

Machine Learning Theory (CS 6783)

Lecture 4 : Linear Betting with Covariates Contd., Learning Frameworks, Examples

1 Linear Betting Game With Covariates

For $t = 1$ to n :

1. Receive instance $x_t \in \mathcal{X}$.
2. Predict $\hat{y}_t \in \mathbb{R}$.
3. Receive label $y_t \in \{\pm 1\}$ and pay loss $\hat{y}_t \cdot y_t$.

End For.

This is a variant of the betting game where we get covariate or side information on every round and we want to perform as well as some benchmark Φ that uses knowledge of this side information. You can think of this as an extension of Cover's result when we have covariates.

We say that the adaptive bound $\Phi : \mathcal{X}^n \times \{\pm 1\}^n \rightarrow \mathbb{R}$ is *achievable* if there exists a strategy for the learner that ensures that

$$\sum_{t=1}^n \hat{y}_t \cdot y_t \leq \Phi(x_1, \dots, x_n, y_1, \dots, y_n). \quad (1)$$

We want to answer the question of when a performance bound Φ is achievable and, when it is achievable, what the algorithm for the learner should be.

Meet the Trees

Definition 1. The sequence of functions $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ with $\mathbf{x}_t : \{\pm 1\}^{t-1} \rightarrow \mathcal{X}$ will be called an \mathcal{X} -valued tree. Here $\mathbf{x}_1 \in \mathcal{X}$ is a constant.

We are now ready to provide the main result characterizing which Φ 's are achievable.

Lemma 1. A necessary and sufficient condition for Φ to be achievable in the sense of (??) is that

$$\inf_{\mathbf{x}} \mathbb{E}_{\epsilon} \Phi(\mathbf{x}_1, \mathbf{x}_2(\epsilon_1), \dots, \mathbf{x}_n(\epsilon_{1:n-1}); \epsilon) \geq 0, \quad (2)$$

where the infimum is taken over all \mathcal{X} -valued trees \mathbf{x} , and $\epsilon_1, \dots, \epsilon_n$ are i.i.d. Rademacher random variables.

Proof you will work out in Assignment 1

Remark 1.1. We will sometimes write either \mathbf{x}_t or $\mathbf{x}_t(\epsilon)$ instead of the more precise but longer expression $\mathbf{x}_t(\epsilon_1, \dots, \epsilon_{t-1})$ whenever this does not cause confusion.

Example 1.1. *Let*

$$\Phi(x_1, \dots, x_n, y_1, \dots, y_n) = \inf_{f \in \mathcal{F}} \sum_{t=1}^n y_t \cdot f(x_t) + \text{Complex}_n(\mathcal{F}).$$

In this case, the smallest value of $\text{Complex}_n(\mathcal{F})$ that makes the above Φ achievable is given by:

$$\text{Complex}_n(\mathcal{F}) = \sup_{\mathbf{x}} \mathbb{E}_{\epsilon} \sup_{f \in \mathcal{F}} \sum_{t=1}^n \epsilon_t \cdot f(\mathbf{x}_t(\epsilon_{1:t-1})),$$

which we refer to as the sequential Rademacher complexity. Note that if the tree \mathbf{x} has the same value for all its nodes on level t (i.e., $\mathbf{x}_t(\epsilon_1, \dots, \epsilon_{t-1}) = x_t$ for all ϵ), then the above recovers the worst-case statistical Rademacher complexity. The crucial difference here is that we allow an arbitrary tree, which can make the sequential Rademacher complexity possible larger than the statistical one in some settings.

Supervised Learning With Convex, L -Lipschitz Losses The linear betting loss $-\hat{y}_t \cot y_t$ seems restrictive. But it turns out that via reduction to such linear betting games one can actually consider more complex losses and prediction settings. Specifically, consider the following learning problem: For $t = 1$ to n :

1. Receive instance $x_t \in \mathcal{X}$.
2. Predict $\hat{y}_t \in \mathbb{R}$.
3. Receive $y_t \in \mathcal{Y}$ and pay loss $\ell(\hat{y}_t, y_t)$.

End For.

