

1. Recall the definition of the expected risk $R(f)$ of a predictor f . While this cannot be computed in general note that here we defined $P_{X \times Y}$. Which function f^* is an obvious Bayes predictor? Make sure to explain why the risk $R(f^*)$ is minimum at f^* .

Answer:

$f^* = a_0 + a_1x + a_2x^2$ is an obvious Bayes predictor

In this case, the values of x and y are linked as $y = g(x) = a_0 + a_1x + a_2x^2$, which fixes the joint distribution $P_{X \times Y}$

Hence, for all $x \in [0,1]$. $\ell(\hat{y}, y) = \frac{1}{2}(f^*(x) - g(x))^2 = 0$
 $R(f^*) = E[\ell(f^*(x), g(x))] = 0$
i.e. $R(f^*)$ is minimum at f^*

2. Using \mathcal{H}_2 as your hypothesis class, which function $f_{\mathcal{H}_2}^*$ is a risk minimizer in \mathcal{H}_2 ? Recall the definition of the approximation error. What is the approximation error achieved by $f_{\mathcal{H}_2}^*$?

Answer:

Using \mathcal{H}_2 as a hypothesis set, $f_{\mathcal{H}_2}^* = a_0 + a_1x + a_2x^2$ is a risk minimizer in \mathcal{H}_2 .

This can be explained by the fact that the true function $g(x)$ fits into the hypothesis class. The approximation is actually zero achieved by $f_{\mathcal{H}_2}^*$

i.e. $R(f_{\mathcal{H}_2}^*) - R(f^*) = 0$

3. Considering now \mathcal{H}_d , with $d > 2$. Justify an inequality between $R(f_{\mathcal{H}_2}^*)$ and $R(f_{\mathcal{H}_d}^*)$. Which function $f_{\mathcal{H}_d}^*$ is a risk minimizer in \mathcal{H}_d ? What is the approximation error achieved by $f_{\mathcal{H}_d}^*$?

Answer:

When $d > 2$, the hypothesis set tends to be larger, which introduces more chances for the hypothesis set to cover the Bayesian predictor.

$$\text{Hence, } R(f_{\mathcal{H}_2}^*) \geq R(f_{\mathcal{H}_d}^*) \geq 0.$$

$$\text{It's known that } R(f_{\mathcal{H}_2}^*) = 0. \text{ thus } R(f_{\mathcal{H}_d}^*) = 0$$

$$\text{We can deduce that: } f_{\mathcal{H}_d}^* = a_0 + a_1 x + a_2 x^2$$

$$\text{The approximation error: } R(f_{\mathcal{H}_d}^*) - R(f^*) = 0$$

4. For this question we assume $a_0 = 0$. Considering $\mathcal{H} = \{f : x \rightarrow b_1 x; b_1 \in \mathbb{R}\}$, which function $f_{\mathcal{H}}^*$ is a risk minimizer in \mathcal{H} ? What is the approximation error achieved by $f_{\mathcal{H}}^*$? In particular what is the approximation error achieved if furthermore $a_2 = 0$ in the definition of true underlying relation $g(x)$ above?

Answer:

$$b_1^* \in \underset{b_1}{\operatorname{argmin}} E[l(f_{\mathcal{H}}^*(x), g(x))]$$

$$b_1^* \in \underset{b_1}{\operatorname{argmin}} \frac{1}{2} \int_0^1 (b_1 x - a_1 x - a_2 x^2)^2 dx$$

$$b_1^* \in \underset{b_1}{\operatorname{argmin}} \frac{1}{2} \int_0^1 (b_1 - a_1)^2 x^2 - 2(a_2 b_1 - a_1 a_2) x^3 + a_2^2 x^4 dx$$

$$b_1^* \in \underset{b_1}{\operatorname{argmin}} \frac{1}{2} \left[\frac{1}{3} (b_1 - a_1)^2 x^3 - \frac{1}{2} (a_2 b_1 - a_1 a_2) x^4 + \frac{1}{5} a_2^2 x^5 \Big|_0^1 \right]$$

