



Open in app

Get started



Published in Towards Data Science

You have **2** free member-only stories left this month.

[Sign up for Medium and get an extra one](#)



Kaushik Katari

Follow

Oct 9, 2020 · 8 min read · ✨ · 🎧 Listen

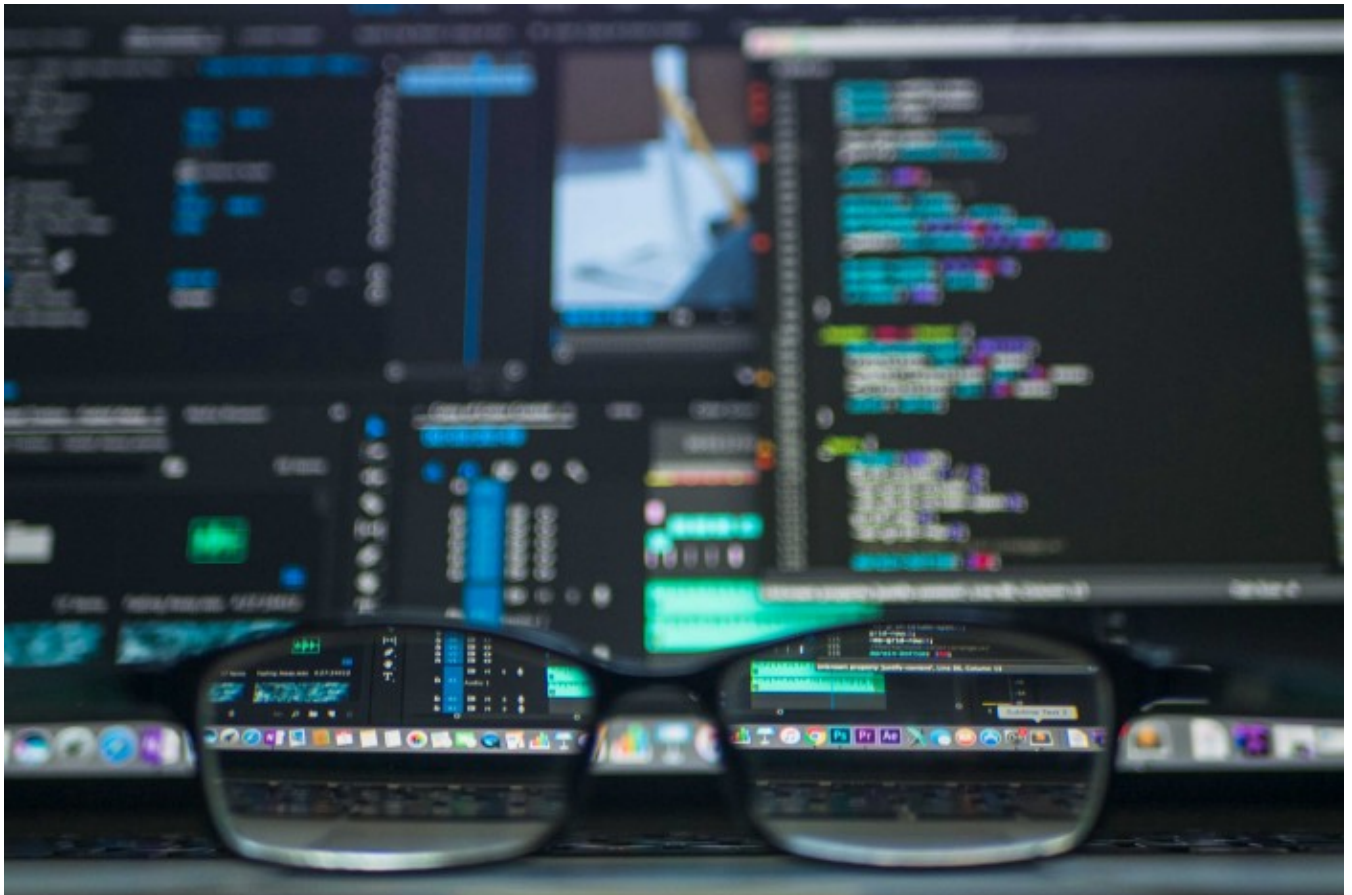


Save



Simple Linear Regression Model using Python: Machine Learning

Learning how to build a simple linear regression model in machine learning using Jupyter notebook in Python



[Open in app](#)[Get started](#)

In the previous article, [the Linear Regression Model](#), we have seen how the linear regression model works theoretically using Microsoft Excel. This article will see how we can build a linear regression model using Python in the Jupyter notebook.

