# Simple Linear Regression

Rafiq Islam

2024-08-29

## Table of contents

## 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 $\beta_0, \beta_1$. 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

$\beta_0, \beta_1$ we can minimize the error term as much as possible. So, we define the residual sum of squares (RSS) as:

$$RSS = \epsilon_1^2 + \epsilon_2^2 + \cdots + \epsilon_{10}^2 \tag{1}$$

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

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

$$\tag{4}$$

Using multivariate calculus we see

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

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

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

$$\frac{\partial l}{\partial \beta_0} = 0$$

$$\implies \sum_{i=1}^{10} y_i - 10\hat{\beta}_0 - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i \right) = 0$$

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

and,

$$\frac{\partial l}{\partial \beta_1} = 0$$

$$\implies \sum_{i=1}^{10} 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-x_i) = 0$$

$$\implies \sum_{i=1}^{10} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(x_i) = 0$$

$$\implies \sum_{i=1}^{10} x_i y_i - \hat{\beta}_0 \left( \sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i^2 \right) = 0$$

$$\implies \sum_{i=1}^{10} x_i y_i - \left( \bar{y} - \hat{\beta}_1 \bar{x} \right) \left( \sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i^2 \right) = 0$$

$$\implies \sum_{i=1}^{10} x_i y_i - \bar{y} \left( \sum_{i=1}^{10} x_i \right) + \hat{\beta}_1 \bar{x} \left( \sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i^2 \right) = 0$$

$$\implies \sum_{i=1}^{10} x_i y_i - \bar{y} \left( \sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i^2 - \bar{x} \sum_{i=1}^{10} x_i \right) = 0$$

$$\implies \sum_{i=1}^{10} x_i y_i - \bar{y} \left( \sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i^2 - 10\bar{x}^2 \right) = 0$$

$$\implies \sum_{i=1}^{10} x_i y_i - \bar{y} \left( \sum_{i=1}^{10} x_i \right) - \hat{\beta}_1 \left( \sum_{i=1}^{10} x_i^2 - 2 \times 10 \times \bar{x}^2 + 10\bar{x}^2 \right) = 0$$

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

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

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

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

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

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 contineous 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.
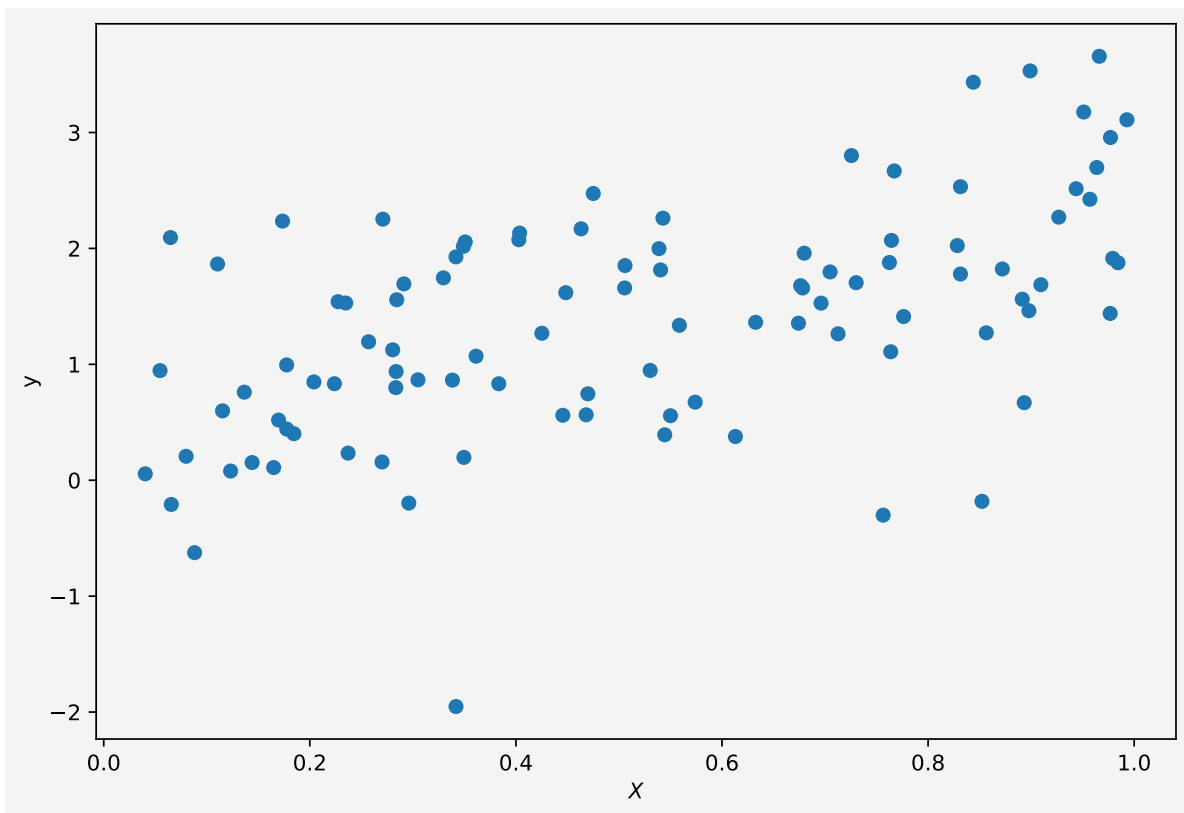
## Synthetic Data

