# 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,

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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.

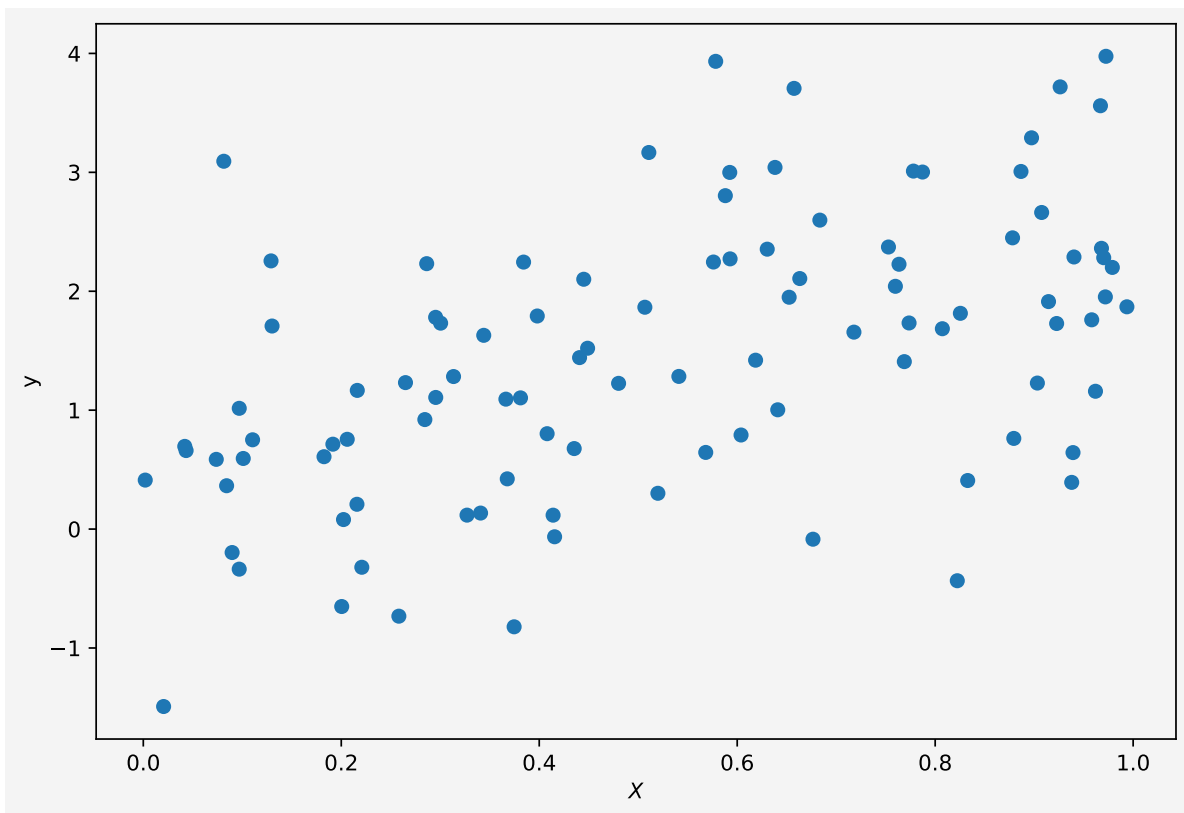- **Normality:** The distribution of the dependent variable $y$ has to be normal.

## 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([0.34073022, 0.82545378, 0.28629932, 0.90324716, 0.1915238 ,
       0.99368257, 0.51072085, 0.57826318, 0.20602109, 0.71791889,
       0.26487518, 0.10097791, 0.93792785, 0.8794269 , 0.66320887,
       0.22075434, 0.38107681, 0.90756416, 0.20226318, 0.43531347,
       0.07372428, 0.65744166, 0.34393001, 0.87816616, 0.6381873 ,
       0.31341946, 0.12904323, 0.08973382, 0.09692669, 0.48016135,
       0.88668101, 0.94021852, 0.97192483, 0.36769188, 0.60388062,
       0.28444115, 0.02051743, 0.51986921, 0.18256894, 0.58800576,
       0.09689651, 0.30037675, 0.9789682 , 0.77802365, 0.3663152 ,
       0.41400016, 0.11041266, 0.13007519, 0.57599606, 0.80725472,
       0.93923743, 0.75980163, 0.64095221, 0.7688229 , 0.89734706,
       0.44497093, 0.5410491 , 0.59247356, 0.39801035, 0.95814563,
       0.25817888, 0.77362097, 0.97254363, 0.08133239, 0.96798539,
       0.21584876, 0.78720903, 0.40796142, 0.44074641, 0.67654785,
       0.82233594, 0.20044603, 0.32703435, 0.29531859, 0.41552352,
       0.75280922, 0.08417087, 0.50670612, 0.63032882, 0.76345822,
       0.65247277, 0.83283659, 0.38418209, 0.92270509, 0.59280774,
       0.61857119, 0.92624467, 0.37462584, 0.96193966, 0.97026701,
       0.21617806, 0.96697119, 0.04324291, 0.29544354, 0.04189044,
       0.44886396, 0.68349321, 0.91450247, 0.56819968, 0.00196069])
```

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([[0.34073022],
       [0.82545378],
       [0.28629932],
       [0.90324716],
       [0.1915238 ],
       [0.99368257],
       [0.51072085],
       [0.57826318],
       [0.20602109],
       [0.71791889]])
```