$$b_1^* \in \underset{b_1}{\operatorname{argmin}} \frac{1}{2} \left[\frac{1}{3} (b_1^2 - 2a_1 b_1 + a_1^2) - \frac{1}{2} (a_2 b_1 - a_1 a_2) + \frac{1}{5} a_2^2 \right]$$

$$\text{let } h(b) = \frac{1}{3} b_1^2 - \frac{2}{3} a_1 b_1 - \frac{1}{2} a_2 b_1 + (\#) \quad \frac{\partial h(b)}{\partial b} = \frac{2}{3} b_1 - \frac{2}{3} a_1 - \frac{1}{2} a_2 = 0$$

$$\text{Hence: } b_1^* = a_1 + \frac{3}{4} a_2 \quad f_{\mathcal{H}}^* = b_1^* x = (a_1 + \frac{3}{4} a_2) x$$

$$R(f_{\mathcal{H}}^*) - R(f^*) = \int_0^1 \frac{1}{2} \left(\frac{3}{4} a_2 x - a_2 x^2 \right)^2 dx = \frac{1}{160} a_2^2$$

$$\text{Furthermore, if } a_2 = 0. \quad f_{\mathcal{H}}^* = a_1 x,$$

$$\text{the approximation error } R(f_{\mathcal{H}}^*) - R(f^*) = 0.$$

5. Show that the empirical risk minimizer (ERM) \hat{b} is given by the following minimization

$$\hat{b} = \arg \min_b \|Xb - y\|_2^2.$$

Answer:

$$\begin{aligned}\hat{b} &= \arg \min_b \|Xb - y\|_2^2 = \arg \min_b (Xb - y)^T (Xb - y) \\&= \arg \min_b \begin{pmatrix} b_0 + b_1 x_1 + \dots + b_d x_1^d - y_0 \\ b_0 + b_1 x_2 + \dots + b_d x_2^d - y_1 \\ \vdots \\ b_0 + b_1 x_N + \dots + b_d x_N^d - y_d \end{pmatrix}^T \begin{pmatrix} b_0 + b_1 x_1 + \dots + b_d x_1^d - y_0 \\ b_0 + b_1 x_2 + \dots + b_d x_2^d - y_1 \\ \vdots \\ b_0 + b_1 x_N + \dots + b_d x_N^d - y_d \end{pmatrix} \\&= \arg \min_b \sum_{i=1}^N (b_0 + b_1 x_i + \dots + b_d x_i^d - y_{N-i})^2 \\&= \arg \min_b \sum_{i=1}^N (Xb - y)_i^2 \quad \text{where } (Xb - y)_i \text{ denotes} \\&\quad \text{the } i\text{th row of } (Xb - y) \\&= \arg \min_b \sum_{i=1}^N (f_b(x_i) - y_{N-i})^2\end{aligned}$$

Consider N is a constant, we can rewrite this equation as:

$$\hat{b} = \arg \min_b \sum_{i=1}^N \frac{1}{N} (f_b(x_i) - y_{N-i})^2 \quad \textcircled{*}$$

$\textcircled{*}$ can be regarded as the empirical minimizer given the loss function: $l(x) = (f(x) - y)^2$

Hence, "the empirical risk minimizer (ERM) is given by the following minimization $\hat{b} = \arg \min_b \|Xb - y\|_2^2$ " is proved.

6. If $N > d$ and X is full rank, show that $\hat{b} = (X^\top X)^{-1} X^\top \mathbf{y}$. (Hint: you should take the gradients of the loss above with respect to b). Why do we need to use the conditions $N > d$ and X full rank?

