# Polynomial Regression and Model Selection Project Report

Lecture: CS454/554 - Introduction to Machine Learning and Artificial Neural Networks

Instructor: Prof. Ethem Alpaydın

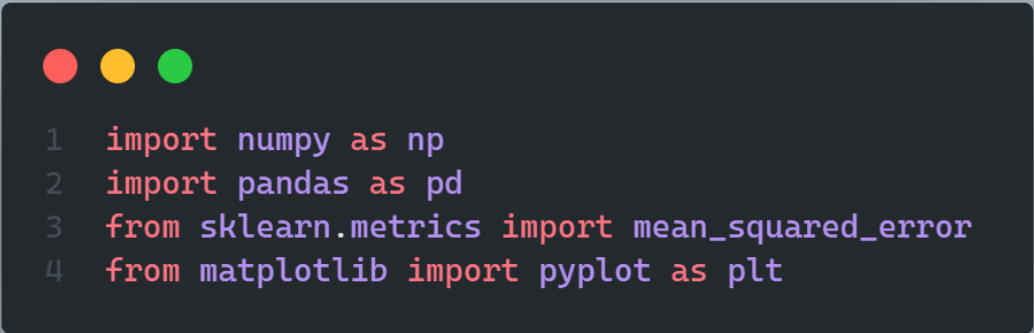Tuna Tuncer - S018474

# 1. Introduction

I was given one test data set with 100 instances and ten training data sets, each having 25 instances, for this assignment. I was expected to create polynomial regression models using the training data sets, and then measure and compare the mean squared errors of the models using the test data set. It was desired that the polynomial regression models I would develop had degrees ranging from 1 to 6. The data was generated using a hidden f(x) function and made more random by adding gaussian noise to it.

# 2. Methodology

For each x degree, 10 different models were created by using each of the 10 training data sets one by one. The mean squared errors of each of these 10 models were calculated using the test data set. Afterwards, a single value was obtained by averaging 10 different mean squared error values. When this process was repeated for each degree, a total of 6 different average mean square error values were obtained. Thus, it has been empirically found that the models with which degree explain the data more successfully.

# 3. Implementation Details

Python is used as the programming language for this assignment. In addition, Jupyter Notebook, which is widely used for machine learning related tasks, is also used as an editor.

```
1  import numpy as np
2  import pandas as pd
3  from sklearn.metrics import mean_squared_error
4  from matplotlib import pyplot as plt
```

Figure [1]: Imported Libraries

As seen in Figure 1, necessary libraries are imported. Pandas is used to read data from CSV files and save it in data frames. Numpy is used to create polynomial regression models. Sklearn is used to calculate the mean squared error of the models. Matplotlib is used to visualize the whole process.

```python
1  train_dfs = []
2  for idx in range(1, 11):
3      train_dfs.append(pd.read_csv("sample{}.csv".format(idx), header=None))
```

```python
1  test_df = pd.read_csv("test.csv", header=None)
```

Figure [2, 3]: Reading CSV Files and Saving into Data Frames

As seen in Figures 2 and 3, the data were read from the 10 training and 1 test data set CSV files and saved in data frames.

```python
1  train_xs = []
2  train_ys = []
3  for df in train_dfs:
4      train_xs.append(df.iloc[:, 0].values)
5      train_ys.append(df.iloc[:, 1].values)
```

```
1  test_x = test_df.iloc[:, 0].values
2  test_y = test_df.iloc[:, 1].values
```

Figure [4, 5]: Splitting Input and Output Data

As seen in Figures 4 and 5, the input and output columns in the given training and test data sets are divided into different data frames.

```
1  model_dict = {
2      1: [],
3      2: [],
4      3: [],
5      4: [],
6      5: [],
7      6: []
8  }
9  for idx in range(1, 7):
10     for train_x, train_y in zip(train_xs, train_ys):
11         model_dict[idx].append(np.polyfit(train_x, train_y, idx))
```

Figure [6]: Creating the Models

In Figure 6, it can be seen how models are created for each degree with the help of the **np.polyfit** function.

```
1   def predict(x, model):
2       arr = []
3       for data in x:
4           dim = len(model) - 1
5           pred = 0
6           for i in model:
7               if model[-1] == i:
8                   pred += i
9                   break
10              pred += (data ** dim) * i
11              dim -= 1
12          arr.append(pred)
13      return arr
```

Figure [7]: Defining a Predict Function

In Figure 7, the predict method to be used when calculating the mean squared error is created. This method takes a list of inputs and one of the models created above in its parameters. Then, using this model, the output value mapped by each input is calculated. Finally, a list of the output values is returned.

```
1   mse_dict = {
2       1: [],
3       2: [],
4       3: [],
5       4: [],
6       5: [],
7       6: []
8   }
9   for idx in range(1, 7):
10      for model in model_dict[idx]:
11          mse_dict[idx].append(mean_squared_error(test_y, predict(test_x, model)))
```

