

Linear Regression in Python

Jose Chavez

July 26, 2023

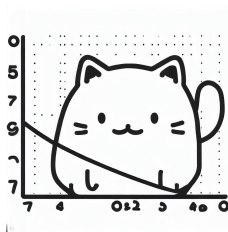


Table of Contents

- 1 Introduction
- 2 Data Set: Fake Twitter Data
- 3 Python and Linear Regression
- 4 An Example
 - Data Collection

Recap and Introduction

- Linear Regressions: Dependent variable (y) and one or more independent variables (x_1, x_2, \dots, x_n).

Fitting a linear equation to the observed data.

Minimize differences between. real and predicted

The linear regression model is represented by the equation:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

where:

- y is the dependent variable (the one we want to predict),
 - b_0 is the intercept (the value of y when all x values are zero),
 - b_1, b_2, \dots, b_n are the coefficients (slope) for each independent variable x_1, x_2, \dots, x_n .
- Simple yet instructive case. . . good old $y = mx + b$

Hypothetical Project

Question:

Are older Twitter users less likely to use the app throughout the week? What's the relationship between age usage?

Data: Fake Twitter Data

- Hypothetical Data Scientist scenario:
You are asked to investigate the connection between age and tweet rates.
You are given an access token to an API.
- Dataset for our demonstration:
Python scripts you can find at: [github](#) were used to generate hundreds of records of consisting of records that include: age, date, tweet (a string);
- Why not the real thing?
Privacy: Twitter does not release even that user information.

```
1 const express = require('express');  
2 const app = express();  
3  
4 const fs = require('fs');  
5 const path = require('path');  
6  
7 app.get('/download', function(req, res){  
8   const token = req.query.token;  
9  
10  // Check if the token is valid  
11  if(token === 'playtoken'){  
12    res.status(401).send('Unauthorized: Invalid token');  
13    return;  
14  }  
15  
16  // If the token is valid, serve the CSV file  
17  const file = path.join(__dirname, 'linear_data.csv');  
18  res.download(file);  
19  
20 });  
21  
22 app.listen(3000, function(){  
23   console.log('Server is running on port 3000');  
24 });  
25
```

The Data Science Process

- **Understand the Problem:** Define problem and objective.
Here: The problem is finding a relationship between the ages and weekly tweet rates.
- **Data Collection:** Gather data.
- **Exploratory Data Analysis (EDA):** Explore the data to get a feel for its structure.
Here: We're going to view a scatter plot to get a rough idea whether or not we should expect any success.

The Data Science Processs (2)

- **Feature Selection:** Choosing relevant features (given in our problem).
Here: This is straightforward in our scenario.
- **Data Splitting:** split to train and to test.
- **Model Building:** Use data to create a model.
- **Model Evaluation:** Assess the model's performance using evaluation metrics like Mean Squared Error (MSE), R-squared, or Root Mean Squared Error (RMSE).

The Data Science Process(3)

- **Model Validation:** Validate the final model on the testing data to ensure it generalizes well to new, unseen data.
- **Interpretation:** Interpret your findings.
- **Deployment:** Deploy the trained model to make predictions on new data or integrate it into a larger application.
- **Monitoring and Maintenance:** Monitor the model's performance, retrain, tune, etc.

Python and Linear Regression

- Idea: Use Python's linear regression library

- Python Libraries:

Importing and cleaning the data:

- pandas: data structures to hydrate;
- python scripts for detecting and removing outliers and null values;

Exploring the data:

- matplotlib

Model training and evaluation:

- `import sklearn as sk`
- `sk.LinearRegression;`
-

Visualizing the model's results

- 1 matplotlib

Collect Data from our API

How does text interact with this.
Our training data we set aside
seems to “align” with the line...

ages	rates
54	53.855321
67	59.560101
44	58.634860
30	86.417711
58	46.608163
23	77.859235
36	48.900427

Collect Data from our API (2)

The JavaScript Node.js
Application serving the csv file we
generated with Python.

```
ranger (Python)  ● 3x1  ranger (nvim)  ● 3x2

2 const fs = require('fs');
3 const path = require('path');
4
5 app.get('/download', function(req, res){
6   const token = req.query.token;
7
8   // Check if the token is valid
9   if(token !== 'playtoken'){
10     res.status(401).send('Unauthorized: Invalid token');
11     return;
12   }
13
14   // If the token is valid, serve the CSV file
15   const file = path.join(__dirname, 'linear_data.csv');
16   res.download(file);
17 });
18
19 app.listen(3000, function(){
20   console.log('Server is running on port 3000');
21 });
22
~
~
~
```

To access: `curl -i http://csv.bluewifjc.`

`online/download?token=playtoken`



Hydrating a DataFrame

We load the csv into a pandas DataFrame.

```
1 import pandas as pd
2 df = pd.read_csv('filename.csv')
```

Examine Visuals of Data

Plotting ages versus rates per week Does it look linear?

