtosis:                2.593    Cond. No.         328.
=====

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

```
In [21]: plt.scatter(X_train, y_train)
plt.plot(X_train, 6.948 + 0.054*X_train, 'r')
plt.show()
```

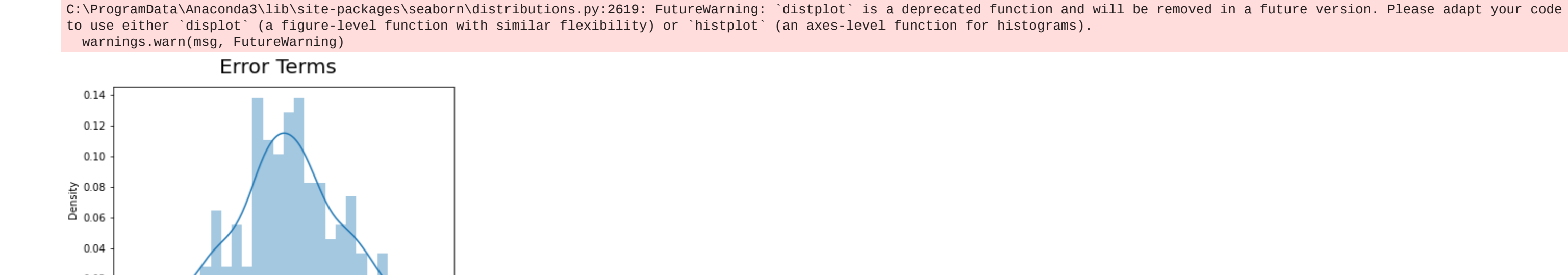


```
In [22]: y_train_pred = lr.predict(X_train_sm)
res = (y_train - y_train_pred)
```

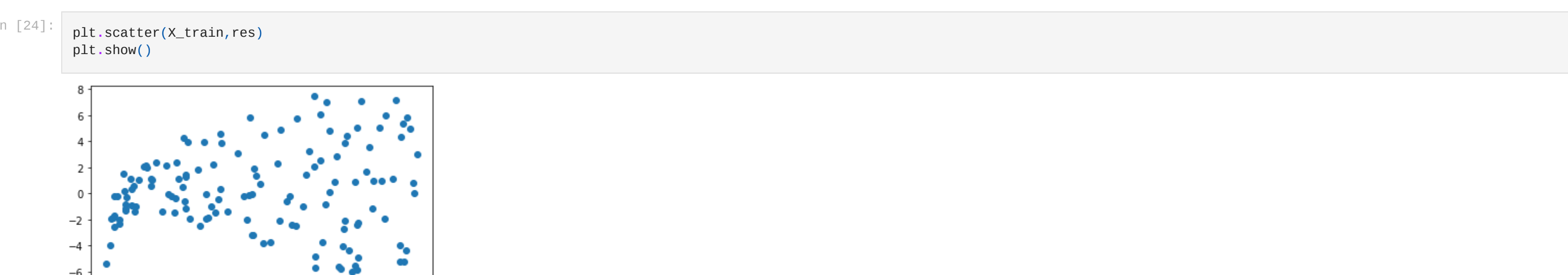
```
In [23]: fig = plt.figure()
sns.distplot(res, bins = 20)
fig.suptitle('Error Terms', fontsize = 20) # Plot heading
plt.xlabel('y_train - y_train_pred', fontsize = 20) # x-label
plt.show()
```

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 flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)



```
In [24]: plt.scatter(X_train, res)
plt.show()
```



```
In [30]: # Adding 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 [31]: # Predicting the y values corresponding to X_test_sm
y_pred = lr.predict(X_test_sm)
```

```
In [32]: y_pred.head()
```

```
Out[32]: 126     7.352345
104    18.005337
99     13.276109
92     17.112141
111    18.228077
dtype: float64
```

```
In [33]: from sklearn.metrics import mean_squared_error
from sklearn.metrics import r2_score
```

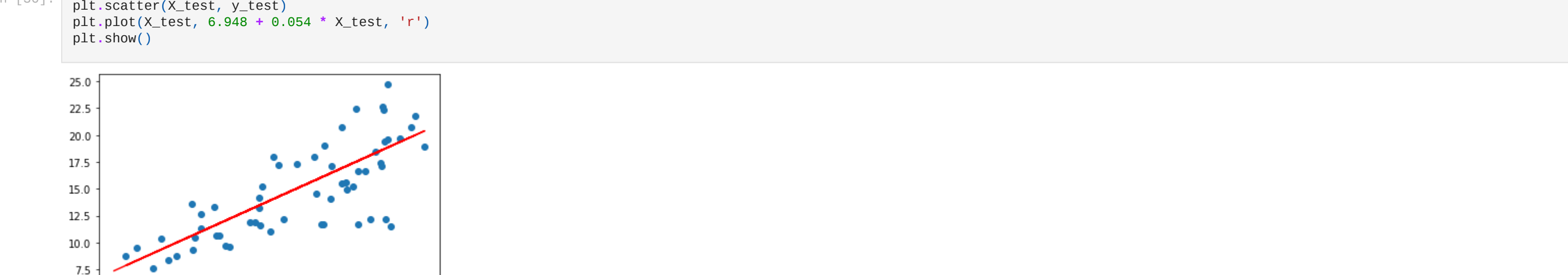
```
In [34]: np.sqrt(mean_squared_error(y_test, y_pred))
```

```
Out[34]: 2.8241456288327816
```

```
In [35]: # Checking R-squared on the test set
r_squared = r2_score(y_test, y_pred)
r_squared
```

```
Out[35]: 0.59429872677833
```

```
In [36]: plt.scatter(X_test, y_test)
plt.plot(X_test, 6.948 + 0.054 * X_test, 'r')
plt.show()
```



```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```