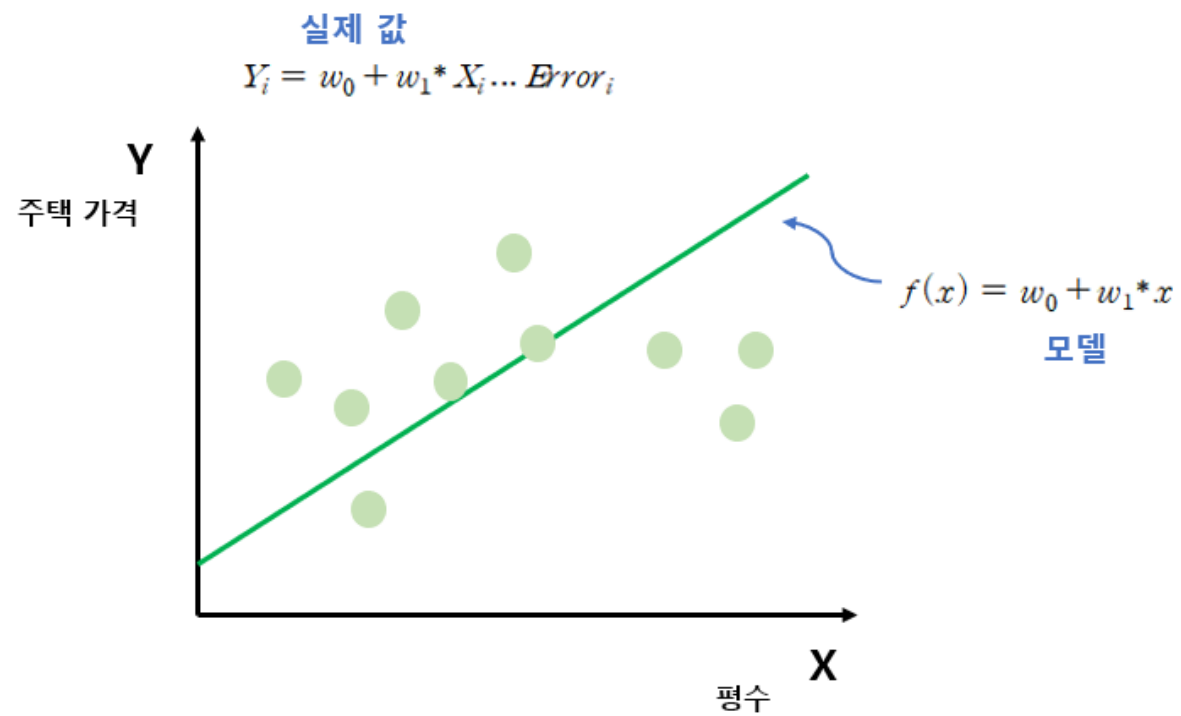
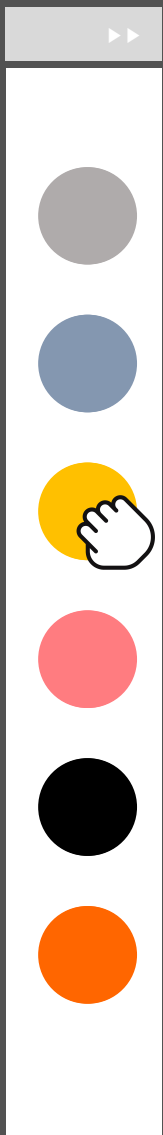



