

Machine Learning Theory (CS 6783)

Lecture 3: Cover's Result and Rademacher Complexity

1 Bit Prediction

We covered the result of Thomas Cover given below in the first couple classes.

Lemma 1 (T. Cover'65). *Let $\phi : \{\pm 1\}^n \mapsto \mathbb{R}$ be a function such that, for any i , and any $y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n$,*

$$|\phi(y_1, \dots, y_{i-1}, +1, y_{i+1}, \dots, y_n) - \phi(y_1, \dots, y_{i-1}, -1, y_{i+1}, \dots, y_n)| \leq \frac{1}{n}, \text{ (stability condition)}$$

then, there exists a randomized strategy such that for any sequence of bits,

$$\frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\hat{y}_t \sim q_t} [\mathbf{1}\{\hat{y}_t \neq y_t\}] \leq \phi(y_1, \dots, y_n)$$

if and only if,

$$\mathbb{E}_{\epsilon} \phi(\epsilon_1, \dots, \epsilon_n) \geq \frac{1}{2}$$

and further, the strategy achieving this bound on expected error is given by:

$$q_t = \frac{1}{2} + \frac{n}{2} \mathbb{E}_{\epsilon_{t+1}, \dots, \epsilon_n} [\phi(y_1, \dots, y_{t-1}, -1, \epsilon_{t+1}, \dots, \epsilon_n) - \phi(y_1, \dots, y_{t-1}, +1, \epsilon_{t+1}, \dots, \epsilon_n)]$$

Using the above lemma one can conclude the following claim easily.

Claim 2. *There exists a randomized prediction strategy that ensures that*

$$\mathbb{E} [\text{Reg}_n] \leq \frac{1}{2n} \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \sum_{t=1}^n f_t \epsilon_t \right]$$

where $\text{Reg}_n = \frac{1}{n} \sum_{t=1}^n \mathbf{1}\{\hat{y}_t \neq y_t\} - \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \mathbf{1}\{f_t \neq y_t\}$ and $\mathcal{F} \subseteq \{\pm 1\}^n$

2 Rademacher Complexity

Given a set $\mathcal{F} \subseteq \mathbb{R}^n$, we define the Rademacher complexity of this set as

$$\mathcal{R}_n(\mathcal{F}) = \frac{1}{n} \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \sum_{t=1}^n f_t \epsilon_t \right]$$

While we have already seen the Rademacher complexity as coming from cover's result, it turns out that this quantity or rather complexity measure is a key tool in Statistical learning theory.

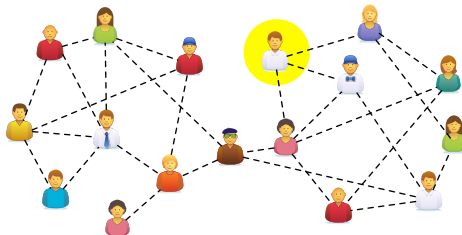
Hence lets try to see what the quantity represents. Note that if \mathcal{F} was binary labels, then for any vector $f \in \mathcal{F}$, $\|f\|_2 = \sqrt{n}$ and $\|\epsilon\|_2 = \sqrt{n}$. Hence we can interpret,

$$\frac{1}{n} \sum_{t=1}^n f_t \epsilon_t = \frac{1}{n} f^\top \epsilon = \frac{f^\top \epsilon}{\|f\|_2 \|\epsilon\|_2} = \cos(\epsilon, f)$$

Hence, we can think of $\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_\epsilon [\max_{f \in \mathcal{F}} \cos(\epsilon, f)]$, that is, how well we can correlate with random draw of labels using set \mathcal{F} .

Now before we go into statistical learning, let us get back to our bit prediction problem.

3 Application: Binary Node Classification



Let $G = (V, E)$ be a known undirected graph representing a social network. At each time step t , a user in the network opens her Facebook page, and the system needs to decide whether to classify the user as type “−1” or “+1”, say, in order to decide on an advertisement to display. We assume here that the feedback on the “correct” type is revealed to the system after the prediction is made. Suppose we have a hunch that the type of the user (+1 or −1) is correlated with the community to which she belongs. For simplicity, suppose there are two communities, more densely connected within than across. To capture the idea of correlating communities and labels, we set ϕ to be small on labelings that assign homogenous values within each community. We make the following simplifying assumptions: (i) $|V| = n$, (ii) we only predict the label of each node once, and (iii) the order in which the nodes are presented is fixed (this assumption is easily removed). Smoothness of a labeling $f \in \{\pm 1\}^n$ with respect to the graph may be computed via

$$\text{Cut}(f) = \sum_{(u,v) \in E} \mathbf{1}_{\{f_u \neq f_v\}} = \frac{1}{4} \sum_{(u,v) \in E} (f_u - f_v)^2 = f^\top L f \quad (1)$$

where $L = D - A$, the diagonal matrix D contains degrees of the nodes, and A is the adjacency matrix and $f_v \in \{\pm 1\}$ is the label in f that corresponds to vertex $v \in V$. This function in (1) counts the number of disagreements in labels at the endpoints of each edge. The value is also known as the size of the cut induced by f (the smallest possible being MinCut). As desired, the function in (1) gives a smaller value to the labelings that are homogenous within the communities.

