

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.

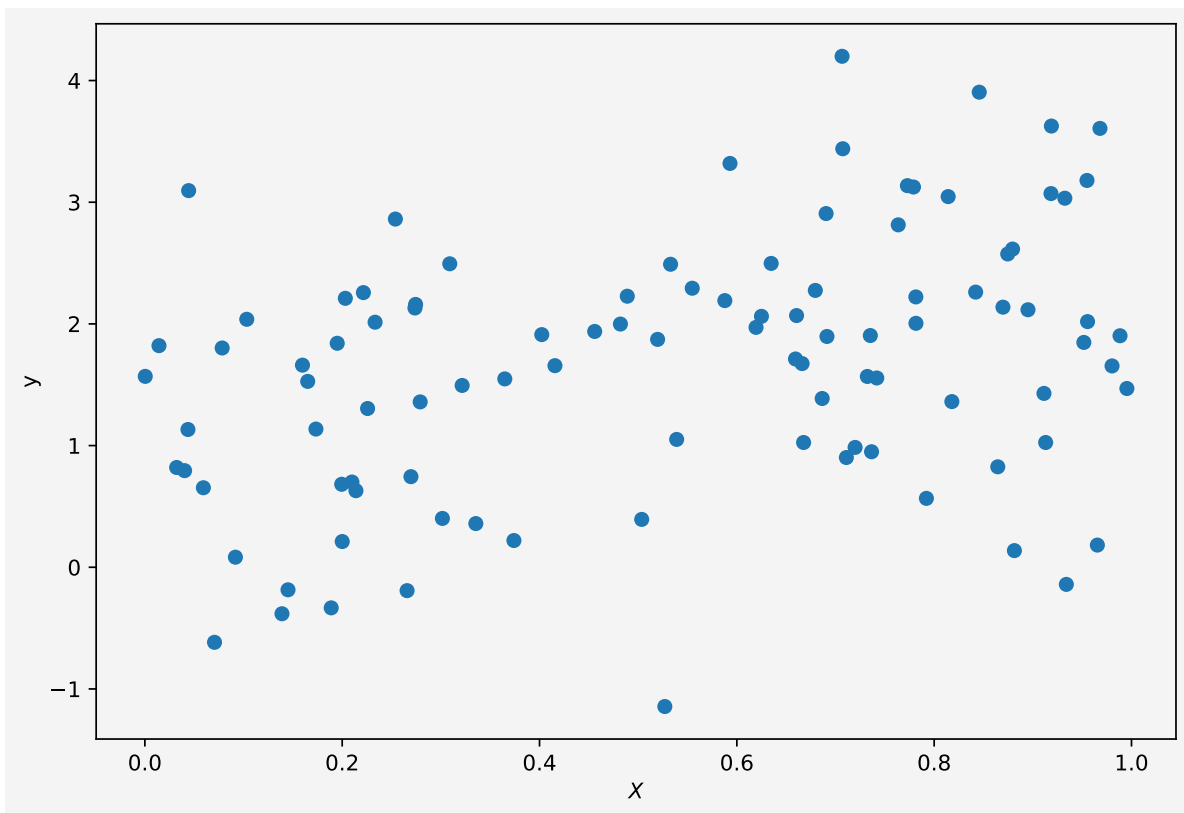
Synthetic Data

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([5.27003731e-01, 7.36523722e-01, 8.64474770e-01, 3.64802584e-01,
       1.45070094e-01, 9.51708911e-01, 2.53905382e-01, 6.24927801e-01,
       1.59723243e-01, 5.03580300e-01, 7.79015901e-01, 9.68031258e-01,
       6.66128255e-01, 3.35435355e-01, 2.03169892e-01, 9.11276386e-01,
       4.02568084e-02, 4.43644055e-02, 5.92947285e-02, 5.93057526e-01,
       7.19849019e-01, 1.99501201e-01, 2.09892325e-01, 6.67642162e-01,
       8.14174673e-01, 1.38958932e-01, 9.95336129e-01, 1.88862021e-01,
       2.33331464e-01, 1.42986458e-02, 3.01564731e-01, 4.55915477e-01,
       7.72856637e-01, 2.79029569e-01, 4.15642755e-01, 7.81434601e-01,
       2.65773134e-01, 9.18402833e-01, 6.86495731e-01, 9.18811513e-01,
       8.74544118e-01, 3.74047645e-01, 2.74402296e-01, 2.21480080e-01,
       7.63575592e-01, 7.11005155e-01, 6.79530168e-01, 3.22071858e-02,
       7.41827001e-01, 9.80277420e-01, 7.05884630e-02, 9.55443654e-01,
       3.21576557e-01, 5.54775778e-01, 4.88940639e-01, 8.17930829e-01,
       5.32822124e-01, 7.32273790e-01, 2.00036602e-01, 2.69678498e-01,
       9.32434470e-01, 9.33952260e-01, 7.92115148e-01, 9.17693866e-02,
       6.91328448e-01, 8.42044973e-01, 1.73377318e-01, 9.54882403e-01,
       7.81532149e-01, 2.13925789e-01, 8.79452186e-01, 1.03274532e-01,
       7.06645940e-01, 8.45693051e-01, 4.36615838e-02, 9.13002625e-01,
       6.34789369e-01, 6.19442551e-01, 2.73732843e-01, 3.09024018e-01,
       9.88373125e-01, 6.60531765e-01, 5.19674441e-01, 8.69703288e-01,
       4.81992399e-01, 3.71086675e-04, 6.59451250e-01, 9.65407586e-01,
       4.02277967e-01, 7.83492578e-02, 7.07351040e-01, 8.81315343e-01,
       1.94948481e-01, 8.95073943e-01, 1.65038183e-01, 2.25729740e-01,
       5.39003054e-01, 6.90452839e-01, 7.35295273e-01, 5.87818935e-01])
```

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([[0.52700373],
       [0.73652372],
       [0.86447477],
       [0.36480258],
       [0.14507009],
       [0.95170891],
       [0.25390538],
```

```
[0.6249278 ],
[0.15972324],
[0.5035803 ]])
```

