

Simple Linear Regression

Lecture 14

Centre for Data Science, ITER
Siksha 'O' Anusandhan (Deemed to be University), Bhubaneswar, Odisha, India.



Contents

- 1 Introduction
- 2 Simple Linear Regression model
- 3 Gradient Descent
- 4 Maximum Likelihood Estimation

Correlation:

- Correlation function to measure the strength and direction of the linear linear relationship between two variables.
- The value of the correlation between two variable lies between -1 to 1.
- If the two variables move in the same direction, then those variables are said to have a positive correlation. If they move in opposite directions, then they have a negative correlation.

What is the need of linear regression ?

Simple Linear Regression model

- The model for linear regression with one predictor variable can be stated as follows:

$$y_i = \beta x_i + \alpha + \epsilon_i, \quad \text{where}$$

- y_i is the value of the response variable in the i th trial.
- α and β are parameters.
- x_i is known, it is the value of the predictor variable in the i th trial.
- ϵ_i is a random error term, normally distributed with mean, $E(\epsilon_i) = 0$ and variance, $V(\epsilon_i) = \sigma^2$
- $\text{Cov}(\epsilon_i, \epsilon_j) = 0$ if $i \neq j$ for $i=1, \dots, n$ and $j=1, \dots, n$.

The above regression model is said to be simple linear regression because it is linear in terms of parameters.

Simple Linear Regression (Contd.)

Examples of Simple Linear Regression model:



Figure 1: first figure

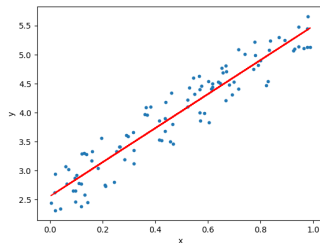


Figure 2: second figure

Simple Linear Regression (Contd.)

Program to plot regression line and scatter plot :

```
1 import numpy as np
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 x=[1,2,3,4,9]
6 y=[5,2,9,7,3]
7
8 plt.figure()
9 sns.regplot(x,y,fit_reg = True)
10 plt.scatter(np.mean(x),np.mean(y),color='green')
```

Simple Linear Regression (Contd.)

Least square Estimation of α and β :

$$\hat{\beta} = \frac{\text{Correlation}(x, y) * \sigma_y}{\sigma_x}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}, \text{ where}$$

\bar{y} = mean of y

\bar{x} = mean of x

σ_y = standard deviation of y

σ_x = standard deviation of x

Simple Linear Regression (Cont.)

Program to find the least square estimates of α and β :

1.By Formula:

```
1 import numpy as np
2 x = np.array([5, 15, 25, 35, 45, 55])
3 y = np.array([5, 20, 14, 32, 22, 38])
4
5 from typing import Tuple
6 from scratch.linear_algebra import Vector
7 from scratch.statistics import correlation, standard_deviation,
   mean
8 """Given two vectors x and y,
9 find the least-squares values of alpha and beta"""
10
11 def least_squares_fit(x: Vector, y: Vector) -> Tuple[float,
   float]:
12     beta = correlation(x, y) * standard_deviation(y) /
   standard_deviation(x)
13     alpha = mean(y) - beta * mean(x)
14     return alpha, beta
```


Simple Linear Regression (Cont.)

Program to predict for any value of x and to find the error :

```
1 def predict(alpha: float, beta: float, x_i: float) -> float:
2     return beta * x_i + alpha
3     """ The error from predicting beta * x_i + alpha
4     when the actual value is y_i """
5 def error(alpha: float, beta: float, x_i: float, y_i: float) ->
6     float:
7     return predict(alpha, beta, x_i) - y_i
8
9 from scratch.linear_algebra import Vector
10 def sum_of_sqerrors(alpha: float, beta: float, x: Vector, y:
11     Vector) -> float:
12     return sum(error(alpha, beta, x_i, y_i) ** 2
13         for x_i, y_i in zip(x, y))
14
15 m1=least_squares_fit(x,y)
16 predict(m1[0], m1[1], 20)
17 E=error(m1[0], m1[1],x,y)
```

Simple Linear Regression (Cont.)

Program to find the least square estimates of α and β :

2.By Inbuilt module of python:

```
1 import statsmodels.api as s
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import statsmodels.formula.api as sm
6
7 x = np.array([5, 15, 25, 35, 45, 55]).reshape((-1, 1))
8 y = np.array([5, 20, 14, 32, 22, 38])
9
10 x= s.add_constant(x)
11
12 model1=s.OLS(y, x)
13 result1=model1.fit()
14 print(result1.summary())
```

Simple Linear Regression (Cont.)