Simple Linear Regression

To predict the relationship between two variables, we'll use a simple linear regression model.

In a simple linear regression model, we'll predict the outcome of a variable known as the dependent variable using only one independent variable.

We'll directly dive into building the model in this article. More about the linear regression model and the factors we have to consider are explained in detail [here](#).

Building a linear regression model

To build a linear regression model in python, we'll follow five steps:

1. Reading and understanding the data
2. Visualizing the data
3. Performing simple linear regression
4. Residual analysis
5. Predictions on the test set

Reading and understanding the data

In this step, first, we'll import the necessary libraries to import the data. After that, we'll perform some basic commands to understand the structure of the data.

We can download the sample dataset, which we'll be



[Open in app](#)[Get started](#)

Import libraries

We'll import the `numpy` and `pandas` library in the Jupyter notebook and read the data using `pandas`.

```
1 # Importing numpy and pandas libraries to read the data
2
3 # Supress Warnings
4 import warnings
5 warnings.filterwarnings('ignore')
6
7 # Import the numpy and pandas package
8 import numpy as np
9 import pandas as pd
10
11 # Read the given CSV file, and view some sample records
12 advertising = pd.read_csv("Company_data.csv")
13 advertising
```

lrm_import_libraries.py hosted with ❤ by GitHub

[view raw](#)

The dataset looks like this. Here our target variable is the Sales column.

	TV	Radio	Newspaper	Sales
0	230.1	37.8	69.2	22.1
1	44.5	39.3	45.1	10.4
2	17.2	45.9	69.3	12.0
3	151.5	41.3	58.5	16.5
4	180.8	10.8	58.4	17.9
...
195	38.2	3.7	13.8	7.6
196	94.2	4.9	8.1	14.0
197	177.0	9.3	6.4	14.8
198	283.6	42.0	66.2	25.5
199	232.1	8.6	8.7	18.4



[Open in app](#)[Get started](#)

```
1 # Shape of our dataset
2 advertising.shape
3
4 # Info our dataset
5 advertising.info()
6
7 # Describe our dataset
8 advertising.describe()
```

lrm_shape.py hosted with ❤ by GitHub

[view raw](#)