Now we fit the data to our model

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

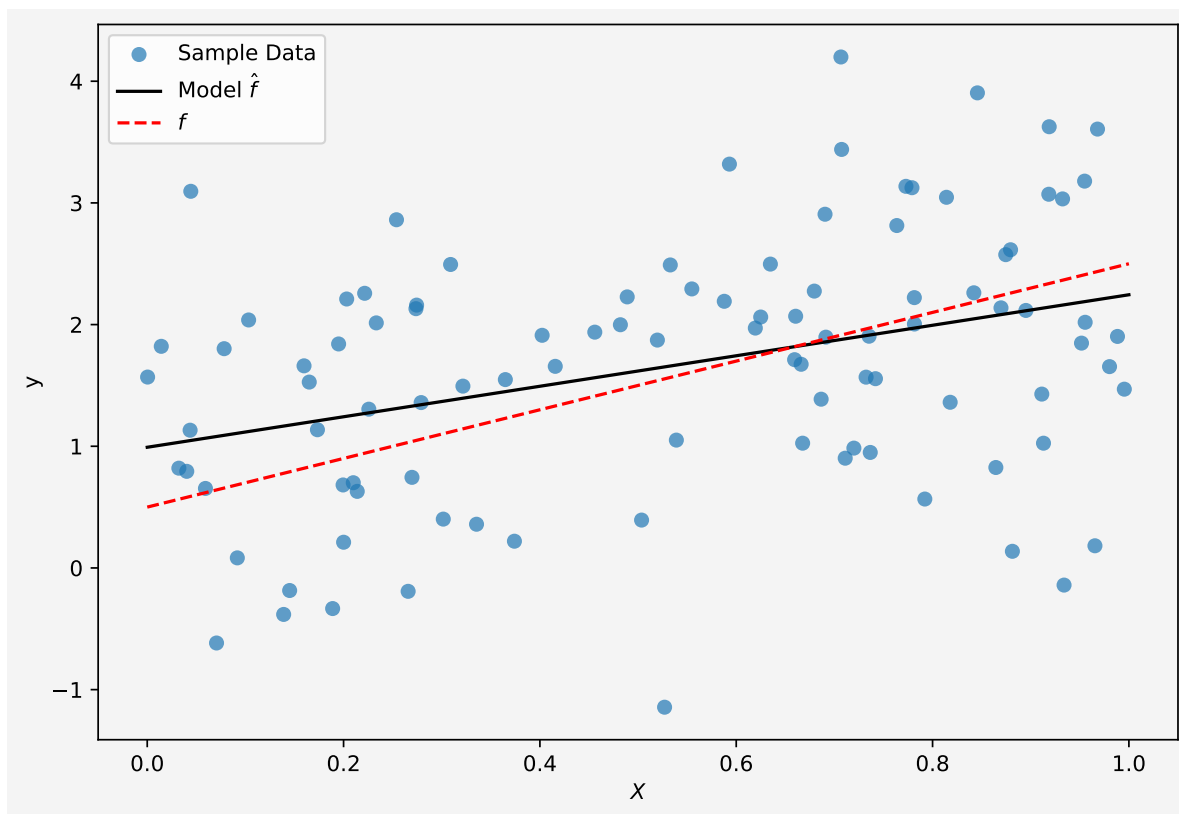
```
array([3.49844208, 4.75176355])
```

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 $\beta_0 = 0.9918$ and slope $\beta_1 = \text{array}([1.25332147])$.

```
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