Unfortunately, the function $\text{Cut}(f)$ is not stable. Further, the cut size is $n - 1$ for a star graph, where $n - 1$ nodes, labeled as +1, are connected to the center node, labeled as −1. The large value of the cut does not capture the simplicity of this labeling, which is only one bit away from being a constant +1. Instead, we opt for the indirect definition:

$$F_\kappa = \left\{ f \in \{\pm 1\}^n : f^\top L f \leq \kappa \right\} \quad (2)$$

for $\kappa \geq 0$, and then set

$$\phi(y_1, \dots, y_n) = \inf_{f \in \mathcal{F}_\kappa} \frac{1}{n} \sum_{t=1}^n \mathbf{1}_{\{f_t \neq y_t\}} + \frac{1}{2n} \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}_\kappa} \sum_{t=1}^n f_t \epsilon_t \right] \quad (3)$$

Parameter κ should be larger than the value of MinCut , for otherwise the set F_κ is empty. This gives an interesting algorithm for the prediction problem What does this look like?

Well we want to use the strategy

$$\begin{aligned} q_t &= \frac{1}{2} + \frac{n}{2} \mathbb{E}_{\epsilon_{t+1}, \dots, \epsilon_n} [\phi(y_1, \dots, y_{t-1}, -1, \epsilon_{t+1}, \dots, \epsilon_n) - \phi(y_1, \dots, y_{t-1}, +1, \epsilon_{t+1}, \dots, \epsilon_n)] \\ &= \frac{1}{2} + \frac{n}{2} \mathbb{E}_{\epsilon_{t+1}, \dots, \epsilon_n} \left[\inf_{f \in \mathcal{F}_\kappa} \left\{ \frac{1}{n} \sum_{j=1}^{t-1} \mathbf{1}_{\{f_j \neq y_j\}} + \mathbf{1}_{\{f_t \neq -1\}} + \sum_{j=t+1}^n \mathbf{1}_{\{f_j \neq \epsilon_j\}} \right\} \right. \\ &\quad \left. - \inf_{f \in \mathcal{F}_\kappa} \left\{ \frac{1}{n} \sum_{j=1}^{t-1} \mathbf{1}_{\{f_j \neq y_j\}} + \mathbf{1}_{\{f_t \neq +1\}} + \sum_{j=t+1}^n \mathbf{1}_{\{f_j \neq \epsilon_j\}} \right\} \right] \end{aligned}$$

It turns out that by concentration inequalities, it even suffices to take a single new sample of $\epsilon_{t+1}, \dots, \epsilon_n$ for round t to compute q_t above. In this case the underlying strategy is peculiar: At time t , to predict label for vertex v_t , we fill seen entries by labels, unseen entries by random ϵ_v 's and solve two optimization problems. One with labels set as mentioned and with label of v_t set to -1 we solve for $\inf_{f \in \mathcal{F}_\kappa} \left\{ \frac{1}{n} \sum_{j=1}^{t-1} \mathbf{1}_{\{f_j \neq y_j\}} + \mathbf{1}_{\{f_t \neq -1\}} + \sum_{j=t+1}^n \mathbf{1}_{\{f_j \neq \epsilon_j\}} \right\}$. Now we do the optimization with only changing the label of v_t to a $+1$. We can then set q_t by equation above. Here once can view the random signs we draw as a kind of regularization or protection against worst case adversarial future.

4 A Game of Betting

In the previous section, we assumed ϕ was stable. While stable ϕ 's consist of a large number of benchmark, it might not be expressive enough for some problems. Unfortunately, if we want to do classification, it is not easy to get rid of such an assumption easily. Instead below we consider a slightly different betting game on binary outcomes where we are allowed to bet arbitrary amounts on outcomes. In such game, the same idea as above can be used without requiring stability.

Consider a gambler who bets on the outcomes of games one every round. Specifically, on any round t , the gambler can choose an amount $|\hat{y}_t|$ to bet on the outcome of game between two players or teams A and B . The gambler can choose to place this bet of $|\hat{y}_t|$ on either team A to win or on team B . If the chosen team wins, the gambler gains an additional amount of \hat{y}_t and if the chosen team loses the gambler loses the bet amount of \hat{y}_t . This game of betting can be formalized as the following linear game between the gambler and the house. Specifically, we can view the choice of the gambler at round t as a real number \hat{y}_t . The magnitude \hat{y}_t denotes the bet amount and the sign of \hat{y}_t denotes whether the bet is placed on team A or team B . The corresponding outcome of the game is encoded by the variable $y_t \in \{\pm 1\}$ which indicates whether team A won or team B . At time t , $-\hat{y}_t \cdot y_t$ denotes the loss of the gambler. That is if the gambler guessed the outcome right, that is if $\text{sign}(\hat{y}_t) = y_t$, then the loss is the negative value of $|\hat{y}_t|$ (or in other words the gambler gains) and if the outcome is guessed in correctly the gambler loses the amount of $|\hat{y}_t|$.