The `shape` of our dataset is,

```
(200, 4)
```

Using the `info`, we can see whether there are any null values in the data. If yes, then we have to do some data manipulation.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 4 columns):
TV                200 non-null float64
Radio             200 non-null float64
Newspaper         200 non-null float64
Sales             200 non-null float64
dtypes: float64(4)
memory usage: 6.4 KB
```

Info of the dataset

As we can observe, there are no null values present in the data.

Using `describe`, we'll see whether there is any sudden jump in the data's values.



[Open in app](#)[Get started](#)

	TV	Radio	Newspaper	Sales
count	200.000000	200.000000	200.000000	200.000000
mean	147.042500	23.264000	30.554000	15.130500
std	85.854236	14.846809	21.778621	5.283892
min	0.700000	0.000000	0.300000	1.600000
25%	74.375000	9.975000	12.750000	11.000000
50%	149.750000	22.900000	25.750000	16.000000
75%	218.825000	36.525000	45.100000	19.050000
max	296.400000	49.600000	114.000000	27.000000

Describing the dataset

The values present in the columns are pretty consistent throughout the data.

Visualizing the data

Let's now visualize the data using the `matplotlib` and `seaborn` library. We'll make a pairplot of all the columns and see which columns are the most correlated to `Sales`.

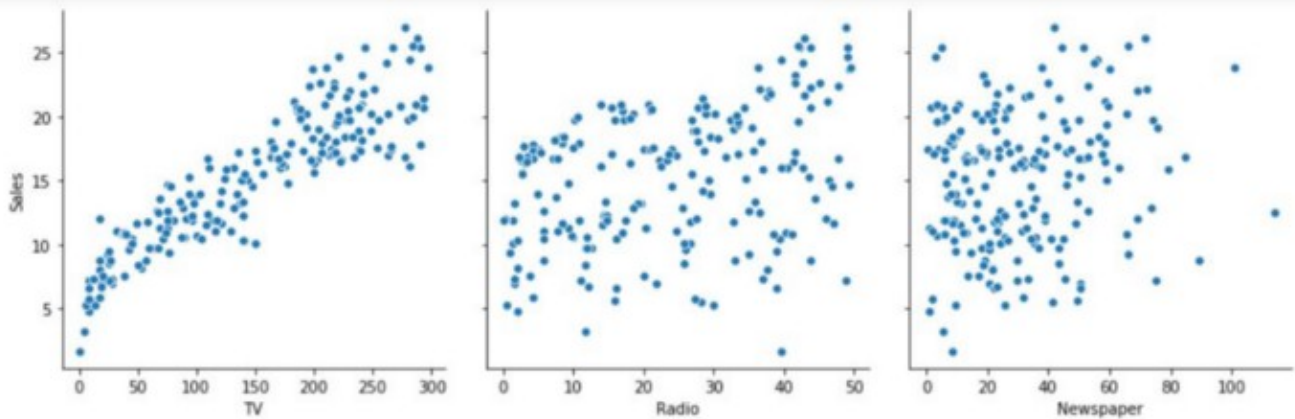
```
1 # Import matplotlib and seaborn libraries to visualize the data
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Using pairplot we'll visualize the data for correlation
6 sns.pairplot(advertising, x_vars=['TV', 'Radio', 'Newspaper'],
7             y_vars='Sales', size=4, aspect=1, kind='scatter')
8 plt.show()
```

Irm_visualize.py hosted with ❤ by GitHub

[view raw](#)

It is always better to use a scatter plot between two numeric variables. The pairplot for the above code looks like,



[Open in app](#)[Get started](#)

Pairplot of each Column w.r.t. Sales column

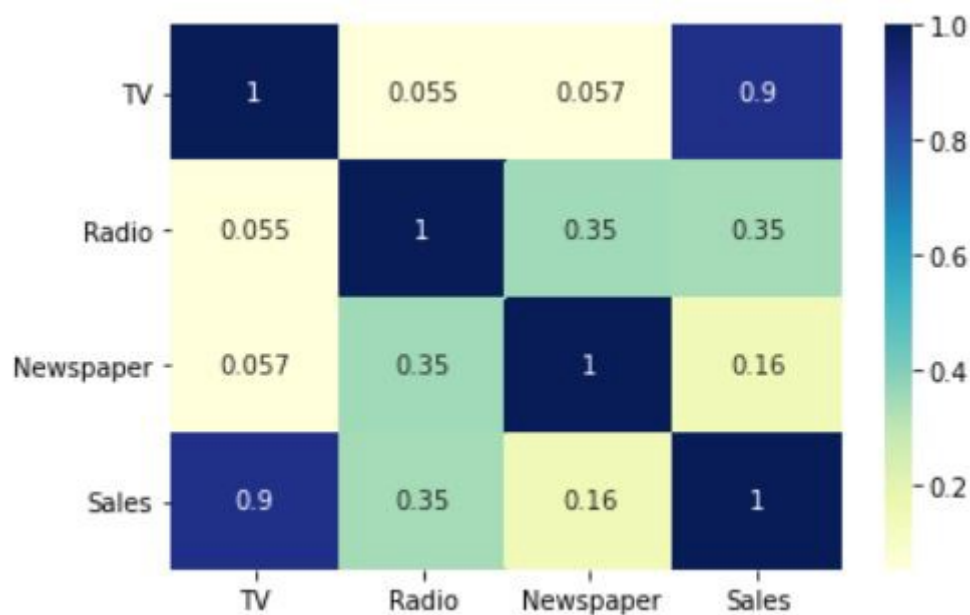
If we cannot determine the correlation using a scatter plot, we can use the seaborn heatmap to visualize the data.

