

# SDS 384 11: Theoretical Statistics Lecture 14: Uniform Law of Large Numbers- Covering number

Purnamrita Sarkar Department of Statistics and Data Science The University of Texas at Austin

- Recall that a metric space (*T*, *ρ*) consists of a nonempty set *T* and a mapping *ρ* : *T* × *T* → ℝ that satisfies:
  - Non-negative:  $\rho(\theta, \theta') \ge 0$  for all  $(\theta, \theta')$  with equality iff  $\theta = \theta'$ .
  - Symmetric:  $\rho(\theta, \theta') = \rho(\theta', \theta)$  for all pairs  $(\theta', \theta)$ , and
  - Triangle ineq holds:  $\rho(\theta, \theta') + \rho(\theta', \theta'') \ge \rho(\theta, \theta'')$
- Examples:

• 
$$\mathcal{T} = \mathbb{R}^d$$
,  $\rho(\theta, \theta') = \|\theta - \theta'\|_2$   
•  $\mathcal{T} = \{0, 1\}^d$  with  $\rho(\theta, \theta') = \frac{1}{d} \sum_i \mathbb{1}(\theta_i \neq \theta'_i)$ 

#### Definition

A  $\delta$  cover of a set  $\mathcal{T}$  w.r.t to a metric  $\rho$  is a set  $\{\theta^1, \ldots, \theta^N\} \subset \mathcal{T}$  such that for every  $\theta \in \mathcal{T}$ ,  $\exists i \in [N]$ , s.t.  $\rho(\theta, \theta^i) \leq \delta$ . The  $\delta$  covering number  $N(\delta; \mathcal{T}, \rho)$  is the cardinality of the smallest  $\delta$  cover.

- We will consider metric spaces which are totally bounded, i.e.
   N(δ; T, ρ) < ∞ for all δ > 0.
- The covering number is non-increasing in  $\delta$ , i.e.  $N(\delta) \ge N(\delta')$  for all  $\delta < \delta'$
- We are interested in something called Metric entropy, which is the logarithm of the covering number.

# Picture



Figure 1: [courtesy: Martin Wainwright's book]

• A  $\delta$  covering can be thought of as a union of balls with radius  $\delta$ .

#### Example

Consider the interval [-1, 1] with  $\rho(\theta, \theta') = |\theta - \theta'|$ . We have  $N(\delta; [-1, 1], |.|) \leq \frac{1}{\delta} + 1$ 

- Divide the interval into L sub-intervals centered at  $\theta^i := -1 + (2i - 1)\delta$  for  $i \in [L]$  and each of length at most  $2\delta$ .
- By construction this is a  $\delta$  covering.
- So  $L \leq 1 + 1/\delta$

#### Example

Consider a *d* dimensional binary hypercube  $T = \{0, 1\}^d$  with the Hamming metric defined before.

$$rac{\log \mathsf{N}(\delta;\mathcal{T},
ho)}{\log 2} \leq \lceil \mathsf{d}(1-\delta) 
ceil$$

- Let  $S = \{1, 2, \dots, \lceil \delta d \rceil\}$
- Consider the set of binary vectors S(δ) := {θ ∈ T : θ<sub>j</sub> = 0, j ∈ S(δ)}.
- By construction, for every binary vector  $\theta' \in \mathcal{T}$ , we can find a vector  $\theta \in S(\delta)$  such that  $\rho(\theta, \theta') \leq \delta$
- $N(\delta; \mathcal{T}, \rho) \leq |\mathcal{S}(\delta)| = 2^{\lceil d(1-\delta) \rceil}$

- Let  $\delta \in (0, 1/2)$
- If {θ<sup>1</sup>,...,θ<sup>N</sup>} is a δ covering, then the (unrescaled) Hamming balls of radius s = δd around each θ<sup>ℓ</sup> must contain all 2<sup>d</sup> vectors.
- Let  $s = \lfloor \delta d \rfloor$

• For each 
$$\theta^i$$
 there are exactly  $\sum_{j=0}^{s} \binom{d}{j}$  vectors within  $\delta d$  distance.

• So 
$$N \sum_{j=0}^{s} \binom{d}{j} \ge 2^{d}$$

- Let  $\delta \in (0, 1/2)$
- So  $N \sum_{j=0}^{s} \binom{d}{j} \ge 2^{d}$
- Now take a Binomial (d, 1/2) random variable X.

