# 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.

```python
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.

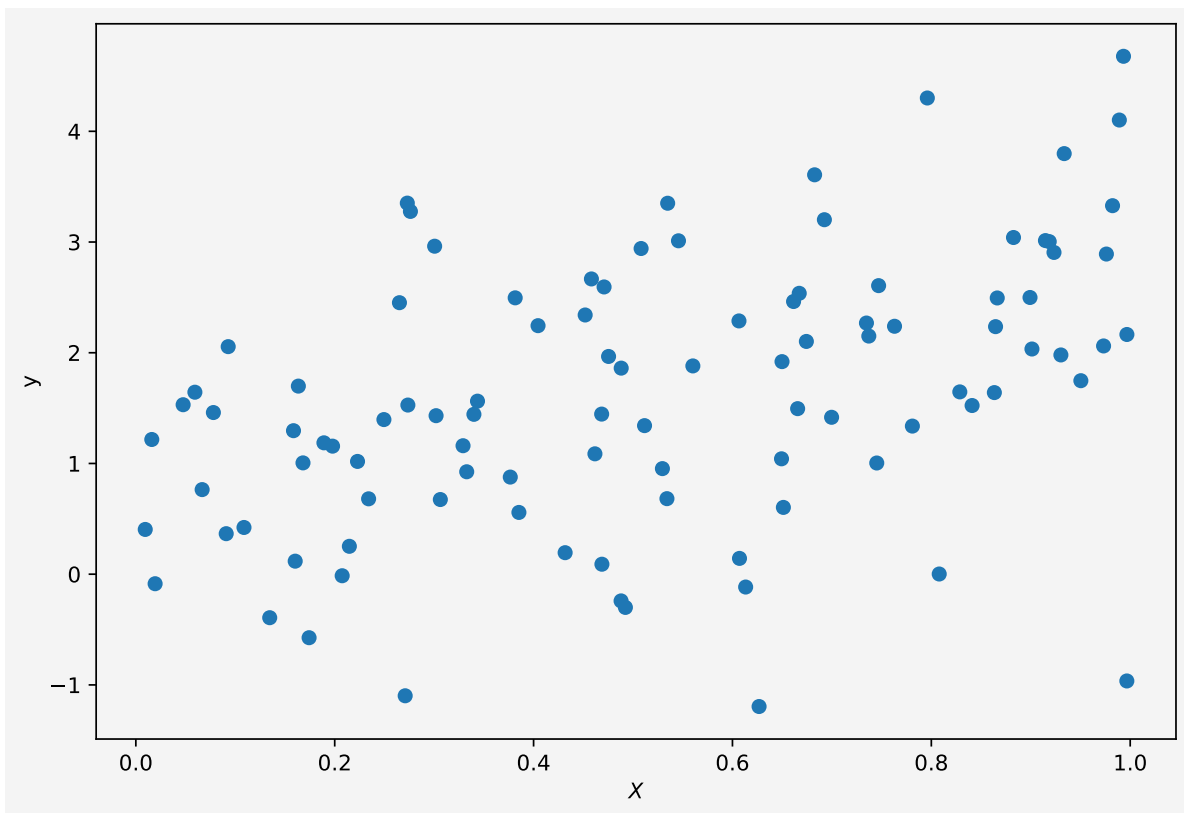- **Normality:** The distribution of the dependent variable $y$ has to be normal.

### 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 $\xi \sim \mathbf{N}(0, 1)$

```python
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([0.98907093, 0.45187652, 0.66554759, 0.99661672, 0.52955205,
       0.90119618, 0.22294027, 0.65120945, 0.92335664, 0.06668294,
       0.20750548, 0.54572072, 0.88263001, 0.99325433, 0.73712474,
       0.10875866, 0.38150137, 0.99667558, 0.78094622, 0.91476824,
       0.5081332 , 0.68261442, 0.00945387, 0.17422965, 0.07798732,
       0.97321964, 0.01938885, 0.64934061, 0.6498437 , 0.26510027,
       0.69245038, 0.97605869, 0.0593869 , 0.18915494, 0.33274342,
       0.09088217, 0.74511371, 0.98223257, 0.43172891, 0.46845941,
       0.30194763, 0.19783473, 0.32905317, 0.34359824, 0.27301398,
       0.24948218, 0.61326598, 0.76298609, 0.27357904, 0.79597725,
       0.13463345, 0.8410004 , 0.93362478, 0.80793408, 0.47540998,
       0.66145251, 0.60654765, 0.23409486, 0.93029632, 0.53411553,
       0.86338524, 0.62677276, 0.67432942, 0.09288274, 0.51158008,
       0.69969889, 0.48798909, 0.89924337, 0.45816974, 0.56018206,
       0.15854677, 0.40448149, 0.27613497, 0.16026923, 0.38524776,
       0.49234027, 0.047561  , 0.95038148, 0.30049301, 0.47088721,
       0.37655915, 0.86458979, 0.16343701, 0.46872092, 0.73475334,
       0.91856918, 0.01604243, 0.16809838, 0.66715195, 0.46169156,
       0.82864409, 0.27086898, 0.33998957, 0.2147328 , 0.60707983,
       0.86629109, 0.48811632, 0.53484162, 0.30622177, 0.74698352])
```

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([[0.98907093],
       [0.45187652],
       [0.66554759],
       [0.99661672],
       [0.52955205],
       [0.90119618],
       [0.22294027],
       [0.65120945],
       [0.92335664],
       [0.06668294]])
```