In the above assume that the loss ℓ is convex and L -Lipschitz in its first argument. Say our goal was to minimize regret w.r.t. some class of predictors $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$ given by

$$\text{Reg}_n = \sum_{t=1}^n \ell(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^n \ell(f(x_t), y_t)$$

Now we claim that if one has an algorithm for the linear betting game with $\Phi(x_1, \dots, x_n, y_1, \dots, y_n) = \inf_{f \in \mathcal{F}} \sum_{t=1}^n y_t \cdot f(x_t) + \text{Complex}_n(\mathcal{F})$, then, one can reduce the above game to linear betting game and use that algorithm and obtain a bound on regret of the form:

$$\text{Reg}_n \leq L \text{Complex}_n(\mathcal{F})$$

Proof. Note that for any $f \in \mathcal{F}$, by convexity: $\ell(\hat{y}_t, y_t) - \ell(f(x_t), y_t) \leq \partial \ell(\hat{y}_t, y_t) \cdot (\hat{y}_t - f(x_t))$. Hence we can conclude that:

$$\sum_{t=1}^n \ell(\hat{y}_t, y_t) - \sum_{t=1}^n \ell(f(x_t), y_t) \leq \sum_{t=1}^n \partial \ell(\hat{y}_t, y_t) \cdot (\hat{y}_t - f(x_t))$$

We use the fact that by L -Lipschitzness $\partial\ell(\hat{y}_t, y_t)$ is a number between $-L$ and L . Next we interpret this number as L times the expected value of $b_t \in \{\pm 1\}$ drawn such that probability of b_t is one is given by $\frac{\partial\ell(\hat{y}_t, y_t)+1}{2}$

$$= L \sum_{t=1}^n \mathbb{E}_{b_t \sim \frac{\partial\ell(\hat{y}_t, y_t)+1}{2}} [b_t \cdot (\hat{y}_t - f(x_t))]$$

Hence

$$\begin{aligned} \text{Reg}_n &\leq L \sup_{f \in \mathcal{F}} \sum_{t=1}^n \mathbb{E}_{b_t \sim \frac{\partial\ell(\hat{y}_t, y_t)+1}{2}} [b_t \cdot (\hat{y}_t - f(x_t))] \\ &\leq L \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sum_{t=1}^n b_t \cdot (\hat{y}_t - f(x_t)) \right] \\ &= L \mathbb{E} \left[\sum_{t=1}^n b_t \cdot \hat{y}_t - \inf_{f \in \mathcal{F}} \sum_{t=1}^n b_t \cdot f(x_t) \right] \end{aligned}$$

In other words, the way we can use the linear betting algorithm is as follows. On every iteration predict \hat{y}_t as suggested by the linear betting algorithm based on what it is fed so far as outcomes. (at round $t = 1$ there is no input and so we just get a recommendation and subsequently we can proceed as stated here). Next, after playing \hat{y}_t $y_t \in \mathcal{Y}$ is revealed to us. We then compute $\partial\ell(\hat{y}_t, y_t)$ and then draw $b_t \in \{\pm 1\}$ as described in the proof and feed this as outcome of the game for round t of linear betting game to the linear betting algorithm. Thus on the next round when we provide x_{t+1} to the linear betting algorithm it will in turn suggest the next \hat{y}_{t+1} to play and so on. \square

2 Setting up learning problems

1. \mathbf{x} : instance space or input space

Examples:

- Computer Vision: Raw $M \times N$ image vectorized $\mathcal{X} = [0, 255]^{M \times N}$, SIFT features (typically $\mathcal{X} \subseteq \mathbb{R}^d$)
- Speech recognition: Mel Cepstral co-efficients $\mathcal{X} \subset \mathbb{R}^{12 \times \text{length}}$
- Natural Language Processing: Bag-of-words features ($\mathcal{X} \subset \mathbb{N}^{\text{document size}}$), n-grams

2. \mathcal{Y} : Outcome space, label space

Examples: Binary classification $\mathcal{Y} = \{\pm 1\}$, multiclass classification $\mathcal{Y} = \{1, \dots, K\}$, regression $\mathcal{Y} \subset \mathbb{R}$