Figure [8]: Calculating Mean Squared Error Values

In Figure 8, a mean squared error value was calculated using the test data set for each model created.

```
mse_dict1 [2.567272528135012, 2.609548148096624464, 3.4449716507017644, 2.786426817331451, 2.5836743224409223, 3.4455544764623103, 2.6158264871464136, 3.0777750790671536, 2.7511853969409197, 2.588870387820492]
mse_dict2 [1.9268225266724486, 2.013260422264308, 2.3447657982100414, 2.113691075554605, 2.2107624869713383, 2.133199658671474745, 2.007064044199733, 2.20899710223812, 2.0692536648931, 2.0895040591667096]
mse_dict3 [1.5497249787356444, 1.914317046598652, 1.9885742712429306, 1.5816195995297704, 1.67141915352443175, 1.95202744442040078, 1.63047027109046, 1.639318538530181089, 1.963922707714174, 1.633888000445399974]
mse_dict4 [1.6018397980686216, 1.9407904534903055, 2.1377027970042897, 1.5831023434000844, 1.639973722450097, 1.9575776282353987, 1.8094072474835605, 1.6869870519995696, 1.88934834056969, 1.5818635254672428]
mse_dict5 [1.7282666326675555, 1.9852435955958276, 2.3371592607999627, 1.591852886057904, 1.6617003939839352, 2.00260988657892, 1.9507288342183324, 1.79630501691287111, 3.4535966495613133, 1.7483061385547247]
mse_dict6 [1.7368982626137008, 2.0602243875256594, 2.3389434877101046, 1.6186765173254736, 2.04027401425899995, 3.6910090331239074, 3.036933051199549, 2.0784809330596836, 14.835643970673473, 1.8919263383224896]
```

Figure [9]: Mean Squared Errors for All Models of All Degrees

In Figure 9, the calculated mean squared errors of each model created for each degree can be seen.

# 4. Results

Below you can see the plots of total 60 different models created for each degree.

- Plots for Degree 1:

0.05*x^1+0.11

0.79*x^1+0.25

0.52*x^1+(-0.03)

0.23*x^1+0.56

0.45*x^1+(-0.08)

0.86*x^1+0.48

0.37*x^1+0.08

0.62*x^1+0.04

- Plots for Degree 2:


(-0.35)*x^2+0.72*x^1+0.8


(-0.28)*x^2+0.46*x^1+0.74

(-0.38)*x^2+0.3*x^1+1.16

(-0.33)*x^2+0.64*x^1+1.23

(-0.17)*x^2+0.45*x^1+0.62

(-0.34)*x^2+0.56*x^1+1.25

(-0.27)*x^2+0.46*x^1+0.69

(-0.27)*x^2+0.73*x^1+1.14

(-0.45)*x^2+0.64*x^1+1.36

(-0.17)*x^2+0.54*x^1+0.51

- Plots for Degree 3:



(-0.15)*x^3+(-0.32)*x^2+1.45*x^1+0.85



(-0.02)*x^3+(-0.26)*x^2+0.59*x^1+0.72

$(-0.14)*x^3+(-0.3)*x^2+1.02*x^1+1.11$

$(-0.2)*x^3+(-0.33)*x^2+1.68*x^1+0.98$

(-0.2)*x^3+(-0.33)*x^2+1.54*x^1+1.18

(-0.11)*x^3+(-0.26)*x^2+1.02*x^1+1.17

(-0.17)*x^3+(-0.25)*x^2+1.39*x^1+0.56

(-0.17)*x^3+(-0.3)*x^2+1.65*x^1+1.04

(-0.02)*x^3+(-0.44)*x^2+0.73*x^1+1.32

(-0.23)*x^3+(-0.28)*x^2+1.75*x^1+0.78

- Plots for Degree 4:



0.06*x^4+(-0.14)*x^3+(-0.75)*x^2+1.37*x^1+1.28



0.03*x^4+(-0.03)*x^3+(-0.48)*x^2+0.56*x^1+0.86

0.07*x^4+(-0.17)*x^3+(-0.82)*x^2+1.1*x^1+1.54

(-0.0)*x^4+(-0.2)*x^3+(-0.32)*x^2+1.68*x^1+0.97

0.02*x^4+(-0.18)*x^3+(-0.46)*x^2+1.46*x^1+1.29

0.02*x^4+(-0.13)*x^3+(-0.4)*x^2+1.11*x^1+1.27

0.08*x^4+(-0.24)*x^3+(-0.8)*x^2+1.74*x^1+0.93

0.06*x^4+(-0.17)*x^3+(-0.77)*x^2+1.64*x^1+1.41

0.09*x^4+(-0.12)*x^3+(-0.98)*x^2+1.1*x^1+1.69

0.03*x^4+(-0.21)*x^3+(-0.52)*x^2+1.67*x^1+0.98

- Plots for Degree 5:

(-0.03)*x^5+0.06*x^4+0.18*x^3+(-0.82)*x^2+0.75*x^1+1.37
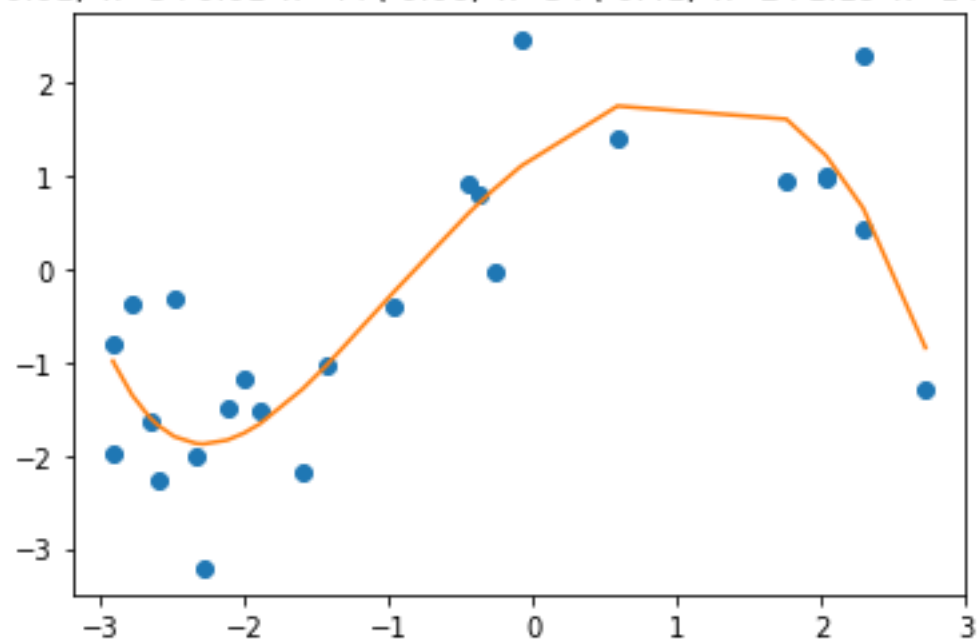


