

# Lecture 4: Regression

## ENEE 691 – Machine Learning and Photonics @ UMBC

Ergun Simsek, Ph.D. and Masoud Soroush, Ph.D.

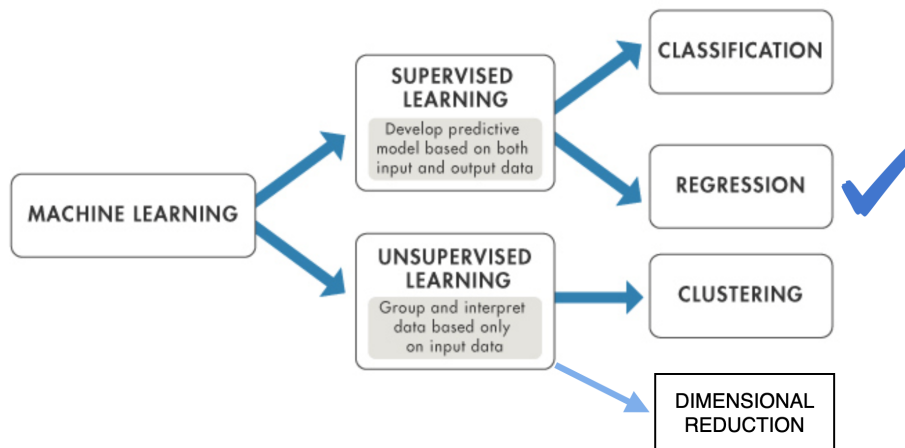
February 14, 2023

### 1 Regression Analysis

In today's lecture<sup>1</sup>, we introduce one of the simplest and most common machine learning algorithms, namely the regression. Despite its simplicity, it is still widely used in statistical learning.

#### 1.1 Background

In the context of *supervised learning*, regression is used to predict a quantitative continuous target variable. Moreover, in this lecture we will assume that all features (*i.e.* explanatory variables) are all continuous too.



**Note:** Be aware that regression methods have been generalized to cases where features include categorical variables as well. However, analyzing those cases require a deeper knowledge of statistics, and are beyond the scope of our current course.

<sup>1</sup> This product is a property of UMBC, and no distribution is allowed. These notes are solely for your own use, as a registered student in this class.

## 1.2 Formulating the Regression Problem

Let us first start with a simple case. Assume that we have a continuous target variable  $y$  explained by one single continuous feature  $x$ . We have recorded the values of  $x$  and  $y$  for  $n$  observations in a dataset. The dataset appears to have the following format:

Observation	$x$	$y$
1	$x_1$	$y_1$
2	$x_2$	$y_2$
3	$x_3$	$y_3$
$\vdots$	$\vdots$	$\vdots$
$n$	$x_n$	$y_n$

We assume a linear relation between the feature and the target variables

$$\hat{y} = \omega_1 x + \omega_0, \quad (1)$$

where  $\hat{y}$  denotes the predicted value for the target (Note that  $\hat{y}$  is the predicted value of the target whereas  $y$  is the true value of the target). In equation (1),  $\omega_0$  and  $\omega_1$  are the intercept and the slope of the line, respectively. Note that in the literature, the constant term  $\omega_0$  is sometimes called the *bias term*. The bias term  $\omega_0$  has nothing to do with the concept of bias (appearing in Bias-Variance trade off for instance) in machine learning.

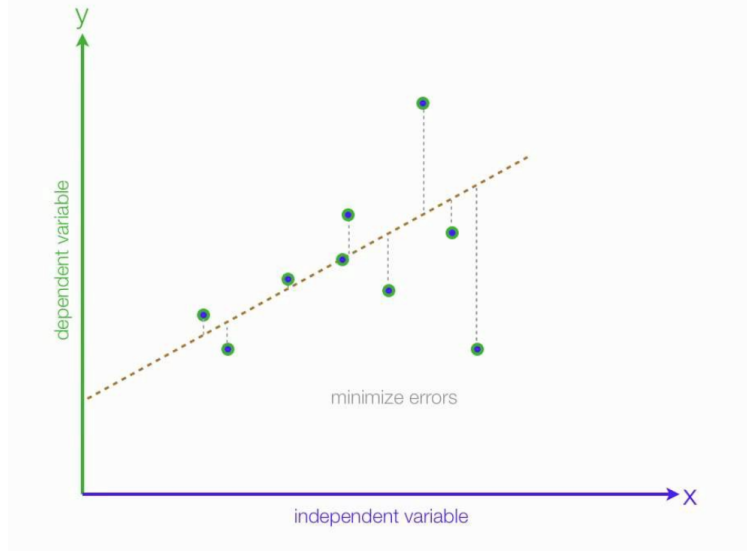
**Goal:** To determine the *best fit line*, we have to determine  $\omega_1$  (i.e. the slope of the line) and  $\omega_0$  (i.e. the intercept of the line) in equation (1). But before that, we need a criterion by which we can compare different lines.

To achieve the above goal, we define an error term  $e_i = (y_i - \hat{y}_i)$  for each observation. We then define the *cost function*  $J(\omega_0, \omega_1)$  associated with the regression algorithm to be

$$J(\omega_0, \omega_1) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \omega_1 x_i - \omega_0)^2. \quad (2)$$

The best fit line is the one whose slope and intercept (i.e.  $\omega_1, \omega_0$ ) minimizes the cost function in equation (2). This implies that, to minimize the cost function, one has to solve

$$\frac{\partial J}{\partial \omega_0} = 0 \quad , \quad \frac{\partial J}{\partial \omega_1} = 0. \quad (3)$$



Before we solve equation (3), we have to think about the generalization of the above setup. In reality, we rarely encounter a situation where the feature is a single variable  $x$ . The features of a given model typically consists of several variables. We can employ the matrix notation to easily generalize the above picture to the case where we have a set of explanatory variables. Now, suppose we collect the following dataset

Observation	$\vec{x} \in \mathbb{R}^d$	$y$
1	$\vec{x}^{(1)}$	$y_1$
2	$\vec{x}^{(2)}$	$y_2$
3	$\vec{x}^{(3)}$	$y_3$
$\vdots$	$\vdots$	$\vdots$
$n$	$\vec{x}^{(n)}$	$y_n$

We now construct an  $n \times (d + 1)$  matrix  $\mathbb{X}$  associated with the above dataset as follows:

$$\mathbb{X} = (\mathbb{1}, x_1, x_2, \dots, x_d) = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{pmatrix}. \quad (4)$$

There are a couple of comments in order:

- Note that the subscript  $i$  in  $x_i^{(j)}$  denotes the  $i$ -th feature among  $d$  features, whereas the superscript  $(j)$  in  $x_i^{(j)}$  denotes the  $j$ -th observation for the  $i$ -th feature.
- A column of 1's has been inserted in matrix  $\mathbb{X}$  (the first column in (4)) to capture the intercepts (*i.e.* bias term  $\omega_0$  in (1)) in a convenient way in matrix notation.