다항회귀

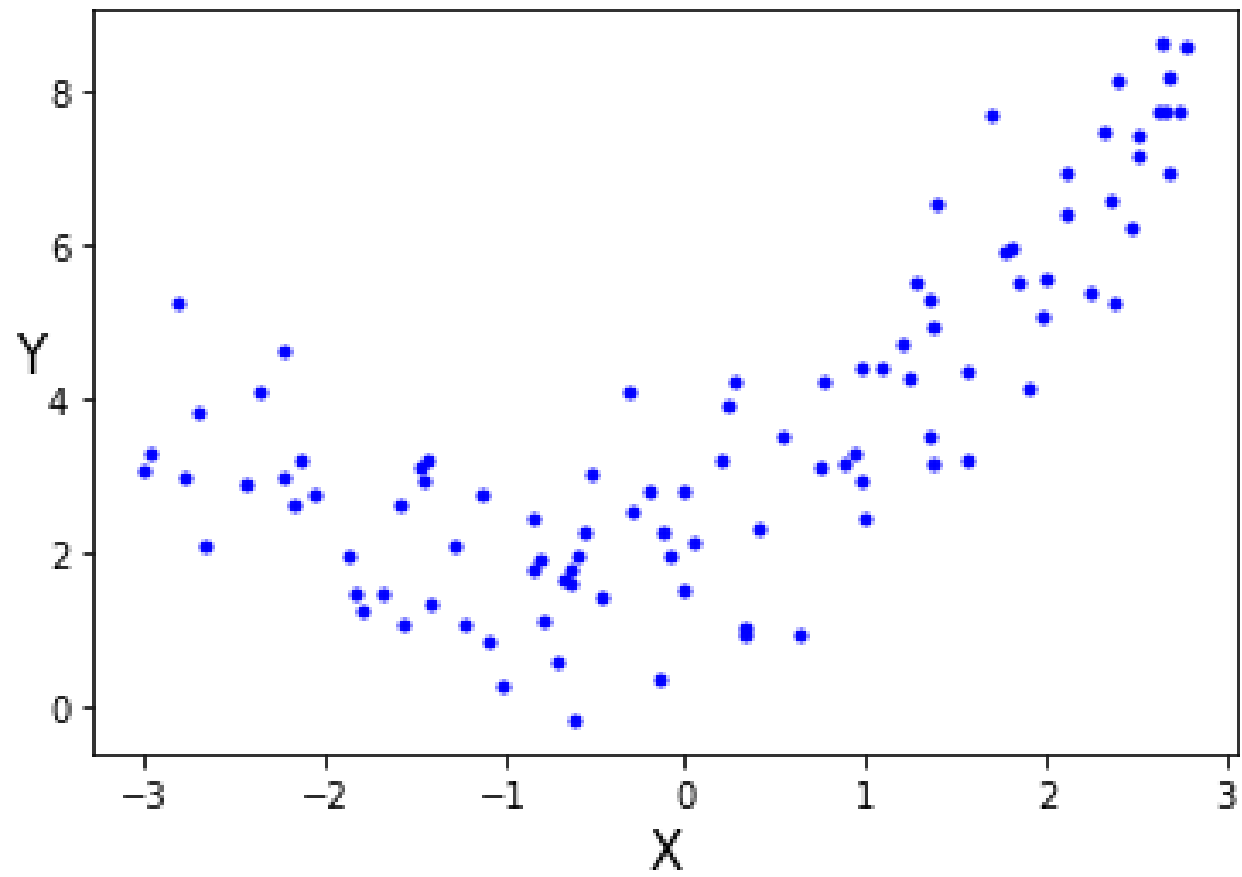
목차

1. 다항회귀
2. 다항회귀 Overfitting




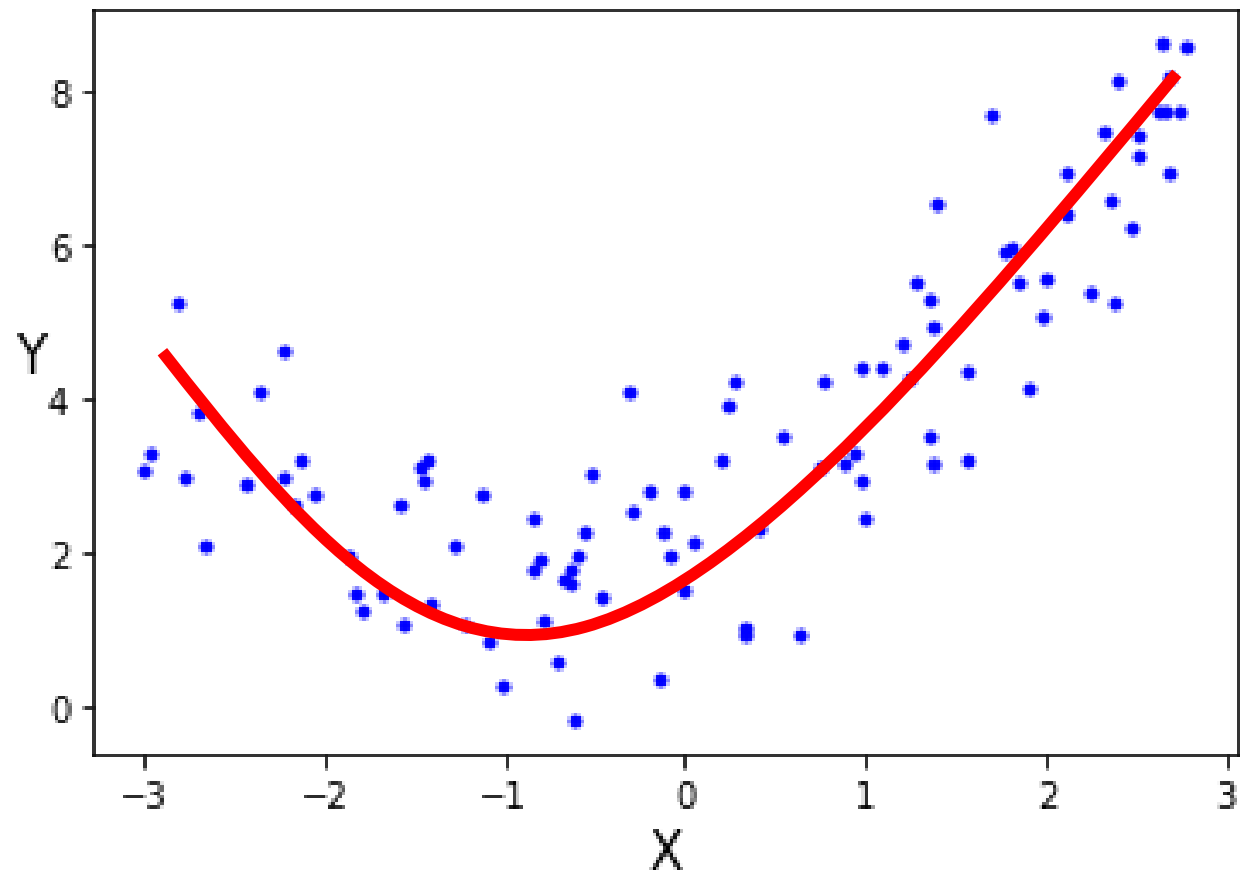
▶▶





▶▶





다항회귀

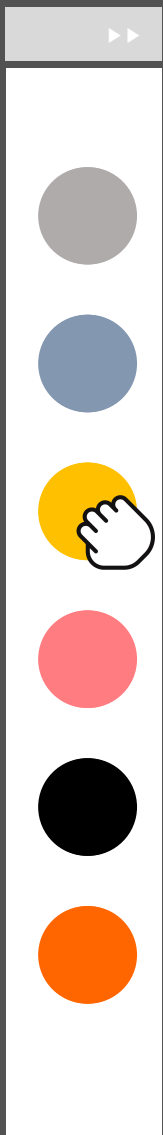
$$Y = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4x_1^2 + w_5x_2^2$$

(w_n 은 상수, x_n 은 변수)