Answer:

According to the previous proof in 5.

$$\text{we get that } \hat{b} = \underset{b}{\operatorname{argmin}} \|Xb - \mathbf{y}\|_2^2$$

$$= \underset{b}{\operatorname{argmin}} (Xb - \mathbf{y})^\top (Xb - \mathbf{y})$$

$$= \underset{b}{\operatorname{argmin}} (b^\top X^\top X b - \mathbf{y}^\top X b - b^\top X \mathbf{y} + \mathbf{y}^\top \mathbf{y})$$

$$\text{let } f(b) = b^\top X^\top X b - \mathbf{y}^\top X b - b^\top X \mathbf{y} + \mathbf{y}^\top \mathbf{y}$$

$$\frac{\partial f(b)}{\partial b} = (X^\top X + X^\top X)b - X^\top \mathbf{y} - X^\top \mathbf{y} = 2X^\top X b - 2X^\top \mathbf{y} = 0$$

$$\text{let } g(b) = 2X^\top X b - 2X^\top \mathbf{y}; \quad \frac{\partial g(b)}{\partial b} = 2X^\top X \geq 0$$

$$\text{Hence } X^\top X b = X^\top \mathbf{y}$$

With $N > d$ and X being full rank, we can obtain that $X^\top X$ is invertible. thus: $\hat{b} = (X^\top X)^{-1} X^\top \mathbf{y}$.

To state why we need to use the condition:

$$\left\{ \begin{array}{l} N > d \\ X \text{ being full rank} \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \operatorname{rank}(X) = d \\ X \text{ is invertible} \end{array} \right. \Rightarrow X^\top X \text{ is invertible}$$

Once these two restrictions are disobeyed:

then $\ker(X) \neq \{0\}$. then $\ker(X^\top X) \neq \{0\}$.

$X^\top X$ isn't invertible anymore.

7. Write a function called `least_square_estimator` taking as input a design matrix $X \in \mathbb{R}^{N \times (d+1)}$ and the corresponding vector $y \in \mathbb{R}^N$ returning $\hat{b} \in \mathbb{R}^{(d+1)}$. Your function should handle any value of N and d , and in particular return an error if $N \leq d$. (Drawing x at random from the uniform distribution makes it almost certain that any design matrix X with $d \geq 1$ we generate is full rank).

ANSWER TO Q_7

```

1 def least_square_estimator(x, y):
2     """
3     array: N*(d+1)
4     array: d+1
5     """
6     N, d_plus_1 = X.shape
7     if N > (d_plus_1-1):
8         return np.linalg.inv(X.T @ X) @ X.T @ y
9     else:
10        raise ValueError("N must be larger than d")

```

8. Recall the definition of the empirical risk $\hat{R}(\hat{f})$ on a sample $\{(x_i, y_i)\}_{i=1}^N$ for a prediction function \hat{f} . Write a function `empirical_risk` to compute the empirical risk of f_b taking as input a design matrix $X \in \mathbb{R}^{N \times (d+1)}$, a vector $y \in \mathbb{R}^N$ and the vector $b \in \mathbb{R}^{(d+1)}$ parametrizing the predictor.

ANSWER TO Q_8

```

1 def empirical_risk(X, y, b):
2     N, _ = X.shape
3     emp_risk = 1/(2*N) * np.linalg.norm(X @ b - y) **2
4     return emp_risk

```

9. Use your code to estimate $\hat{\mathbf{b}}$ from $\mathbf{x_train}$, $\mathbf{y_train}$ using $d = 5$. Compare $\hat{\mathbf{b}}$ and \mathbf{a} . Make a single plot (Plot 1) of the plan (x, y) displaying the points in the training set, values of the true underlying function $g(x)$ in $[0, 1]$ and values of the estimated function $f_{\hat{\mathbf{b}}}(x)$ in $[0, 1]$. Make sure to include a legend to your plot.