```
1 # Visualizing the data using heatmap
2 sns.heatmap(advertising.corr(), cmap="YlGnBu", annot = True)
3 plt.show()
```

lrm_heatmap.py hosted with ❤ by GitHub

[view raw](#)

The heatmap looks like this,



Heatmap of all the columns in the data



[Open in app](#)[Get started](#)

Performing Simple Linear Regression

Equation of simple linear regression

$$y = c + mX$$

In our case:

$$y = c + m * TV$$

The m values are known as **model coefficients** or **model parameters**.

We'll perform simple linear regression in four steps.

1. Create X and y
2. Create Train and Test set
3. Train your model
4. Evaluate the model

Create X and y First, we'll assign our feature variable/column TV as x and our target variable $Sales$ as y .

To generalize,

The independent variable represents x , and y represents the target variable in a simple linear regression model.

```
1 # Creating X and y
2 X = advertising['TV']
3 y = advertising['Sales']
```

lrm_X_y.py hosted with ❤ by GitHub

[view raw](#)

Create Train and Test sets

We need to split our variables into training and testing sets. Using the training set, we'll build the model and perform the model on the testing set. We'll divide the training and testing sets into a 7:3 ratio, respectively.

We'll split the data by using the `train_test_split` function from the `sklearn` library.





Open in app

Get started

```

2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.7,
4                                                    test_size = 0.3, random_state = 100)

```

lrm_train_test.py hosted with ❤ by GitHub

view raw

Let's take a look at the training dataset,

```

1 # Take a look at the train dataset
2 X_train
3 y_train

```

lrm_X_y_train.py hosted with ❤ by GitHub

view raw

X_train data looks like this after splitting.

```

74      213.4
3       151.5
185     205.0
26      142.9
90      134.3
...
87      110.7
103     187.9
67      139.3
24       62.3
8         8.6
Name: TV, Length: 140, dtype: float64

```

X_train data after splitting

y_train data looks like this after splitting.

```

74      17.0
3       16.5
185     22.6
26      15.0
90      14.0
...
87      16.0
103     19.7

```



[Open in app](#)[Get started](#)

y_train data after splitting

Building and training the model

Using the following two packages, we can build a simple linear regression model.

- `statsmodel`
- `sklearn`

First, we'll build the model using the `statsmodel` package. To do that, we need to import the `statsmodel.api` library to perform linear regression.

By default, the `statsmodel` library fits a line that passes through the origin. But if we observe the simple linear regression equation $y = c + mX$, it has an intercept value as c . So, to have an intercept, we need to add the `add_constant` attribute manually.

```
1 # Importing Statsmodels.api library from Stamodel package
2 import statsmodels.api as sm
3
4 # Adding a constant to get an intercept
5 X_train_sm = sm.add_constant(X_train)
```

lrm_train_model.py hosted with ❤ by GitHub

[view raw](#)

Once we've added constant, we can fit the regression line using `OLS` (Ordinary Least Square) method present in the `statsmodel`. After that, we'll see the parameters, i.e., c and m of the straight line.

```
1 # Fitting the regression line using 'OLS'
2 lr = sm.OLS(y_train, X_train_sm).fit()
3
4 # Printing the parameters
5 lr.params
```

lrm_OLS.py hosted with ❤ by GitHub

[view raw](#)

The output is,





Open in app

Get started

Intercept and Slope of the line

Let's see the summary of all the different parameters of the regression line fitted like R^2 , probability of F-statistic, and p-value.