We can now represent the target variable by the  $n \times 1$  matrix  $\mathbb{Y}$ . The linear regression model assumes a linear relation between the target  $\mathbb{Y}$  and the features  $\mathbb{X}$

$$\hat{\mathbb{Y}} = \mathbb{X} \cdot \omega^\top, \quad (5)$$

where matrix  $\omega$  is a  $1 \times (d + 1)$  matrix, namely  $\omega = (\omega_0, \omega_1, \dots, \omega_d)$ , that encapsulates the coefficients of the linear relation. As in equation (1),  $\hat{\mathbb{Y}}$  indicates the *predicted values* of the target as opposed to the true value of the target,  $\mathbb{Y}$ . Moreover, note that equation (5) is the analog of equation (1) in the case where we have more than one feature.

We can now define the cost function in a similar manner and proceed with its minimization. The cost function is given by

$$J(\omega) = \frac{1}{2}(\mathbb{Y} - \hat{\mathbb{Y}})^\top \cdot (\mathbb{Y} - \hat{\mathbb{Y}}) \quad (6)$$

Equation (6) should be thought as the analog of equation (2). The difference between the two is that in equation (6), the cost function  $J$  is the function of  $d + 1$  omega's  $\omega = (\omega_0, \omega_1, \dots, \omega_d)$ . Notice that  $\omega_0$  denotes the bias term (*i.e.* the constant term) as usual. We now want to minimize the cost function  $J(\omega)$ . Replacing equation (5) for the predicted target  $\hat{\mathbb{Y}}$  in equation (6), we obtain

$$J(\omega) = \frac{1}{2}(\mathbb{Y}^\top - \omega \mathbb{X}^\top) \cdot (\mathbb{Y} - \mathbb{X} \omega^\top) = \frac{1}{2}(\mathbb{Y}^\top \mathbb{Y} - 2 \omega \mathbb{X}^\top \mathbb{Y} + \omega \mathbb{X}^\top \mathbb{X} \omega^\top). \quad (7)$$

Differentiating with respect to  $\omega$  and setting  $\frac{\partial J}{\partial \omega} = 0$ , we obtain

$$-\mathbb{X}^\top \mathbb{Y} + \mathbb{X}^\top \mathbb{X} \omega^\top = 0 \quad \implies \quad \boxed{\omega^\top = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbb{Y}} \quad (8)$$

The boxed equation in (8) is the **exact solution** of the linear regression problem! Also note that the superscript  $-1$  in equation (8) indicates the *inverse matrix*.

**Question:** Does scikit-learn library use equation (8) to find the coefficients of linear regression?

**Answer:** No! It is computationally very costly to calculate the inverse of big matrices. The scikit-learn library finds the coefficients of regression (*i.e.* matrix  $\omega$ ) by approximation through an iterative process.

### 1.3 Assessing the Accuracy of the Model

Now we would like to see how we should assess the performance of the linear regression model. In order to assess the performance of the linear regression, we first need to define metrics by which the accuracy of the model can be measured. These measures are not exclusive to linear regression model, and can be used to assess the performance of any model whose job is the prediction of a continuous variable. There are several different metrics defined for assessing the performance of a given model, and in here we discuss three important metrics, namely the *RSE*, the *R<sup>2</sup> score*, and the *F-statistic*.

### 1.3.1 Residual Standard Error (RSE)

The residual standard error is one metric to assess the performance of the linear regression. It is denoted by  $RSE$  and is defined by

$$RSE = \sqrt{\frac{1}{n-2}J(\omega)} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}, \quad (9)$$

where  $J(\omega)$  is the cost function defined in (6). Minimization of the cost function through linear regression implies that  $RSE$  will assume its minimum value. In other words, any other linear model will have a higher value of  $RSE$ . Among linear models, the model with the smallest value of  $RSE$  is most favorable. If you construct two linear models, and they happen to have different values for  $RSE$ , the model with smaller  $RSE$  should be preferred over the other. Note that  $RSE$  carries the same unit as the target variable.

### 1.3.2 $R^2$ Statistic

One problem with  $RSE$  metric is that it is dimensionful, and it is not always clear what constitutes a good  $RSE$ . To address this issue, an alternative measure of the fit (*i.e.*  $R^2$ ) is defined that takes the form of proportion, and hence, is dimensionless.  $R^2$  statistic measures *what portion of variance is explained*, and  $0 \leq R^2 \leq 1$ . The closer  $R^2$  is to 1, the more successful the linear model is in predicting the target variable. To define  $R^2$ , we need to define the Mean Standard Error  $MSE$ . First, let us recall the definition of variance from your statistics knowledge:

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mu_y)^2, \quad (10)$$

where  $\mu_y$  is the mean of  $y$ -variable (*i.e.*  $\mu_y = \frac{1}{n} \sum_{i=1}^n y_i$ ). The mean standard error,  $MSE$ , is then defined by

$$MSE = \frac{1}{n}J(\omega) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2. \quad (11)$$

Note that  $MSE$  is different from the cost function  $J$  only by a factor of  $\frac{1}{n}$ . The  $R^2$  score/statistic is then defined by

$$R^2 = 1 - \frac{MSE}{\sigma^2}. \quad (12)$$

Note that for the perfect linear model ( $RSE = 0$ ), we have  $R^2 = 1$ . If  $R^2 = 0.6$  for instance, then this implies that 60% of the variance is explained by the model.

### 1.3.3 F-Statistic

The  $F$ -statistic is used to verify if there is a relationship between the response/target variable and the features. In other words, the  $F$ -statistic performs a hypothesis testing with the following null and alternative hypotheses. The **null hypothesis** states that *there is no relationship between the target and features*

$$H_0 : \forall i \in \{1, 2, \dots, d\} \quad \omega_i = 0, \quad (13)$$

while the **alternative hypothesis** states that *there is a relationship between the target and features*

$$H_1 : \exists i \in \{1, 2, \dots, d\} \quad \text{such that} \quad \omega_i \neq 0. \quad (14)$$

Note that the above hypotheses,  $H_0$  and  $H_1$ , do not concern the bias term  $\omega_0$ . The hypothesis testing is performed by computing the  $F$ -statistic

$$F = \frac{(n - d - 1)(\sigma^2 - \text{MSE})}{d \text{MSE}} = \frac{n - d - 1}{d} \frac{R^2}{1 - R^2}. \quad (15)$$

Note that just like  $R^2$  statistic, the  $F$ -statistic is a dimensionless quantity. Moreover, the  $F$ -statistic is always greater than 1 (*i.e.*  $F > 1$ ). If the  $F$ -statistic is a number close to 1, you should reject the alternative hypothesis (This means that  $H_0$  is true and there is no relationship between the target and the features). On the other hand, if  $F$ -statistic is greater than 1, then you should reject the null hypothesis and accept the alternative hypothesis (This means that there is a relationship between the target and the features).

