

Simple Linear Regression

Rafiq Islam

2024-08-29

Table of contents

| | |
|---|----------|
| Simple Linear Regression | 1 |
| Assumptions of Linear Regressions | 4 |
| Synthetic Data | 4 |
| Model | 5 |

Simple Linear Regression

A simple linear regression in multiple predictors/input variables/features/independent variables/explanatory variables/regressors/ covariates (many names) often takes the form

$$y = f(\mathbf{x}) + \epsilon = \beta\mathbf{x} + \epsilon$$

where $\beta \in \mathbb{R}^d$ are regression parameters or constant values that we aim to estimate and $\epsilon \sim \mathcal{N}(0, 1)$ is a normally distributed error term independent of x or also called the white noise.

In this case, the model:

$$y = f(x) + \epsilon = \beta_0 + \beta_1 x + \epsilon$$

Therefore, in our model we need to estimate the parameters β_0, β_1 . The true relationship between the explanatory variables and the dependent variable is $y = f(x)$. But our model is $y = f(x) + \epsilon$. Here, this $f(x)$ is the working model with the data. In other words, $\hat{y} = f(x) = \hat{\beta}_0 + \hat{\beta}_1 x$. Therefore, there should be some error in the model prediction which we are calling $\epsilon = \|y - \hat{y}\|$ where y is the true value and \hat{y} is the predicted value. This error term is normally distributed with mean 0 and variance 1. To get the best estimate of the parameters

β_0, β_1 we can minimize the error term as much as possible. So, we define the residual sum of squares (RSS) as:

$$RSS = \epsilon_1^2 + \epsilon_2^2 + \cdots + \epsilon_{10}^2 \quad (1)$$

$$= \sum_{i=1}^{10} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (2)$$

$$\hat{\Downarrow}(\bar{\beta}) = \sum_{i=1}^{10} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (3)$$

$$(4)$$

Using multivariate calculus we see

$$\frac{\partial l}{\partial \beta_0} = \sum_{i=1}^{10} 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-1) \quad (5)$$

$$\frac{\partial l}{\partial \beta_1} = \sum_{i=1}^{10} 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-x_i) \quad (6)$$

Setting the partial derivatives to zero we solve for $\hat{\beta}_0, \hat{\beta}_1$ as follows

$$\begin{aligned} \frac{\partial l}{\partial \beta_0} &= 0 \\ \Rightarrow \sum_{i=1}^{10} y_i - 10\hat{\beta}_0 - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i \right) &= 0 \\ \Rightarrow \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

and,

$$\begin{aligned}
& \frac{\partial l}{\partial \beta_1} = 0 \\
& \Rightarrow \sum_{i=1}^{10} 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-x_i) = 0 \\
& \Rightarrow \sum_{i=1}^{10} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(x_i) = 0 \\
& \Rightarrow \sum_{i=1}^{10} x_i y_i - \hat{\beta}_0 \left(\sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i^2 \right) = 0 \\
& \Rightarrow \sum_{i=1}^{10} x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \left(\sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i^2 \right) = 0 \\
& \Rightarrow \sum_{i=1}^{10} x_i y_i - \bar{y} \left(\sum_{i=1}^{10} x_i \right) + \hat{\beta}_1 \bar{x} \left(\sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i^2 \right) = 0 \\
& \Rightarrow \sum_{i=1}^{10} x_i y_i - \bar{y} \left(\sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i^2 - \bar{x} \sum_{i=1}^{10} x_i \right) = 0 \\
& \Rightarrow \sum_{i=1}^{10} x_i y_i - \bar{y} \left(\sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i^2 - 10\bar{x}^2 \right) = 0 \\
& \Rightarrow \sum_{i=1}^{10} x_i y_i - \bar{y} \left(\sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left(\sum_{i=1}^{10} x_i^2 - 2 \times 10 \times \bar{x}^2 + 10\bar{x}^2 \right) = 0 \\
& \Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^{10} x_i y_i - 10\bar{x}\bar{y}}{\sum_{i=1}^{10} x_i^2 - 2 \times 10 \times \bar{x}^2 + 10\bar{x}^2} \\
& \Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^{10} x_i y_i - 10\bar{x}\bar{y} - 10\bar{x}\bar{y} + 10\bar{x}\bar{y}}{\sum_{i=1}^{10} x_i^2 - 2\bar{x} \times 10 \times \frac{1}{10} \sum_{i=1}^{10} x_i + 10\bar{x}^2} \\
& \Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^{10} x_i y_i - \bar{y} \left(\sum_{i=1}^{10} x_i \right) - \bar{x} \left(\sum_{i=1}^{10} y_i \right) + 10\bar{x}\bar{y}}{\sum_{i=1}^{10} (x_i - \bar{x})^2} \\
& \Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^{10} (x_i y_i - x_i \bar{y} - \bar{x} y_i + \bar{x} \bar{y})}{\sum_{i=1}^{10} (x_i - \bar{x})^2} \\
& \Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^{10} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{10} (x_i - \bar{x})^2}
\end{aligned}$$