```
1 # Performing a summary to list out all the different parameters of the regression line fitted
2 lr.summary()
```

lr_summary.py hosted with ❤ by GitHub

[view raw](#)

The statistics for the above regression line is,

OLS Regression Results

Dep. Variable:	Sales	R-squared:	0.816
Model:	OLS	Adj. R-squared:	0.814
Method:	Least Squares	F-statistic:	611.2
Date:	Fri, 09 Oct 2020	Prob (F-statistic):	1.52e-52
Time:	14:33:49	Log-Likelihood:	-321.12
No. Observations:	140	AIC:	646.2
Df Residuals:	138	BIC:	652.1
Df Model:	1		
Covariance Type:	nonrobust		
	coef	std err	t P> t [0.025 0.975]
const	6.9487	0.385	18.068 0.000 6.188 7.709
TV	0.0545	0.002	24.722 0.000 0.050 0.059
Omnibus:	0.027	Durbin-Watson:	2.196
Prob(Omnibus):	0.987	Jarque-Bera (JB):	0.150
Skew:	-0.006	Prob(JB):	0.928
Kurtosis:	2.840	Cond. No.	328.

All the statistics for the above best-fit line





Open in app

Get started

1. The `coefficients` and its `p-value` (significance)
2. `R-squared` value
3. `F-statistic` and its significance

OLS Regression Results

Dep. Variable:	Sales	R-squared:	0.816
Model:	OLS	Adj. R-squared:	0.814
Method:	Least Squares	F-statistic:	611.2
Date:	Fri, 09 Oct 2020	Prob (F-statistic):	1.52e-52
Time:	14:33:49	Log-Likelihood:	-321.12
No. Observations:	140	AIC:	646.2
Df Residuals:	138	BIC:	652.1
Df Model:	1		
Covariance Type:	nonrobust		

	coef	std err	t	P> t 	[0.025	0.975]
const	6.9487	0.385	18.068	0.000	6.188	7.709
TV	0.0545	0.002	24.722	0.000	0.050	0.059

Omnibus:	0.027	Durbin-Watson:	2.196
Prob(Omnibus):	0.987	Jarque-Bera (JB):	0.150
Skew:	-0.006	Prob(JB):	0.928
Kurtosis:	2.840	Cond. No.	328.

Statistics we need to look

1. The `coefficient` for TV is 0.054, and its corresponding `p-value` is very low, almost 0. That means the `coefficient` is statistically significant.

We have to make sure that the `p-value` should



[Open in app](#)[Get started](#)

2. `R-squared` value is 0.816, which means that 81.6% of the `Sales` variance can be explained by the `TV` column using this line.

3. `Prob F-statistic` has a very low `p-value`, practically zero, which gives us that the model fit is statistically significant.

Since the fit is significant, let's go ahead and visualize how well the straight-line fits the scatter plot between `TV` and `Sales` columns.

From the parameters, we got the values of the `intercept` and the `slope` for the straight line. The equation of the line is,

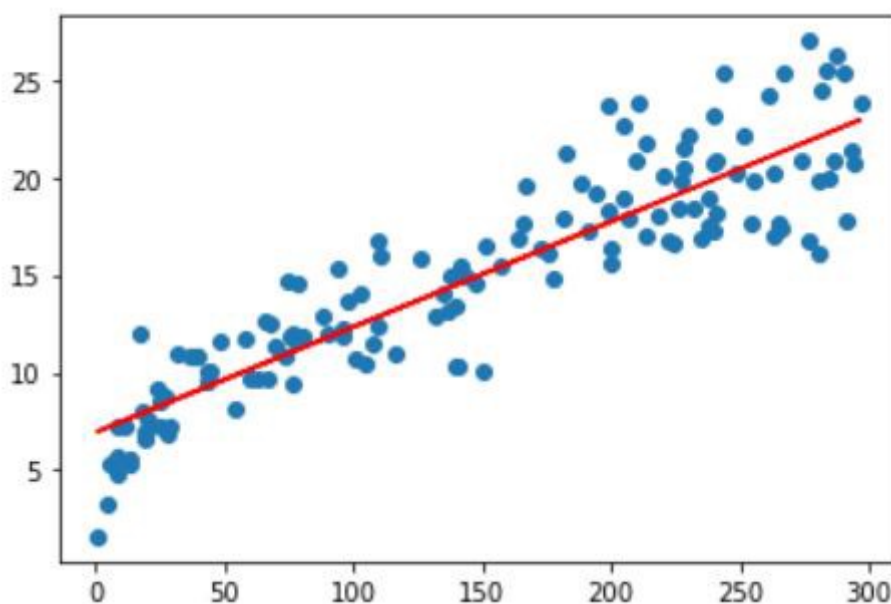
$$\text{Sales} = 6.948 + 0.054 * \text{TV}$$

```
1 # Visualizing the regression line
2 plt.scatter(X_train, y_train)
3 plt.plot(X_train, 6.948 + 0.054*X_train, 'r')
4 plt.show()
```

lrm_visualize_line.py hosted with ❤ by GitHub

[view raw](#)

The graph looks like this,



Best-fit regression line



[Open in app](#)[Get started](#)

Residual Analysis

One of the major assumptions of the linear regression model is the error terms are normally distributed.

Error = Actual y value - y predicted value

Now from the dataset,

We have to predict the y value from the training dataset of X using the `predict` attribute. After that, we'll create the error terms(Residuals) from the predicted data.

