

CS/ECE/STAT-861: Theoretical Foundations of Machine Learning

University of Wisconsin–Madison, Fall 2024

Instructor: Kirthivasan Kandasamy

Homework 0.

Due 09/14/2024, 11.59 pm

Homework 0 serves as a diagnostic tool to assess your background knowledge and mathematical maturity for this course. It also familiarizes you with course expectations as you advance to more complex topics. While you are not expected to know all the solutions right away, you should be able to solve most of the questions with reasonable effort, using references as needed.

Instructions:

1. Homework is due on Canvas by 11.59 pm on the due date. Please plan to submit well before the deadline. Refer to the course website for policies on late submission.
2. Homework must be typeset using appropriate software; handwritten and scanned submissions will **not** be accepted. If you typeset your homework using \LaTeX , you will receive 5 percent extra credit.
3. Your solutions will be evaluated on correctness, clarity, and conciseness.
4. Unless otherwise specified, you may use any result we have already proved in class. Clearly state which result you are using.
5. Solutions to some of the problems may be found in the recommended textbook or other resources. Unless stated otherwise, you should try the problems on your own instead of searching for answers. If you used any external references, please cite them in your submission.
6. **Collaboration:** You are allowed to collaborate on problem 2 of this homework in groups of size up to 3. If you do so, please indicate your collaborators at the top of your solution. You may *not* collaborate on problem 1.

1 Estimating the mean of a normal distribution

You may **not** collaborate on problem 1.

We are given n independent samples $S = \{X_1, \dots, X_n\}$, where each X_i is sampled from a normal distribution $\mathcal{N}(\mu, \sigma^2)$ with *unknown* mean μ , but known variance σ^2 . We wish to estimate the mean μ . An estimator $\hat{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function of the data collected. We will quantify the performance of an estimator via the risk R , which is simply the expected squared error of the estimator,

$$R(\hat{\mu}, \mu) = \mathbb{E}[(\hat{\mu}(S) - \mu)^2].$$

Here, the expectation is with respect to the randomness in the data.

1. [2 pts] (*Bias–variance decomposition*) First, show that the following holds for any estimator $\hat{\mu}$,

$$R(\hat{\mu}, \mu) = \underbrace{(\mathbb{E}[\hat{\mu}(S)] - \mu)^2}_{\text{bias}} + \underbrace{\mathbb{E}[(\hat{\mu}(S) - \mathbb{E}[\hat{\mu}(S)])^2]}_{\text{variance}}.$$

2. [2 pts] (*Sample mean*) The most natural estimator for μ is the *sample mean* $\hat{\mu}_{\text{SM}}$, defined below.

$$\hat{\mu}_{\text{SM}}(S) = \frac{1}{n} \sum_{i=1}^n X_i.$$

Using the result of part 1, show that $R(\hat{\mu}_{\text{SM}}, \mu) = \sigma^2/n$.

3. [4 pts] (*Concentration*) While the risk measures how well an estimator does in expectation, sometimes we also wish to know that $\hat{\mu}_{\text{SM}}$ is within some margin of error ϵ of the true mean μ with high probability. Prove the following result for any $\epsilon > 0$:

$$\mathbb{P}(|\hat{\mu}_{\text{SM}}(S) - \mu| > \epsilon) \leq 2 \exp\left(\frac{-n\epsilon^2}{2\sigma^2}\right).$$

where the probability \mathbb{P} is with respect to the randomness in the data.

You may use the following facts about normal random variables:

- If X_1, \dots, X_n are normal, then so is $\sum_{i=1}^n X_i$. (You will need to compute the mean and variance.)
 - If X are normal, then so is aX for any $a \in \mathbb{R}$. (You will need to compute the mean and variance.)
 - If $Z \sim \mathcal{N}(0, 1)$ is a standard normal random variable, then $\mathbb{P}(|Z| > \epsilon) \leq 2e^{-\epsilon^2/2}$.
4. [2 pts] (*Sample complexity*) Suppose you are given some $\epsilon > 0$ and $\delta \in (0, 1)$. You wish to collect enough samples so that your estimator is within an ϵ margin of error with probability at least δ . Show that if $n \geq \frac{2\sigma^2}{\epsilon^2} \log(2/\delta)$, we will have the following guarantee: $\mathbb{P}(|\hat{\mu}_{\text{SM}} - \mu| > \epsilon) \leq \delta$.
5. [2 pts] (*An alternative estimator*) The sample mean is just one of several possible estimators for μ . Student A proposes the following alternative estimator $\hat{\mu}_\alpha$ with some parameter $\alpha \in (0, 1)$,