## 1.4 Pearson Correlation

Although  $F$ -statistic can be used to decide whether there is a relationship between a target variable and a set of features, the  $F$ -statistic is unable to determine the role of dominant features when there is a relationship between the target and features. In other words, we need a **feature selection criterion**. To compare the significance of different features, one needs a more effective statistical tool. It turns out that the Pearson correlation enables us to compare the significance of different variables in an effective way.

To introduce the notion of correlation, we do not take a rigorous statistical approach. Rather we introduce Pearson correlation in a practical manner. Suppose  $X$  and  $Y$  are two random variables (Roughly speaking, this implies that they are associated with probability distributions). Then, the (Pearson) correlation between  $X$  and  $Y$  is defined by

$$\text{corr}[x, y] = \frac{1}{n} \frac{\text{cov}[x, y]}{\sigma_x \sigma_y} \in [-1, 1], \quad (16)$$

where  $\sigma_x$  and  $\sigma_y$  refer to the standard deviations of variables  $X$  and  $Y$ , respectively. Moreover, the covariance between variables  $X$  and  $Y$  in (16) is defined by

$$\text{cov}[x, y] = \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y). \quad (17)$$

There are a couple of important comments in order:

- Correlation  $\text{corr}[x, y]$  between two (random) variables  $X$  and  $Y$  is always a number between  $-1$  and  $1$  (i.e.  $-1 \leq \text{corr}[x, y] \leq 1$ ).
- A value close to  $0$  shows that there is little correlation between  $X$  and  $Y$ . The extreme case  $\text{corr}[x, y] = 0$  shows that there is absolutely no correlation between  $X$  and  $Y$ , and the two variables are independent.
- A positive value for the correlation ( $\sim \text{corr}[x, y] > 0.2$ ) shows a positive correlation between  $X$  and  $Y$  (i.e.  $Y$  increases as  $X$  increases and vice versa). The higher  $\text{corr}[x, y]$  is, the stronger the correlation between  $X$  and  $Y$  is.
- A negative value for the correlation ( $\sim \text{corr}[x, y] < -0.2$ ) shows a negative correlation between  $X$  and  $Y$  (i.e.  $Y$  increases as  $X$  decreases and vice versa). The smaller  $\text{corr}[x, y]$  is, the stronger the negative correlation between  $X$  and  $Y$  is.
- $\text{corr}[x, y] = +1$  and  $\text{corr}[x, y] = -1$  correspond to perfect linear correlation between  $X$  and  $Y$ .

We can employ the (Pearson) correlation to find the dominant features that affect the target variable of the model. Moreover, (Pearson) correlation can be used to detect strong correlations among features themselves.

## 1.5 Spearman Rank Correlation Coefficient

The Pearson correlation coefficient determines to what extent a linear model of the form  $y = mx + b$  between the target and features can fit the observed data. This *generally* does a good job in measuring the similarity between the variables, but it is possible to construct pathological examples where the correlation coefficient between  $X$  and  $Y$  is zero, yet  $Y$  is completely dependent on (and hence perfectly predictable from)  $X$ . Therefore, it is essential to measure correlation between  $X$  and  $Y$  through different benchmarks.

Another important quantity that measures correlation between two variables is the **Spearman correlation coefficient** (Note that there exists a very similar notation, called the **Kendall correlation coefficient**, to Spearman correlation. The two notions are very similar, and in hence, we only talk about one of them). The Spearman correlation coefficient essentially counts the number of pairs of data points which are *out of order*. Assume that our data set contains points  $(x_1, y_1)$  and  $(x_2, y_2)$  where  $x_1 < x_2$  and  $y_1 < y_2$ . This is a vote that the values are positively correlated, whereas the vote would be for a negative correlation if  $y_2 < y_1$ . More specifically, the Spearman rank correlation coefficient  $\rho$  is defined by

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}, \quad (18)$$

where  $d_i = \text{rank}(x_i) - \text{rank}(y_i)$ . To see how the above formula works, let us consider the following baby example. Consider the following dataset.

$x_i$	$y_i$	$\text{rank}(x_i)$	$\text{rank}(y_i)$	$d_i$	$d_i^2$
5	6	3	2	1	1

$x_i$	$y_i$	$\text{rank}(x_i)$	$\text{rank}(y_i)$	$d_i$	$d_i^2$
2	14	1	3	-2	4
4	-2	2	1	1	1
8	17	4	4	0	0

We have basically ranked  $x$  and  $y$  columns, and calculated  $d_i$  for each row. Now we can use equation (18) to calculate the correlation coefficient.

$$\rho = 1 - \frac{6(1 + 4 + 1 + 0)}{4(16 - 1)} = 1 - \frac{36}{60} = \frac{2}{5} = 0.4. \quad (19)$$

We can easily calculate the Spearman (and the Kendall) correlation coefficients through the scipy library. Let us quickly see how this is done.

```
[1]: # Import the necessary libraries

from numpy.random import rand
from scipy.stats import spearmanr # For Spearman correlation
from scipy.stats import kendalltau # For Kendall correlation

# Create a toy dataset

x = rand(1000) * 20 # create a feature x as numpy array

y = x + (rand(1000) * 5) # creat a target y by adding some noise to x

# Calculate the coefficient rho and its associated p-value

rho, p = spearmanr(x, y)

print('Spearman correlation coefficient rho = %.4f' % rho, '\n')
print('Spearman correlation p-value = %.5f' % p, '\n')

# If p-value < alpha (level of significance) H_0 is rejected --> There is
    ↳ correlation between x and y.

# Note H_0: There is no correlation between x and y.
# Note H_1: x and y are correlated.

# Similarly, you can calculate the Kendall correlation

tau, p = kendalltau(x, y)

print('Kendall correlation coefficient tau = %.4f' % tau, '\n')
print('Kendall correlation p-value = %.5f' % p)
```