• 
$$P(X \le \delta d) = \sum_{j=0}^{s} {d \choose j} / 2^{d}$$

• So 
$$N \ge \frac{1}{P(X \le \delta d)}$$

- Using the Hoeffding bound gives:  $N \ge \exp(\frac{d}{2}(1/2 \delta)^2)$
- Using the refined version in your homework gives:  $N \ge \exp(d \times KL(\delta||1/2))$

#### Definition

An  $\delta$ -packing of  $\mathcal{T}$  w.r.t a metric  $\rho$  is a set  $\{\theta^1, \ldots, \theta^M\} \subset \mathcal{T}$  such that  $\rho(\theta^i, \theta^j) > \delta$  for every distinct pair  $i, j \in [M]$ . The  $\delta$  packing number  $M(\delta; \mathcal{T}, \rho)$  is the cardinality of the largest  $\delta$  packing.

# Picture



Figure 2: [courtesy: Martin Wainwright's book]

• A  $2\delta$  packing can be thought of as a union of balls with radius  $\delta$  such that no two balls touch.

#### Theorem

For all  $\delta > 0$ ,

#### $M(2\delta; \mathcal{T}, \rho) \leq N(\delta; \mathcal{T}, \rho) \leq M(\delta; \mathcal{T}, \rho)$

• This is saying that packing and covering numbers exhibit the same scaling behavior as  $\delta \rightarrow 0$ .

# Proof

- Upper bound: Let V = {x<sub>1</sub>,..., x<sub>N</sub>} be a δ packing of T. So for each y ∈ T \ V, ∃i, ||y x<sub>i</sub>|| ≤ δ. Otherwise we could have added this point and increased the packing number. So, V is also a δ cover. But since the covering number is the size of the smallest δ covering, the lower bound holds.
- Lower bound: Say there is a  $2\delta$  packing  $\{y_1, \ldots, y_M\}$  and a  $\delta$  covering  $\{v_1, \ldots, v_n\}$  with M > n. Now by pigeonhole, there must be two  $y_i, y_j$  who both are in the  $\delta$  ball around some  $v_k$ . But using triangle, we will have  $|y_i y_j| \le 2\delta$ , which is a contradiction. So we must have  $m \le n$ .

# Covering and Packing numbers-example

#### Theorem

Let  $\rho$  be the Euclidean norm on  $\mathbb{R}^d$ . Let  $B_1(0)$  be the unit ball centered at the origin (WLOG).

$$\frac{1}{\epsilon^d} \le N(\epsilon, B_1, \rho) \le (1 + 2/\epsilon)^d$$

• Consider an 
$$\epsilon$$
 cover  $\{\theta^1, \ldots, \theta^N\}$ . Now,

$$B_{1} \subseteq \bigcup_{i=1}^{N} B_{\epsilon}(\theta^{i})$$
$$\mathsf{vol}(B_{1}) \le N\mathsf{vol}(B_{\epsilon}(\theta^{i})) = N\epsilon^{d}\mathsf{vol}(B_{1})$$
$$N \ge 1/\epsilon^{d}$$

- Consider a  $\epsilon$  packing  $\{\theta^1, \dots, \theta^M\}$
- This is an union of disjoint balls of radius  $\epsilon/2$

$$\bigcup_{i} B_{\epsilon/2}(\theta^{i}) \subseteq B_{1+\epsilon/2}$$
$$M \operatorname{vol}(B_{\epsilon/2}(\theta^{i})) \leq (1+\epsilon/2) \operatorname{vol}(B_{1+\epsilon/2})$$
$$M(\epsilon/2)^{d} \operatorname{vol}(B_{1}) \leq (1+\epsilon/2)^{d} \operatorname{vol}(B_{1})$$
$$M \leq (1+2/\epsilon)^{d}$$

#### Theorem

Consider a d dimensional vector of independent sub $G(\sigma^2)$  random variables. Let  $B_d$  be the unit ball in  $\|.\|_2$  norm. Then the following holds:

$$E[\sup_{\theta \in B_d} \theta^T X] \le 4\sigma \sqrt{d}$$

Also, for  $\delta \in (0,1)$ , with probability  $1 - \delta$ ,

$$\sup_{\theta \in B_d} \theta^T X \le 4\sigma \sqrt{d} + \sqrt{2\sigma \log(1/\delta)}.$$

# [Recall] Size of a function class $\mathcal{F}$

• Let 
$$\mathcal{F}(X) = \{(f(X_1), \dots, f(X_n)) : f \in \mathcal{F}\}$$
  
•  $\mathcal{R}_{\mathcal{F}} = E \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i} \epsilon_i f(X_i) \right| = E \left[ E \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i} \epsilon_i f(X_i) \right| \left| X_1, \dots, X_n \right] \right]$ 

• In the next slide we will bound this using the cardinality of  $\mathcal{F}(X)$ 

#### Theorem

Let 
$$A \subseteq \mathbb{R}^{n}$$
,  $R = \max_{a \in A} ||a||$ ,  
 $E \sup_{a \in A} \langle \epsilon, a \rangle \leq \sqrt{2R^{2} \log |A|}$ .  
And,  
 $E \sup_{a \in A} |\langle \epsilon, a \rangle| \leq \sqrt{2R^{2} \log |2A|}$ .

