

Homework 1

Due date: 11:59pm on Wednesday Oct. 12th

1. Bias-variance decomposition (15 points)

Prove the claim in class:

$$\begin{aligned} \mathbb{E}_D \mathbb{E}_{X,Y} [(Y - \hat{h}_D(X))^2] &= \mathbb{E}_X [(\mathbb{E}_D[\hat{h}_D(X)] - h^*(X))^2] \\ &\quad + \mathbb{E}_X \mathbb{E}_D [(\hat{h}_D(X) - \mathbb{E}_D[\hat{h}_D(X)])^2] \\ &\quad + \mathbb{E}_{X,Y} [(Y - h^*(X))^2]. \end{aligned}$$

2. MLE (15 points) Let \mathbf{x} be a sample from some distribution with parameter θ . Let $L(\theta|\mathbf{x})$ be the likelihood function for $\theta \in \Theta$. Let $g: \Theta \rightarrow \mathbb{R}$ be a function and $\tilde{L}(\xi|\mathbf{x}) := \sup_{\theta: g(\theta)=\xi} L(\theta|\mathbf{x})$ be the induced likelihood function. Suppose $\hat{\theta}$ is the MLE of θ , prove that $g(\hat{\theta})$ is the MLE for $g(\theta)$.

3. Maximum likelihood estimation (30 points)

Suppose we have n i.i.d. samples X_1, \dots, X_n from the following distributions. Find the MLE of the required parameters.

a. (5 points) Unif(0, θ). Find MLE for θ .

b. (5 points) $N(\mu, \sigma^2)$ with both parameters unknown. Find MLE for σ^2 .

c. (10 points) $N(0, \Sigma)$. Find the MLE for Σ .

d. (10 points) $N(0, \Sigma)$. Find the MLE for $\Theta = \Sigma^{-1}$. Please also provide conditions that such the MLE exists.

4. Convex optimization (10 points)

Suppose that $f: \mathbb{R}^d \mapsto \mathbb{R}$ is a convex and differentiable function. Show that x is a minimizer of f if and only if $\nabla f(x) = 0$.

5. Programming assignment: Empirical risk minimization (30 points) Given a sample data set D and loss function $\ell(y, \hat{y})$, the *empirical risk* (or empirical loss) of a hypothesis h is defined as the sample mean of the loss on $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$:

$$\hat{R}_D(h) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i))$$

where $h(x_i)$ is the predicted label for x_i . In *empirical risk minimization*, we use the empirical risk $\hat{R}_D(h)$ as an estimator of the minimal true risk, a.k.a. the *Bayes risk*, defined as

$$R^* := \min_h \mathbb{E}_{X,Y} \ell(Y, h(X)).$$

Given hypothesis class \mathcal{H} (e.g. a collection of predictors), denote the empirical risk estimator as

$$\hat{h}_D := \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_D(h).$$

The expected error of the empirical risk estimator \hat{h}_D is

$$\mathbb{E}_D \hat{R}_D(\hat{h}_D) - R^* = \underbrace{\mathbb{E}_D \hat{R}_D(\hat{h}_D) - \min_{h \in \mathcal{H}} R(h)}_{\text{estimation error}} + \underbrace{\min_{h \in \mathcal{H}} R(h) - R^*}_{\text{approximation error}}$$

Hereby, the *estimation error* is due to error caused by n training samples instead of full knowledge of the joint distribution of X, Y , and the *approximation error* is due to restricting our attention to model class \mathcal{H} .

In this question, we will investigate the trade-off between *estimation error* and *approximation error* given different collections of predictors.

Use the following code to generate a dataset:

```
#python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error

def data_generator(n_samples):
    x = np.random.uniform(-10, 10, n_samples)
    y = np.cos(0.5 + np.exp(-x)) + 1/(1 + np.exp(-x))
    noise = np.random.normal(0, 0.01, n_samples)
    y += noise
    return x, y

complete_X, complete_Y = data_generator(5000)
train_X, train_Y = complete_X[:100], complete_Y[:100]
large_X, large_Y = complete_X[100:], complete_Y[100:]

loss_func = mean_squared_error
```

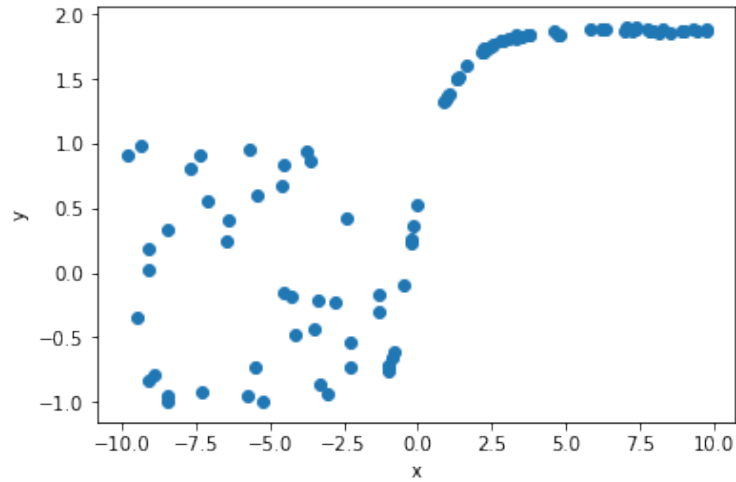
Use `train_X` and `train_Y` as training samples for ERM. `large_X` and `large_Y` are for the approximation of true data distribution of \mathbf{X} and Y , in order to estimate true risk. A plot of a small portion of the dataset.

Use mean squared error as loss function in this problem (where $\ell(y_i, h(x_i)) = (y_i - h(x_i))^2$), unless noted otherwise.

a.(10 points) Let's first define \mathcal{H}_k to be a collection of all possible polynomial functions of degree k . Implement the ERM process to select the predictor $\hat{h} \in \mathcal{H}_k$ with the lowest empirical risk. (Hint: `polyfit` function in numpy could be useful.)

b.(10 points) Experiment with the ERM you built with k of \mathcal{H}_k ranging from 0 to 30. Report empirical loss and plot a graph of empirical risk v.s. k .

c.(10 points) We will further explore the approximation vs. estimation trade-off. First, use the noise-free distribution in `data_generator` to estimate R^* with the **complete dataset**. i.e., the complete dataset



has noisy (x, y) pairs and we know that without noise

$$y^* := \mathbb{E}[Y \mid \mathbf{X} = x] = \cos(0.5 + e^{-x}) + \frac{1}{1 + e^{-x}}. \quad (1)$$

(Note that in real applications, you normally do not have access to the true distribution of \mathbf{X} and Y .) Now use `large_X`, `large_Y` to estimate the Bayes risk R^* , the risk of the ERM for each k , the estimation error, and the approximation error. Plot graphs of these errors v.s. k . Experiment with k range from 0 to 25. (Note: these different errors might not be in the same scale. You can plot one graph for each error v.s. k)