Therefore, we have the following

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{10} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{10} (x_i - \bar{x})^2}$$

Simple Linear Regression `slr` is applicable for a single feature data set with continuous response variable.

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
```

Assumptions of Linear Regressions

- **Linearity:** The relationship between the feature set and the target variable has to be linear.
- **Homoscedasticity:** The variance of the residuals has to be constant.
- **Independence:** All the observations are independent of each other.
- **Normality:** The distribution of the dependent variable y has to be normal.

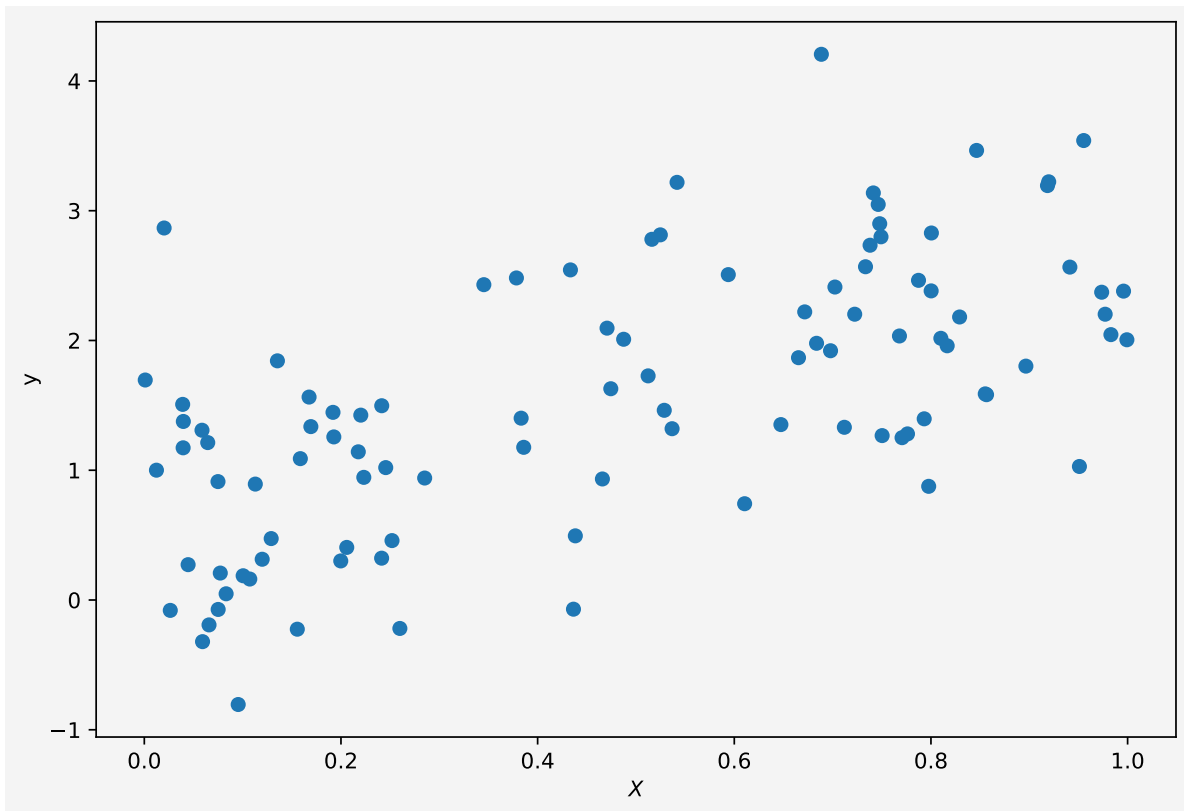
Synthetic Data

To implement the algorithm, we need some synthetic data. To generate the synthetic data we use the linear equation $y(x) = 2x + \frac{1}{2} + \xi$ where $\xi \sim \mathbf{N}(0, 1)$

```
X=np.random.random(100)
y=2*X+0.5+np.random.randn(100)
```

Note that we used two random number generators, `np.random.random(n)` and `np.random.randn(n)`. The first one generates n random numbers of values from the range (0,1) and the second one generates values from the standard normal distribution with mean 0 and variance or standard deviation 1.