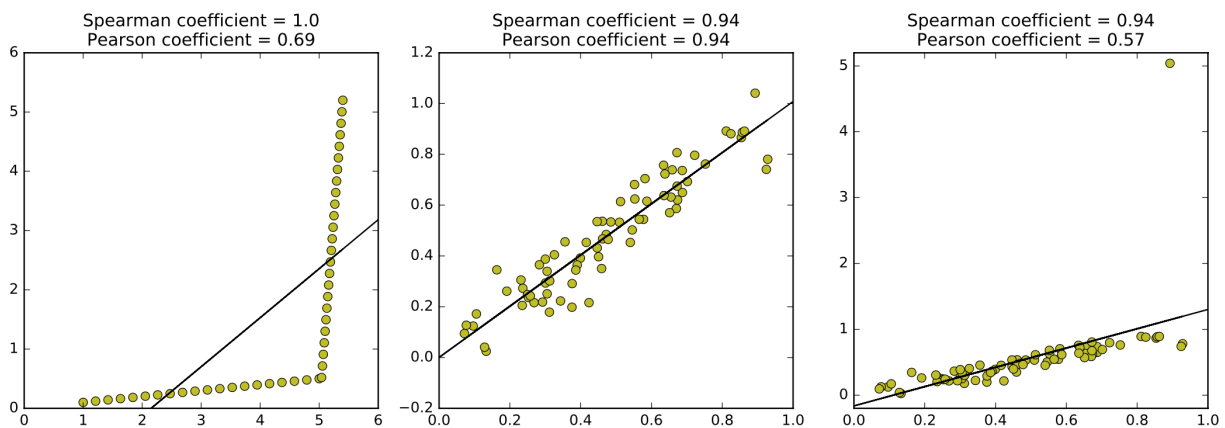
Spearman's correlation coefficient  $\rho = 0.9708$

Spearman's correlation p-value = 0.00000

Kendall correlation coefficient  $\tau = 0.8418$

Kendall correlation p-value = 0.00000

In practice, one should consider both the Pearson and the Spearman (or Kendall) correlations. The combination of the two will provide more insight about your dataset. Consider the following plots.



**Question:** As observed in above plots, the middle plot receives high Pearson correlation and Spearman correlation scores. Why are the Pearson and Spearman correlation coefficients for the left and right plots very different?

We finish the discussion on correlation by reminding you an important fact from statistics:

**Warning: Correlation does not imply causation!**

You should keep the above fact in mind whenever you work with the concept of correlation.

## 1.6 Bias-Variance Trade-Off

In statistics, the contrast between complexity and performance is known as the **bias-variance trade-off**. Let us review the two concepts briefly:

- **Bias** is an error that is caused by *assuming incorrect assumptions* that built into the model. For instance, assuming that a nonlinear relation between a target and features is linear. Bias can be associated with *prejudice*. In this case, models perform on both training and test datasets lousy, and they are *underfit*.
- **Variance** is a sensitivity that is caused by fluctuations in the training dataset. If the training set contains sampling or measurement errors, this noise induces variance in the resulting

model. Errors of variance lead to *overfit* models. The search of overfit models for accuracy causes them to mistake noise for actual signal. They adjust so well to the training data that noise leads them astray. Models that *perform much better on the training dataset than the testing dataset are overfit*.

## 2 Coding for Regression Analysis

```
[1]: # Importing the PIN Dataset

import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

df = pd.read_csv("https://raw.githubusercontent.com/simsekergun/MLPdatasets/main/
→pin_dataset1.csv")
df.head()
```

```
[1]:
```

	t1	t2	t3	d1	d2	d3	PN \
0	1625	363	138	1.780000e+18	1.340000e+16	3.250000e+18	-161.121081
1	1813	406	119	7.500000e+18	5.470000e+16	8.290000e+18	-161.420393
2	1938	494	306	4.220000e+19	7.310000e+15	5.700000e+18	-162.376610
3	1500	450	250	1.000000e+19	2.000000e+16	2.240000e+18	-162.087268
4	1432	268	473	1.870000e+18	8.480000e+14	5.020000e+17	-160.281601

	Qeff	DT	BW
0	0.127847	119.739190	6.426738
1	0.142941	146.779380	5.777518
2	0.146864	113.419730	8.953805
3	0.181376	119.203478	7.439119
4	0.195884	125.800230	5.969470

```
[2]: df.shape # to know no. of rows and columns
```

```
[2]: (200, 10)
```

```
[3]: df.isnull().sum() # to check if there is any missing value
```

```
[3]: t1      0
      t2      0
      t3      0
      d1      0
      d2      0
      d3      0
      PN      0
      Qeff    0
      DT      0
```

```
BW      0
dtype: int64
```

```
[4]: df.dtypes
```

```
[4]: t1      int64
      t2      int64
      t3      int64
      d1     float64
      d2     float64
      d3     float64
      PN     float64
      Qeff   float64
      DT     float64
      BW     float64
      dtype: object
```

Predict the quantum efficiency using phase noise

```
[5]: # Defining the feature and the target of the model

X = df[['PN']]          # Feature
y = df[['Qeff']]        # Target
```

The linear model is now  $y = \omega_1 X_1 + \omega_0$ , where  $X_1$  refers to PN.  $\omega_1$  is the coefficient of regression, and the constant  $\omega_0$  is the intercept (*i.e.* bias term).

```
[6]: # Breaking the data into train and test subsets

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
→random_state=3)
```

## 2.1 Performing Regression through scikit-learn

```
[7]: # Importing 'LinearRegression' through linear_model module

from sklearn.linear_model import LinearRegression

reg = LinearRegression()          # Instantiate
reg.fit(X_train, y_train)         # Fit the train data

r2_train_score = reg.score(X_train, y_train) # Calculating R^2 score for train

print('R^2 score for train dataset = ', round(r2_train_score, 4), '\n')
print('Coefficients of Linear Model:', reg.coef_, '\n')
```

```
print('Intercept:', reg.intercept_)
```

R<sup>2</sup> score for train dataset = 0.6462

Coefficients of Linear Model:  $\begin{bmatrix} -0.0290857 \end{bmatrix}$

Intercept:  $\begin{bmatrix} -4.52400226 \end{bmatrix}$

```
[8]: # Finding the predictions of the model for test dataset
```

```
y_pred = reg.predict(X_test)
```

```
y_pred[:10] # Representing the Qeff prediction for the first 10 data points in  
            ↳ test dataset
```

```
[8]: array([[0.29433772],  
            [0.31366907],  
            [0.42008194],  
            [0.44967163],  
            [0.43908034],  
            [0.3398191 ],  
            [0.43807007],  
            [0.3034589 ],  
            [0.35224422],  
            [0.3565935 ]])
```

```
[9]: # Evaluating the performance of the model on the test dataset
```

```
r2_test_score = reg.score(X_test, y_test) # Calculating R2 score for test
```

```
print('R2 score for test dataset = ', round(r2_test_score, 4), '\n')
```

