

STA3000F: Final Exam

December 12, 2024

- **Submission:** Please submit your solutions by noon at Dec 12th EST through Quercus. In an emergent situation, you may also submit the solution directly to wenlong.mou@utoronto.ca, and you will typically have a one-hour grace period.
- **Policy:** Please work on the problem set by yourself. Collaboration or resorting to external help is not allowed. On the other hand, please feel free to refer to any textbooks, papers, and online materials (or even ChatGPT, if you trust it).
- **Grading:** Each question is worth 20% of the final exam. All these questions can be solved using results from the lectures and the homeworks. You are also welcome to use ideas from other resources (books, papers, etc.). However, you are required to provide self-contained solutions to the problems using only the results from lectures and homeworks. Citing existing results directly as a black box may lead to deductions in the points depending on the nature of these results.
- **Hints:** The difficulties of problems are *not* in ascending (or descending) order. Try to allocate your time wisely. Besides, partially-solved questions may get partial credits.
- **Have fun!**

Q1. Asymptotics can be misleading

Let $X_1, X_2, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1)$ for some unknown $\theta \in \mathbb{R}$. Consider the following estimator

$$\hat{\theta}_n := \bar{X}_n \mathbf{1}_{|\bar{X}_n| \geq n^{-1/4}}, \quad \text{where } \bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i.$$

- (1/2 credit) For any fixed $\theta \in \mathbb{R}$, compute

$$R_{\text{asym}}(\theta) := \lim_{n \rightarrow +\infty} \left\{ n \cdot \mathbb{E}_\theta [|\hat{\theta}_n - \theta|^2] \right\}.$$

- (1/2 credit) Show that

$$\lim_{n \rightarrow +\infty} \sup_{\theta \in \mathbb{R}} \left\{ n \cdot \mathbb{E}_\theta [|\hat{\theta}_n - \theta|^2] \right\} = +\infty.$$

Q2. Gaussian mixture model

Suppose that we observe d -dimensional data from a mixture of two normal distributions

$$X_1, X_2, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \frac{1}{2} \mathcal{N}(\theta^*, I_d) + \frac{1}{2} \mathcal{N}(-\theta^*, I_d)$$

where $\theta^* \in \mathbb{R}^d$ is the underlying unknown parameter.

Prove that there exist universal constants $c, c_1 > 0$, such that for $n \geq c_1 d$, we have

$$\inf_{\hat{\theta}_n} \sup_{\theta^* \in \mathbb{R}^d} \mathbb{E} \left[\left\| \hat{\theta}_n - \theta^* \right\|_2^2 \right] \geq c \sqrt{\frac{d}{n}},$$

Q3. convergence rate for k -means clustering

Let \mathbb{P} be a probability distribution in \mathbb{R}^d , satisfying $\mathbb{E}_{\mathbb{P}}[\|X\|_2^2] < +\infty$. Given i.i.d. samples from \mathbb{P} , we try to divide them into K clusters of symmetric locations. In doing so, we consider the following constrained M -estimator within a constrained set $\Theta \subseteq \mathbb{R}^d$.

$$(\hat{\theta}_n^{(1)}, \hat{\theta}_n^{(2)}, \dots, \hat{\theta}_n^{(K)}) := \arg \min_{\theta_1, \theta_2, \dots, \theta_K \in \Theta} \frac{1}{n} \sum_{i=1}^n \min_{1 \leq j \leq K} \|X_i - \theta_j\|_2^2. \quad (1)$$

Clearly, the minimizer may not be unique in general due to symmetry in the problem. In such a case, we take any possible minimizer. We define the population-level loss function

$$F(\theta_1, \theta_2, \dots, \theta_K) := \mathbb{E} \left[\min_{1 \leq j \leq K} \|X_i - \theta_j\|_2^2 \right].$$

Assume that Θ is compact. Show that there exists a constant $c > 0$ that depends only on the constraint set Θ and the distribution of X (and of course, the dimension d), such that

$$F(\hat{\theta}_n^{(1)}, \hat{\theta}_n^{(2)}, \dots, \hat{\theta}_n^{(K)}) - \min_{\theta_1, \theta_2, \dots, \theta_K \in \Theta} F(\theta_1, \theta_2, \dots, \theta_K) \leq c\sqrt{K/n}$$

with probability 9/10.

Q4. Rademacher complexity of ℓ_1 ball

Given d basis functions $\phi_1, \phi_2, \dots, \phi_d$ satisfying $\|\phi_j\|_\infty \leq 1$ for $j \in [d]$, define the function class

$$\mathcal{F} := \{x \mapsto \theta^\top \phi(x) \mid \theta \in \mathbb{R}^d, \|\theta\|_1 \leq 1\}.$$

Let X_1, X_2, \dots, X_n be i.i.d. data from the underlying distribution \mathbb{P} . Show that

- (1/2 credit) For any distribution \mathbb{P} and any set of basis functions under above assumptions, we have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i f(X_i) \right] \leq c \sqrt{\log d}.$$

- (1/2 credit) There exists a distribution \mathbb{P} and a class of basis functions satisfying above assumptions, such that

$$\liminf_{n \rightarrow +\infty} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i f(X_i) \right] \geq c' \sqrt{\log d}.$$

Q5. sup norm convergence rate

Let p be a probability density function supported on the interval $[0, 1]$, satisfying the uniform Lipschitz condition

$$|p(x) - p(y)| \leq |x - y|, \quad \text{for any pair } x, y \in [0, 1].$$

Given samples $X_1, X_2, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} p$, construct a kernel density estimator \hat{p}_n (with your choice of kernel and bandwidth) satisfying

$$\mathbb{E} \left[\sup_{x \in [0, 1]} |\hat{p}_n(x) - p(x)|^2 \right] \leq c \left(\frac{\log n}{n} \right)^{2/3}.$$

[Hint: you may use the Bernstein inequality as given.]

$$\mathbb{P} \left(\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z]) \right| \geq t \right) \leq 2 \exp \left(\frac{-t^2/2}{n \text{var}(Z) + t/3} \right),$$

for i.i.d. random variables Z_1, Z_2, \dots, Z_n that takes value in $[0, 1]$.