0.02*x^5+0.01*x^4+(-0.16)*x^3+(-0.39)*x^2+0.77*x^1+0.78
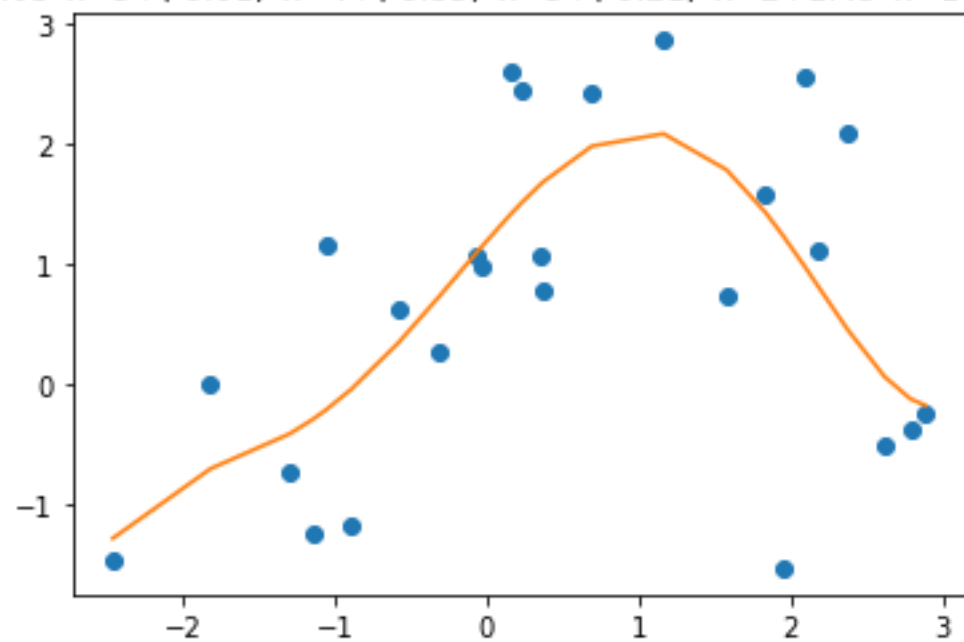
(-0.04)*x^5+0.09*x^4+0.25*x^3+(-0.93)*x^2+0.34*x^1+1.57

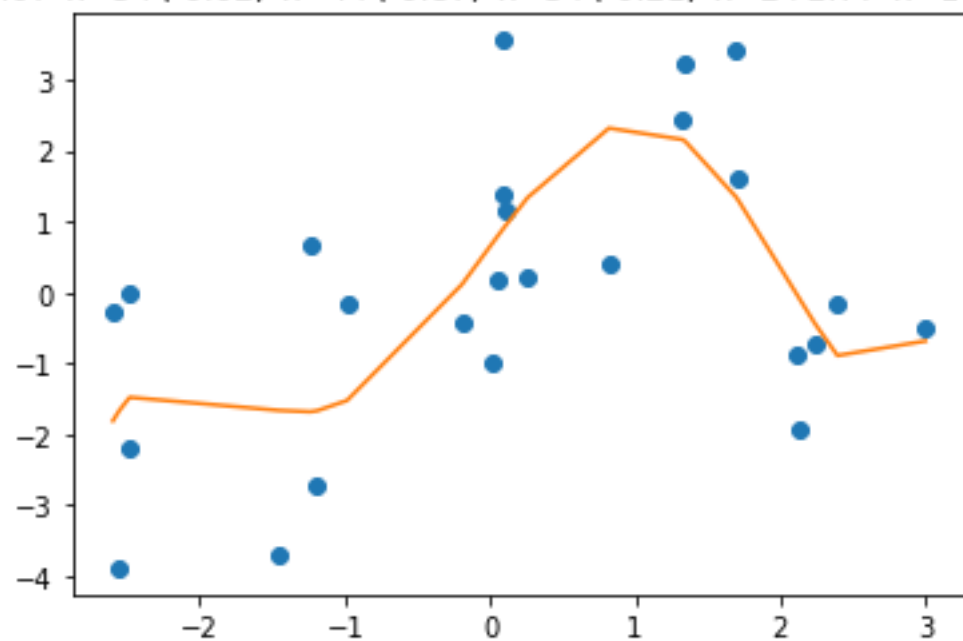(-0.0)*x^5+(-0.0)*x^4+(-0.18)*x^3+(-0.32)*x^2+1.64*x^1+0.97
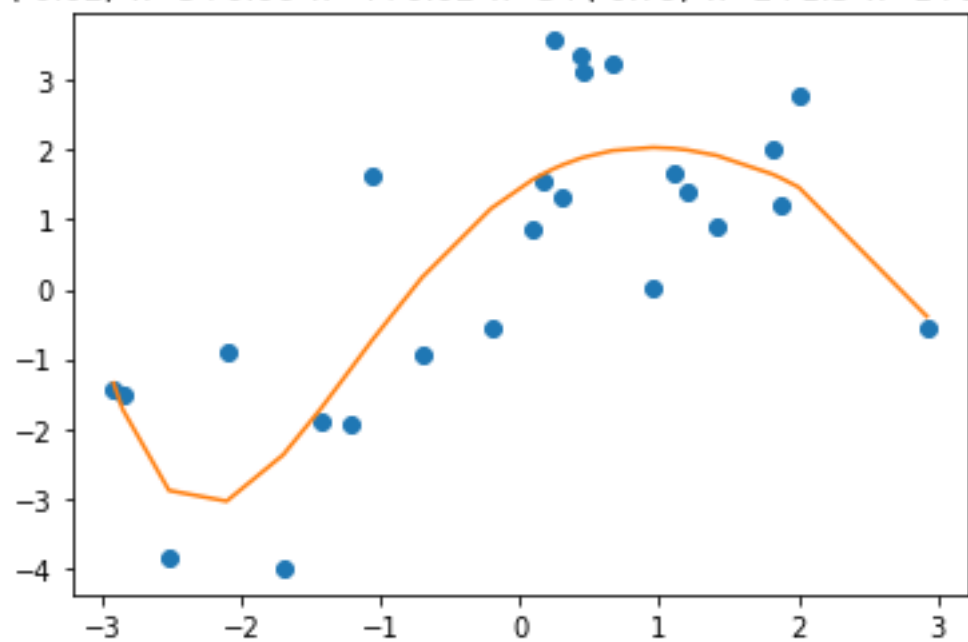
(-0.01)*x^5+0.01*x^4+(-0.06)*x^3+(-0.41)*x^2+1.19*x^1+1.21