R<sup>2</sup> score for test dataset = 0.6863

```
[10]: # Let's plot predictions vs ground truth for 'Qeff'
```

```
fig, ax = plt.subplots(figsize=(8,5))
```

```
ax.scatter(y_test, y_pred, edgecolors=(0.3, 0.2, 0.7)) # Scatter plot for  
            ↳ predictions vs truth
```

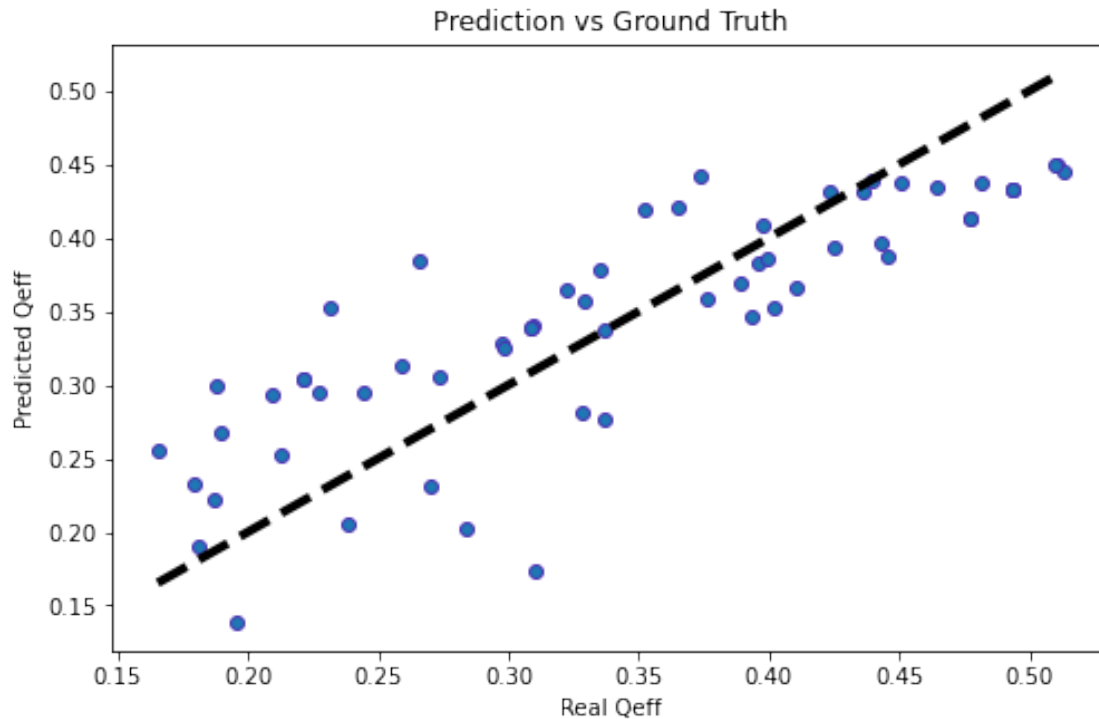
```
ax.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=4)  
            ↳ # Draw line y=x
```

```
ax.set_xlabel('Real Qeff')
```

```
ax.set_ylabel('Predicted Qeff')
```

```
plt.title('Prediction vs Ground Truth', fontdict=None, loc='center')
```

```
plt.show()
```



The closer a point in above plot is to the  $y = x$  line, the more realistic the prediction for the Qeff is.

We can see here that not all the points are close to the line which means the prediction is not very good.

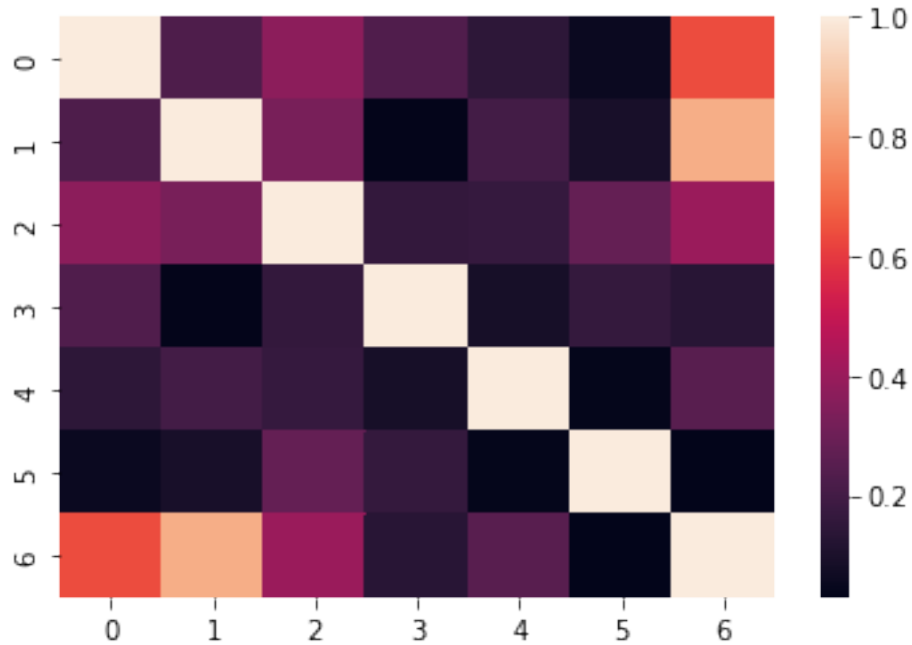
### 2.1.1 Predict the phase noise

#### Calculating Pearson Correlation between Variables

```
[11]: import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import kendalltau    # For Kendall correlation

cols = df.columns[:7]                # List of columns of the dataframe

cm = np.corrcoef(df[cols].values.T)  # Calculate Pearson correlation
hm = sns.heatmap(abs(cm))
plt.show()
```



The above matrix shows that t1 and t2 have the highest correlations with the target variable PN. Therefore, as our first step in our regression, we identify t1 and t2 as the two features to explain the target PN. Now, let's check the Spearman correlation. This is a cross check to make sure that we have not missed any influential variable (on PN) that Pearson correlation hasn't been able to detect.

```
[12]: # Calculate Spearman correlation coefficient between 'PN' and each feature
from scipy.stats import spearmanr      # For Spearman correlation

for col in cols[:-1]:
    rho, p = spearmanr(df[col].values, df['PN'].values)
    print('Spearman correlation between PN and %s is %s' %(col, round(rho, 4)))
```

```
Spearman correlation between PN and t1 is 0.6802
Spearman correlation between PN and t2 is -0.7705
Spearman correlation between PN and t3 is -0.4798
Spearman correlation between PN and d1 is 0.0351
Spearman correlation between PN and d2 is 0.2031
Spearman correlation between PN and d3 is 0.2953
```

It is again observed that the dominant effect comes from the two variables t1 (with  $\rho = 0.6802$ ) and t2 (with  $\rho = -0.7705$ ). The Pearson and Spearman correlations agree on the role of the dominant variables.