3. $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$: loss function (measures prediction error)

Examples: Classification $\ell(y', y) = \mathbf{1}_{\{y' \neq y\}}$, Support vector machines $\ell(y', y) = \max\{0, 1 - y' \cdot y\}$, regression $\ell(y', y) = (y - y')^2$

4. $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$: Model/ Hypothesis class (set of functions from input space to outcome space)

Examples:

- Linear classifier: $\mathcal{F} = \{x \mapsto \text{sign}(f^\top x) : f \in \mathbb{R}^d\}$

- Linear SVM: $\mathcal{F} = \{x \mapsto f^\top x : f \in \mathbb{R}^d, \|f\|_2 \leq R\}$
- Neural Networks (deep learning): $\mathcal{F} = \{x \mapsto \sigma(W_{out}\sigma(W_K\sigma(\dots\sigma(W_2(W_1\sigma(W_{in}x))))))\}$ where σ is some non-linear transformation (Eg. ReLU)

Learner observes sample: $S = (x_1, y_1), \dots, (x_n, y_n)$

Learning Algorithm : (forecasting strategy, estimation procedure)

$$\hat{\mathbf{y}} : \mathcal{X} \times \bigcup_{t=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^t \mapsto \mathcal{Y}$$

Given new input instance x the learning algorithm predicts $\hat{\mathbf{y}}(x, S)$. When context is clear (ie. sample S is understood) we will fudge notation and simply use notation $\hat{\mathbf{y}}(\cdot) = \hat{\mathbf{y}}(\cdot, S)$. $\hat{\mathbf{y}}$ is the predictor returned by the learning algorithm.

Example: linear SVM Learning algorithm solves the optimization problem:

$$\mathbf{w}_{\text{SVM}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{t=1}^n \max\{0, 1 - y_t \mathbf{w}^\top x_t\} + \lambda \|\mathbf{w}\|$$

and the predictor is $\hat{\mathbf{y}}(x) = \hat{\mathbf{y}}(x, S) = \mathbf{w}_{\text{SVM}}^\top x$

2.1 PAC framework

$$\mathcal{Y} = \{\pm 1\}, \quad \ell(y', y) = \mathbf{1}_{\{y' \neq y\}}$$

Input instances generated as $x_1, \dots, x_n \sim D_X$ where D_X is some unknown distribution over input space. The labels are generated as

$$y_t = f^*(x_t)$$

where target function $f^* \in \mathcal{F}$. Learning algorithm only gets sample S and does not know f^* or D_X .

Goal: Find $\hat{\mathbf{y}}$ that minimizes

$$\mathbb{P}_{x \sim D_X} (\hat{\mathbf{y}}(x) \neq f^*(x))$$

2.2 Non-parametric Regression

$$\mathcal{Y} \subseteq \mathbb{R}, \quad \ell(y', y) = (y' - y)^2$$

Input instances generated as $x_1, \dots, x_n \sim D_X$ where D_X is some unknown distribution over input space. The labels are generated as

$$y_t = f^*(x_t) + \varepsilon_t \quad \text{where } \varepsilon_t \sim N(0, \sigma)$$

where target function $f^* \in \mathcal{F}$. Learning algorithm only gets sample S and does not know f^* or D_X .

Goal: Find $\hat{\mathbf{y}}$ that minimizes

$$\mathbb{E}_{x \sim D_X} [(\hat{\mathbf{y}}(x) - f^*(x))^2] =: \|\hat{\mathbf{y}} - f^*\|_{L_2(D_X)}$$

2.3 Statistical Learning (Agnostic PAC)

Generic \mathcal{X} , \mathcal{Y} , ℓ and \mathcal{F}

Samples generated as $(x_1, y_1), \dots, (x_n, y_n) \sim D$ where D is some unknown distribution over $\mathcal{X} \times \mathcal{Y}$.