0.03*x^5+(-0.01)*x^4+(-0.35)*x^3+(-0.21)*x^2+1.45*x^1+1.19

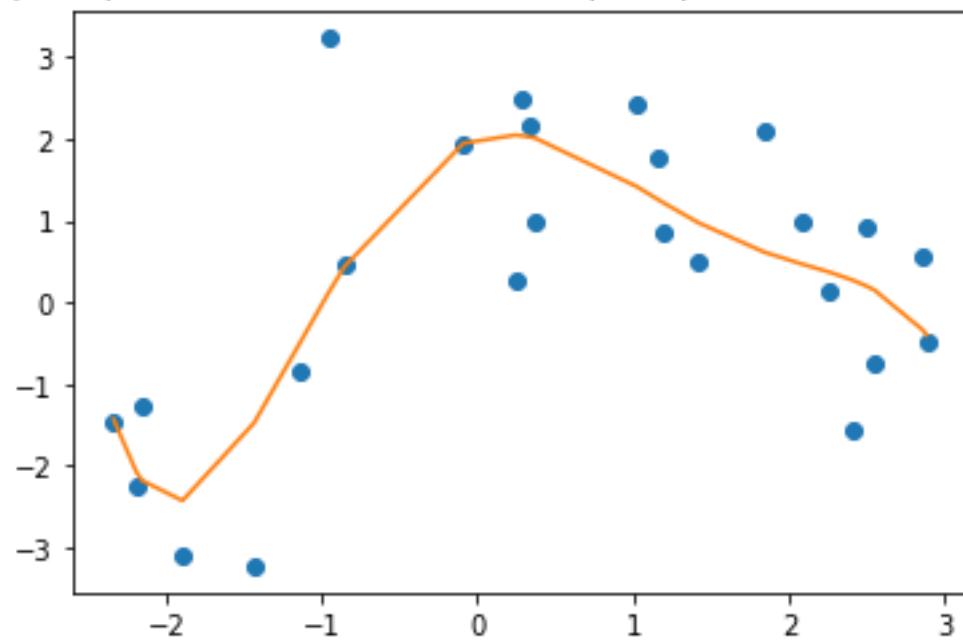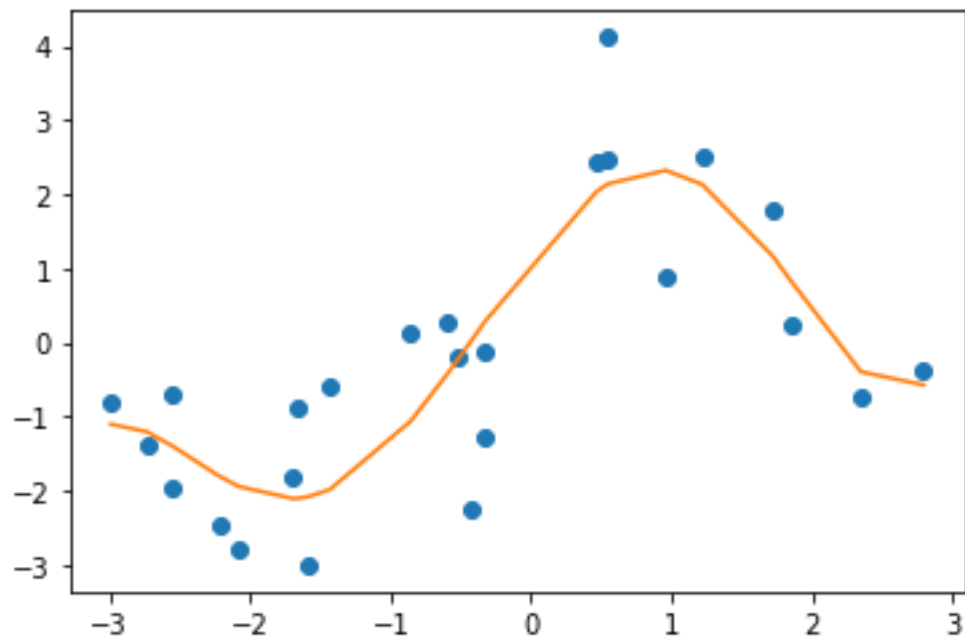0.07*x^5+(-0.02)*x^4+(-0.87)*x^3+(-0.21)*x^2+2.77*x^1+0.65

