

TP 2 : underfitting, overfitting

Definition 1

Your model is **underfitting** the training data when the model performs poorly on the training data.

Your model is **overfitting** your training data when you see that the model performs well on the training data but does not perform well on the evaluation data

Generalization is the ability of a model to make correct predictions on new data, which were not used to build it.

Example 1 The data set is generated as follows. The sample data $(x_i, y_i)_{i \in \{1, \dots, n\}}$ is a realization of the random variables (X, Y)

- The random variable X follows a uniform distribution over $[0, 1]$ ($X \sim \mathcal{U}[0, 1]$)
- The noise $\epsilon \in \mathbb{R}$ is independent from X and $\epsilon \sim \mathcal{N}(0, 1)$
- $Y = \cos(\frac{3\pi}{2}X) + \epsilon$

The data set is splitted into training data (X_{train}, y_{train}) and test data (X_{test}, y_{test}) (see Figure 1). For any $d \in \{1, \dots, 19\}$, we choose $\mathcal{F}_d = \{f_\theta : x \rightarrow \varphi_d(x).\theta, \theta \in \mathbb{R}^{d+1}\}$ as approximate set where for any $x \in \mathcal{X} = \mathbb{R}$, $\varphi_d(x) = (1, x, x^2, \dots, x^d) \in \mathbb{R}^{d+1}$.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 !pip install scikit-learn if you use JupyterHub
4 from sklearn.pipeline import Pipeline
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.linear_model import LinearRegression
7 from sklearn.model_selection import train_test_split
8 from numpy import linalg as LA
```

```
1 def true_fun(X):
2     return np.cos(1.5 * np.pi * X)
3
4 np.random.seed(0)
5
6 n_samples = 30
7 degrees = [1, 3, 20]
8 X = np.sort(np.random.rand(n_samples))
9 y = true_fun(X) + np.random.randn(n_samples) * 0.1
10
11 X_train, X_test, y_train, y_test=train_test_split(X, y, test_size=0.20,
12                                                    random_state=42)
```

```

1 plt.figure
2 X_true = np.linspace(0, 1, 100)
3 plt.scatter(X_train, y_train, edgecolor="b", s=20, label="Training data")
4 plt.scatter(X_test, y_test, edgecolor="r", s=20, label="Test data")
5 plt.plot(X_true, true_fun(X_true), color='orange', label="True function")
6 plt.legend(loc="best")
7 plt.xlabel("x")
8 plt.ylabel("y")
9 plt.xlim((0, 1))
10 plt.ylim((-2, 2))
11 plt.show()

```

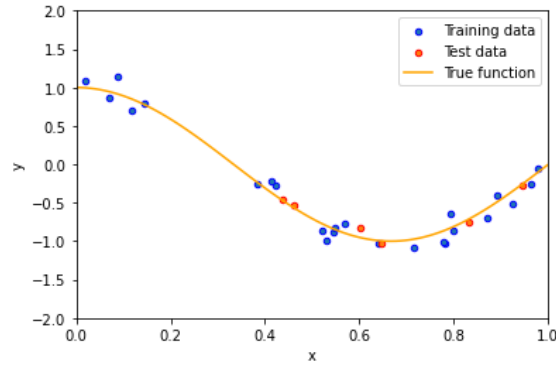


Figure 1: Splitting data

For any degree d , a model $f_{\theta_d^*}$ is created only using (X_{train}, y_{train}) . More precisely,

$$\theta_d^* = \underset{\theta}{\operatorname{argmin}} L(\theta) := \frac{1}{n} \sum_{i=0}^{n-1} (y_{train}[i] - f_{\theta}(X_{train}[i]))^2$$

where n is the length of X_{train} . Then the training Mean Squared Error (briefly MSE) and the test MSE is computed as follows :

$$\text{Training MSE} = \frac{1}{n} \sum_{i=0}^{n-1} (y_{train}[i] - f_{\theta_d^*}(X_{train}[i]))^2$$

$$\text{Test MSE} = \frac{1}{p} \sum_{i=0}^{p-1} (y_{test}[i] - f_{\theta_d^*}(X_{test}[i]))^2$$

where n is the length of X_{train} and p is the length of X_{test}

```

1 Test_MSE=[]
2 Training_MSE=[]
3 for i in range(1,20) :
4
5     polynomial_features = PolynomialFeatures(degree=i, include_bias=False)
6     linear_regression = LinearRegression()
7     pipeline = Pipeline([("polynomial_features", polynomial_features),
8                           ("linear_regression", linear_regression), ])

```

```

9 )
10 model=pipeline.fit(X_train.reshape(-1, 1), y_train.reshape(-1, 1))
11 model.predict(X_test.reshape(-1, 1))
12 Test_MSE.append(LA.norm(model.predict(X_test.reshape(-1, 1))-y_test.reshape
13 (-1,1),2))
    Training_MSE.append(LA.norm(model.predict(X_train.reshape(-1, 1))-y_train.
        reshape(-1,1),2))

```