### Proof of first half

- Let  $\mathcal{N}_{1/2}$  be a half covering of  $B_1$ . So  $N(1/2, B_1, |||_2) \le 5^d$
- So for each  $heta\in B_d$ ,  $\exists z_ heta\in \mathcal{N}_{1/2}$  such that

$$\theta = z_{\theta} + x, \qquad \|x\| \le 1/2$$

• So,  
$$Y := \sup_{\theta \in B_1} \theta^T X \le \max_{z_{\theta} \in \mathcal{N}_{1/2}} z_{\theta}^T X + \underbrace{\sup_{x \in B_{1/2}} x^T X}_{Y/2}$$

• Thus, we have:

$$EY \leq 2E \left[ \max_{z_{\theta} \in \mathcal{N}_{1/2}} z_{\theta}^{\mathsf{T}} X \right] \leq 2\sigma \sqrt{2\log|\mathcal{N}_{1/2}|} \leq \sigma \sqrt{8d\log 5} \leq 4\sigma \sqrt{d}$$

• We used the same result as last time.

$$P(Y \ge t) \le P(\max_{z \in \mathcal{N}_{1/2}} z^T X \ge t/2)$$
  
$$\le |\mathcal{N}_{1/2}||P(z^T X \ge t/2)$$
  
$$\le 5^d \exp(-t^2/8\sigma^2 ||z||^2) \le 5^d \exp(-t^2/4\sigma^2) = \delta$$

Solving for t gives,  $\sqrt{\log 5 + \log(1/\delta)} = t/2\sigma$ . In fact, we can get an upper bound on t as follows.

$$t = 2\sigma\sqrt{d\log 5 + \log(1/\delta)} \le 2\sigma\sqrt{d\log 5} + 2\sigma\sqrt{\log(1/\delta)} =: t_0$$

Thus,  $P(Y \ge t_0) \le P(Y \ge t) \le \delta$ 

• Consider the following function class parametrized by  $\theta \in \Theta$ .

$$\mathcal{F} := \{ f_{\theta}(.) : \theta \in \Theta \}$$

• Let  $\|.\|_{\Theta}$  be the norm for  $\theta$  and  $\|.\|_{\mathcal{F}}$  be the norm for  $\mathcal{F}$ .

• Say 
$$\|f_{\theta}(.) - f_{\theta'}(.)\|_{\mathcal{F}} \le L \|\theta - \theta'\|_{\Theta}$$

• Then  $N(\epsilon; \mathcal{F}, \|.\|_F) \leq N(\epsilon/L; \Theta, \|.\|_{\Theta})$ 

- A Lipschtiz parametrization allows us to go from cover of the Θ space to cover of the f<sub>θ</sub> space with a loss of L.
- If  $\mathcal{F}$  is parametrized by a compact set of d parameters then  $N(\epsilon, \mathcal{F}) = O(1/\epsilon^d)$

#### Example

For any fixed  $\theta$ , define the real-valued function  $f_{\theta}(x) := \exp(-\theta x)$ , and consider the function class

$$\mathcal{F} = \{ f_{\theta} : [0,1] \rightarrow \mathbb{R} | \theta \in [0,1] \}$$

Using the uniform norm as a metric, i.e.

$$\|f - g\|_{\infty} := \sup_{x \in [0,1]} |f(x) - g(x)|.$$
 Prove that  
 $\lfloor \frac{1 - 1/e}{2\delta} 
floor + 1 \le N(\delta; \mathcal{F}, \|.\|_{\infty}) \le \frac{1}{2\delta} + 2.$ 

### Proof-upper bound

- First note that  $\|f_{\theta} f_{\theta'}\|_{\infty} \le |\theta \theta'|$
- For any  $\delta \in (0,1)$ , let  $T = \lfloor \frac{1}{2\delta} \rfloor$
- Consider  $S = \{\theta^0, \dots, \theta^{T+1}\}$  where  $\theta^i = 2\delta i$  for  $i \leq T$  and  $\theta^{T+1} = 1$ .

• 
$$\{f_{\theta^i}: \theta^i \in S\}$$
 is a  $\delta$  cover for  $\mathcal{F}$ .

- For any  $heta \in [0,1]$  we can find  $heta^i \in S$  such that  $| heta^i heta| \leq \delta$
- Indeed we have,

$$\|f_{\theta^{i}} - f_{\theta}\|_{\infty} = \sup_{x \in [0,1]} |\exp(-\theta^{i}x) - \exp(-\theta x)|$$
$$\leq |\theta^{i} - \theta| \leq \delta$$

So  $N(\delta; \mathcal{F}, \|.\|_{\infty}) \leq 2 + T \leq 2 + \frac{1}{2\delta}$ 

• We will do a  $\delta$  packing.