```
1 # Predicting y_value using traingn data of X
2 y_train_pred = lr.predict(X_train_sm)
3
4 # Creating residuals from the y_train data and predicted y_data
5 res = (y_train - y_train_pred)
```

lrm_residual.py hosted with ❤ by GitHub

[view raw](#)

Now, let's plot the histogram of the residuals and see whether it looks like normal distribution or not.

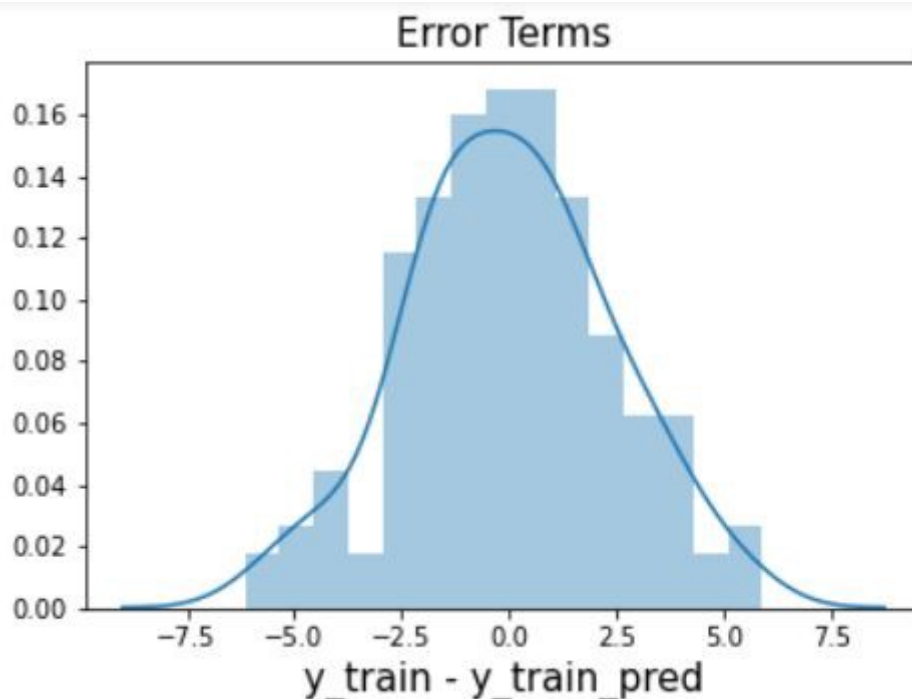
```
1 # Plotting the histogram using the residual values
2 fig = plt.figure()
3 sns.distplot(res, bins = 15)
4 plt.title('Error Terms', fontsize = 15)
5 plt.xlabel('y_train - y_train_pred', fontsize = 15)
6 plt.show()
```

lrm_residual_histogram.py hosted with ❤ by GitHub

[view raw](#)

The histogram of the residuals looks like,



[Open in app](#)[Get started](#)

Residuals distribution

As we can see, the residuals are following the normal distribution graph with a mean 0.

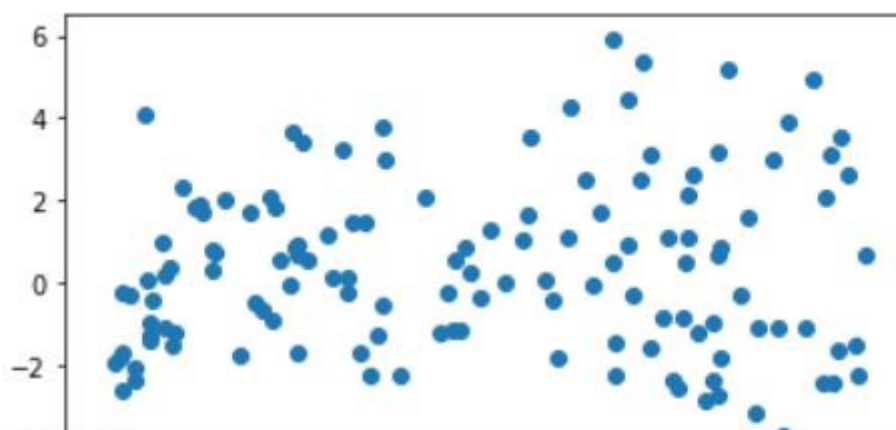
Now, make sure that the residuals are not following any specific pattern.