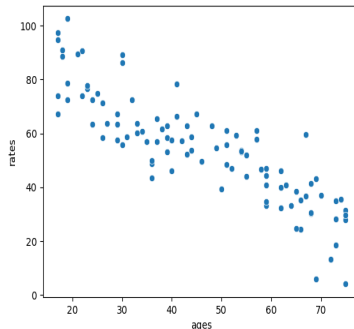


Figure: Scatter Plot of Data
Produced with Matplotlib

Cleanup

Removes negative 'rates' values

Drops rows with null values

Displays initial cleaned data

Calculates quartiles and IQR

Removes outliers from data

Displays data without outliers

```
1  # Drop negative values
2  df = df[df['rates'] >= 0]
3  # Drop null values
4  df = df.dropna()
5
6  # Display the first few rows of the
   dataframe
7  df.head()
8  Q1 = df.quantile(0.25)
9  Q3 = df.quantile(0.75)
10 IQR = Q3 - Q1
11
12 # Remove outliers
13 df_no_outliers = df[~((df < (Q1 - 1.5 *
   IQR)) | (df > (Q3 + 1.5 * IQR))).
   any(axis=1)]
14
15 # Display the first few rows of the
   dataframe without outliers
16 df_no_outliers.head()
```

Figure: Pyt Data Scatter Plot With
Line

Clean Up

Code run and differences (post
IQR Method)

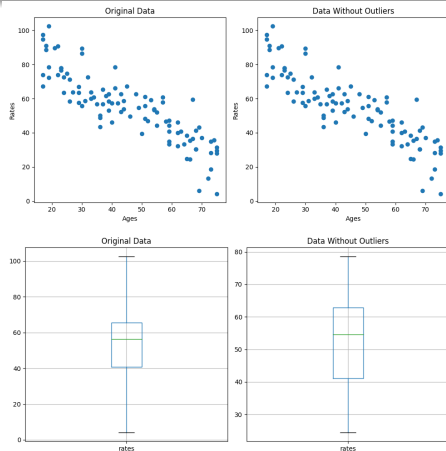


Figure: Box Plot Shows Outliers
Have Been Extracted

Library Imports and Splitting the Data

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error, r2_score
4 # Split the data into training and test sets
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
6     random_state=42)
7 # Create a linear regression object
8 regr = LinearRegression()
9
10 # Train the model using the training sets
11 regr.fit(X_train, y_train)
```


Gathering the Model's Parameters and Testing

```
12 # Make predictions using the testing set
13 y_pred = regr.predict(X_test)
14 # Print the coefficients
15 #print('Coefficients:', regr.coef_)
16 model_cos=regr.coef_
17
18 # Print the mean squared error
19 #print('Mean squared error:', mean_squared_error(y_test, y_pred))
20 mse=mean_squared_error(y_test, y_pred)
21
22 coefficients = regr.coef_[0]
23 intercept = regr.intercept_
24 # Print the coefficient of determination (R^2 score)
25 #print('Coefficient of determination (R^2 score):', r2_score(y_test, y_pred))
26 #r2score=r2_score(y_test,y_pred)
```

The Model We Have Created...

Model Coefficients: Model Intercept (only one):

Mean Squared Error: [-0.74629129]

and MSE (means squared error)

33.912656604887246 -0.7462912894807127

and intercept ...

88.23380257847109

and as a formula:

$$y = -0.75x + 88.23$$

```
kForLACC$ python3 helpers.py  
R2 score: 0.8266009264508798
```

Evaluating the Model Visually

Our training data we set aside seems to “align” with the line. . .

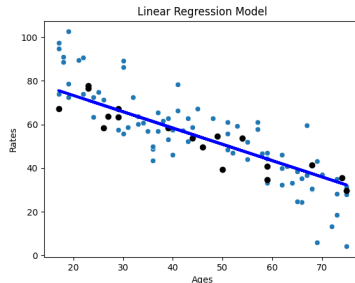


Figure: Training Data Scatter Plot With Line