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 from scipy import stats
6 from sklearn.datasets import load_boston
7 %matplotlib inline
8
9 boston = load_boston()
10 bostonDF = pd.DataFrame(boston.data , columns = boston.feature_names)
11 bostonDF['PRICE'] = boston.target
12 print('Boston 데이터셋 크기 :', bostonDF.shape)
13 bostonDF.head()
```

Boston 데이터셋 크기 : (506, 14)

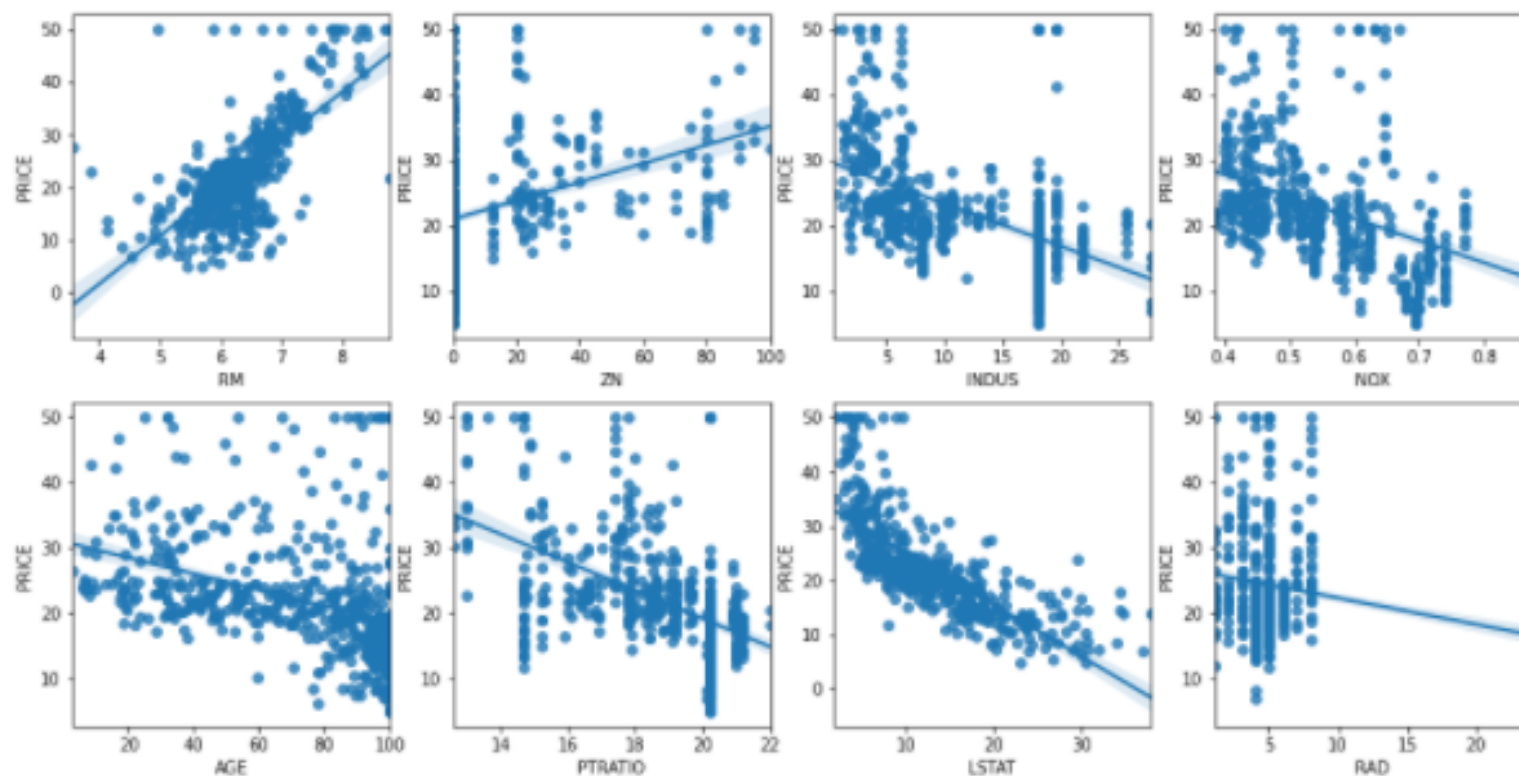


	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LST
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5

- CRIM: 지역별 범죄 발생률
- ZN: 25,000평방피트를 초과하는 거주 지역의 비율
- NDUS: 비상업 지역 넓이 비율
- CHAS: 찰스강에 대한 더미 변수(강의 경계에 위치한 경우는 1, 아니면 0)
- NOX: 일산화질소 농도
- RM: 거주할 수 있는 방 개수
- AGE: 1940년 이전에 건축된 소유 주택의 비율
- DIS: 5개 주요 고용센터까지의 가중 거리
- RAD: 고속도로 접근 용이도
- TAX: 10,000달러당 재산세율
- PTRATIO: 지역의 교사와 학생 수 비율
- B: 지역의 흑인 거주 비율
- LSTAT: 하위 계층의 비율
- MEDV: 본인 소유의 주택 가격(중앙값)

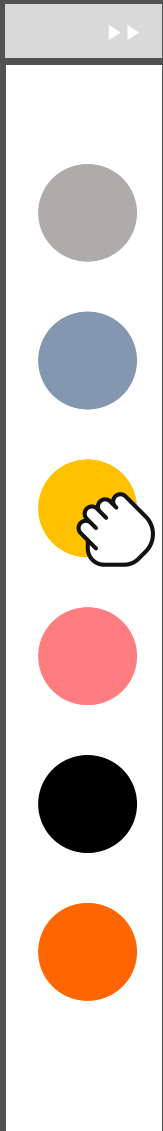