### 2.1.2 Performing Regression through scikit-learn

```
[13]: # Defining the features and the target of the model
```

```
X = df[cols[:2]]          # Features
y = df[['PN']]            # Target
```

```
[14]: # Breaking the data into train and test subsets
```

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
→random_state=3)
```

```
[15]: # Importing 'LinearRegression' through linear_model module
```

```
from sklearn.linear_model import LinearRegression

reg = LinearRegression()      # Instantiate
reg.fit(X_train, y_train)     # Fit the train data

r2_train_score = reg.score(X_train, y_train) # Calculating R^2 score for train

print('R^2 score for train dataset = ', round(r2_train_score, 4), '\n')
print('Coefficients of Linear Model:', reg.coef_, '\n')
print('Intercept:', reg.intercept_)
```

R<sup>2</sup> score for train dataset = 0.9185

Coefficients of Linear Model: [[ 0.00240653 -0.00986512]]

Intercept: [-161.96202025]

```
[16]: # Finding the predictions of the model for test dataset
```

```
y_pred = reg.predict(X_test)

y_pred[:10] # Representing the PN prediction for the first 10 data points in
→test dataset
```

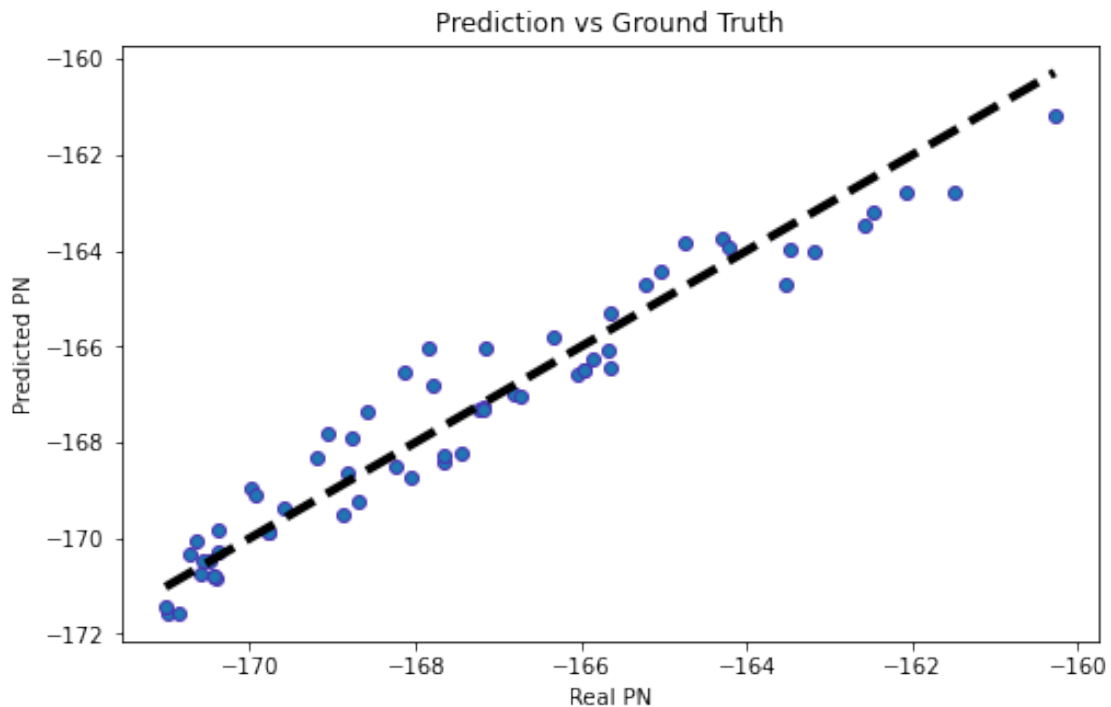
```
[16]: array([[ -166.42703086],
              [ -165.82076695],
              [ -168.96735833],
              [ -171.54703017],
              [ -170.07376956],
              [ -167.30319243],
              [ -170.73319349],
```

```
[-166.49700325],  
[-168.39301892],  
[-166.79319063]])
```

```
[17]: # Evaluating the performance of the model on the test dataset  
  
r2_test_score = reg.score(X_test, y_test) # Calculating R^2 score for test  
  
print('R^2 score for test dataset = ', round(r2_test_score, 4), '\n')
```

R<sup>2</sup> score for test dataset = 0.9328

```
[18]: # Let's plot predictions vs ground truth for 'PN'  
  
fig, ax = plt.subplots(figsize=(8,5))  
  
ax.scatter(y_test, y_pred, edgecolors=(0.3, 0.2, 0.7)) # Scatter plot for  
→ predictions vs truth  
ax.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=4)  
→ # Draw line y=x  
  
ax.set_xlabel('Real PN')  
ax.set_ylabel('Predicted PN')  
plt.title('Prediction vs Ground Truth', fontdict=None, loc='center')  
plt.show()
```





The closer a point in above plot is to the  $y = x$  line, the more realistic the prediction for the PN is. Here the data points are much closer to the line.

### 2.1.3 Predict the quantum efficiency

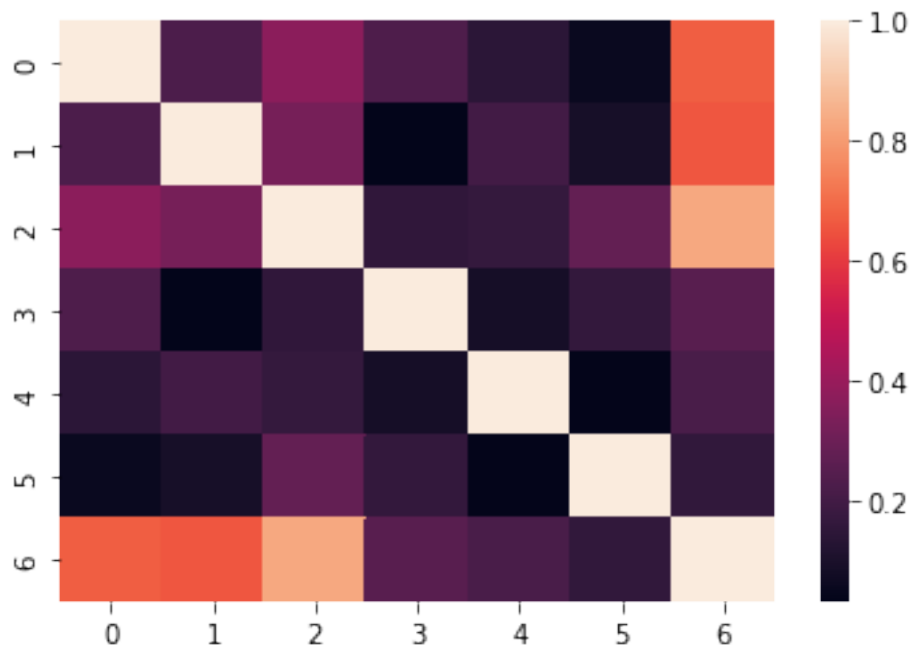
#### Calculating Pearson Correlation between Variables

