# 10-301/601: Introduction to Machine Learning Lecture 15 – Learning Theory (Finite Case)

Matt Gormley & Henry Chai 3/10/25

#### **Front Matter**

- Announcements
  - HW5 released 2/27, due 3/16 at 11:59 PM
  - Exam 1 Exit Poll due 3/10 (today!) at 11:59 PM
  - Peer tutoring information will be posted to Piazza some time this week

### Statistical Learning Theory Model

1. Data points are generated i.i.d. from some *unknown* distribution

$$\mathbf{x}^{(n)} \sim p^*(\mathbf{x})$$

2. Labels are generated from some *unknown* function  $y^{(n)} = c^*(x^{(n)})$ 

- 3. The learning algorithm chooses the hypothesis (or classifier) with lowest training error rate from a specified hypothesis set,  $\mathcal{H}$
- 4. Goal: return a hypothesis (or classifier) with low *true* error rate

### Statistical Learning Theory Model

1. Data points are generated i.i.d. from some *unknown* distribution

$$\mathbf{x}^{(n)} \sim p^*(\mathbf{x})$$

2. Labels are generated from some *unknown* function

$$y^{(n)} = c^*(x^{(n)}) \in \{-1, +1\}$$

- 3. The learning algorithm chooses the hypothesis (or classifier) with lowest training error rate from a specified hypothesis set,  $\mathcal{H}$
- 4. Goal: return a hypothesis (or classifier) with low *true* error rate

#### Types of Error

- True error rate
  - Actual quantity of interest in machine learning
  - How well your hypothesis will perform on average across all possible data points
- Test error rate
  - Used to evaluate hypothesis performance
  - Good estimate of your hypothesis's true error
- Validation error rate
  - Used to set hypothesis hyperparameters
  - Slightly "optimistic" estimate of your hypothesis's true error
- Training error rate
  - Used to set model parameters
  - Very "optimistic" estimate of your hypothesis's true error

### Types of Risk (a.k.a. Error)

Expected risk of a hypothesis h (a.k.a. true error)

$$R(h) = P_{\boldsymbol{x} \sim p^*} (c^*(\boldsymbol{x}) \neq h(\boldsymbol{x}))$$

• Empirical risk of a hypothesis h (a.k.a. training error)

$$\widehat{R}(h) = P_{x \sim \mathcal{D}} \left( c^*(x) \neq h(x) \right)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \mathbb{1} \left( c^*(x^{(n)}) \neq h(x^{(n)}) \right)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \mathbb{1} \left( y^{(n)} \neq h(x^{(n)}) \right)$$

where  $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^{N}$  is the training data set and  $x \sim \mathcal{D}$  denotes a point sampled uniformly at random from  $\mathcal{D}$ 

### Three Hypotheses of Interest

1. The true function,  $c^*$ 

2. The expected risk minimizer,

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R(h)$$

3. The *empirical risk minimizer*,

$$\hat{h} = \operatorname*{argmin}_{h \in \mathcal{H}} \hat{R}(h)$$

## Poll Question 1: Which of the following are *always* true?

A. 
$$c^* = h^*$$
  
B.  $c^* = \hat{h}$   
C.  $h^* = \hat{h}$   
D.  $c^* = h^* = \hat{h}$   
E. None of the above  
F. **TOXIC**

• The *true function*, *c*\*

The expected risk minimizer,

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R(h)$$

• The empirical risk minimizer,

$$\hat{h} = \operatorname*{argmin}_{h \in \mathcal{H}} \hat{R}(h)$$

#### **Key Question**

• Given a hypothesis with zero/low training error, what can we say about its true error?

#### **PAC Learning**

• PAC = Probably Approximately Correct

• PAC Criterion:

$$P(|R(h) - \hat{R}(h)| \le \epsilon) \ge 1 - \delta \ \forall \ h \in \mathcal{H}$$

for some  $\epsilon$  (difference between expected and empirical risk) and  $\delta$  (probability of "failure")

• We want the PAC criterion to be satisfied for  ${\cal H}$  with small values of  $\epsilon$  and  $\delta$ 

### Sample Complexity

- The sample complexity of an algorithm/hypothesis set,  $\mathcal H$ , is the number of labelled training data points needed to satisfy the PAC criterion for some  $\delta$  and  $\epsilon$
- Four cases
  - Realizable vs. Agnostic
    - Realizable  $\rightarrow c^* \in \mathcal{H}$
    - Agnostic  $\rightarrow c^*$  might or might not be in  $\mathcal{H}$
  - Finite vs. Infinite
    - Finite  $\rightarrow |\mathcal{H}| < \infty$
    - Infinite  $\rightarrow |\mathcal{H}| = \infty$

### Theorem 1: Finite, Realizable Case

• For a finite hypothesis set  $\mathcal{H}$  s.t.  $c^* \in \mathcal{H}$  and arbitrary distribution  $p^*$ , if the number of labelled training data points satisfies

$$M \ge \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

then with probability at least  $1-\delta$ , all  $h\in\mathcal{H}$  with  $\widehat{R}(h)=0$  have  $R(h)\leq\epsilon$ 

- 1. Assume there are K "bad" hypotheses in  $\mathcal{H}$ , i.e.,  $h_1, h_2, \dots, h_K$  that all have  $R(h_k) > \epsilon$
- 2. Pick one bad hypothesis,  $h_k$ 
  - A. Probability that  $h_k$  correctly classifies the first training data point  $< 1 \epsilon$
  - B. Probability that  $h_k$  correctly classifies all M training data points  $< (1 \epsilon)^M$
- 3. Probability that at least one bad hypothesis correctly classifies all M training data points =
   P(h<sub>1</sub> correctly classifies all M training data points U h<sub>2</sub> correctly classifies all M training data points U :
  - $\cup$   $h_K$  correctly classifies all M training data points)