```

1 plt.figure
2 plt.plot(range(1,20),Test_MSE)
3 plt.plot(range(1,20),Training_MSE)
4 plt.xlabel("d")
5 plt.ylabel("MSE")
6 plt.xticks(range(1,20,3))
7 plt.legend(loc="best", labels=['Test MSE', 'Training MSE'])
8 plt.show()

```

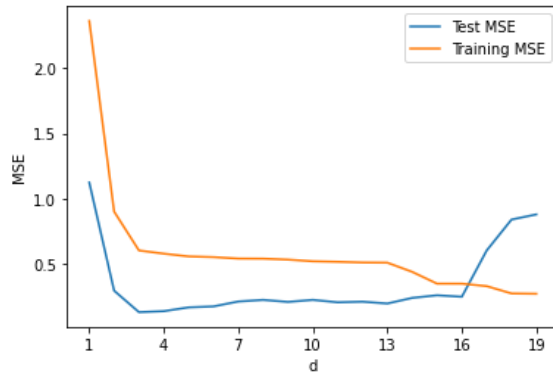


Figure 2: MSE (Mean Squared Error) with respect to the size of \mathcal{F}_d

In Figure 2, the training MSE and the test MSE is plotted with respect to the degree d . Since $\mathcal{F}_1 \subset \dots \subset \mathcal{F}_{19}$, the 'size' of \mathcal{F}_d inscreases with respect to the degree d .

```

1 plt.figure(figsize=(14, 5))
2 for i in range(len(degrees)):
3     ax = plt.subplot(1, len(degrees), i + 1)
4     plt.setp(ax, xticks=(), yticks=())
5
6     polynomial_features = PolynomialFeatures(degree=degrees[i], include_bias=
7     False)
8     linear_regression = LinearRegression()
9     pipeline = Pipeline(
10         [
11             ("polynomial_features", polynomial_features),
12             ("linear_regression", linear_regression),
13         ]
14     )
15     model=pipeline.fit(X_train.reshape(-1, 1), y_train.reshape(-1, 1))
16     model.predict(X_test.reshape(-1, 1))
17     # Evaluate the models using crossvalidation

```

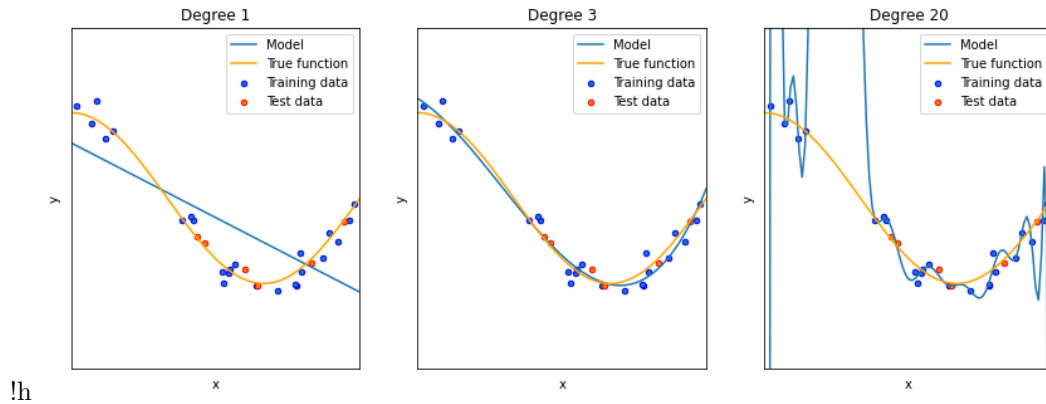


Figure 3: Polynomial model using `sklearn`

```

18 X_true = np.linspace(0, 1, 100)
19 plt.plot(X_true, pipeline.predict(X_true[:, np.newaxis]), label="Model")
20 plt.plot(X_true, true_fun(X_true), color='orange', label="True function")
21 plt.scatter(X_train, y_train, edgecolor="b", s=20, label="Training data")
22 plt.scatter(X_test, y_test, edgecolor="r", s=20, label="Test data")
23 plt.xlabel("x")
24 plt.ylabel("y")
25 plt.xlim((0, 1))
26 plt.ylim((-2, 2))
27 plt.legend(loc="best")
28 plt.title(
29     "Degree {}".format(
30         degrees[i]
31     )
32 )
33 plt.show()
34

```

- When $d = 1$ the model is **underfitting**; the model has “not learned enough” from the training data, resulting in low generalization and unreliable predictions. More precisely, it is not flexible enough to take into account the structure of the data set (see on the left side of Figure 3)
- When d is large enough ($d \geq 16$), the model is **overfitting**; the model learning “too much” from the training data set. More precisely, a model learns the details and noises in the training data that negatively impacts the performance of the model on new data (see on the right side of Figure 3).
- Ideally, you want to select a model between underfitting and overfitting. For example $d = 3$ (see Figure 3).