Goal: Find \hat{y} that minimizes

$$\mathbb{E}_{(x,y) \sim D} [\ell(\hat{y}(x), y)] - \inf_{f \in \mathcal{F}} \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)]$$

For any mapping $g : \mathcal{X} \mapsto \mathcal{Y}$ we shall use the notation $L_D(g) = \mathbb{E}_{(x,y) \sim D} [\ell(g(x), y)]$ and so our goal can be re-written as:

$$L_D(\hat{y}) - \inf_{f \in \mathcal{F}} L_D(f)$$

Remarks:

1. \hat{y} is a random quantity as it depends on the sample
2. Hence formal statements we make will be in high probability over the sample or in expectation over draw of samples

2.4 Online Learning

For $t = 1$ to n

- (a) Input instance $x_t \in \mathcal{X}$ is produced
- (b) Learning algorithm outputs prediction \hat{y}_t
- (c) True outcome y_t is revealed to learner

End For

One can think of $\hat{y}_t = \hat{y}_t(x_t, ((x_1, y_1), \dots, (x_{t-1}, y_{t-1})))$.

Goal: Find learning algorithm \hat{y} that minimizes regret w.r.t. hypothesis class $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ given by,

$$\text{Reg}_n = \sum_{t=1}^n \ell(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^n \ell(f(x_t), y_t)$$

3 Example 1: Classification using Finite Class, Realizable Setting

In this section we consider the classification setting where $\mathcal{Y} = \{\pm 1\}$ and $\ell(y', y) = \mathbf{1}\{y' \neq y\}$. We further make the realizability assumption meaning $y_t = f^*(x_t)$ where f^* is obviously not known to the learner.

3.1 Online Framework

The online framework is just as described earlier with the realizability assumption added in. That is, at every round the true label y_t revealed to us is set as $y_t = f^*(x_t)$ for some fixed f^* not known to the learning algorithm. However x_t 's can be presented to us arbitrarily. First note that under the realizability assumption, we have that

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) = \frac{1}{n} \sum_{t=1}^n \mathbf{1}\{f^*(x_t) \neq y_t\} = 0$$

Hence the aim in such a framework is to simply minimize number of mistakes $\sum_{t=1}^n \ell(\hat{y}_t, y_t)$ and prove mistake bounds.

Now say $\mathcal{F} = \{f_1, \dots, f_N\}$, a finite set of hypothesis. What strategy can we provide for this problem? How well does it work?

If we simply pick some hypothesis that has not made a mistake so far, such an algorithm can make a large number of mistakes (Eg. as many as N). A simple strategy that works in this scenario is the following. At any point t , we have observed x_1, \dots, x_{t-1} and labels y_1, \dots, y_{t-1} . Now say

$$\mathcal{F}_t = \{f \in \mathcal{F} : \forall i \in [t-1], f(x_i) = y_i\}.$$

Now given x_t , we pick $\hat{y}_t = \text{sign}(\sum_{f \in \mathcal{F}_t} f(x_t))$. That is we go with the majority of predictions by hypothesis in \mathcal{F}_t . How well does this algorithm work?

Claim 2. *For any sequence x_1, \dots, x_n , the above algorithm makes at most $\lceil \log_2 N \rceil$ number of mistakes.*

Proof. Notice that each time we make a mistake, ie. $\text{sign}(\sum_{f \in \mathcal{F}_t} f(x_t)) \neq y_t$, then we know that at least half the number of functions in \mathcal{F}_t are wrong and so each time we make a mistake, $|\mathcal{F}_{t+1}| \leq |\mathcal{F}_t|/2$ and hence, we can make at most $\log_2 N$ number of mistakes. \square

That is the average error is $\frac{\log_2 N}{n}$.

3.2 PAC Framework

In the PAC framework, x_1, \dots, x_n are drawn iid from some fixed distribution $D_{\mathcal{X}}$ and our goal is to minimize $P_{x \sim D_x}(\hat{y}(x) \neq f^*(x))$ either in expectation or high probability over sample $\{x_1, \dots, x_n\}$. Unlike the online setting, in the PAC setting one can simply pick any hypothesis that has not made any mistakes on training sample. That is,

$$\hat{y}(\cdot, S) = \underset{f \in \mathcal{F}}{\text{argmin}} \sum_{(x_t, y_t) \in S} \mathbf{1}\{f(x_t) \neq y_t\}.$$

How well does this algorithm work? How should we analyze this?