(-0.02)*x^5+0.06*x^4+0.02*x^3+(-0.79)*x^2+1.3*x^1+1.46
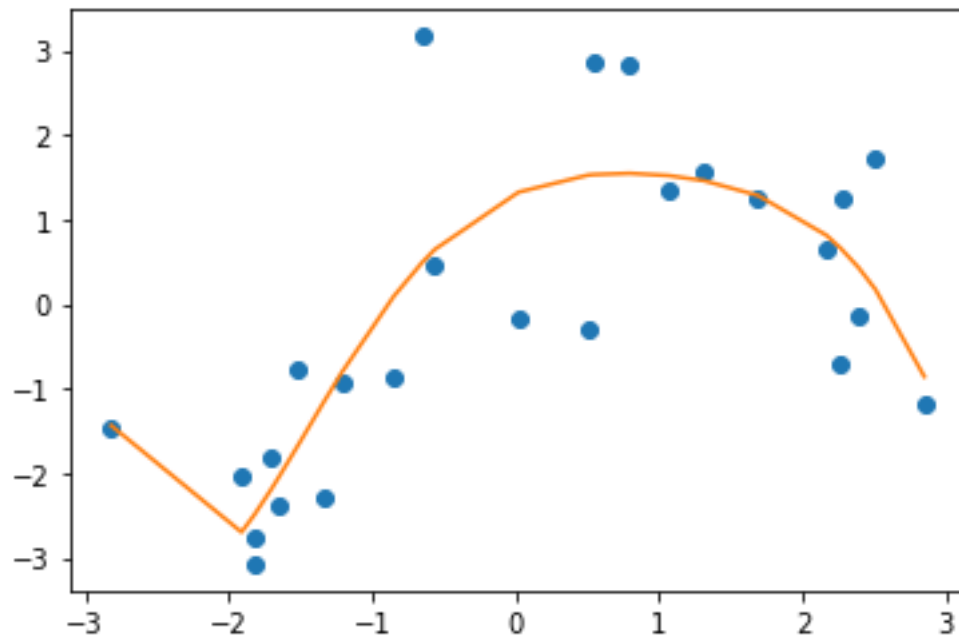
(-0.06)*x^5+0.18*x^4+0.28*x^3+(-1.48)*x^2+0.51*x^1+2.0


0.04*x^5+0.06*x^4+(-0.61)*x^3+(-0.71)*x^2+2.41*x^1+1.13

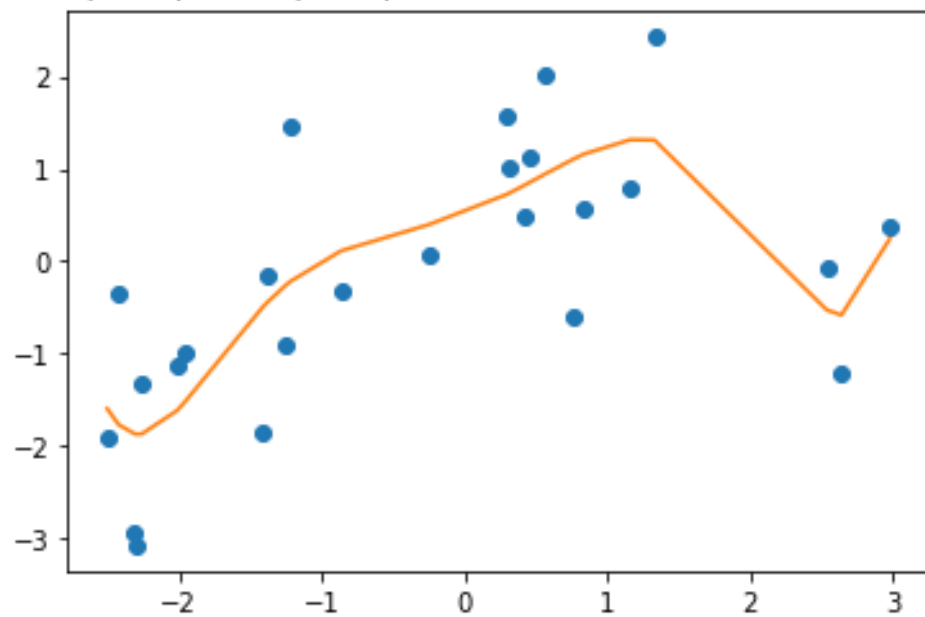- Plots for Degree 6:
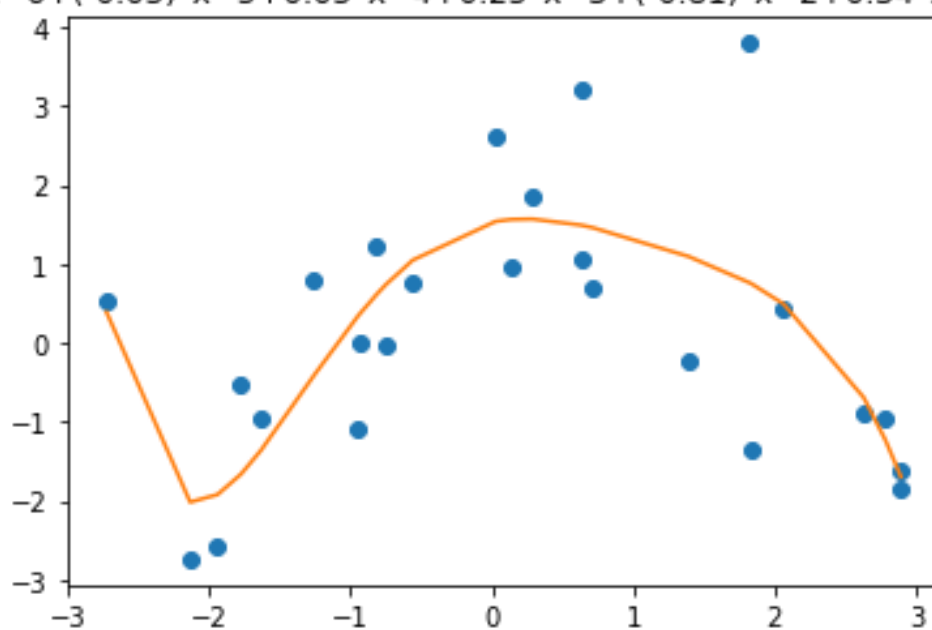
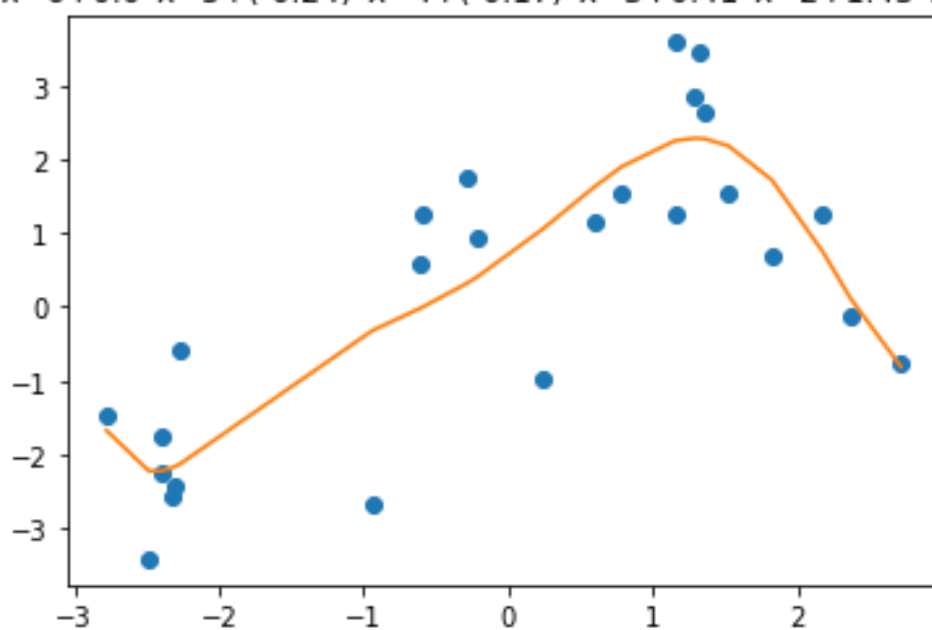0.0*x^6+(-0.03)*x^5+0.02*x^4+0.2*x^3+(-0.7)*x^2+0.73*x^1+1.31