• Let 
$$\theta' = -\log(1-i\delta)$$
 for  $i = 0, \dots, T$ 

•  $-\log(1 - T\delta) = 1$ , and so the largest integral value is  $T = \lfloor \frac{1 - 1/e}{\delta} \rfloor$ 

• So 
$$M(\delta; \mathcal{F}, \|.\|_{\infty}) \ge 1 + \lfloor \frac{1 - 1/e}{\delta} \rfloor$$

•  $N(\delta; \mathcal{F}, \|.\|_{\infty}) \ge M(2\delta; \mathcal{F}, \|.\|_{\infty}) \ge 1 + \lfloor \frac{1 - 1/e}{2\delta} \rfloor$ 



**Figure 3:**  $\exp(-\theta^{i}x)$  where  $\theta^{i} = -\log(1-i\delta)$ 

#### **Example**

 $\mathcal{F}_{L} = \{g: [0,1] \to \mathbb{R} | g(0) = 0, |g(x) - g(y)| \le L | x - x'|, \forall x, x' \in [0,1] \}$ 

Metric entropy scales as log  $N(\delta; \mathcal{F}_L, \|.\|_{\infty}) \simeq L/\delta$  for small enough  $\delta > 0$ .

# Proof

- Its sufficient to consider a sufficiently large packing of  $\mathcal{F}_L$
- For a given  $\epsilon$  define  $M = \lfloor \frac{1}{\epsilon} \rfloor$

• Let 
$$x_i = (i-1)\epsilon$$
 for  $i = 1, \dots, M+1$ 

$$\phi(x) := \begin{cases} 0 & x < 0 \\ x & x \in [0, 1] \\ 1 & x > 1 \end{cases}$$
(1)

• Define 
$$f_{\beta}(x) = \sum_{i=1} \beta_i L \epsilon \phi\left(\frac{x-x_i}{\epsilon}\right)$$
 for  $\beta \in \{-1,1\}^M$ 

#### Picture



Figure 5-2. The function class  $\{f_{\beta}, \beta \in \{-1, +1\}^M\}$  used to construct a packing of the Lipschitz class  $\mathscr{F}_L$ . Each function is piecewise linear over the intervals  $[0, \epsilon], [\epsilon, 2\epsilon], \ldots, [(M-1)\epsilon, M\epsilon]$  with slope either +L or -L. There are  $2^M$  functions in total, where  $M = \lfloor 1/\epsilon \rfloor$ .

- Note that the functions in the packing are all continuous, although they are not differentiable.
- Lipschitz continuous functions are continuous, but they do not have to be necessarily differentiable.
- Lipschitz continuous functions are differentiable almost everywhere.
- Rademacher's theorem basically says that the number of discontinuities of a Lipschitz continuous function are of measure zero.

# example

- For any pair  $\beta \neq \beta' \in \{-1, 1\}^M$  there is at least one interval where they have the same starting point.
- So  $||f_{\beta}(x) f_{\beta}'(x)||_{\infty} \ge 2L\epsilon$
- $f_{\beta} \in \mathcal{F}_L$  for all  $\beta \in \{-1, 1\}^M$
- So  $f_{\beta}$  forms a  $2L\epsilon$  packing.
- Making  $\epsilon L = \delta$  we see

$$N(\delta; \mathcal{F}_L, \|.\|_{\infty}) \geq M(2L\epsilon; \mathcal{F}_L, \|.\|_{\infty}) = 2^{\lfloor \frac{1}{\epsilon} \rfloor} = 2^{\lfloor \frac{L}{\delta} \rfloor}$$

• Also the set  $f_{\beta}$  also form a suitable covering of the original functions, and this gives the upper bound.

• The last example can be extended to Lipschitz functions on the Unit cube in higher dimensions, i.e.

$$|f(x) - f(y)| \le ||x - y||_{\infty}$$
 for all  $x, y \in [0, 1]^d$ 

• The same method can be used to show that the metric entropy for this class is the same order as  $(L/\delta)^d$ 

# Make a comparison

- Recall that for a L Lipschitz continuous functions supported on [0,1] with f(0) = 0, the metric entropy was L/δ
- Also recall that for a *L* Lipschitz continuous functions supported on  $[0,1]^d$  with f(0) = 0, the metric entropy was  $(L/\delta)^d$
- However for a given function class like the last one the metric entropy is  $\log(1/\delta)$
- Recall that for Unit hypercubes in *d* dimensions the metric entropy is  $d \log(1 + 1/\delta)$
- Note that for Lipschitz continuous functions the dependence on *d* is exponential. This is a much richer class of functions, so the size is considerably larger and scales poorly with *d*.

This lecture was very much based on Martin Wainwright's book.