```
1 fig, axs = plt.subplots(figsize=(16,8) , ncols=4 , nrows=2)
2 lm_features = ['RM', 'ZN', 'INDUS', 'NOX', 'AGE', 'PTRATIO', 'LSTAT', 'RAD']
3 for i , feature in enumerate(lm_features):
4     row = int(i/4)
5     col = i%4
6     sns.regplot(x=feature , y='PRICE', data=bostonDF , ax=axs[row][col])
```





```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error , r2_score
4
5 y_target = bostonDF['PRICE']
6 X_data = bostonDF.drop(['PRICE'],axis=1,inplace=False)
7
8 X_train , X_test , y_train , y_test = train_test_split(X_data , y_target ,test_si
9 |
10 lr = LinearRegression()
11 lr.fit(X_train ,y_train )
12
13
```



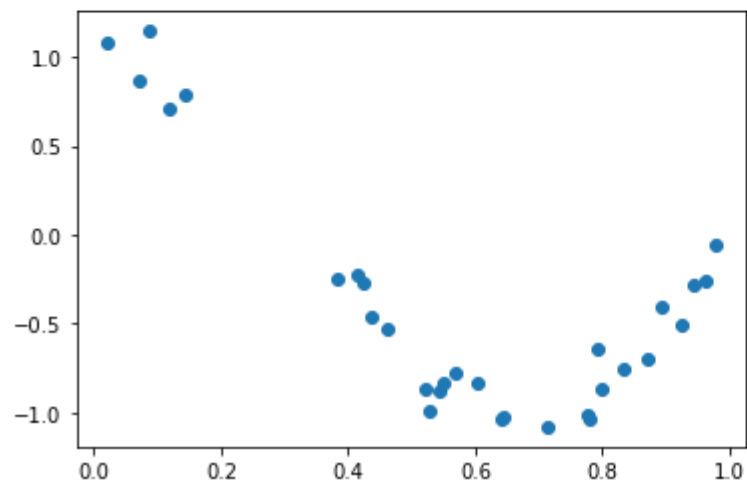
```
1 |coeff = pd.Series(data=np.round(lr.coef_, 1), index=X_data.columns )  
2 |coeff.sort_values(ascending=False)
```

```
RM          3.4  
CHAS        3.0  
RAD         0.4  
ZN          0.1  
B           0.0  
TAX        -0.0  
AGE         0.0  
INDUS       0.0  
CRIM       -0.1  
LSTAT      -0.6  
PTRATIO    -0.9  
DIS        -1.7  
NOX       -19.8  
dtype: float64
```



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.pipeline import Pipeline
4 from sklearn.preprocessing import PolynomialFeatures
5 from sklearn.linear_model import LinearRegression
6 from sklearn.model_selection import cross_val_score
7 %matplotlib inline
8
9 def true_fun(X):
10     return np.cos(1.5 * np.pi * X)
11
12 np.random.seed(0)
13 n_samples = 30
14 X = np.sort(np.random.rand(n_samples))
15 y = true_fun(X) + np.random.randn(n_samples) * 0.1
16 plt.scatter(X, y)
```