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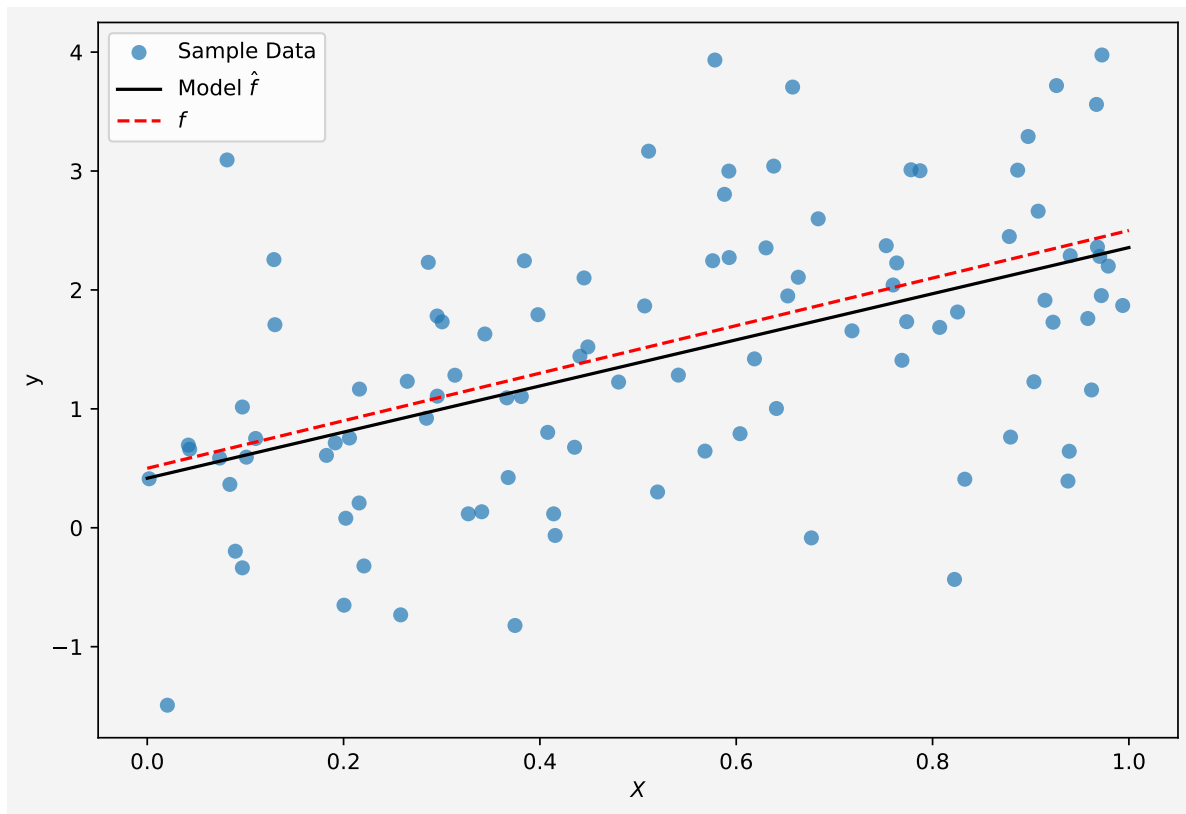
```
slr.fit(X,y)
slr.predict([[2],[3]])
```

```
array([4.29611973, 6.23592738])
```

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 $\beta_0 = 0.4165$ and slope $\beta_1 = \text{array}([1.93980765])$.

```
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**