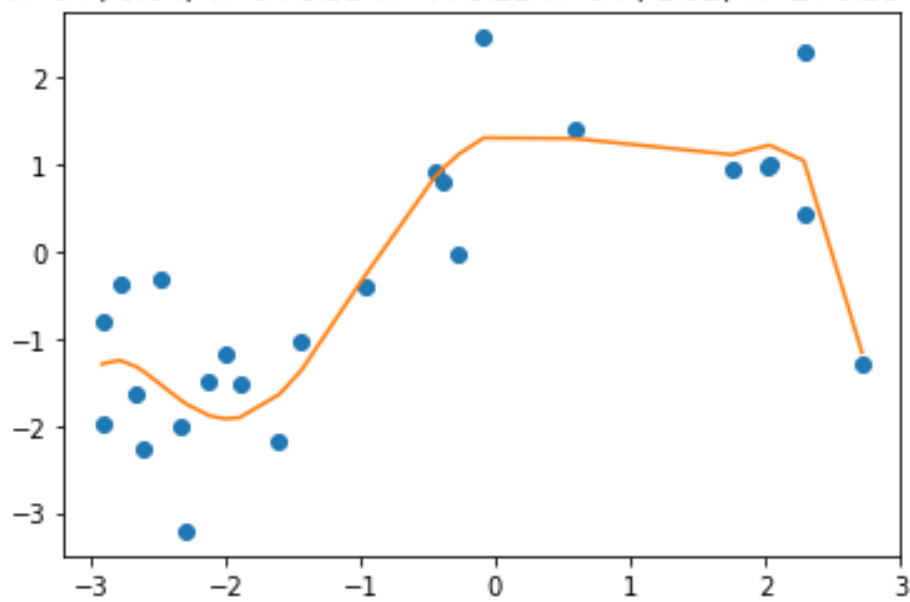0.03*x^6+(-0.02)*x^5+(-0.25)*x^4+0.06*x^3+0.34*x^2+0.58*x^1+0.52