<matplotlib.collections.PathCollection at 0x187ae2cdd30>



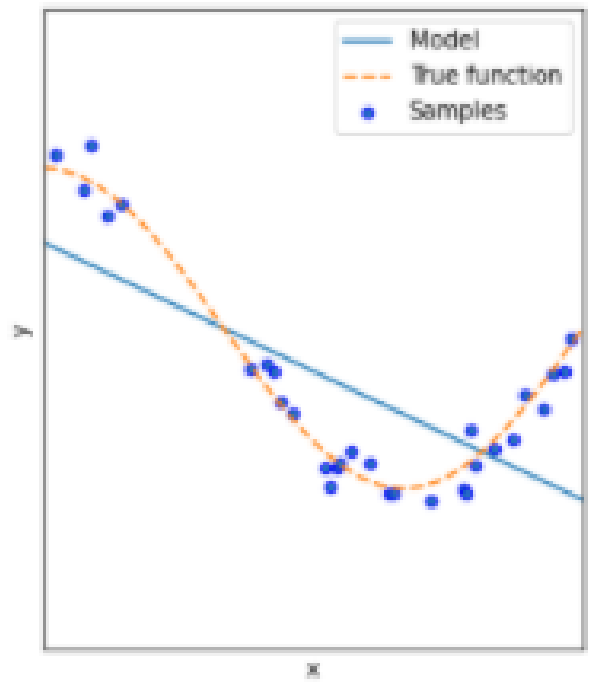


```
1 plt.figure(figsize = (14,5))
2 degrees = [1,4,15]
3
4 for i in range(len(degrees)):
5     ax = plt.subplot(1, len(degrees), i + 1)
6     plt.setp(ax, xticks=(), yticks=())
7
8     polynomial_features = PolynomialFeatures(degree = degrees[i]
9                                             , include_bias = False)
10    linear_regression = LinearRegression()
11    pipeline = Pipeline([("polynomial_features", polynomial_features),
12                        ("linear_regression", linear_regression)])
13    pipeline.fit(X.reshape(-1, 1), y)
14
15    X_test = np.linspace(0, 1, 100)
16
17    plt.plot(X_test, pipeline.predict(X_test[:, np.newaxis]), label="Model")
18
19    plt.plot(X_test, true_fun(X_test), '--', label="True function")
20    plt.scatter(X, y, edgecolor='b', s=20, label="Samples")
21
22    plt.xlabel("x"); plt.ylabel("y"); plt.xlim((0,1)); plt.ylim((-2,2)); plt.legend()
23    plt.title("Degree {} \nMSE = {:.2e} (+/- {:.2e})".format(degrees[i], -scores.me
24
```

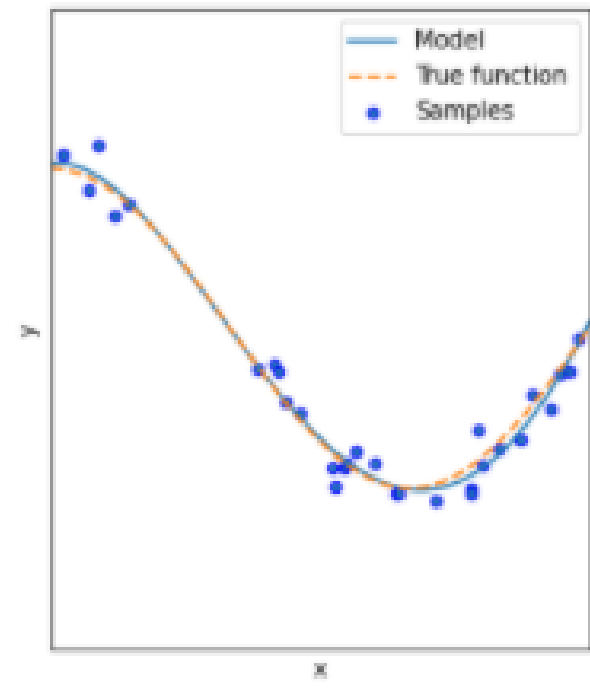
▶▶



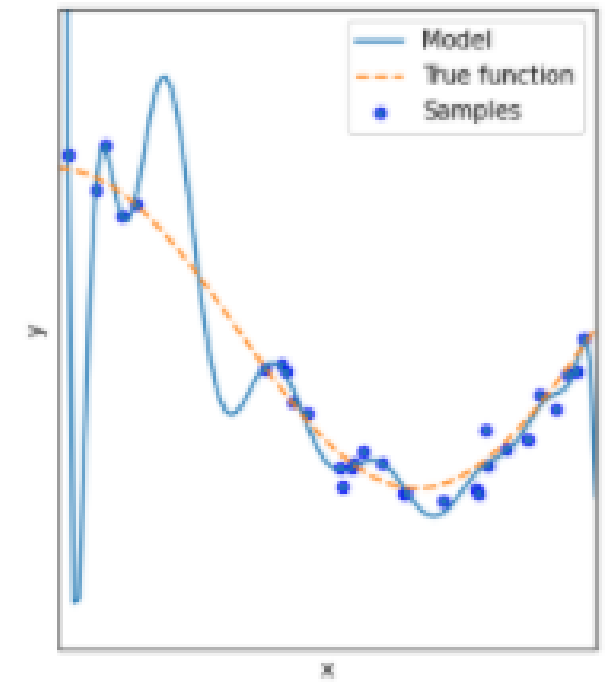
Degree 1



Degree 4



Degree 15



Any Question??

감사합니다