```
[19]: import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import kendalltau # For Kendall correlation

c = ['t1', 't2', 't3', 'd1', 'd2', 'd3', 'Qeff']

cols = df[c].columns # List of columns of the dataframe

cm = np.corrcoef(df[cols].values.T) # Calculate Pearson correlation
hm = sns.heatmap(abs(cm)) # Represent correlation by a heat map
plt.show()
```



The above matrix shows that  $t_1$ ,  $t_2$  and  $t_3$  have the highest correlations with the target variable  $Q_{eff}$ . Therefore, as our first step in our regression, we identify  $t_1$ ,  $t_2$  and  $t_3$  as the three features to explain the target  $Q_{eff}$ . Now, let's check the Spearman correlation. This is a cross check to make sure that we have not missed any influential variable (on  $Q_{eff}$ ) that Pearson correlation hasn't been

able to detect.

```
[20]: # Calculate Spearman correlation coefficient between 'Qeff' and each feature
from scipy.stats import spearmanr      # For Spearman correlation

for col in cols[:-1]:
    rho, p = spearmanr(df[col].values, df['Qeff'].values)
    print('Spearman correlation between Qeff and %s is %s' %(col, round(rho, 4)))
```

```
Spearman correlation between Qeff and t1 is -0.6429
Spearman correlation between Qeff and t2 is 0.6707
Spearman correlation between Qeff and t3 is 0.8374
Spearman correlation between Qeff and d1 is -0.0406
Spearman correlation between Qeff and d2 is -0.1754
Spearman correlation between Qeff and d3 is -0.4013
```

It is again observed that the dominant effect comes from the three variables  $t_1$  (with  $\rho = -0.6429$ ),  $t_2$  (with  $\rho = 0.6707$ ) and  $t_3$  (with  $\rho = 0.8374$ ). The Pearson and Spearman correlations agree on the role of the dominant variables.

#### 2.1.4 Performing Regression through scikit-learn

```
[21]: # Defining the features and the target of the model
```

```
X = df[cols[:3]]          # Features
y = df[['Qeff']]          # Target
```

```
[22]: # Breaking the data into train and test subsets
```

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
→random_state=3)
```

```
[23]: # Importing 'LinearRegression' through linear_model module
```

```
from sklearn.linear_model import LinearRegression

reg = LinearRegression()      # Instantiate
reg.fit(X_train, y_train)     # Fit the train data

r2_train_score = reg.score(X_train, y_train)  # Calculating R^2 score for train

print('R^2 score for train dataset = ', round(r2_train_score, 4), '\n')
print('Coefficients of Linear Model:', reg.coef_, '\n')
print('Intercept:', reg.intercept_)
```

```
R^2 score for train dataset = 0.9801
```

Coefficients of Linear Model:  $\begin{bmatrix} -7.21293744 \times 10^{-5} & 1.95226816 \times 10^{-4} \\ 1.82729618 \times 10^{-4} \end{bmatrix}$

Intercept: [0.14996192]

```
[24]: # Finding the predictions of the model for test dataset

y_pred = reg.predict(X_test)

y_pred[:10] # Representing the Qeff prediction for the first 10 data points in
            ↳ test dataset
```

```
[24]: array([[0.25414513],
             [0.24957537],
             [0.35306449],
             [0.51992451],
             [0.43019417],
             [0.30277808],
             [0.43881116],
             [0.24078971],
             [0.26053308],
             [0.31316042]])
```

```
[25]: # Evaluating the performance of the model on the test dataset

r2_test_score = reg.score(X_test, y_test) # Calculating R^2 score for test

print('R^2 score for test dataset = ', round(r2_test_score, 4), '\n')
```

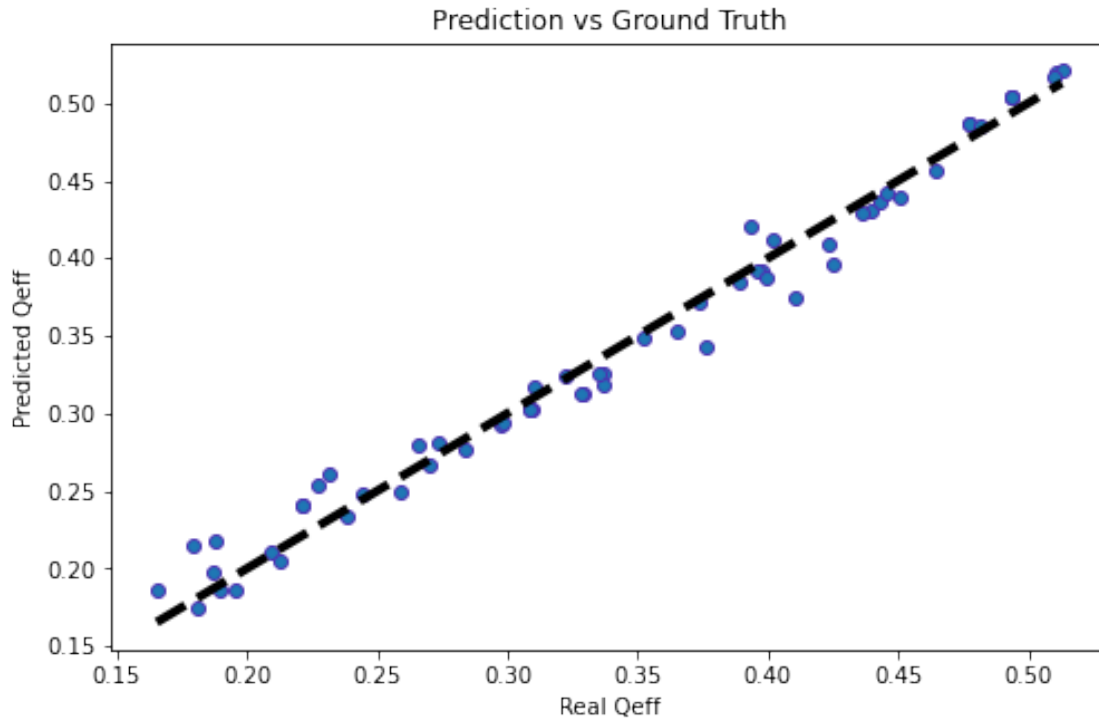
R<sup>2</sup> score for test dataset = 0.9799

```
[26]: # Let's plot predictions vs ground truth for 'Qeff'

fig, ax = plt.subplots(figsize=(8,5))

ax.scatter(y_test, y_pred, edgecolors=(0.3, 0.2, 0.7)) # Scatter plot for
↳ predictions vs truth
ax.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=4)
↳ # Draw line y=x

ax.set_xlabel('Real Qeff')
ax.set_ylabel('Predicted Qeff')
plt.title('Prediction vs Ground Truth', fontdict=None, loc='center')
plt.show()
```