```
1 # Looking for any patterns in the residuals
2 plt.scatter(X_train, res)
3 plt.show()
```

lrm_residuals_pattern.py hosted with ❤ by GitHub

[view raw](#)

The scatter plot looks like,



[Open in app](#)[Get started](#)

Scatter plot of Residual values

Since the Residuals follow a normal distribution and do not follow any specific pattern, we can use the linear regression model we have built to evaluate test data.

Predictions on the Test data or Evaluating the model

Now that we have fitted the regression line on our train dataset, we can make some predictions to the test data. Similar to the training dataset, we have to `add_constant` to the test data and predict the y values using the `predict` attribute present in the `statsmodel`.

```
1 # Adding a constant to X_test
2 X_test_sm = sm.add_constant(X_test)
3
4 # Predicting the y values corresponding to X_test_sm
5 y_test_pred = lr.predict(X_test_sm)
6
7 # Printing the first 15 predicted values
8 y_test_pred
```

lrm_test_data.py hosted with ❤ by GitHub

[view raw](#)

The predicted y-values on test data are,

126	7.374140
104	19.941482
99	14.323269
92	18.823294
111	20.132392
167	18.228745
116	14.541452
96	17.726924
52	18.752384
69	18.774202
164	13.341445
124	19.466933
182	10.014155
154	17.192376
125	11.705073



[Open in app](#)[Get started](#)

Now, let's calculate the R^2 value for the above-predicted y-values. We can do that by merely importing the `r2_score` library from `sklearn.metrics` package.

```
1 # Importing r2_square
2 from sklearn.metrics import r2_score
3
4 # Checking the R-squared value
5 r_squared = r2_score(y_test, y_test_pred)
6 r_squared
```

lrm_r_2.py hosted with ❤ by GitHub

[view raw](#)

The R^2 value by using the above code = 0.792

If we can remember from the training data, the R^2 value = 0.815

Since the R^2 value on test data is within 5% of the R^2 value on training data, we can conclude that the model is pretty stable. Which means, what the model has learned on the training set can generalize on the unseen test set.

Let's visualize the line on the test data.

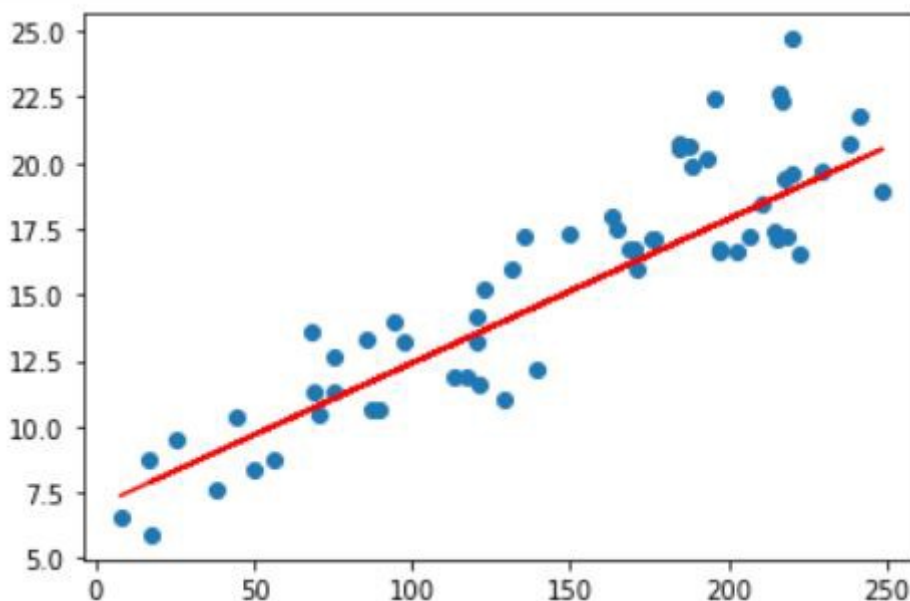
```
1 # Visualize the line on the test set
2 plt.scatter(X_test, y_test)
3 plt.plot(X_test, y_test_pred, 'r')
4 plt.show()
```

lrm_line_test.py hosted with ❤ by GitHub

[view raw](#)

The scatter-plot with best-fit line looks like,



[Open in app](#)[Get started](#)

Best-fit line on test data

This is how we build a linear regression model using the `statsmodel` package.

Apart from the `statsmodel`, we can build a linear regression model using `sklearn`. Using the `linear_model` library from `sklearn`, we can make the model.

Similar to `statsmodel`, we'll split the data into `train` and `test`.