At time $t = 1, \dots, n$, the forecaster chooses $\hat{y}_t \in \mathbb{R}$ based on the history y_1, \dots, y_{t-1} and then observes the value $y_t \in \{\pm 1\}$.

Given some benchmark function $\phi : \{\pm 1\}^n \rightarrow \mathbb{R}_{\geq 0}$, the goal of the gambler is to ensure that the loss of the gambler is smaller than this benchmark. In other words, the gambler would like to ensure that,

$$\forall \mathbf{y}, \quad \mathbb{E} \left[\frac{1}{n} \sum_{t=1}^n -\hat{y}_t y_t \right] \leq \phi(\mathbf{y}) \quad (4)$$

Lemma 3. ϕ is achievable if and only if $\mathbb{E}[\phi(\epsilon)] \geq 0$. Further, in this case, the strategy for the gambler is given by: $\hat{y}_t = \frac{n}{2} \cdot \mathbb{E}[\phi(y_{1:t-1}, -1, \epsilon_{t+1:n}) - \phi(y_{1:t-1}, +1, \epsilon_{t+1:n})]$.

Remark: stability is not required.

Example 4.1. We have a gambler who likes to bet on games played between m teams. Assume that the information about which pairs of teams play each other for the n matches is announced in advance. Specifically, say we know that on round t , teams i_t and j_t play each other. Let us further denote by n_i the number of games played by player i . This game of betting can be formalized in the linear betting games framework above. As specific benchmark a gambler might consider is the one where each of the m team is given a score represented by an m dimensional vector \mathbf{w} . Further, when team i plays team j , a bet of amount of $|w[i] - w[j]|$ on the team with the larger score is placed. Further, assume that the largest bet amount is restricted to B . The goal of the gambler is to do as well as the best scoring of the teams selected in hindsight. This example, can be represented by the benchmark $\phi : \{\pm 1\}^n \mapsto \mathbb{R}$ as follows:

$$\phi(y_1, \dots, y_n) = \inf_{\mathbf{w} \in \mathbb{R}^m : \max_{i,j} |w[i] - w[j]| \leq B} \frac{1}{n} \sum_{t=1}^n y_t \cdot (\mathbf{w}[i_t] - \mathbf{w}[j_t]) + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \quad (5)$$

$$\leq \inf_{\mathbf{w} \in \mathbb{R}^m : \max_{i,j} |w[i] - w[j]| \leq B} \frac{1}{n} \sum_{t=1}^n y_t \cdot (\mathbf{w}[i_t] - \mathbf{w}[j_t]) + \frac{B}{2} \sqrt{\frac{m}{n}} \quad (6)$$

This benchmark satisfies the property that $\mathbb{E}[\phi(\epsilon)] \geq 0$. This is because

$$\begin{aligned}
\mathbb{E}[\phi(\epsilon)] &= \mathbb{E} \left[\inf_{\mathbf{w} \in \mathbb{R}^m: \max_{i,j} \mathbf{w}[i] - \mathbf{w}[j] \leq B} \frac{1}{n} \sum_{t=1}^n y_t \cdot (\mathbf{w}[i_t] - \mathbf{w}[j_t]) \right] + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \\
&= \mathbb{E} \left[\inf_{\mathbf{w} \in [0,B]^m} \frac{1}{n} \sum_{t=1}^n \epsilon_t (\mathbf{w}[i_t] - \mathbf{w}[j_t]) \right] + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \\
&= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^m \min_{\mathbf{w}[i] \in [0,B]} \sum_{t=1}^n \mathbf{w}[i] \epsilon_t (\mathbf{1}_{\{i_t=i\}} - \mathbf{1}_{\{j_t=i\}}) \right] + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \\
&= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^m \min \left\{ B \sum_{t=1}^n \epsilon_t (\mathbf{1}_{\{i_t=i\}} - \mathbf{1}_{\{j_t=i\}}), 0 \right\} \right] + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \\
&= \frac{B}{n} \sum_{i=1}^m \mathbb{E} \left[\min \left\{ \sum_{j=1}^{n_i} \epsilon_j, 0 \right\} \right] + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \\
&\geq -\frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} + \frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} = 0
\end{aligned}$$

where in the last line we used the fact that for any integer N , $\mathbb{E} \left[\min \left\{ \sum_{j=1}^N \epsilon_j, 0 \right\} \right] \geq -\sqrt{N}/2$. Hence, from Lemma 3 this benchmark is achievable by the gambler using the strategy $\hat{y}_t = n \cdot \mathbb{E}[\phi(y_{1:t-1}, -1, \varepsilon_{t+1:n}) - \phi(y_{1:t-1}, +1, \varepsilon_{t+1:n})]$. Finally, noting that square-root is a concave function and applying Jensen's inequality, yields that $\frac{B}{2n} \sum_{i=1}^m \sqrt{n_i} \leq \frac{B}{2} \sqrt{\frac{m}{n}}$.