ANSWER TO Q_9

```

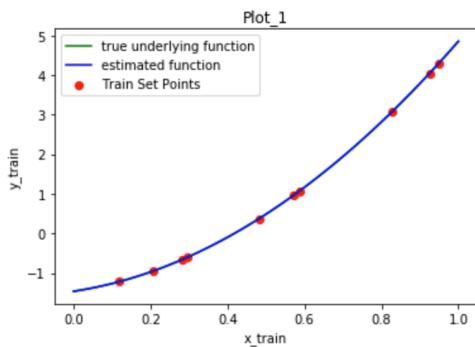
1 #X_train: N_train * (d+1)
2 #coeff: (d+1)*1
3 #y_train: N_train * 1
4 N_train, N_test = 10, 1000
5 true_d, esti_d = 2, 5
6 a = get_a(deg_true=true_d)
7
8 # Design X_train by changing true_d to esti_d
9 # With the esti_d, use estimated_b to represent the latent relationship between x_train & y_train
10 x_train, y_train = draw_sample(deg_true=true_d, a=a, N=N_train)
11 X_train = get_design_mat(x=x_train, deg=esti_d)
12 esti_b = least_square_estimator(X_train, y_train)
13 print(f'The generated a is {a}')
14 print(f'The estimated b is {esti_b}')
15 print('*'*100)
16 print('From the Plot_1, we can see that the scatter nearly perfect fits.')
17
18 gx = get_design_mat(np.linspace(0,1),2) @ a
19 fb = get_design_mat(np.linspace(0,1),esti_d) @ esti_b
20
21 # Plot the True underlying functions
22 # Plot the Estimated function
23 plt.scatter(x_train, y_train, label='Train Set Points', color='red')
24 plt.plot(np.linspace(0,1), gx, label='true underlying function', color='green')
25 plt.plot(np.linspace(0,1), fb, label='estimated function', color='blue')
26 plt.title('Plot_1')
27 plt.xlabel('x_train')
28 plt.ylabel('y_train')
29 plt.legend()
30 plt.show()

```

```

The generated a is [-1.46027496  1.476993   4.84103081]
The estimated b is [-1.46027496e+00  1.47699300e+00  4.84103081e+00  2.04789785e-09
 1.30665967e-09  3.17754711e-11]
*****
From the Plot_1, we can see that the scatter nearly perfect fits.

```



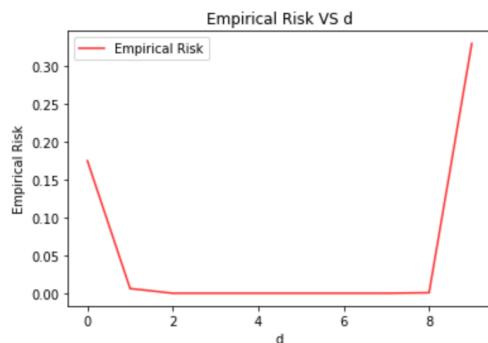
10. Now you can adjust d . What is the minimum value for which we get a “perfect fit”? How does this result relates with your conclusions on the approximation error above?

ANSWER TO Q_10

```

1 N_train, N_test = 10, 1000
2 emp_risk = []
3 for d in range(10):
4     true_d, esti_d = 2, d
5     x_train, y_train = draw_sample(deg_true=true_d, a=a, N=N_train)
6     X_train = get_design_mat(x=x_train, deg=esti_d)
7     esti_b = least_square_estimator(X_train, y_train)
8     emp_risk.append(empirical_risk(X_train, y_train, esti_b))
9
10 plt.plot(range(10), emp_risk, label='Empirical Risk', color='red', alpha=0.8)
11 plt.title('Empirical Risk VS d')
12 plt.xlabel('d')
13 plt.ylabel('Empirical Risk')
14 plt.legend()
15 plt.show()

```



```
1 emp_risk
```

```
[1.5030719653456406,
 0.04886300621441407,
 9.68326761158792e-28,
 1.9938577708596854e-25,
 1.143015372461334e-22,
 1.0966908714390253e-19,
 1.9132537079401906e-16,
 2.9029778107946162e-11,
 0.01016300044355156,
 8.651662401906895]
```

- The Minimum value that we get a 'perfect fit' is 2, with the empirical risk=9.68326761158792e-28.
- This result is the same as the conclusion on the approximation error above.
- Once the hypothesis space H_d has a degree d that is larger than the degree of the true underlying function, we could get the approximation error 0.