```
1 # Splitting the data into train and test
2 from sklearn.model_selection import train_test_split
3 X_train_lm, X_test_lm, y_train_lm, y_test_lm = train_test_split(X, y, train_size = 0.7,
4                                                                 test_size = 0.3, random_state =
```

lrm_sk_train_test_split.py hosted with ❤ by GitHub

[view raw](#)

For simple linear regression, we need to add a column to perform the regression fit properly.

```
1 # Shape of the train set without adding column
2 X_train_lm.shape
3
4 # Adding additional column to the train and test data
5 X_train_lm = X_train_lm.values.reshape(-1,1)
```



[Open in app](#)[Get started](#)

The shape of X_train before adding a column is (140,).

The shape of X for train and test data is (140, 1).

Now, let's fit the line to the plot importing the `LinearRegression` library from the `sklearn.linear_model`.

```
1 from sklearn.linear_model import LinearRegression
2
3 # Creating an object of Linear Regression
4 lm = LinearRegression()
5
6 # Fit the model using .fit() method
7 lm.fit(X_train_lm, y_train_lm)
```

lrm_sk_fit_model.py hosted with ❤ by GitHub

[view raw](#)

Now, let's find the coefficients of the model.

```
1 # Intercept value
2 print("Intercept :",lm.intercept_)
3
4 # Slope value
5 print('Slope :',lm.coef_)
```

lrm_sk_coeff.py hosted with ❤ by GitHub

[view raw](#)

The value of intercept and slope is,

```
Intercept : 6.948683200001357
Slope : [0.05454575]
```

Coefficient Values

The straight-line equation we get for the above values is,

$\text{Sales} = 6.948 + 0.054 * \text{TV}$

If we observe, the equation we got here is the same as the one we got in the

`statsmodel`.



[Open in app](#)[Get started](#)

```
1 # Making Predictions of y_value
2 y_train_pred = lm.predict(X_train_lm)
3 y_test_pred = lm.predict(X_test_lm)
4
5 # Comparing the r2 value of both train and test data
6 print(r2_score(y_train,y_train_pred))
7 print(r2_score(y_test,y_test_pred))
```

lrm_sk_pred_r2.py hosted with ❤ by GitHub

[view raw](#)

The R^2 values of the train and test data are

R^2 train_data = 0.816

R^2 test_data = 0.792

Same as the `statsmodel`, the R^2 value on test data is within 5% of the R^2 value on training data. We can apply the model to the unseen test set in the future.

Conclusion

As we have seen, we can build a linear regression model using either a `statsmodel` or `sklearn`.

We have to make sure to follow these five steps to build the simple linear regression model:

1. Reading and understanding the data
2. Visualizing the data
3. Performing simple linear regression
4. Residual analysis
5. Predictions on the test set

In the next article, we'll see how the multiple linear regression model works.

Thank you for reading and happy coding!!!

[Check out my previous articles here](#)





Open in app

Get started

- [Central Limit Theorem\(CLT\): Data Science](#)
- [Inferential Statistics: Data Analysis](#)
- [Seaborn: Python](#)
- [Pandas: Python](#)
- [Matplotlib: Python](#)
- [NumPy: Python](#)

References

- **Machine Learning — Linear Regression:**
https://www.w3schools.com/python/python_ml_linear_regression.asp
- **Linear Regression in Python:** <https://realpython.com/linear-regression-in-python/>
- **Linear Regression (Python Implementation):**
<https://www.geeksforgeeks.org/linear-regression-python-implementation/>
- **A Beginner's Guide to Linear Regression in Python with Scikit-Learn:**
<https://www.kdnuggets.com/2019/03/beginners-guide-linear-regression-python-scikit-learn.html>

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)



[Open in app](#)[Get started](#)[Get this newsletter](#)[About](#) [Help](#) [Terms](#) [Privacy](#)

Get the Medium app