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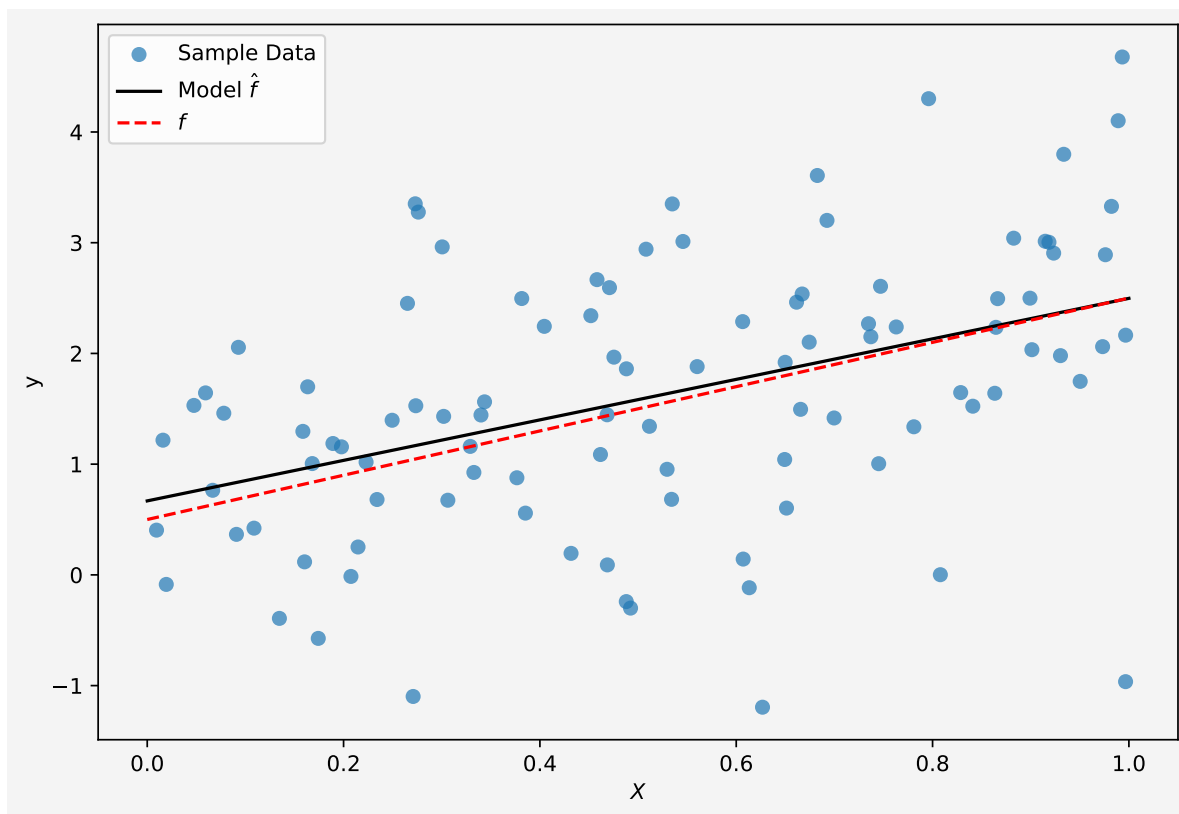
```
slr.fit(X,y)
slr.predict([[2],[3]])
```

```
array([4.32672342, 6.15599091])
```

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 $\beta_0 =$ np.float64(0.6682) and slope $\beta_1 =$ array([1.82926749]).

```
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**