$$\hat{\mu}_\alpha(S) = \frac{\alpha}{n} \sum_{i=1}^n X_i.$$

Using the result from part 1, compute the risk of the estimator $\hat{\mu}_\alpha$. Note that, unlike the sample mean, the risk of $\hat{\mu}_\alpha$ depends on the true mean μ .

6. [2 pts] Show that there exists at least one value for μ such that $\hat{\mu}_\alpha$ is a strictly better estimator than $\hat{\mu}_{\text{SM}}$. That is, there exists $\mu \in \mathbb{R}$, such that, for all $\alpha \in (0, 1)$, we have $R(\hat{\mu}_\alpha, \mu) < R(\hat{\mu}_{\text{SM}}, \mu)$.

7. [4 pts] (*Maximum risk*) Despite the result from part 6, student B is not satisfied with student A's proposition, as an estimator should perform well for all values of μ , and not just for one value of μ . In particular, she argues that the worst-case risk over all μ should be small. She proposes the following criterion, the *maximum risk* \mathcal{R} , as a way to measure how well an estimator performs.

$$\mathcal{R}(\hat{\mu}) = \sup_{\mu \in \mathbb{R}} R(\hat{\mu}, \mu) = \sup_{\mu \in \mathbb{R}} \mathbb{E}[(\hat{\mu} - \mu)^2].$$

- (a) Compute $\mathcal{R}(\hat{\mu}_{\text{SM}})$ and $\mathcal{R}(\hat{\mu}_\alpha)$.
- (b) Based on the above answers, which estimator would you choose?
8. [5 pts] (*Maximum risk over a bounded domain*) Suppose we had prior knowledge that $\mu \in [0, 1]$. While student A agrees with student B's criterion, she argues that we should modify the definition of the maximum risk to incorporate this prior knowledge. She proposes the following definition instead:

$$\mathcal{R}'(\hat{\mu}) = \sup_{\mu \in [0,1]} R(\hat{\mu}, \mu) = \sup_{\mu \in [0,1]} \mathbb{E}[(\hat{\mu} - \mu)^2].$$

- (a) Compute $\mathcal{R}'(\hat{\mu}_{\text{SM}})$ and $\mathcal{R}'(\hat{\mu}_\alpha)$, the maximum risk for the two estimators discussed above.
- (b) Is there any particular value of α (possibly dependent on n and σ) for which $\mathcal{R}'(\hat{\mu}_\alpha) < \mathcal{R}'(\hat{\mu}_{\text{SM}})$?
- (c) Based on the above answer, which estimator would you choose? Intuitively, explain the discrepancy in the conclusions in part 7 and part 8.

N.B. In the first few lectures, we will study *empirical risk minimization* (ERM), a learning method which returns a model which performs well on observed data. We will use concentration tools, similar to part 3 and 4, to establish conditions under which ERM does well on future, unseen data.

Later in class, we will study minimax optimality, where our goal is to design algorithms with the smallest maximum risk across a class of distributions. We will begin with simple mean estimation problems, like the above, and extend these concepts to regression, classification, density estimation, online learning, and bandits.

2 Understanding exploration–exploitation trade-offs

Consider the following game which proceeds over T rounds. You have access to two normal distributions $\nu^{(1)} = \mathcal{N}(\mu_1, \sigma^2)$ and $\nu^{(2)} = \mathcal{N}(\mu_2, \sigma^2)$, where σ^2 is known but $\mu_1, \mu_2 \in [0, 1]$ are not. On each round t , you have the opportunity to earn a monetary reward by drawing a sample from either distribution. You choose $I_t \in \{1, 2\}$, where $I_t = i$ corresponds to drawing a sample X_t from $\nu_i = \mathcal{N}(\mu_i, \sigma^2)$, and earning $\$X_t$. If $X_t < 0$, you should pay that amount instead. Your total cumulative reward, over T rounds is $\sum_{t=1}^T X_t$. We will measure how well we perform via our *average regret*, defined below:

$$R_T = \max\{\mu_1, \mu_2\} - \frac{1}{T} \sum_{t=1}^T X_t.$$

We wish to design an algorithm whose average regret vanishes¹ with T in expectation, i.e $\mathbb{E}[R_T] \rightarrow 0$ as $T \rightarrow \infty$.