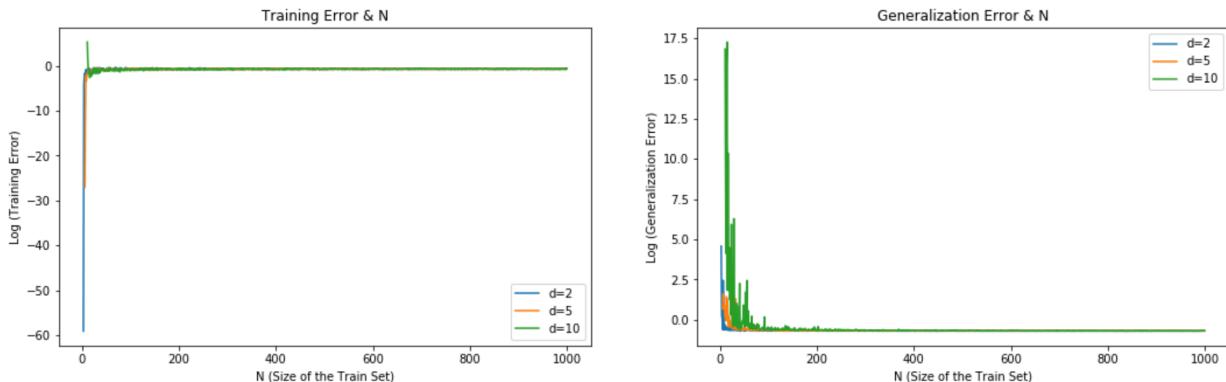
11. Plot e_t and e_g as a function of N for $d < N < 1000$ for $d = 2$, $d = 5$ and $d = 10$ (Plot 2). You may want to use a logarithmic scale in the plot. Include also plots similar to Plot 1 for 2 or 3 different values of N for each value of d .