Let us show a bound of error with high probability over samples. To this end we will use the so called Bernstein concentration bound.

Fact: Consider binary r.v. Z_1, \dots, Z_n drawn iid. Let $\mu = \mathbb{E}[Z]$ be their expectation. We have the following bound on the average of these random variables. (notice that since Z 's are binary their variance is given by $\mu - \mu^2$)

$$P\left(\mu - \frac{1}{n} \sum_{t=1}^n Z_t > \theta\right) \leq \exp\left(-\frac{n\theta^2}{2\mu + \frac{\theta}{3}}\right)$$

Now for any $f \in \mathcal{F}$, let $Z_t^f = \mathbf{1}\{f(x_t) \neq f^*(x_t)\}$ where x_t are drawn from $D_{\mathcal{X}}$. Note that $\mathbb{E}[Z^f] = P_{x \sim D_{\mathcal{X}}}(f(x) \neq f^*(x))$. Hence note that for any single $f \in \mathcal{F}$,

$$P_S\left(P_{x \sim D_{\mathcal{X}}}(f(x) \neq f^*(x)) - \frac{1}{n} \sum_{t=1}^n \mathbf{1}\{f(x_t) \neq f^*(x_t)\} > \theta\right) \leq \exp\left(-\frac{n\theta^2}{2\mu + \frac{\theta}{3}}\right)$$

Let us write the R.H.S. above as δ , and hence, rewriting, we have that with probability at least $1 - \delta$ over sample,

$$P_{x \sim D_{\mathcal{X}}}(f(x) \neq f^*(x)) - \frac{1}{n} \sum_{t=1}^n \mathbf{1}\{f(x_t) \neq f^*(x_t)\} \leq \frac{\log(1/\delta)}{3n} + \sqrt{\frac{P_{x \sim D_{\mathcal{X}}}(f(x) \neq f^*(x)) \log(1/\delta)}{n}}$$

This upon further simplifying (use inequality $\sqrt{ab} \leq a/2 + b/2$) leads to the bound

$$P_{x \sim D_{\mathcal{X}}}(f(x) \neq f^*(x)) - \frac{2}{n} \sum_{t=1}^n \mathbf{1}\{f(x_t) \neq f^*(x_t)\} \leq \frac{2 \log(1/\delta)}{n}$$

Using union bound, we have that for any $\delta > 0$, with probability at least $1 - \delta$ over sample, simultaneously,

$$\forall f \in \mathcal{F} \quad P_{x \sim D_{\mathcal{X}}}(f(x) \neq f^*(x)) - \frac{2}{n} \sum_{t=1}^n \mathbf{1}\{f(x_t) \neq f^*(x_t)\} \leq \frac{2 \log(|\mathcal{F}|/\delta)}{n}$$

Since $\hat{\mathbf{y}} \in \mathcal{F}$, from the above we conclude that, for any $\delta > 0$, with probability at least $1 - \delta$ over sample,

$$P_{x \sim D_{\mathcal{X}}}(\hat{\mathbf{y}}(x) \neq f^*(x)) - \frac{2}{n} \sum_{t=1}^n \mathbf{1}\{\hat{\mathbf{y}}(x_t) \neq f^*(x_t)\} \leq \frac{2 \log(|\mathcal{F}|/\delta)}{n}$$

But note that by realizability assumption and the definition of $\hat{\mathbf{y}}$, we have that

$$\sum_{t=1}^n \mathbf{1}\{\hat{\mathbf{y}} \neq f^*(x_t)\} = \sum_{t=1}^n \mathbf{1}\{\hat{\mathbf{y}} \neq y_t\} = 0$$

and so, with probability at least $1 - \delta$ over sample,

$$P_{x \sim D_{\mathcal{X}}}(\hat{\mathbf{y}}(x) \neq f^*(x)) \leq \frac{2 \log(|\mathcal{F}|/\delta)}{n}$$