CPSC 532D, Fall 2025: Assignment 3 due Friday, 14 November 2025, **11:59 pm**

You can do this assignment, and future ones, with a partner. **Read the website section on academic integrity** here for what you're allowed to do and not do; in particular, cite your sources (including people you talked to!) and don't use ChatGPT/etc for anything content-related. If you're not sure if something is okay, ask.

Prepare your answers to these questions using LATEX; hopefully you're reasonably familiar with it, but if not, try using Overleaf and looking around for tutorials online. Feel free to ask questions if you get stuck on things on Piazza (but remove any details about the actual answers to the questions...make a private post if that's tough). If you prefer, the .tex source for this file is available on the course website, and you can put your answers in \begin{answer} My answer here... \end{answer} environments to make them stand out; feel free to delete whatever boilerplate you want. Or answer in a fresh document if you'd prefer.

Submit your answers as a single PDF on Gradescope: here's the link. You'll be prompted to mark where each question is in your PDF; make sure you mark all relevant pages for each part (which saves us a surprising amount of grading time).

Make sure you submit using Gradescope's group feature if you're submitting a joint assignment, and put both your names on the first page to be safe; if you did the assignment partially together and partially separately, hand in separate PDFs, and put a note on each question where you worked together like *I did this problem with Alice* so we don't think you cheated. :)

On the off chance something goes wrong, you can also email your assignment to me directly (dsuth@cs.ubc.ca).

1 Rademacher complexity of deep networks [45 points]

As promised in class, we're now going to prove a Rademacher complexity bound for deep networks. To do that, we're going to build up our repertoire of Rademacher properties a bit first.

Lemma 1.1. If $V \subseteq W$, then $Rad(V) \leq Rad(W)$.

[1.1] [5 points] Prove Lemma 1.1.

Answer: TODO

Lemma 1.2. Consider finitely many sets V_i such that for all $\sigma \in \{-1,1\}^m$, it holds that $\sup_{v \in V_i} v \cdot \sigma \ge 0$; for instance, this holds if $0 \in V_i$, or if for all $v \in V_i$ we also have $-v \in V_i$. Then $\operatorname{Rad}(\cup_i V_i) \le \sum_i \operatorname{Rad}(V_i)$.

[1.2] [5 points] Prove Lemma 1.2.

Answer: TODO

The convex hull of a set V is the set of all convex combinations of points in V:

$$\operatorname{conv}(V) = \bigcup_{k>1} \left\{ \sum_{i=1}^k \alpha_i v_i : \alpha_i \ge 0; \sum_{i=1}^k \alpha_i = 1; v_1, \dots, v_k \in V \right\}.$$

Lemma 1.3. For any set V, Rad(conv(V)) = Rad(V).

[1.3] [10 points] Prove Lemma 1.3.

Answer: TODO

Lemma 1.4. For any set V, Rad $\left(\left\{\sum_{i=1}^d w_i v_i : w_i \in \mathbb{R}, \sum_{i=1}^d |w_i| \leq B, v_i \in V\right\}\right) \leq B \operatorname{Rad}(V \cup (-V))$.

[1.4] [10 points] Prove Lemma 1.4. Hint: You might want to apply Lemmas 1.2 and 1.3.

Answer: TODO

Now we're ready to bound a class of multilayer perceptrons (without bias terms because it makes things look a little cleaner – in practice, you should use bias terms!). Specifically,

$$\mathcal{H}_D = \{x \mapsto \sigma_D(W_D \sigma_{D-1}(\cdots \sigma_1(W_1 x) \cdots)) : W_1 \in \mathcal{W}_1, \dots, W_D \in \mathcal{W}_D\}.$$

The σ_i are M_i -Lipschitz elementwise activation functions such that $\sigma_i(0) = 0$; for example, ReLU(x) = $[\max(x_i, 0)]$. The W_i are matrices of shape $d_i \times d_{i-1}$, where the input dimension is $d_0 = d$, the output dimension is $d_D = 1$, and the in-between dimensions are some arbitrary, fixed sequence. The constraints are

$$\mathcal{W}_i = \left\{ W \in \mathbb{R}^{d_i \times d_{i-1}} : \forall j \in [d_i], \sum_{k=1}^{d_{i-1}} |W_{jk}| \le B_i \right\}.$$

Since \mathcal{H}_D has a nice recursive form, let's think about "peeling off" a layer at a time: bounding $\operatorname{Rad}(\mathcal{H}_D)$ in terms of $\operatorname{Rad}(\mathcal{H}_{D-1})$. To do this, recall that since we're dealing with a real-valued network, W_D is of shape $1 \times d_{D-1}$, and then notice that for $D \geq 2$,

$$\mathcal{H}_D \subseteq \left\{ x \mapsto \sigma_D \left(\sum_{j=1}^{d_{D-1}} (W_D)_j h_j(x) \right) : h_1, \dots, h_{d_{D-1}} \in \mathcal{H}_{D-1}, W_D \in \mathcal{W}_D \right\}. \tag{1}$$

[1.5] [5 points] Prove that $\operatorname{Rad}(\mathcal{H}_D|_{S_x}) \leq 2M_D B_D \operatorname{Rad}(\mathcal{H}_{D-1}|_{S_x})$.

Answer: TODO

If we define \mathcal{H}_0 in a way so that (1) also makes sense for D=1, this leaves us with a bound of the form $\operatorname{Rad}(\mathcal{H}|_{S_x}) \leq \left(\prod_{i=1}^D (2M_iB_i)\right) \operatorname{Rad}(\mathcal{H}_0|_{S_x}).$

[1.6] [10 points] Give a definition of \mathcal{H}_0 so that (1) makes sense for D=1. Bound $\operatorname{Rad}(\mathcal{H}_0|_{S_x})$ under the assumption that $\max_{x\in S_x}\|x\|_p\leq C$, for some $p\in[1,\infty]$ of your choice. Your bound should be $\mathcal{O}(1/\sqrt{m})$, treating everything but m as a constant.

Answer: TODO

Armed with this bound, we can show generalization bounds for scalar-output MLPs in the same way as for anything else: for example, we can immediately get an expectation bound on $L_{\mathcal{D}}(\mathrm{ERM}_{\mathcal{H}_{\mathcal{D}}})$ for any Lipschitz loss, and if the loss is also bounded (either "naturally" or based on a bound of |h(x)| as for logistic regression) then we can get a high-probability bound too. (The bound won't be very good for very deep networks, though – it's exponential in the depth! It's possible to improve on this somewhat with fancier techniques, but if the W_i are all norm balls, a dependence on the product of those norms is unavoidable.)

2 Threshold functions [20 points]

This question is about the class of threshold functions on \mathbb{R} :

$$\mathcal{H} = \{x \mapsto \mathbb{1}(x \ge \theta) : \theta \in \mathbb{R}\}.$$

We showed in class (notes section 6.4.1.1) that $VCdim(\mathcal{H}) = 1$: it can shatter a single point, but it cannot shatter any set of size two (since it can't label the left point 1 and the right point 0).

[2.1] [5 points] Use Sauer-Shelah (Lemma 6.12), and also the simpler Corollary 6.10, to give two upper bounds on the growth function $\Gamma_{\mathcal{H}}(m)$.

Answer: TODO

[2.2] [5 points] Directly derive the exact value of the growth function $\Pi_{\mathcal{H}}$ from its definition. How tight are the upper bounds from Question [2.1]?

Answer: TODO

[2.3] [5 points] Plug the previous parts in to upper bound $Rad(\mathcal{H}|_{S_x})$ for an S containing m distinct real numbers. You should give multiple bounds here: one for each bound, and one for the exact value of the growth function.

Answer: TODO

[2.4] [5 points] Give the asymptotic value of $\operatorname{Rad}(\mathcal{H}|_{S_x})$ for an S_x containing m distinct real numbers. Your answer might look something like " $\operatorname{Rad}(\mathcal{H}|_{S_x}) = 7m + \mathcal{O}(1)$," with a justification. To be clear, this means that $7m - a_n \leq \operatorname{Rad}(\mathcal{H}|_{S_x}) \leq 7m + a_n$ for some $a_m = \mathcal{O}(1)$. How does it compare to the bound from Question [2.3]?

Hint: Imagine playing a (pretty boring) betting game where you bet \$1 whether a coin I'm flipping comes up heads or tails, with even odds. Since all physical coin flips are unbiased, you have a 50-50 shot of getting it right. The distribution of how much money I owe you is known as a simple random walk. Your expected winnings at any time t are always 0 (it's the sum of a bunch of mean-zero variables). If we play for a while, and then you conveniently "lose" the records of what happened after some time t that just so happens to be the best possible time for you to have forgotten, you'll probably be able to win some money: the expected maximum value achieved at any point during a simple random walk of length m turns out to be $\sqrt{\frac{2m}{\pi}} - \frac{1}{2} + \mathcal{O}(m^{-\frac{1}{2}})$. (This is from equations (4) and (7) of the linked paper, which you don't need to read, just FYI.)

Answer: TODO

3 Piecewise-constant functions [25 points + 10 challenge points]

Let $a = (a_1, a_2, ..., a_k, 0, 0, ...)$ be an eventually-zero sequence with entries $a_i \in \{0, 1\}$. Then define a hypothesis $h_a : \mathbb{R}_{>0} \to \{0, 1\}$ by

$$h_a(x) = a_{\lceil x \rceil} = \begin{cases} a_1 & \text{if } 0 < x \le 1 \\ a_2 & \text{if } 1 < x \le 2 \\ & \vdots \end{cases}$$

Consider the hypothesis class of all such functions: $\mathcal{H} = \{h_a : \forall i \in \mathbb{N}, a_i \in \{0,1\} \text{ and } a \text{ is eventually zero}\}$. We'll use the 0-1 loss in this question.

[3.1] [5 points] Show $VCdim(\mathcal{H}) = \infty$.

Answer: TODO

[3.2] [5 points] Give an example of a continuous distribution \mathcal{D}_x on (a subset of) $\mathbb{R}_{>0}$ where, for some $m < \text{VCdim}(\mathcal{H})$, samples $S_x \sim \mathcal{D}_x^m$ have probability zero of being shattered by \mathcal{H} . Thus prove that, for any \mathcal{D} with this x marginal \mathcal{D}_x , ERM over \mathcal{H} obtains error at most $\inf_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon(m, \delta)$ with probability at least $1 - \delta$, where $\varepsilon(m, \delta)$ is some finite quantity such that $\lim_{m \to \infty} \varepsilon(m, \delta) = 0$ for each δ . By comparison, the VC bound would only show the approximation error is at most ∞ .

Answer: TODO

[3.3] [5 points] Write $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \cdots$, where each \mathcal{H}_k has a finite VC dimension, and write down an explicit SRM algorithm that nonuniformly learns \mathcal{H} .

By "an explicit algorithm," I mean to expand out things like the uniform convergence bound for \mathcal{H}_k . It's okay to write something as an argmin over \mathcal{H} like in equation (8.3) of the notes, if you say what k_h is for a given h and give the value of the corresponding Rademacher complexity. It's also okay to appeal to the SRM algorithm pseudocode from the notes, as long as you say what's in each \mathcal{H}_k , what the ε_k functions are, and how to compute the stopping condition.

Answer: TODO

[3.4] [5 challenge points] Challenge question: Suppose that instead of eventually-zero sequences, we allowed all possible sequences $a \in \{0,1\}^{\mathbb{N}}$, e.g. the a that infinitely alternates between 0 and 1 is now an option. Prove that this bigger \mathcal{H}' is not nonuniformly learnable. This implies a sort of no-free-lunch theorem for nonuniform learnability.

Hint: Try a diagonalization argument.

Answer: TODO

The following result will be useful momentarily:

Proposition 3.1. Let \mathcal{D} be any distribution over the positive integers \mathbb{N} , and $S \sim \mathcal{D}^m$. Define a random variable Q_S to be the number of unique samples seen out of m draws, $Q_S = |\{n : n \in S\}|$. Then $\mathbb{E}Q_S = o(m)$.

(Recall little-o notation in this case is equivalent to saying $\lim_{m\to\infty} \frac{\mathbb{E}Q_S}{m} = 0$.)

[3.5] [5 points] Prove that, for any \mathcal{D}_x , $\mathbb{E}_{S_x \sim \mathcal{D}_x^m} \operatorname{Rad}(\mathcal{H}|_{S_x}) \to 0$ as $m \to \infty$.

Hint: You can use Proposition 3.1, if you reframe the problem slightly.

Answer: TODO

[3.6] [5 challenge points] Challenge question: Prove Proposition 3.1.

Answer: TODO

[3.7] [5 points] An absentminded professor once made the following argument on the final exam for a course: If a hypothesis class has $\mathbb{E}_{S_x \sim \mathcal{D}_x^m} \operatorname{Rad}(\mathcal{H}|_{S_x}) \to 0$ for all \mathcal{D}_x , then for all realizable \mathcal{D} ,

$$L_{\mathcal{D}}(\hat{h}_S) \leq \underset{S_X \sim \mathcal{D}_x^m}{\mathbb{E}} \operatorname{Rad}(\mathcal{H}|_{S_x}) + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \to 0.$$

Thus, by the "fundamental theorem of statistical learning," H must have finite VC dimension.

Clearly this argument is wrong, since it puts Questions [3.1] and [3.5] in contradiction. What was her mistake?

Answer: TODO