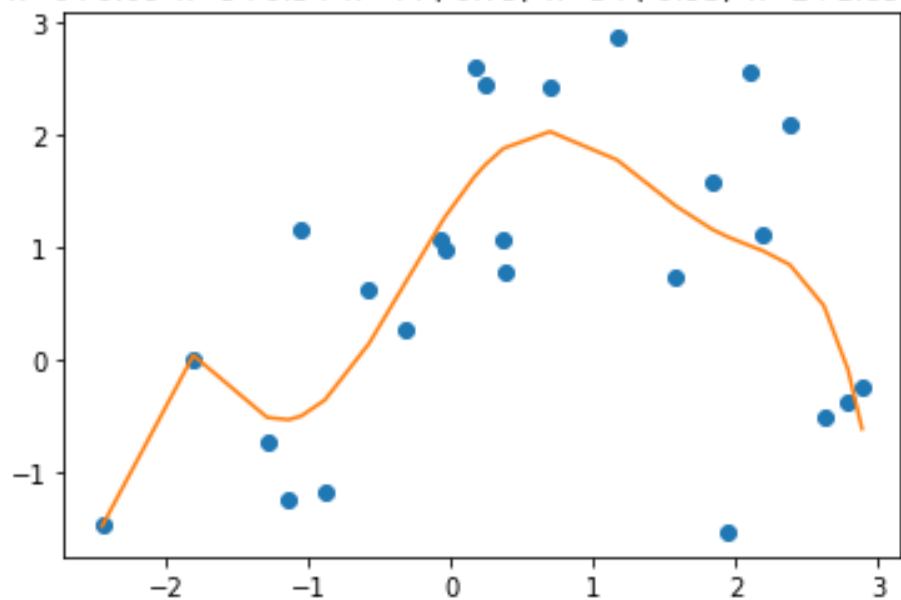0.0*x^6+(-0.05)*x^5+0.05*x^4+0.25*x^3+(-0.81)*x^2+0.34*x^1+1.53
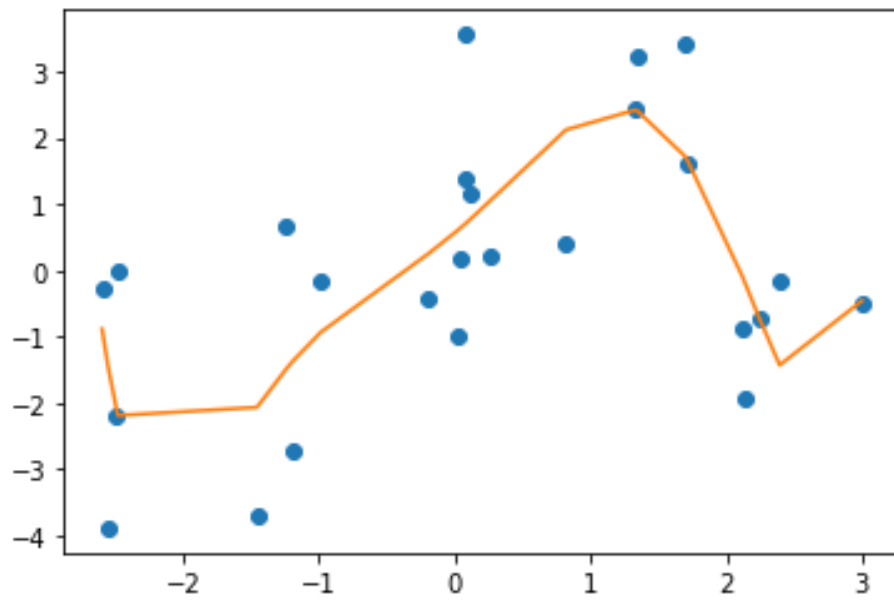
0.02*x^6+0.0*x^5+(-0.24)*x^4+(-0.17)*x^3+0.41*x^2+1.43*x^1+0.69

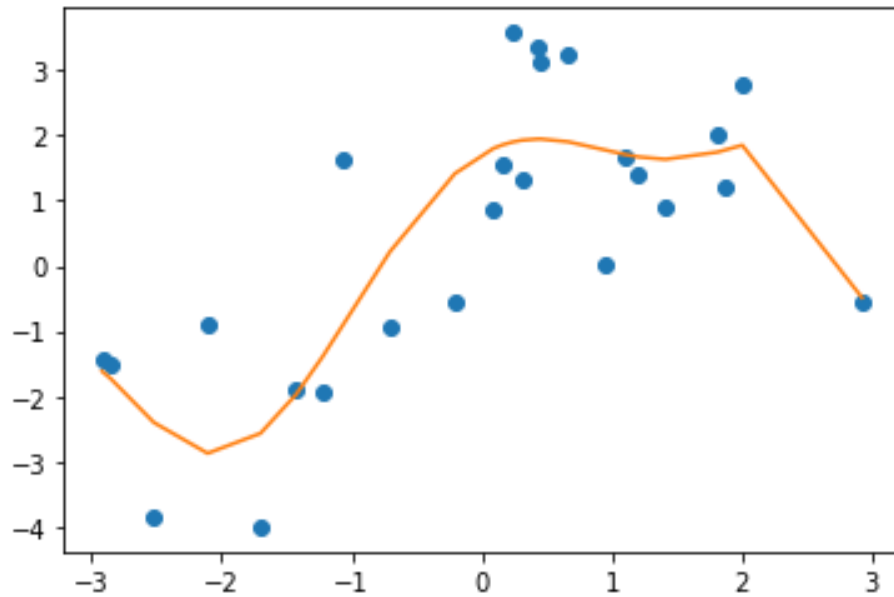(-0.03)*x^6+(-0.04)*x^5+0.33*x^4+0.22*x^3+(-1.32)*x^2+0.53*x^1+1.36