ANSWER TO Q_11

```

1 x_test, y_test = draw_sample_with_noise(deg_true=true_d, a=a, N=N_test)
2
3 def get_et_eg(esti_d):
4     et = []
5     eg = []
6     for n in range(esti_d+1, 1000):
7         x_train, y_train = draw_sample_with_noise(deg_true=true_d, a=a, N=n)
8         X_train = get_design_mat(x_train, deg=esti_d)
9         X_test = get_design_mat(x_test, deg=esti_d)
10        esti_b = least_square_estimator(X_train, y_train)
11        et.append(empirical_risk(X_train, y_train, esti_b))
12        eg.append(empirical_risk(X_test, y_test, esti_b))
13    return et, eg
14
15
16 fig = plt.figure(figsize=(18,5))
17
18 ax1 = fig.add_subplot(1,2,1)
19 ax2 = fig.add_subplot(1,2,2)
20
21 for esti_d in [2, 5, 10]:
22     ax1.plot(range(esti_d+1,1000), np.log(get_et_eg(esti_d)[0]), label=f'd={esti_d}')
23     ax2.plot(range(esti_d+1,1000), np.log(get_et_eg(esti_d)[1]), label=f'd={esti_d}')
24
25 ax1.set_title('Training Error & N')
26 ax2.set_title('Generalization Error & N')
27 ax1.set_xlabel('N (Size of the Train Set)')
28 ax2.set_xlabel('N (Size of the Train Set)')
29 ax1.set_ylabel('Log (Training Error)')
30 ax2.set_ylabel('Log (Generalization Error)')
31 ax1.legend()
32 ax2.legend()
33
34 plt.show()
35

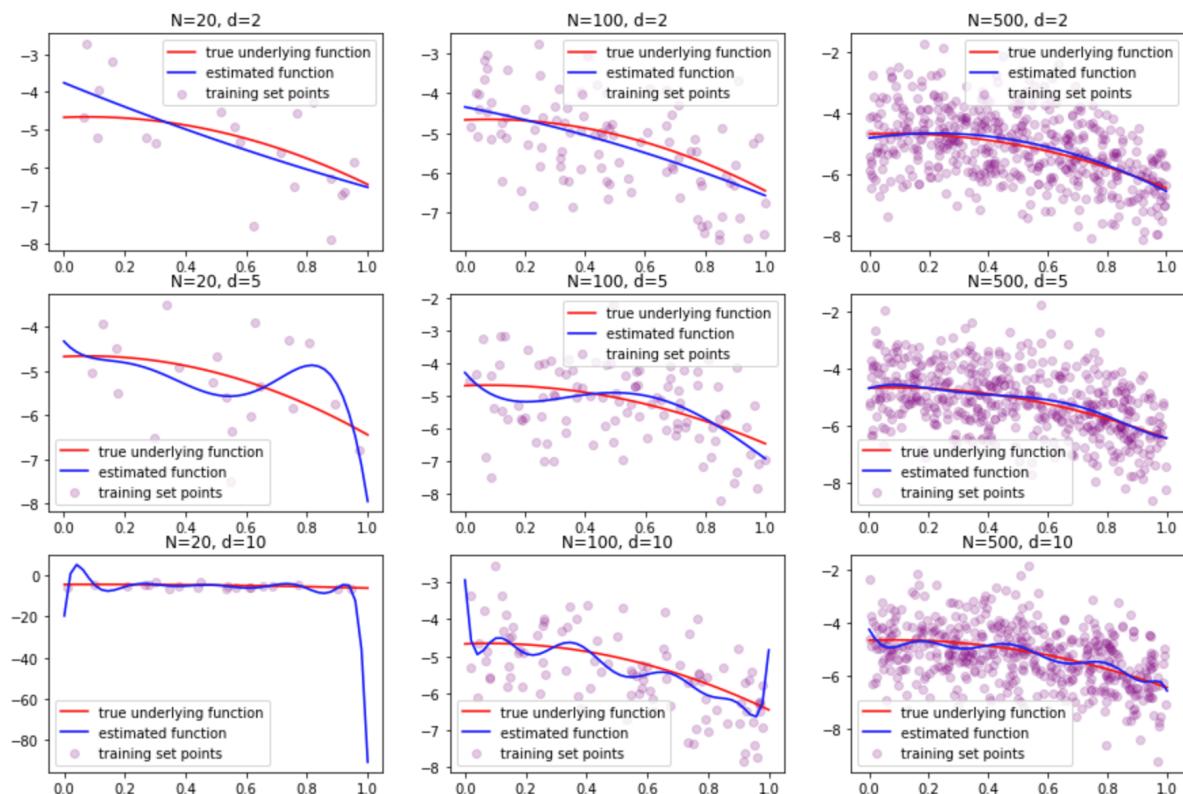
```



```

1 fig = plt.figure(figsize=(15,10))
2 iteration_times = 1
3 N_test = 1000
4
5 for esti_d in [2, 5, 10]:
6     for N_train in [20, 100, 500]:
7         x_train, y_train = draw_sample_with_noise(deg_true=true_d, a=a, N=N_train)
8         X_train = get_design_mat(x=x_train, deg=esti_d)
9         esti_b = least_square_estimator(X_train, y_train)
10
11     gx = get_design_mat(np.linspace(0,1), 2) @ a
12     fb = get_design_mat(np.linspace(0,1), esti_d) @ esti_b
13
14     ax = fig.add_subplot(3,3,iteration_times)
15     iteration_times += 1
16     ax.scatter(x_train, y_train, label='training set points', color='purple', alpha=0.2)
17     ax.plot(np.linspace(0,1), gx, label='true underlying function', color='red')
18     ax.plot(np.linspace(0,1), fb, label='estimated function', color='blue')
19     ax.set_title(f'N={N_train}, d={esti_d}')
20     ax.legend()