Test the least square module:

```
1 x = [i for i in range(-100, 110, 10)]
2 y = [3 * i - 5 for i in x]
3 # Should find that  $y = 3x - 5$ 
4 p = least_squares_fit(x, y)
5 assert least_squares_fit(x, y) == (-5, 3)
```

Use of assert statement in least square method:

```
1 from scratch.statistics import num_friends_good,
   daily_minutes_good
2 alpha, beta = least_squares_fit(num_friends_good,
   daily_minutes_good)
3 assert 22.9 < alpha < 23.0
4 assert 0.9 < beta < 0.905
```

Simple Linear Regression (Cont.)

Calculation of R-square:

```
1 from scratch.statistics import de_mean
2 def total_sum_of_squares(y: Vector) -> float:
3     """the total squared variation of y_i's from their mean"""
4     return sum(v ** 2 for v in de_mean(y))
5 def r_squared(alpha: float, beta: float, x: Vector, y: Vector)
6     -> float:
7     return 1.0 - (sum_of_sqerrors(alpha, beta, x, y) /
8                   total_sum_of_squares(y))
9 """the fraction of variation in y captured by the model, which
10    equals
11    #1 - the fraction of variation in y not captured by the model
12    """
13
14 rsq = r_squared(alpha, beta, num_friends_good,
15                 daily_minutes_good)
16 assert 0.328 < rsq < 0.330
```

R-square:

- R-square value lies between 0 to 1
- The model whose R-squared value close to 1 indicates that the model is a maximum variance can be explained using the model. If the R-square value close to 0 indicates that is the variance of response can't be explained using x .

Simple Linear Regression (Cont.)

R-square comparison between two plots:

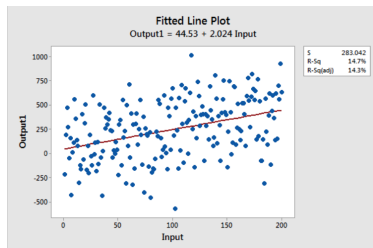


Figure 3: R^2 close to 0

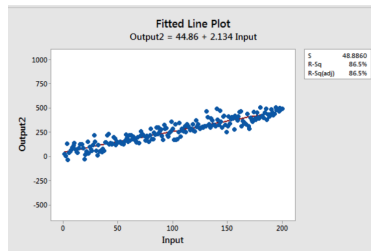


Figure 4: R^2 close to 1

Gradient Descent :

Program to estimate the parameters using gradient and Descent:

```
1 import random
2 import tqdm
3 from scratch.gradient_descent import gradient_step
4 num_epochs = 10000
5 random.seed(0)
6 guess = [random.random(), random.random()] # choose random
       value to start
7 learning_rate = 0.00001
8 with tqdm.trange(num_epochs) as t:
9     for _ in t:
10         alpha, beta = guess
11         # Partial derivative of loss with respect to alpha
12         grad_a = sum(2 * error(alpha, beta, x_i, y_i)
13                      for x_i, y_i in zip(num_friends_good,
14                      daily_minutes_good))
```

Gradient Descent (Cont...) :

Program to estimate the parameters using gradient and Descent:

```
1 # Partial derivative of loss with respect to beta
2     grad_b = sum(2 * error(alpha, beta, x_i, y_i) * x_i
3                 for x_i, y_i in zip(num_friends_good,
4                                     daily_minutes_good))
5 # Compute loss to stick in the tqdm description
6     loss = sum_of_sqerrors(alpha, beta,
7                             num_friends_good, daily_minutes_good
8     )
9     t.set_description(f"loss: {loss:.3f}")
10 # Finally, update the guess
11     guess = gradient_step(guess, [grad_a, grad_b], -
12                             learning_rate)
13 # We should get pretty much the same results:
14 alpha, beta
15 assert 22.9 < alpha < 23.0
16 assert 0.9 < beta < 0.905
```


Maximum Likelihood Estimation :

- Imagine that we have a sample of data v_1, v_2, \dots, v_n that comes from a distribution that depends on some unknown parameter θ .
- If θ is unknown then we can consider the likelihood of θ given the sample that is

$$L(\theta | v_1, \dots, v_n)$$

Under this approach, the most likely is the value that maximizes this likelihood function.

- As per the assumption of the Linear regression model, the errors are normally distributed with mean 0 and standard deviation σ . So the likelihood of α and β based on x and y data set is :

$$L(\alpha, \beta | x_i, y_i, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(- (y_i - \alpha - \beta x_i)^2 / 2\sigma^2\right)$$

Thank You
Any Questions?