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([0.22385393, 0.7048668 , 0.28435213, 0.47505871, 0.67977931,
       0.05464596, 0.44840441, 0.83143578, 0.18458166, 0.53887021,
       0.95705949, 0.87205369, 0.83141033, 0.909443  , 0.89902134,
       0.46327038, 0.5444457 , 0.55000912, 0.22742005, 0.73018549,
       0.6130214 , 0.75643099, 0.42518078, 0.67633813, 0.97700959,
       0.76716119, 0.1776102 , 0.54270048, 0.46815185, 0.92681044,
       0.85647156, 0.23714362, 0.9436864 , 0.28373978, 0.0881827 ,
       0.89140139, 0.63263335, 0.0799422 , 0.72559597, 0.50587889,
       0.67420761, 0.67821368, 0.16966618, 0.55863319, 0.11041683,
       0.34904709, 0.30494733, 0.27009376, 0.34186032, 0.06555454,
       0.89785966, 0.76446226, 0.17750662, 0.57395209, 0.270881  ,
       0.29607447, 0.40284376, 0.97674039, 0.32961657, 0.76246074,
       0.06475832, 0.54039416, 0.34181623, 0.77623075, 0.96613664,
       0.11530654, 0.17353008, 0.35080613, 0.98427141, 0.46979096,
       0.82834317, 0.13643374, 0.33841693, 0.28043549, 0.38336485,
       0.44553404, 0.69611976, 0.12313408, 0.36132081, 0.28344435,
       0.20396632, 0.97919804, 0.40362642, 0.8933104 , 0.29125611,
       0.16492845, 0.53029186, 0.96371704, 0.0402838 , 0.23485863,
       0.99294302, 0.2570391 , 0.9510301 , 0.85230072, 0.76377794,
       0.84386968, 0.71263068, 0.34951971, 0.50554734, 0.14399891])
```

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([[0.22385393],
       [0.7048668 ],
       [0.28435213],
       [0.47505871],
       [0.67977931],
       [0.05464596],
       [0.44840441],
       [0.83143578],
       [0.18458166],
       [0.53887021]])
```

Now we fit the data to our model
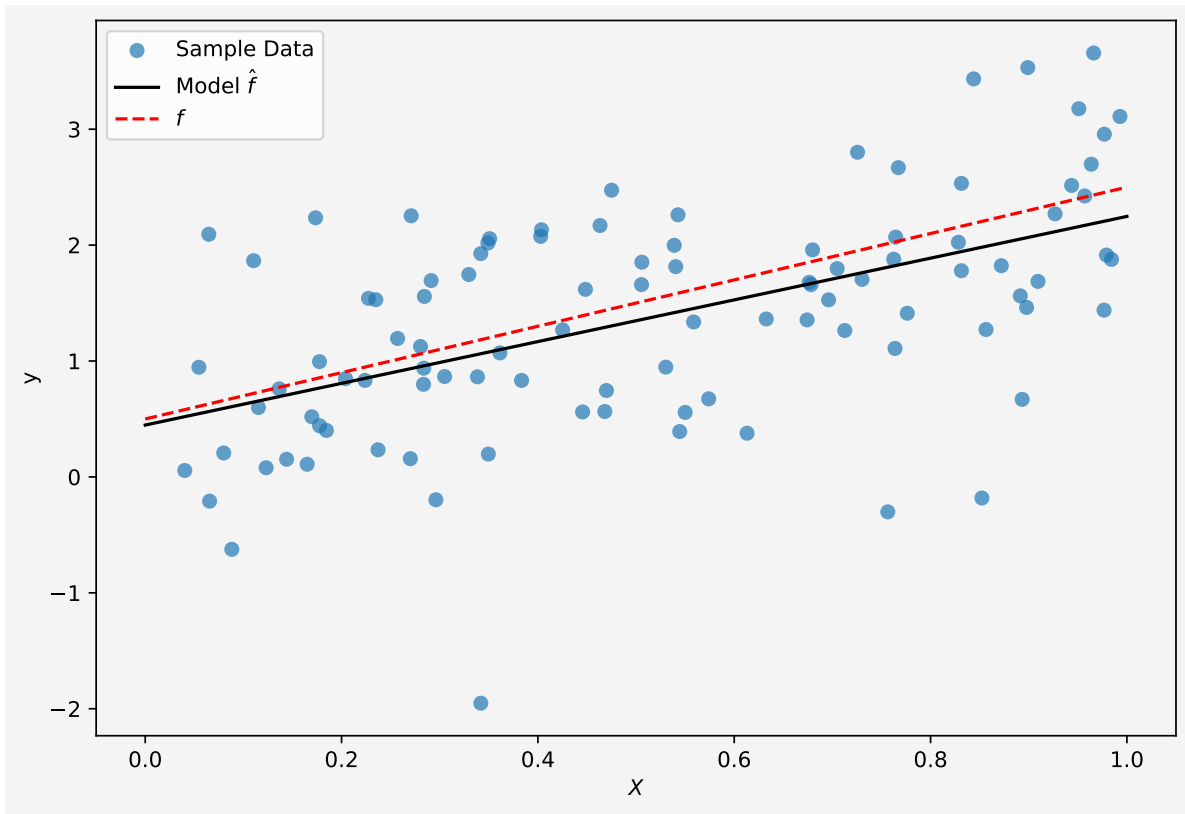
```
slr.fit(X,y)
slr.predict([[2],[3]])
```

```
array([4.04774739, 5.84789178])
```

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 $\beta_0 = 0.4475$ and slope $\beta_1 = \text{array}([1.80014438])$.

```
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**