

# Lecture 13

## Statistical Learning: First Steps

Sasha Rakhlin

Oct 22, 2018

# Outline

Setup

Perceptron

# Outline

Setup

Perceptron

Supervised Learning: data  $\mathcal{S} = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$  are i.i.d. from *unknown* distribution  $\mathbf{P}$ .

*Learning algorithm:* a mapping  $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \mapsto \hat{f}_n$ .

Goals:

- ▶ **Prediction:** small expected loss

$$\mathbf{L}(\hat{f}_n) = \mathbb{E}_{X,Y} \ell(Y, \hat{f}_n(X)).$$

Here  $(X, Y) \sim \mathbf{P}$ . Interpretation: good prediction on a random example from same population.

- ▶ **Estimation:** small  $\|\hat{f}_n - f^*\|$ , or  $\|\hat{\theta} - \theta^*\|$ , where  $f^*$  or  $\theta^*$  are parameters of  $\mathbf{P}$  (e.g. regression function  $f^*(x) = \mathbb{E}[Y|X=x]$ , or  $f^*(x) = \langle \theta^*, x \rangle$ , etc).

In this course, we mostly focus on prediction, but will also outline connections between prediction and estimation.

Why not estimate the underlying distribution  $P$  first?

This is in general a harder problem than prediction. Consider classification. We might be attempting to learn parts/properties of the distribution that are irrelevant, while all we care about is the “boundary” between the two classes.

*Key difficulty:* our goals are in terms of unknown quantities related to unknown  $\mathbf{P}$ . Have to use empirical data instead. Preview of statistics.

For instance, we can calculate the *empirical loss* of  $f: \mathcal{X} \rightarrow \mathcal{Y}$

$$\widehat{\mathbf{L}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i))$$

## Quiz: what is random here?

1.  $\widehat{\mathbf{L}}(\mathbf{f})$  for a given fixed  $\mathbf{f}$
2.  $\widehat{\mathbf{f}}_n$
3.  $\widehat{\mathbf{L}}(\widehat{\mathbf{f}}_n)$
4.  $\mathbf{L}(\widehat{\mathbf{f}}_n)$
5.  $\mathbf{L}(\mathbf{f})$  for a given fixed  $\mathbf{f}$

It is important that these are understood before we proceed further.

Theoretical analysis of performance is typically easier if  $\hat{f}_n$  has closed form (in terms of the training data).

E.g. ordinary least squares  $\hat{f}_n(\mathbf{x}) = \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ .

Unfortunately, most ML and many Statistical procedures are not explicitly defined but arise as

- ▶ solutions to an optimization objective (e.g. logistic regression)
- ▶ as an iterative procedure without an immediately obvious objective function (e.g. AdaBoost, Random Forests, etc)



# The Gold Standard

Within the framework we set up, the smallest expected loss is achieved by the *Bayes optimal* function

$$f^* = \arg \min_f \mathbf{L}(f)$$

where the minimization is over all (measurable) prediction rules  $f: \mathcal{X} \rightarrow \mathcal{Y}$ .

The value of the lowest expected loss is called the *Bayes error*:

$$\mathbf{L}(f^*) = \inf_f \mathbf{L}(f)$$

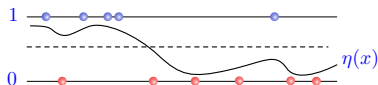
Of course, we cannot calculate any of these quantities since  $\mathbf{P}$  is unknown.

# Bayes Optimal Function

Bayes optimal function  $f^*$  takes on the following forms in these two particular cases:

- Binary classification ( $\mathcal{Y} = \{0, 1\}$ ) with the indicator loss:

$$f^*(x) = \mathbf{I}\{\eta(x) \geq 1/2\}, \quad \text{where} \quad \eta(x) = \mathbb{E}[Y|X = x]$$



- Regression ( $\mathcal{Y} = \mathbb{R}$ ) with squared loss:

$$f^*(x) = \eta(x), \quad \text{where} \quad \eta(x) = \mathbb{E}[Y|X = x]$$

**The big question:** is there a way to construct a learning algorithm with a guarantee that

$$L(\widehat{f}_n) - L(f^*)$$

is small for large enough sample size  $n$ ?

# Consistency

An algorithm that ensures

$$\lim_{n \rightarrow \infty} L(\hat{f}_n) = L(f^*) \quad \text{almost surely}$$

is called *consistent*. Consistency ensures that our algorithm is approaching the best possible prediction performance as the sample size increases.

The good news: consistency is possible to achieve.

- ▶ easy if  $\mathcal{X}$  is a finite or countable set
- ▶ not too hard if  $\mathcal{X}$  is infinite, and the underlying relationship between  $\mathbf{x}$  and  $\mathbf{y}$  is “continuous”

# The bad news...

In general, we cannot prove anything “interesting” about  $\mathbf{L}(\widehat{\mathbf{f}}_n) - \mathbf{L}(\mathbf{f}^*)$ ,  
unless we make further assumptions (incorporate prior knowledge).