```
plt.figure(figsize=(9,6))
plt.scatter(X,y)
plt.xlabel('$X$')
plt.ylabel('y')
plt.gca().set_facecolor('#f4f4f4')
plt.gcf().patch.set_facecolor('#f4f4f4')
plt.show()
```



Model

We want to fit a simple linear regression to the above data.

```
slr=LinearRegression()
```

Now to fit our data X and y we need to reshape the input variable. Because if we look at X ,

X

```
array([7.97697867e-01, 8.10080076e-01, 6.65307198e-01, 9.18514039e-01,
       7.87351209e-01, 3.94812583e-02, 3.78435774e-01, 9.50995789e-01,
       1.99784304e-01, 3.45307666e-01, 5.93899309e-01, 7.50353127e-01,
       2.01567820e-02, 1.67563408e-01, 3.85842386e-01, 5.36742775e-01,
       9.83020514e-01, 7.67966776e-01, 7.49197431e-01, 4.70594104e-01,
       6.97888259e-01, 1.92735153e-01, 7.72605138e-02, 1.58677108e-01,
       7.11934178e-01, 1.91889689e-01, 6.47410465e-01, 5.28775771e-01,
       8.96559759e-01, 3.83230846e-01, 7.38087323e-01, 2.85120631e-01,
       1.07278100e-01, 1.19884764e-01, 8.55222192e-01, 6.71627366e-01,
       7.22459373e-01, 8.46408228e-01, 8.56548902e-01, 7.47896797e-01,
       2.17612054e-01, 1.69380208e-01, 1.23482564e-02, 8.32043150e-02,
       4.36448838e-01, 8.29144362e-01, 4.38347216e-01, 5.16194930e-01,
       9.46824436e-04, 6.58194588e-02, 5.86764882e-02, 7.41408620e-01,
       1.55591996e-01, 2.41350225e-01, 9.54875755e-02, 6.45060259e-02,
       1.29037864e-01, 5.12358645e-01, 5.41790977e-01, 2.20162048e-01,
       2.05875632e-01, 9.95857424e-01, 1.00517262e-01, 7.70616864e-01,
       2.59902086e-01, 9.73625457e-01, 2.23223127e-01, 3.96829056e-02,
       6.88548080e-01, 2.64113344e-02, 8.00448321e-01, 6.83661804e-01,
       7.51150415e-02, 7.33441194e-01, 6.10515112e-01, 8.16505619e-01,
       5.24788329e-01, 3.90160357e-02, 2.41456035e-01, 8.00197797e-01,
       9.41289449e-01, 7.49243932e-02, 2.45522726e-01, 9.19764460e-01,
       9.77178764e-01, 2.51909092e-01, 7.02348747e-01, 4.65876622e-01,
       1.12884544e-01, 4.87477477e-01, 1.35366123e-01, 4.45481778e-02,
       4.74419136e-01, 9.99275073e-01, 7.46366902e-01, 9.55344104e-01,
       7.93133419e-01, 7.76139969e-01, 4.33358840e-01, 5.92181209e-02])
```

It is a one-dimensional array/vector but the `slr` object accepts input variable as matrix or two-dimensional format.