Algorithm: A student proposes the following simple algorithm. First sample each of the distributions N times (where $N < T/2$). Then, for the remaining $T - 2N$ rounds, sample the distribution with the highest observed sample mean

¹Intuitively, if we knew *a priori* which mean was larger, we will always pull the arm with the highest mean and have $\mathbb{E}[R_T] = 0$ as $\frac{1}{T}\mathbb{E}[\sum_t X_t] = \max\{\mu_1, \mu_2\}$. If $\mathbb{E}[R_T] \rightarrow 0$, this means we are able to learn which of the two distributions has a larger mean and converge towards the correct answer as we collect more samples.

using the N samples. That is, I_t is chosen as follows:

$$I_t = \begin{cases} 1 & \text{if } t \leq N, \\ 2 & \text{if } N + 1 \leq t \leq 2N, \\ 1 & \text{if } t > 2N \text{ and } \hat{\mu}_1 \geq \hat{\mu}_2, \\ 2 & \text{if } t > 2N \text{ and } \hat{\mu}_1 < \hat{\mu}_2. \end{cases}$$

$$\text{where, } \hat{\mu}_1 = \frac{1}{N} \sum_{t=1}^N X_t, \quad \hat{\mu}_2 = \frac{1}{N} \sum_{t=N+1}^{2N} X_t,$$

For what follows, let $\Delta = |\mu_1 - \mu_2|$ denote the gap between the two means.

1. **[5 pts]** (*Regret decomposition*) Establish the following identity for the expected average regret:

$$\mathbb{E}[R_T] = \frac{N\Delta}{T} + \frac{(T-2N)\Delta}{T} \Phi\left(-\Delta\sqrt{\frac{N}{2\sigma^2}}\right)$$

Here, $\Phi(x) = \mathbb{P}_{Z \sim \mathcal{N}(0,1)}(Z < x)$ is the CDF of the standard normal distribution.

2. **[2 pts]** Using the result from part 1 and the fact that $\mu_1, \mu_2 \in [0, 1]$, show the following upper bound on the expected average regret.

$$\mathbb{E}[R_T] \leq \frac{N}{T} + \Delta \exp\left(\frac{-N\Delta^2}{4\sigma^2}\right).$$

You may use the following property about standard normals, which is a one-sided version of the inequality given in problem 1. If $Z \sim \mathcal{N}(0, 1)$, then $\mathbb{P}(Z < -\epsilon) \leq e^{-\epsilon^2/2}$.

3. **[3 pts]** Use the result in part 2 to show the following upper bound.

$$\mathbb{E}[R_T] \leq \frac{N}{T} + C \frac{\sigma}{\sqrt{N}}, \quad \text{where, } C = \sqrt{2}e^{-1/2}.$$

Hint: Consider the function $f(x) = \log(x) - \alpha x^2$, where $\alpha > 0$. What is the maximizer of f ?

4. **[2 pts]** (*An optimal choice of N .*) Specify a choice for N , depending only on σ and T , so that the upper bound in part 3 is minimized. Are you able to achieve $\mathbb{E}[R_T] \rightarrow 0$ as $T \rightarrow \infty$? If so, at what rate does it go to zero?
5. **[2 pts]** (*Exploration–exploitation trade-off.*) Let N^* denote the optimal choice in part 4. In words, explain what would happen had we chosen $N \ll N^*$ or $N \gg N^*$.

N.B. Later in class, we will study several models for adaptive decision-making. The model discussed in this question is an example of a *stochastic bandit*, which is one paradigm for decision-making. In bandit settings, we often have to trade-off between *exploration* (learning about the environment) and *exploitation* (leveraging what we have learned to maximize rewards). The above algorithm is a simple, albeit sub-optimal, approach that separates exploration (the first $2N$ rounds) from exploitation (the remaining $T - 2N$ rounds). In class, we will look at better algorithms to manage this trade-off which have faster rates of convergence.