What do we mean by “nothing interesting”? This is the subject of the so-called “No Free Lunch” Theorems. Unless we posit further assumptions,

- ▶ For any algorithm  $\widehat{\mathbf{f}}_n$ , any  $n$  and any  $\epsilon > 0$ , there exists a distribution  $\mathbf{P}$  such that  $\mathbf{L}(\mathbf{f}^*) = 0$  and

$$\mathbb{E}\mathbf{L}(\widehat{\mathbf{f}}_n) \geq \frac{1}{2} - \epsilon$$

- ▶ For any algorithm  $\widehat{\mathbf{f}}_n$ , and any sequence  $\mathbf{a}_n$  that converges to 0, there exists a probability distribution  $\mathbf{P}$  such that  $\mathbf{L}(\mathbf{f}^*) = 0$  and for all  $n$

$$\mathbb{E}\mathbf{L}(\widehat{\mathbf{f}}_n) \geq \mathbf{a}_n$$

Reference: (Devroye, Györfi, Lugosi: *A Probabilistic Theory of Pattern Recognition*),  
(Bousquet, Boucheron, Lugosi, 2004)

# is this really “bad news”?

Not really. We always have some domain knowledge.

Two ways of incorporating prior knowledge:

- ▶ Direct way: assumptions on distribution  $\mathbf{P}$  (e.g. margin)
- ▶ Indirect way: redefine the goal to perform as well as a reference set  $\mathcal{F}$  of predictors:

$$\mathbf{L}(\widehat{\mathbf{f}}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f)$$

$\mathcal{F}$  encapsulates our *inductive bias*.

We often make both of these assumptions.

# Outline

Setup

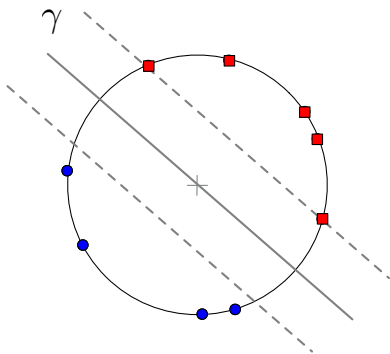
Perceptron

We start our study of Statistical Learning with the classical Perceptron algorithm.

Reason: simplicity. We will give a three-line proof of Perceptron, followed by two interesting consequences with one-line proofs each. These consequences are, perhaps, the easiest nontrivial statistical guarantees I can think of.



# Perceptron



# Perceptron

$(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T) \in \mathcal{X} \times \{\pm 1\}$  ( $T$  may or may not be same as  $n$ )

Maintain a hypothesis  $\mathbf{w}_t \in \mathbb{R}^d$  (initialize  $\mathbf{w}_1 = \mathbf{0}$ ).

On round  $t$ ,

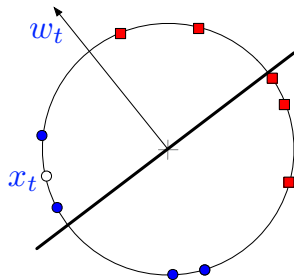
- ▶ Consider  $(\mathbf{x}_t, \mathbf{y}_t)$
- ▶ Form prediction  $\hat{\mathbf{y}}_t = \text{sign}(\langle \mathbf{w}_t, \mathbf{x}_t \rangle)$
- ▶ If  $\hat{\mathbf{y}}_t \neq \mathbf{y}_t$ , update

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{y}_t \mathbf{x}_t$$

else

$$\mathbf{w}_{t+1} = \mathbf{w}_t$$

# Perceptron



For simplicity, suppose all data are in a unit ball,  $\|x_t\| \leq 1$ .

Definition of margin of  $(x_1, y_1), \dots, (x_T, y_T)$ :

$$\gamma = \max_{\|w\|=1} \min_{i \in [T]} y_i \langle w, x_i \rangle$$

or  $\gamma = 0$  if no margin.

**Theorem (Novikoff '62):** Perceptron makes at most  $1/\gamma^2$  mistakes (and corrections) on **any** sequence of examples with margin  $\gamma$ .

**Proof:** Let  $m$  be the number of mistakes after  $T$  iterations. If a mistake is made on round  $t$ ,

$$\|w_{t+1}\|^2 = \|w_t + y_t x_t\|^2 \leq \|w_t\|^2 + 2y_t \langle w_t, x_t \rangle + 1 \leq \|w_t\|^2 + 1.$$

Hence,

$$\|w_T\|^2 \leq m.$$

For optimal hyperplane  $w^*$

$$\gamma \leq \langle w^*, y_t x_t \rangle = \langle w^*, w_{t+1} - w_t \rangle.$$

Hence (adding and canceling),

$$m\gamma \leq \langle w^*, w_T \rangle \leq \|w_T\| \leq \sqrt{m}.$$

More formally, for any  $T$  and  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)$ ,

$$\sum_{t=1}^T \mathbf{I}\{\mathbf{y}_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle \leq 0\} \leq \frac{D^2}{\gamma^2}$$

where  $\gamma = \gamma(\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  is margin and  $D = D(\mathbf{x}_{1:T}, \mathbf{y}_{1:T}) = \max_t \|\mathbf{x}_t\|$ .