

# CSI 5138 Homework Exercise I

This homework is programming based, where **you must use Python** to explore the fitting and generalization of regression models via simulation.

Suppose that  $X$  and  $Y$  are both real valued random variables, where  $X$  takes value in  $(0, 1)$  and  $Y$  depends on  $X$  according to

$$Y = \cos(2\pi X) + Z \quad (1)$$

where  $Z$  is a zero mean Gaussian random variable with variance  $\sigma^2$ , and  $Z$  is independent of  $X$ . But assume that you do not know this dependency of  $Y$  on  $X$  and that you only observe a sample of  $N$   $(X, Y)$  pairs. Based on the observed sample, you must learn a polynomial regression model and examine the fitting and generalization capability of your model in relation to the model complexity and sample size.

Below detailed instructions are given to guide you through this exercise. These instructions only serve as a guideline, which your implementation need not to rigorously follow. You must use Python to write your code. It is fine and encouraged, but NOT compulsory, if you use a Python package that does automatic differentiation.<sup>1</sup> **But you must implement manually gradient-based optimization. That is, the use of the package is only for you to compute the required gradients and you must manually code up the update of model parameters.** You need to submit the following deliverables.

- All Python code
- A concise report explaining your findings.

(A) Write a function `getData` that generates a dataset  $\{(x_i, y_i) : i = 1, 2, \dots, N\}$  of  $N$   $(X, Y)$  pairs for a given value of  $N$  and  $\sigma^2$ . The  $X$  values are drawn uniformly at random from  $(0, 1)$  and the corresponding  $Y$  values are generated according to (1).

The dataset created by `getData` will then be used to fit your regression models. Of course, in the design of your regression model, you should assume no knowledge on how the dataset is generated.

The regression models we consider will be exclusively polynomial models, namely, predicts  $Y$  from  $X$  according to

$$Y = a_0 + a_1X + a_2X^2 + \dots + a_dX^d$$

where  $d$  is the polynomial degree and  $a_i$ 's are coefficients to be estimated.

(B) Write a function `getMSE` which computes the mean square error (MSE) for a given dataset fitted to a specified polynomial.

(C) Write a function `fitData` that estimates the polynomial coefficients by fitting a given dataset to a degree- $d$  polynomial. The function returns the following:

1. The estimated polynomial coefficients. The estimation of the coefficients should be based on GD/SGD/mini-batched SGD.

---

<sup>1</sup>If you are new to such packages, I recommend PyTorch.

2. The MSE of the dataset fitted to the estimated polynomial. This MSE will be denoted by  $E_{\text{in}}$ .
3.  $E_{\text{out}}$ . To obtain this value, your function needs to generate a separate large testing dataset (say, containing 1000 or 2000 data points) using `getData` and under the same setting of  $\sigma^2$ ) and compute the MSE of the testing dataset fitted to the estimated polynomial.

The computation of  $E_{\text{in}}$  and  $E_{\text{out}}$  calls `getMSE`.

**(D)** Write a function `experiment` that takes as input the size  $N$  of training dataset, the degree  $d$  of the model polynomial and noise variance  $\sigma^2$ , and does the following. For the given values of  $N, d$  and  $\sigma^2$ , it loops over  $M$  trials ( $M$  not smaller than 20; say, 50 would be a decent number), where each trial is defined as generating a training dataset of size  $N$  and noise variance  $\sigma^2$  (by calling `getData`) and then fitting the data to a polynomial of degree  $d$  (by calling `fitData`). The computed  $E_{\text{in}}$  and  $E_{\text{out}}$  are respectively averaged over the  $M$  trials, which are denoted by  $\bar{E}_{\text{in}}$  and  $\bar{E}_{\text{out}}$ . The obtained  $M$  polynomials over the  $M$  trials are also averaged. The function then generates another large dataset with noise variance  $\sigma^2$  and computes the average MSE for the dataset fitted to the average polynomial. This MSE will be denoted by  $E_{\text{bias}}$ . The function outputs  $\bar{E}_{\text{in}}, \bar{E}_{\text{out}}$  and  $E_{\text{bias}}$ . These three values are the metrics that you will examine.

**(E)** Run `experiment` for all combinations of  $N$ ,  $d$ , and  $\sigma^2$ , with  $N \in \{2, 5, 10, 20, 50, 100, 200\}$ ,  $d \in \{0, 1, 2, \dots, 20\}$ ,  $\sigma \in \{0.01, 0.1, 1\}$ . Organize the results in plots and comment on what you observe regarding fitting and generalization of your models, in relation to model complexity ( $d$ ), sample size ( $N$ ), and noise level ( $\sigma$ ). Note:

- Plots under comparison should display the same  $Y$ -range.
- You do not need to plot all the results. Rather, you should carefully design the ways you plot your results so as to be the most illustrative and to fully reflect your understanding.

**(F)** Revise your code to include weight decay regularization, and redo **(E)**. Plot the results and comment on what you observe.