

Practice Final

SDS384

May 6, 2019

This exam three short and 4 long questions. You will have to answer all short questions, and 3 out of 4 long questions. The assigned points are noted next to each question; the total number of points is 50. You have 180 minutes to answer the questions.

Please answer all problems in the space provided on the exam. Use extra pages if needed. Of course, please put your name on extra pages.

Read each question carefully, show your work and clearly present your answers. Note, the exam is printed two-sided - please don't forget the problems on the even pages!

Good Luck!

Name: _____

UTeid: _____

1 Short questions (17 points)

1. (5 pts) Let X be the set of binary strings of length 4 (i.e. the set of 4-character strings consisting only of 0 and 1). Let the class H be the set of schemas over X , where a schema consists of the symbols 1, 0, and *, where * matches either 0 or 1. For example, $h = 1***$ returns true for any string that starts with 1, and false for everything else. Similarly, $h = ** **$ returns true for all strings. Is X shattered by H ?

2. (6 pts) Suppose X_1, \dots, X_n are i.i.d random variables with mean μ and variance σ^2 . Let $T_n = \sum_{j=1}^n z_{nj} X_j$ where z_{nj} are given numbers. Let $\mu_n = E[T_n]$ and $\sigma_n^2 = \text{var}(T_n)$. Using the Lindeberg Feller theorem, show that

$$\frac{T_n - \mu_n}{\sigma_n} \xrightarrow{d} N(0, 1),$$

provided $\max_{j \leq n} \frac{z_{nj}^2}{\sum_{j=1}^n z_{nj}^2} \rightarrow 0$ as $n \rightarrow \infty$

3. (6 pts) Let X_1, \dots, X_n be independent and suppose that $X_n = \sqrt{n}$ with probability $1/2$ and $-\sqrt{n}$ with probability $1/2$, for $n = 1, 2, \dots$. Find the asymptotic distribution of \bar{X}_n .

2 Long questions (33 points)

There are 4 long questions. Please answer any three of them.

1. (11 pts) Let X_1, \dots, X_n be iid random variables. Consider the V statistic $V_n = \frac{\sum_{i=1}^n \sum_{j=1}^n h(x_i, x_j)}{n^2}$, where h is a symmetric kernel such that $E[h^2(X_i, X_j)] < \infty$.

(a) (5 pts) Let U_n be the corresponding U statistic. Show that $V_n - U_n \xrightarrow{P} 0$.

(b) (2 pts) What is the asymptotic distribution of V_n ? Why?

(c) (4 pts) Can we write V_n as a U statistic U_n ? That is, can we find a symmetric function $g(x_i, x_j)$ such that

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n h(x_i, x_j) = \frac{1}{\binom{n}{2}} \sum_{1 \leq i < j \leq n} g(x_i, x_j)$$

. Why or why not?

2. (11 pts) Let X_1, \dots, X_n be i.i.d $N(0, 1)$ random variables. Let $X_{(n)}$ be the largest order statistic.
- (a) (6 pts) Show that $X_{(n)}$ concentrates around $E[X_{(n)}]$. Obtain the tail bound explicitly. If you use Lipschitz continuity, please provide a proof of that.

(b) (4 pts) Obtain an upper bound on $E[X_{(n)}]$.

(c) (1 pts) Using the above two answers, can you upper bound $X_{(n)}$?

3. (11 pts) Let \mathcal{F}_d be the set of L Lipschitz functions

$$\mathcal{F}_d = \{g : [0, 1] \rightarrow [-1, 1] : g(0) = 0, \text{ and } |g(x) - g(y)| \leq L|x - y|, \forall x, y \in [0, 1]\}.$$

Recall the fact that $\log N(\epsilon; \mathcal{F}_d, \|\cdot\|_\infty) = O((L/\epsilon)^d)$.

- (a) (1 pts) Consider the process $Y_f := \langle \epsilon, f(X_1^n) \rangle$, where ϵ are n iid Rademacher random variables and $X_1^n = (X_1, \dots, X_n)$ are n iid random variables in $[0, 1]$. Is $Y_f - Y_g$ a subgaussian process? If so, under what distance metric?

- (b) (2 pts) What is the diameter of the function class \mathcal{F}_d under this distance metric?

(c) (6 pts) Using the one step discretization bound, prove that

$$\mathcal{R}_{\mathcal{F}_d} \leq c_{d,L} \left(\frac{1}{n} \right)^{C_d},$$

where $c_{d,L}$ is a constant which depends only on d and L and C_d is a constant which depends on d . Derive C_d . Bonus points for deriving $c_{d,L}$.

- (d) (2 pts) Recall the smoothly parametrized function class we studied in class. To remind you, in this case, \mathcal{F} is a class of parametric functions $\mathcal{F}_L := \{f(\theta, \cdot) : \theta \in B_2\}$, where B_2 is the unit L_2 ball in \mathbb{R}^d . Assume that \mathcal{F}_L is closed under negation. f is L Lipschitz w.r.t. the Euclidean distance on Θ , i.e. $|f(\theta, \cdot) - f(\theta', \cdot)| \leq L\|\theta - \theta'\|_2$. Briefly compare the Rademacher complexity you obtained in the last part with $\mathcal{R}_{\mathcal{F}_L}$. You can cite your lecture notes to get $\mathcal{R}_{\mathcal{F}_L}$.

4. (11 pts) Consider a random undirected network, where $A_{ij} = A_{ji} \stackrel{iid}{\sim} \text{Bernoulli}(p_n)$ for $1 \leq i < j \leq n$. $A_{ii} = 0$ for $1 \leq i \leq n$. The degree of a node is defined as $d_i = \sum_j A_{ij}$. Consider the regime where $np_n/\log n \rightarrow \infty$. *Hint: remember, not all concentration inequalities work in this regime.*
- (a) (5 pts) Show that the degree of a fixed node concentrates around its expectation $(n-1)p_n$. Obtain the tail bound explicitly.

(b) (4 pts) Can you obtain a uniform error bound on the degrees? That is, can you show that $\max_i \frac{|d_i - (n-1)p_n|}{(n-1)p_n}$ goes to zero in probability? If not, show why not. If yes, obtain the tail bound.

(c) (2 pts) Denote by $d_{(n)}$ the maximum degree. Using the last two questions, show that the maximum degree also concentrates. Obtain the tail bound explicitly.