```



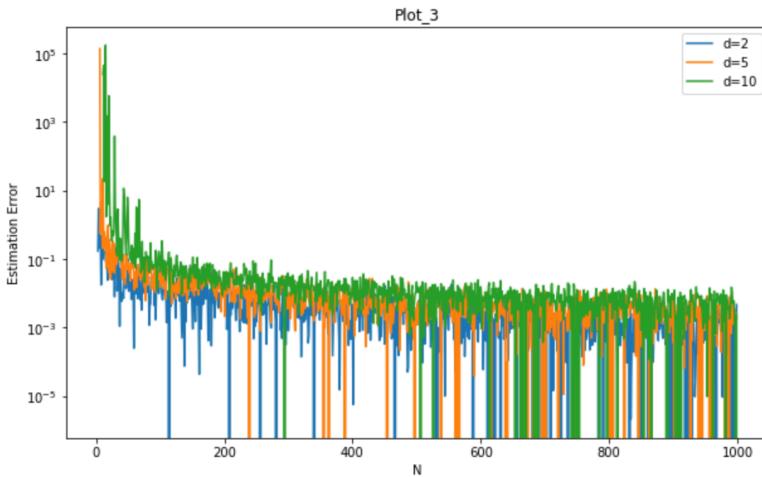
12. Recall the definition of the estimation error. Using the test set, (which we intentionally chose large so as to take advantage of the law of large numbers) give an empirical estimator of the estimation error. For the same values of N and d above plot the estimation error as a function of N (Plot 3).

ANSWER TO Q_12

Recall the definition of the estimation error. Using the test set, (which we intentionally chose large so as to take advantage of the law of large numbers) give an empirical estimator of the estimation error. For the same values of N and d above plot the estimation error as a function of N (Plot 3).

```

1 plt.figure(figsize=(10, 6))
2
3 x_test, y_test = draw_sample_with_noise(deg_true=true_d, a=a, N=N_test)
4 X_test_ = get_design_mat(x=x_test, deg=true_d)
5 risk = empirical_risk(X_test_, y_test, a)
6
7 def test_estimation_error(estimated_d):
8     esti_err = []
9     for N_train in range(estimated_d+1, 1000):
10         x_train, y_train = draw_sample_with_noise(deg_true=true_d, a=a, N=N_train)
11         X_train = get_design_mat(x=x_train, deg=estimated_d)
12         estimated_b = least_square_estimator(X_train, y_train)
13
14         X_test = get_design_mat(x=x_test, deg=estimated_d)
15
16         emp_risk = empirical_risk(X_test, y_test, estimated_b)
17         esti_err.append(emp_risk-risk)
18
19     return esti_err
20
21 for esti_d in [2, 5, 10]:
22     plt.plot(range(esti_d+1, 1000), test_estimation_error(esti_d), label=f'd={esti_d}')
23     plt.yscale('log')
24
25 plt.legend();
26 plt.title('Plot_3');
27 plt.xlabel('N');
28 plt.ylabel('Estimation Error');
```



13. The generalization error gives in practice an information related to the estimation error. Comment on the results of (Plot 2 and 3). What is the effect of increasing N ? What is the effect of increasing d ?

From Plot_2 and Plot_3:

- By increasing N , the estimation error and generalization error both decrease. With more data available, the law of large numbers can be fully taken use of, thus our estimation about the expected risk using the average behavior could be more accurate.
- By increasing d , the estimation error and generalization error both increase. With the same amount of data and an estimated function with greater degree, chances are that we may choose a much more complex function and lead to overfit easily.
- To sum up, the estimation error and generalization are highly correlated.

14. Besides from the approximation and estimation there is a last source of error we have not discussed here. Can you comment on the optimization error of the algorithm we are implementing?

- In this case, there is no optimization function.
- Since the empirical risk function here is convex, we could calculate the global minimum directly, rather than keep wandering around the true empirical risk minimizer.

Answer to Question 15:

