

Computational and Statistical Learning Theory

Problem sets 5 and 6

Due: November 11th

Please send your solutions to learning-submissions@ttic.edu

Notations/Definitions

Recall the definition of sample based Rademacher complexity :

$$\widehat{\mathcal{R}}_S(\mathcal{F}) := \mathbb{E}_{\epsilon \sim \{\pm 1\}^n} \left[\frac{1}{m} \sup_{f \in \mathcal{F}} \sum_{i=1}^m \epsilon_i f(x_i) \right]$$

Definition 1. Given a sample $S = \{x_1, \dots, x_m\}$, and any $\alpha > 0$, a set $V \subset \mathbb{R}^m$ is said to be an α -cover (in ℓ_p) of function class \mathcal{F} on sample S if

$$\forall f \in \mathcal{F}, \exists v \in V \text{ s.t. } \left(\frac{1}{m} \sum_{i=1}^m |f(x_i) - v_i|^p \right)^{1/p} \leq \alpha$$

Specifically for $p = \infty$ $\left(\frac{1}{m} \sum_{i=1}^m |f(x_i) - v_i|^p \right)^{1/p}$ is replaced by $\max_{i \in [m]} |f(x_i) - v_i|$.
Also define

$$\mathcal{N}_p(\mathcal{F}, \alpha, S) := \min\{|V| : V \text{ is an } \alpha\text{-cover of (in } \ell_p) \text{ of } \mathcal{F} \text{ on sample } S\}$$

and

$$\mathcal{N}_p(\mathcal{F}, \alpha, n) := \sup_{x_1, \dots, x_m} \mathcal{N}_p(\mathcal{F}, \alpha, \{x_1, \dots, x_m\})$$

Definition 2. A function \mathcal{F} is said to α -shatter a sample $S = \{x_1, \dots, x_m\}$ if there exists a sequence of thresholds, $\theta_1, \dots, \theta_m \in \mathbb{R}$ such that

$$\forall \xi \in \{\pm 1\}^m, \exists f \in \mathcal{F} \text{ s.t. } \forall i \in [m], \quad \xi_i(f(x_i) - \theta_i) \geq \alpha$$

Problems

1. VC Lemma for Real-valued Function classes :

We shall prove that for any function class \mathcal{F} (assume functions in \mathcal{F} are bounded by 1) and scale $\alpha > 0$, the ℓ_∞ covering number at scale α can be bounded using fat shattering dimension at that scale by proving a statement analogous to VC lemma. We shall proceed by first extending the statement to finite (specifically $\{0, \dots, k\}$) valued function classes and then using this to prove the final bound of form

$$\mathcal{N}_\infty(\mathcal{F}, \alpha, n) \leq \sum_{i=1}^{\text{fat}_{\alpha/2}} \binom{n}{i} \left(\frac{1}{\alpha}\right)^i$$

(a) Let $\mathcal{F}_k \subset \{0, \dots, k\}^{\mathcal{X}}$ be a function class with $\text{fat}_1(\mathcal{F}_k) = d$, show that

$$\mathcal{N}_\infty(\mathcal{F}_k, 1/2, m) \leq \sum_{i=1}^d \binom{m}{i} k^i$$

Show the above statement using induction on $n + d$ (very similar to first problem on Assignment 2). Hint : In Assignment 2 problem 1 where we used \mathcal{H}_S^+ and \mathcal{H}_S^- use instead, for all $i \in \{0, \dots, k\}$, $\mathcal{F}_i = \{f \in \mathcal{F} : f(x') = i\}$ (note that this is a simple multilable extension and for $k = 1$, $\mathcal{F}_0, \mathcal{F}_1$ are identical to $\mathcal{H}^+, \mathcal{H}^-$). Use the notion of 1-shattering instead of shattering for the VC case and use 1/2-cover instead of growth function.

(b) Using the idea of α -discretizing the output of function class \mathcal{F} we shall conclude the required statement. Do the following :

i. Create a $\{0, \dots, k\}$ -valued class \mathcal{G} where k is of order $1/\alpha$. Show that covering \mathcal{G} at scale 1/2 implies we can cover \mathcal{F} at scale α and hence conclude that we can bound $\mathcal{N}_\infty(\mathcal{F}, \alpha, m)$ in terms of covering number at scale 1/2 for \mathcal{G} .

ii. Show that $\text{fat}_1(\mathcal{G}) \leq \text{fat}_{\alpha/2}(\mathcal{F})$

iii. Combine with the bound on $\mathcal{N}_\infty(\mathcal{G}, 1/2, m)$ from previous sub-problem and conclude that

$$\mathcal{N}_\infty(\mathcal{F}, \alpha, n) \leq \sum_{i=1}^{\text{fat}_{\alpha/2}} \binom{n}{i} \left(\frac{1}{\alpha}\right)^i$$

2. Dudley Vs Pollards' Bounds :

In class we saw that Rademacher complexity of a function class \mathcal{F} (assume functions in \mathcal{F} are bounded by 1) can be bounded in terms of covering numbers using Pollard's bound, Dudley integral bound and the slightly modified version of Dudley integral bound as follows

:

$$\widehat{\mathcal{R}}_S(\mathcal{F}) \leq \inf_{\alpha \geq 0} \left\{ \alpha + \sqrt{\frac{2 \log \mathcal{N}_1(\mathcal{F}, \alpha, m)}{m}} \right\} \leq \inf_{\alpha \geq 0} \left\{ \alpha + \sqrt{\frac{2 \log \mathcal{N}_2(\mathcal{F}, \alpha, m)}{m}} \right\} \quad (\text{Pollard})$$

$$\widehat{\mathcal{R}}_S(\mathcal{F}) \leq \inf_{\alpha \geq 0} \left\{ 4\alpha + 12 \int_{\alpha}^1 \sqrt{\frac{\log \mathcal{N}_2(\mathcal{F}, \tau, m)}{m}} d\tau \right\} \quad (\text{Refined Dudley})$$

In this problem using some examples we shall compare these bounds.

(a) Class with finite VC subgraph-dimension :

Assume that the VC subgraph-dimension of function class \mathcal{F} is bounded by D . In this case result in problem 1 can be used to bound the covering number of \mathcal{F} in terms of D . Use this bound on covering number and compare Pollard's bound with refined Dudley integral bound by writing down the bounds implied by each one.

(b) Linear class with bounded norm : Linear classes in high dimensional spaces is probably one of the most important and most used function class in machine learning. Consider the specific example where

$$\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\} \quad \text{and} \quad \mathcal{F} = \{\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} : \|\mathbf{w}\|_2 \leq 1\}$$

In class we saw that for any $\epsilon > 0$, $\text{fat}_\epsilon(\mathcal{F}) \leq \frac{1}{\epsilon^2}$. Using this with the result in problem 1 we have that :

$$\mathcal{N}_2(\mathcal{F}, \alpha, m) \leq \mathcal{N}_\infty(\mathcal{F}, \alpha, m) \leq \left(\frac{en}{\epsilon}\right)^{\frac{1}{\epsilon^2}}$$

Use the above bound on the covering number and write down the bound on Rademacher complexity implied by Pollard's bound. Write down the bound on Rademacher complexity implied by the refined version of the Dudley integral bound.

3. Data Dependent Bound :

Recall the Rademacher complexity bound we proved in class for functions \mathcal{F} bounded by 1. For any $\delta > 0$ with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(x)] - \widehat{\mathbb{E}}_S[f(x)] \right) \leq 2\mathbb{E}_{S \sim \mathcal{D}^m} \left[\widehat{\mathcal{R}}_S(\mathcal{F}) \right] + \sqrt{\frac{\log(1/\delta)}{m}}$$

Note that we don't know the distribution \mathcal{D} . One way we used the above bound was by providing upper bounds on $\widehat{\mathcal{R}}_S(\mathcal{F})$ for any sample of size m and using this instead of $\mathbb{E}_{S \sim \mathcal{D}^m} \left[\widehat{\mathcal{R}}_S(\mathcal{F}) \right]$. But ideally we would like to get tight bounds when the distribution we are faced with is nicer. The aim of this problem is to do this.

Prove that, for any $\delta > 0$ with probability at least $1 - \delta$, over draw of sample $S \sim \mathcal{D}^m$,

$$\sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(x)] - \widehat{\mathbb{E}}_S[f(x)] \right) \leq 2\widehat{\mathcal{R}}_S(\mathcal{F}) + K \sqrt{\frac{\log(2/\delta)}{m}}$$

(provide explicit value of constant K above). Notice that in the above bound the expected Rademacher complexity is replaced by sample based one which can be calculated from the training sample.

Hint : Use McDiarmid's inequality on the expected Rademacher complexity.

4. Learnability and Fat-shattering Dimension:

Recall the setting of stochastic optimization problem where objective function is mapping $r : \mathcal{H} \times \mathcal{Z} \mapsto \mathbb{R}$. Sample $S = \{z_1, \dots, z_m\}$ drawn iid from unknown distribution \mathcal{D} is provided to the learner and the aim of the learner is to output $\hat{h} \in \mathcal{H}$ based on sample that has low expected objective $\mathbb{E}[r(h, z)]$.

- Consider the stochastic optimization problem with r bounded by a , i.e. $|r(h, z)| < a < \infty$ for all $h \in \mathcal{H}$ and $z \in \mathcal{Z}$. If function class $\mathcal{F} := \{z \mapsto r(h, z) | h \in \mathcal{H}\}$ has finite fat_α for all $\alpha > 0$, then show that the problem is learnable.
- Conclude that for a supervised learning problem with bounded hypothesis class \mathcal{H} (ie. $\forall x \in \mathcal{X}, |h(x)| < a$), and loss $\phi : \hat{\mathcal{Y}} \times \mathcal{Y} \mapsto \mathbb{R}$ that is L -Lipschitz (in first argument), if \mathcal{H} has finite fat_α for all $\alpha > 0$, then the problem is learnable.
- Show a stochastic optimization problem that is learnable even though it has infinite fat_α for all $\alpha \leq 0.1$ (or any other constant of your choice). Explicitly write down the hypothesis class, and the learning rule which learns the class, argue that the problem is learnable, and explain why the fat_α is infinite.

Hint : You can make the learning rule that is successful to even be ERM.

- Prove that for a supervised learning problem with the absolute loss $\phi(\hat{y}, y) = |\hat{y} - y|$, if the fat_α is infinite for some $\alpha > 0$, then the problem is not learnable.

Hint: as with the binary case, for every m , construct a distribution which is concentrated on a set of points that can be fat-shattered.

5. Massart's Lemma:

In this problem, we want to prove Massart's lemma, i.e. for finite $\mathcal{F} \subseteq \mathbb{R}^m$ the following inequality holds:

$$\mathcal{R}_m(\mathcal{F}) \leq a \sqrt{\frac{2 \log |\mathcal{F}|}{m}} \quad (1)$$

where $a = \sup_{f \in \mathcal{F}} \|f\|_2$.

- Use Jensen's inequality and union bound to prove that for any $\lambda > 0$

$$e^{\lambda m \mathcal{R}_m} \leq \sum_{f \in \mathcal{F}} \prod_{i=1}^m \mathbb{E}_{\epsilon_i \sim \{\pm 1\}} [e^{\lambda \epsilon_i f_i}] \quad (2)$$

- (b) Find an upper bound for the expected value in the above inequality using the fact that $e^x + e^{-x} \leq 2e^{x^2/2}$.
- (c) Show that for any $\lambda > 0$

$$m\mathcal{R}_m(\mathcal{F}) \leq \frac{\log |\mathcal{F}|}{\lambda} + \frac{\lambda a^2}{2} \quad (3)$$

- (d) Find λ that minimize the above bound.

6. Refined Dudley integral bound on the Rademacher complexity:

Given a sequence of samples $S = (x_1, \dots, x_m)$ we shall prove that for any function class \mathcal{F} containing functions $f : \mathcal{X} \mapsto \mathbb{R}$

$$\widehat{\mathcal{R}}_S(\mathcal{F}) \leq \inf_{\alpha > 0} \left\{ 4\alpha + 10 \int_{\alpha}^{\sup_{f \in \mathcal{F}} \sqrt{\widehat{\mathbb{E}}[f^2]}} \sqrt{\frac{\mathcal{N}_2(\mathcal{F}, \tau, S)}{m}} d\tau \right\} \quad (4)$$

where $\widehat{\mathbb{E}}[f^2] = \frac{1}{m} \sum_{i=1}^m f^2(x_i)$.

- (a) Let $\beta_0 = \sup_{f \in \mathcal{F}} \sqrt{\widehat{\mathbb{E}}[f^2]}$ and for any $j \in \mathbb{Z}_+$ let $\beta_j = \beta_0/2^j$ and C_j be β_j -cover (in ℓ_2) of function class F on sample S . For each $f \in F$ we pick an $\hat{f}_j \in C_j$ such that $\sqrt{\widehat{\mathbb{E}}[(f - \hat{f}_j)^2]} \leq \beta_j$. The main idea of the proof is that for any $N > 0$ we can express f by chaining as

$$f = f - \hat{f}_N + \sum_{j=1}^N (\hat{f}_j - \hat{f}_{j-1}) \quad (5)$$

where $\hat{f}_0 = 0$.

Show that under the same conditions stated above, for any $N > 0$ the following inequality holds:

$$\widehat{\mathcal{R}}_S(\mathcal{F}) \leq \beta_N + \sum_{j=1}^N \widehat{\mathcal{R}}_S(\mathcal{G}_j). \quad (6)$$

where $\mathcal{G}_j = \{\hat{f} - \hat{f}' \mid \hat{f} \in T_j, \hat{f}' \in T_{j-1}\}$

- (b) Show that $|\mathcal{G}_j| \leq \mathcal{N}_2(\mathcal{F}, \beta_j, S)$.
- (c) Use Massart's lemma to find the following upper bound for Rademacher complexity of class \mathcal{G}_j :

$$\widehat{\mathcal{R}}_S(\mathcal{G}_j) \leq 5\beta_j \sqrt{\frac{\log \mathcal{N}_2(\mathcal{F}, \beta_j, S)}{m}} \quad (7)$$

- (d) Show that the following inequality holds:

$$\widehat{\mathcal{R}}_S(\mathcal{F}) \leq \beta_N + 10 \sum_{j=1}^N (\beta_j - \beta_{j-1}) \sqrt{\frac{\log \mathcal{N}_2(\mathcal{F}, \beta_j, S)}{m}} \quad (8)$$

(e) Prove the following bound for any $j > 0$:

$$(\beta_j - \beta_{j-1}) \sqrt{\frac{\log \mathcal{N}_2(\mathcal{F}, \beta_j, S)}{m}} \leq \int_{\beta_{j-1}}^{\beta_j} \sqrt{\log \frac{\mathcal{N}_2(\mathcal{F}, \tau, S)}{m}} d\tau \quad (9)$$

Find the upper bound for the sum term as a single integral.

(f) Show that for any $\alpha > 0$ we can choose N such that

$$\widehat{\mathcal{R}}_S(\mathcal{F}) \leq 4\alpha + 10 \int_{\alpha}^{\sup_{f \in \mathcal{F}} \sqrt{\mathbb{E}[f^2]}} \sqrt{\frac{\log \mathcal{N}_2(\mathcal{F}, \tau, S)}{m}} d\tau \quad (10)$$

Conclude that the Dudley integral bound holds.

7. A general bound for all norms:

In class we studied the hypothesis class $H_B = \{w \mid \|w\| \leq B\}$ and the corresponding rule ERM_B , and showed that for every B , w.h.p over S , for all $w \in H_B$ we had

$$|\hat{L}(w) - L(w)| \leq O\left(\sqrt{\frac{B^2 + \log(1/\delta)}{m}}\right) \quad (11)$$

We then used this to get a bound on $L_{\{0,1\}}$ in terms of \hat{L}_γ , for each γ separately. We will now change the order of quantifiers, and get a bound that holds, w.h.p., over all γ concurrently, with only a very mild (double logarithmic) additional term.

(a) Prove that, for $\mathcal{H} \subseteq [-a, a]^X$, for any $D(X, Y)$, w.h.p. over $S \sim D^m$, for all $h \in \mathcal{H}$ and all $\gamma > 0$:

$$L_{\{0,1\}}(h) \leq \hat{L}_\gamma(h) + O\left(\frac{\mathcal{R}_m(\mathcal{H})}{\gamma} + \sqrt{\frac{\log(\log(a/\gamma)) + \log(1/\delta)}{m}}\right) \quad (12)$$

Rewrite the above bound with proper constants and without big- O notation.

(b) In particular, linear separators with geometric margin γ correspond to using $H_1 = \{w \mid \|w\| \leq 1\}$, in which case show that we have, w.h.p. over $S \sim D^m$, for all $\|w\| \leq 1$ and $\gamma > 0$:

$$L_{\{0,1\}}(w) \leq \hat{L}_\gamma(h) + O\left(\sqrt{\frac{1/\gamma^2 + \log(1/\delta)}{m}}\right) \quad (13)$$

Rewrite the above bound with proper constants and without big- O notation.

Challenge Problems

1. We saw that for any distribution \mathcal{D} , the expected Rademacher complexity provided an upper bound on the maximum deviation between mean and average uniformly over function class, specifically we saw that

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[\sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(x)] - \hat{\mathbb{E}}[f(x)] \right) \right] \leq 2 \mathbb{E}_{S \sim \mathcal{D}^m} \left[\widehat{\mathcal{R}}_S(\mathcal{F}) \right]$$

Prove the (almost) converse that

$$\frac{1}{2} \mathbb{E}_{S \sim \mathcal{D}^m} \left[\widehat{\mathcal{R}}_S(\mathcal{F}) \right] \leq \mathbb{E}_{S \sim \mathcal{D}^m} \left[\sup_{f \in \mathcal{F}} \left(\mathbb{E} [f(x)] - \widehat{\mathbb{E}} [f(x)] \right) \right]$$

This basically establishes that Rademacher complexity tightly bounds the uniform maximal deviation for every distribution.

2. The worst case Rademacher complexity is defined as

$$\widehat{\mathcal{R}}_m(\mathcal{F}) = \sup_{S=\{x_1, \dots, x_m\}} \widehat{\mathcal{R}}_S(\mathcal{F})$$

(ie. supremum over samples of size m).

(a) Prove that for any function class \mathcal{F} and any $\tau > \widehat{\mathcal{R}}_m(\mathcal{F})$, we have that

$$\text{fat}_\tau(\mathcal{F}) \leq \frac{m \widehat{\mathcal{R}}_m(\mathcal{F})^2}{\tau^2}$$

(b) Combine the above with the refined version of Dudley integral bound to prove that

$$\inf_{\alpha \geq 0} \left\{ 4\alpha + 12 \int_\alpha^1 \sqrt{\frac{\mathcal{N}_2(\mathcal{F}, \tau, m)}{m}} d\tau \right\} \leq \widehat{\mathcal{R}}_m(\mathcal{F}) O(\log^{3/2} m)$$

This shows that the refined dudley integral bound is tight to within log factors of the Rademacher complexity. Thus we have established that in the worst case all the complexity measures for function class like Rademacher complexity, covering numbers and fat shattering dimension all tightly govern the rate of uniform maximal deviation for the function class (all to within log factor).

3. Bounded Difference Inequality, Stability and Generalization :

Recall that a function $G : \mathcal{X}^m \mapsto \mathbb{R}$ is said to satisfy the bounded difference inequality if for all $i \in [m]$ and all $x_1, \dots, x_m, x'_i \in \mathcal{X}$,

$$|G(x_1, \dots, x_i, \dots, x_m) - G(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_m)| \leq c$$

for some $c \geq 0$. In this case the McDiarmid's inequality gave us that for any $\delta > 0$, with probability at least $1 - \delta$,

$$G(x_1, \dots, x_m) \leq \mathbb{E} [G(x_1, \dots, x_m)] + \sqrt{\frac{(mc)^2 \log(1/\delta)}{m}}$$

The bounded difference property turns out to be quiet useful to analyze learning algorithms directly (instead of looking at the uniform deviation over function class).

A proper learning algorithm is $A : \bigcup_{m=1}^{\infty} \mathcal{X}^m \mapsto \mathcal{F}$ is said to be a uniformly β stable is for all $i \in [m]$, and any $x_1, \dots, x_m, x'_i \in \mathcal{X}$,

$$\sup_x |A(x_1, \dots, x_i, \dots, x_m)(x) - A(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_m)(x)| \leq \beta$$

Assuming functions in \mathcal{F} are bounded by 1 we shall prove that the learning algorithm generalizes well (expected loss is close to empirical loss of the algorithm). Specifically we shall prove that for any $\delta > 0$, with probability at least $1 - \delta$,

$$L(A(S)) \leq \widehat{L}(A(S)) + \beta + 2(m\beta + 1) \sqrt{\frac{2 \log(1/\delta)}{m}}$$

where $L(A(S)) = \mathbb{E}_x [A(S)(x)]$ and $\widehat{L}(A(S)) = \frac{1}{m} \sum_{i=1}^m A(S)(x_i)$.

(a) First show that $\mathbb{E}_S [L(A(S)) - \widehat{L}(A(S))] \leq \beta$.

Hint : Use renaming of variables to first show that for any $i \in [m]$,

$$\mathbb{E}_S [L(A(S))] = \mathbb{E}_{S, x'_i} [A(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_m)(x_i)]$$

(b) Show that the function $G(S) = L(A(S)) - \widehat{L}(A(S))$ satisfies bounded difference property with $c \leq 2\beta + \frac{2}{m}$. Conclude the required statement using McDiarmid's inequality.

(c) Consider the stochastic convex optimization problem where sample $z = (x, y)$ where y is real valued and x 's are from the unit ball in some Hilbert space and hypothesis is weight vectors w from the same Hilbert space with objective

$$r(w, (x, y)) = |\langle w, x \rangle - y| + \lambda \|w\|^2$$

Show that the ERM algorithm is stable for this problem and thus provide a bound for this algorithm.

4. L_1 Neural Network :

A k -layer 1-norm neural network is given by function class \mathcal{F}_k which is in turn defined recursively as follows.

$$\mathcal{F}_1 = \left\{ x \mapsto \sum_{j=1}^d w_j^1 x_j \mid \|w^1\|_1 \leq B_1 \right\}$$

and further for each $2 \leq i \leq k$,

$$\mathcal{F}_i = \left\{ x \mapsto \sum_{j=1}^{d_i} w_j^i \sigma(f_j(x)) \mid \forall j \in [d_i], f_j \in \mathcal{F}_{i-1}, \|w^i\|_1 \leq B_1 \right\}$$

where d_i is the number of nodes in the i th layer of the network. Function $\sigma : \mathbb{R} \mapsto [-1, 1]$ is called the squash function and is generally a smooth monotonic non-decreasing function (typical example is the tanh function). Assume that input space $\mathcal{X} = [0, 1]^d$ and that σ is L -Lipschitz. Prove that

$$\widehat{\mathcal{R}}_S(\mathcal{F}_k) \leq \left(\prod_{i=1}^k 2B_i \right) L^{k-1} \sqrt{2T \log d}$$

Notice that the above bound the d_i 's don't appear in the bound indicating the number of nodes in intermediate layers don't affect the upper bound on Rademacher complexity.

Hint : prove bound on Rademacher complexity of \mathcal{F}_i recursively in terms of rademacher complexity of \mathcal{F}_{i-1} .