```
X=X.reshape(-1,1)
X[:10]
```

```
array([[0.79769787],
       [0.81008008],
       [0.6653072 ],
       [0.91851404],
       [0.78735121],
       [0.03948126],
       [0.37843577],
```

```
[0.95099579],  
[0.1997843 ],  
[0.34530767]])
```

Now we fit the data to our model

```
slr.fit(X,y)  
slr.predict([[2],[3]])
```

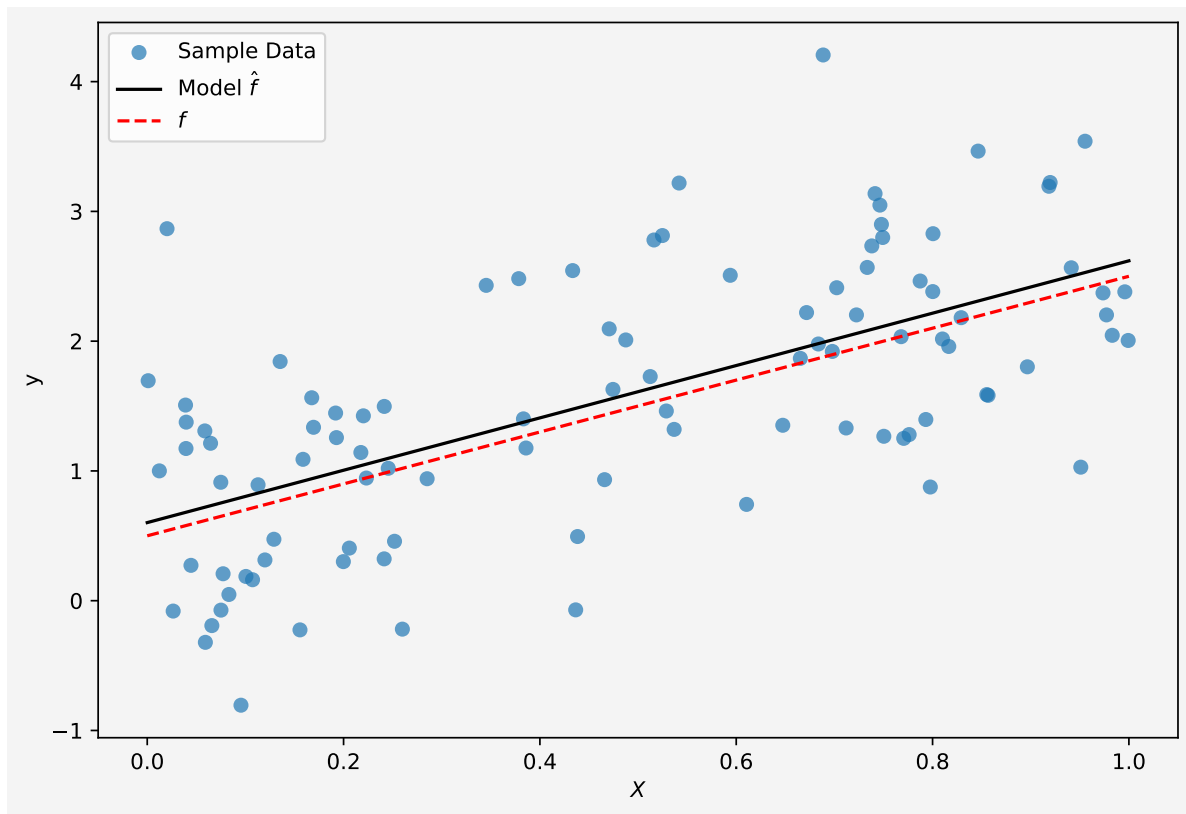
```
array([4.63646575, 6.65380348])
```

We have our $X = 2, 3$ and the corresponding y values are from the above cell output, which are pretty close to the model $y = 2x + \frac{1}{2}$.

```
intercept = round(slr.intercept_,4)  
slope = slr.coef_
```

Now our model parameters are: intercept $\beta_0 = \text{np.float64}(0.6018)$ and slope $\beta_1 = \text{array}([2.01733774])$.

```
plt.figure(figsize=(9,6))  
plt.scatter(X,y, alpha=0.7,label="Sample Data")  
plt.plot(np.linspace(0,1,100),  
         slr.predict(np.linspace(0,1,100).reshape(-1,1)),  
         'k',  
         label='Model  $\hat{f}$ ')  
)  
plt.plot(np.linspace(0,1,100),  
         2*np.linspace(0,1,100)+0.5,  
         'r--',  
         label=' $f$ ')  
)  
plt.xlabel('$X$')  
plt.ylabel('$y$')  
plt.legend(fontsize=10)  
plt.gca().set_facecolor('#f4f4f4')  
plt.gcf().patch.set_facecolor('#f4f4f4')  
plt.show()
```



So the model fits the data almost perfectly.

Up next [multiple linear regression](#).

Share on



You may also like