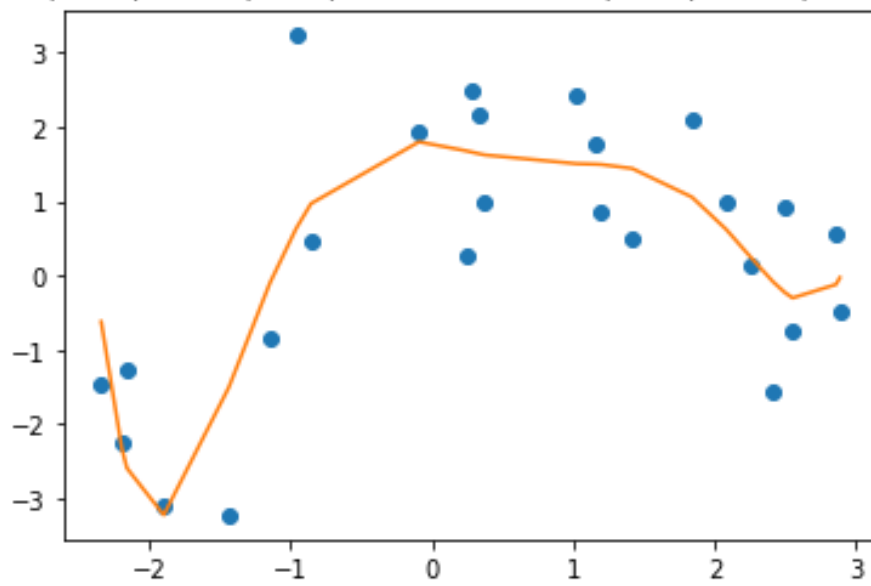(-0.04)*x^6+0.09*x^5+0.34*x^4+(-0.75)*x^3+(-0.93)*x^2+1.85*x^1+1.35

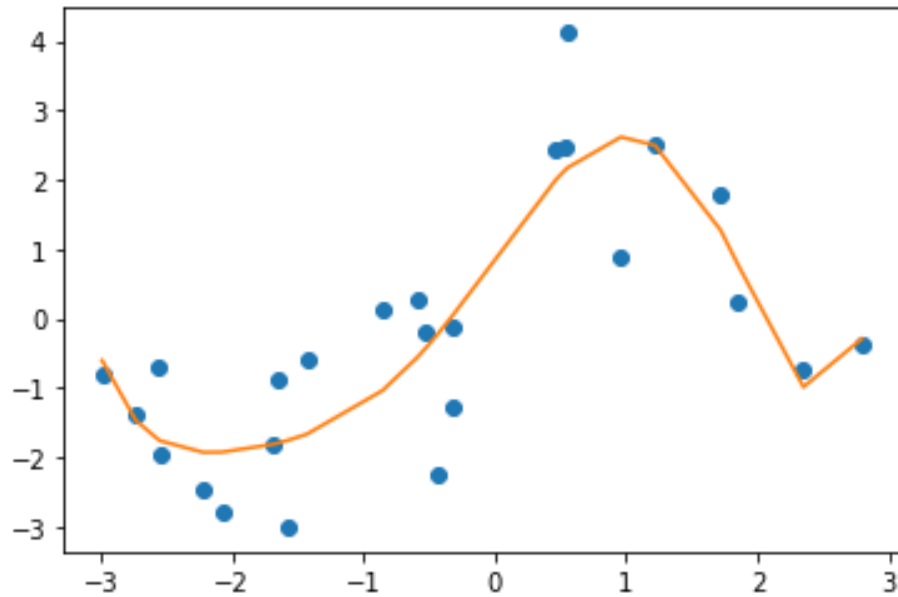0.05*x^6+(-0.04)*x^5+(-0.47)*x^4+(-0.03)*x^3+0.55*x^2+1.72*x^1+0.57

(-0.02)*x^6+(-0.03)*x^5+0.31*x^4+0.12*x^3+(-1.47)*x^2+1.13*x^1+1.71

0.05*x^6+(-0.17)*x^5+(-0.26)*x^4+0.99*x^3+(-0.57)*x^2+(-0.32)*x^1+1.77



0.02*x^6+0.06*x^5+(-0.19)*x^4+(-0.75)*x^3+(-0.01)*x^2+2.63*x^1+0.87



Finally, below you can find the average mean squared error values for each degree.
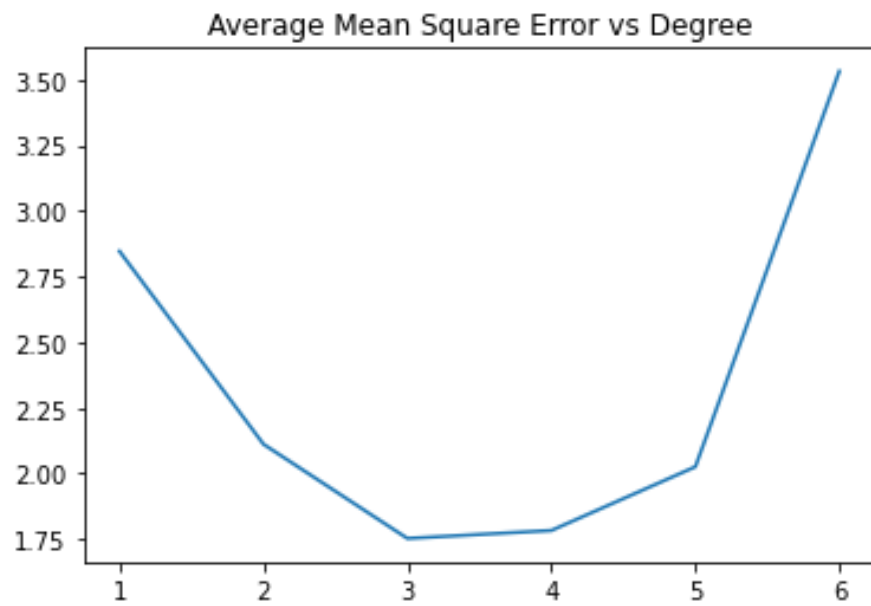
Figure [10]: Average Mean Square Error of Models with Different Degrees

As seen in Figure 10, the lowest error rate was found in the 3rd degree. This shows that 3rd degree models should be used in model selection. In addition, the hidden f(x) function that we mentioned at the beginning is probably a 3rd degree function.