If we compare the above scatter plot with the scatter plot we got when we were using only PN as our feature variable, we will see that the data points in above scatter plot are very close to the line as compared to the first plot. The predictions for Qeff are much more accurate when we are using  $t_1$ ,  $t_2$  and  $t_3$  as target variables instead of using only PN.

The  $R^2$  score was around 65% when only PN was considered where the  $R^2$  score increased to 98% when  $t_1$ ,  $t_2$  and  $t_3$  are used.

## 2.2 Performing Regression through StatsModels

An alternative to scikit-learn library is StatsModels. We now solve the same problem using StatsModels.

```
[27]: import statsmodels.api as sm
```

```
[28]: XC_train = sm.add_constant(X_train) # Creating the bias term omega_0
      model = sm.OLS(y_train, XC_train) # Defining model and the data into the model
      results = model.fit() # Fitting the data into the model
      print('Summary:\n', results.summary()) # Printing the summary of results
```

Summary:

OLS Regression Results			
=====			
Dep. Variable:	Qeff	R-squared:	0.980
Model:	OLS	Adj. R-squared:	0.980

```

Method:                Least Squares    F-statistic:                2228.
Date:                  Tue, 14 Feb 2023  Prob (F-statistic):        2.23e-115
Time:                  21:32:14         Log-Likelihood:             401.31
No. Observations:      140             AIC:                        -794.6
Df Residuals:          136             BIC:                        -782.9
Df Model:               3
Covariance Type:       nonrobust

```

```

=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const          0.1500        0.006     25.575      0.000        0.138        0.162
t1          -7.213e-05        2.4e-06    -29.996      0.000     -7.69e-05     -6.74e-05
t2           0.0002        6.19e-06     31.558      0.000         0.000         0.000
t3           0.0002        4.42e-06     41.338      0.000         0.000         0.000
=====
Omnibus:                 10.494    Durbin-Watson:                 1.993
Prob(Omnibus):            0.005    Jarque-Bera (JB):            13.227
Skew:                    -0.462    Prob(JB):                     0.00134
Kurtosis:                 4.189    Cond. No.                      6.63e+03
=====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 6.63e+03. This might indicate that there are strong multicollinearity or other numerical problems.

/usr/local/lib/python3.8/dist-packages/statsmodels/tsa/tsatools.py:142:

FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only

```
x = pd.concat(x[:, :order], 1)
```

```

[29]: XC_test = sm.add_constant(X_test) # Adding bias for the test set
      y_pred_sm = results.predict(XC_test) # Finding the predictions for the test set

      print(y_pred_sm.values[:10], '\n') # Printing the predictions for the test set

```

```

[0.25414513 0.24957537 0.35306449 0.51992451 0.43019417 0.30277808
 0.43881116 0.24078971 0.26053308 0.31316042]

```

/usr/local/lib/python3.8/dist-packages/statsmodels/tsa/tsatools.py:142:

FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only

```
x = pd.concat(x[:, :order], 1)
```

```

[30]: # Let us compare the predictions of StatsModels with Scikit-learn
      rd_array = np.vectorize(round) # Rounding numpy arrays

```

```
print(rd_array(y_pred[:20].T - y_pred_sm.values[:20], 6)) # Subtract predictions
↳ of StatModels from Scikit-learn
```

```
[[-0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0. -0.
  -0. -0.]]
```

The above results show that the predictions of the StatsModels are in perfect agreement with those of scikit-learn, as they should!

## 2.3 Performing Regression using the Exact Formula

As you noticed, this dataset is not that big in size. Hence, we may find the *exact result* of regression through the exact formula we derived in the lecture:  $\omega^T = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$ . We can easily implement this formula into code using numpy.

```
[31]: # Defining a function to calculate omega from X and Y
from numpy.linalg import inv # Import 'inv' to calculate inverse matrix

def exact_reg(X, Y):
    X = np.insert(X, 0, np.ones(X.shape[0]), axis=1) # Adding a column of 1's
    OmegaT = np.matmul(inv(np.matmul(X.T, X)), np.matmul(X.T, Y)) # Using the
↳ exact formula
    Y_hat = np.matmul(X, OmegaT) # Predicted values of y
    return OmegaT, Y_hat
```

```
[32]: # Let's compare the exact results for coefficients and intercept with the
↳ results of scikit-learn
print('The intercept and coefficients from exact formula:', exact_reg(X_train.
↳ values, y_train.values)[0], '\n')
print('The intercept and coefficients from scikit-learn:', np.insert(reg.
↳ coef_, 0, reg.intercept_))
```

```
The intercept and coefficients from exact formula: [[ 1.49961919e-01]
[-7.21293744e-05]
[ 1.95226816e-04]
[ 1.82729618e-04]]
```

```
The intercept and coefficients from scikit-learn: [ 1.49961919e-01
-7.21293744e-05  1.95226816e-04  1.82729618e-04]
```

```
[33]: # Calculate the exact R^2 scores
y_hat_train = exact_reg(X_train.values, y_train.values)[1] # calculating
↳ predictions for train data
y_hat_test = exact_reg(X_test.values, y_test.values)[1] # calculating
↳ predictions for test data

from sklearn.metrics import r2_score
```

```
print('Exact train R^2 score = ', round(r2_score(y_train.values,   
→y_hat_train),6), '\n')  
print('Exact test R^2 score = ', round(r2_score(y_test.values, y_hat_test), 6))
```

Exact train R^2 score = 0.98006

Exact test R^2 score = 0.982368

The above R2 scores are in very good agreement with those obtained by scikit-learn!