 $P(h_1 \text{ correctly classifies all } M \text{ training data points } \cup h_2 \text{ correctly classifies all } M \text{ training data points } \cup \vdots$ 

 $\cup$   $h_K$  correctly classifies all M training data points)

$$\leq \sum_{k=1}^{K} P(h_k \text{ correctly classifies all } M \text{ training data points})$$

by the union bound: 
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$
  
  $\leq P(A) + P(B)$ 

3/10/25 **14** 

$$\sum_{k=1}^{K} P(h_k \text{ correctly classifies all } M \text{ training data points})$$

$$< k(1 - \epsilon)^M \le |\mathcal{H}|(1 - \epsilon)^M$$

because  $k \leq |\mathcal{H}|$ 

- 3. Probability that at least one bad hypothesis correctly classifies all M training data points  $\leq |\mathcal{H}|(1-\epsilon)^{M}$
- 4. Using the fact that  $1 x \le \exp(-x) \ \forall x$ ,  $|\mathcal{H}|(1 \epsilon)^M \le |\mathcal{H}| \exp(-\epsilon)^M = |\mathcal{H}| \exp(-M\epsilon)$
- 5. Probability that at least one bad hypothesis correctly classifies all M training data points  $\leq |\mathcal{H}| \exp(-M\epsilon)$ , which we want to be low, i.e.,  $|\mathcal{H}| \exp(-M\epsilon) \leq \delta$

3/10/25 **15** 

$$|\mathcal{H}| \exp(-M\epsilon) \le \delta \to \exp(-M\epsilon) \le \frac{\delta}{|\mathcal{H}|}$$

$$\to -M\epsilon \le \ln\left(\frac{\delta}{|\mathcal{H}|}\right)$$

$$\to M \ge \frac{1}{\epsilon} \left(-\ln\left(\frac{\delta}{|\mathcal{H}|}\right)\right)$$

$$\to M \ge \frac{1}{\epsilon} \left(\ln\left(\frac{|\mathcal{H}|}{\delta}\right)\right)$$

$$\to M \ge \frac{1}{\epsilon} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$$

6. Given  $M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$  labelled training data points, the probability that  $\exists$  a bad hypothesis  $h_k \in \mathcal{H}$  with  $R(h_k) > \epsilon$  and  $\hat{R}(h_k) = 0$  is  $\leq \delta$ 

Given  $M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$  labelled training data points, the probability that all hypotheses  $h_k \in \mathcal{H}$  with  $R(h_k) > \epsilon$  have  $\hat{R}(h_k) > 0$  is  $\geq 1 - \delta$ 

### Aside: Proof by Contrapositive

- The contrapositive of a statement  $A \Rightarrow B$  is  $\neg B \Rightarrow \neg A$
- A statement and its contrapositive are logically equivalent, i.e.,  $A \Rightarrow B$  means that  $\neg B \Rightarrow \neg A$
- Example: "it's raining ⇒ Henry brings am umbrella"

is the same as saying

"Henry didn't bring an umbrella ⇒ it's not raining "

6. Given  $M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$  labelled training data points, the probability that  $\exists$  a bad hypothesis  $h_k \in \mathcal{H}$  with  $R(h_k) > \epsilon$  and  $\hat{R}(h_k) = 0$  is  $\leq \delta$ 

Given  $M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$  labelled training data points, the probability that all hypotheses  $h_k \in \mathcal{H}$  with  $R(h_k) > \epsilon$  have  $\hat{R}(h_k) > 0$  is  $\geq 1 - \delta$ 

6. Given  $M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$  labelled training data points, the probability that all hypotheses  $h_k \in \mathcal{H}$  with  $R(h_k) > \epsilon$  have  $\hat{R}(h_k) > 0$  is  $\geq 1 - \delta$ 

Given  $M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$  labelled training data points, the probability that all hypotheses  $h_k \in \mathcal{H}$  with  $\hat{R}(h_k) = 0$  have  $R(h_k) \leq \epsilon$  is  $\geq 1 - \delta$  (proof by contrapositive)

#### Poll Question 2:

### Hint - Recall

M

$$\geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

• Let  $\mathcal{H}$  be the set of all *conjunctions* over M Boolean variables,  $\mathbf{x} \in \{0,1\}^M$ ; examples of conjunctions are

$$h(x) = x_1(1 - x_2)x_4x_{10}$$

• 
$$h(\mathbf{x}) = (1 - x_3)(1 - x_4)x_8$$

- Assuming  $c^* \in \mathcal{H}$ , if M=10,  $\epsilon=0.1$ , and  $\delta=0.01$ , at least how many labelled examples do we need to satisfy the PAC criterion using Theorem 1?
- A. 1 (TOXIC)
- B.  $10(2 \ln 10 + \ln 100) \approx 92$  F.  $100(2 \ln 10 + \ln 10) \approx 691$
- C.  $10(3 \ln 10 + \ln 100) \approx 116$  G.  $100(3 \ln 10 + \ln 10) \approx 922$
- D.  $10(10 \ln 2 + \ln 100) \approx 116$  H.  $100(10 \ln 2 + \ln 10) \approx 924$
- E.  $10(10 \ln 3 + \ln 100) \approx 156$  I.  $100(10 \ln 3 + \ln 10) \approx 1329$

### Theorem 1: Finite, Realizable Case

• For a finite hypothesis set  $\mathcal{H}$  s.t.  $c^* \in \mathcal{H}$  and arbitrary distribution  $p^*$ , if the number of labelled training data points satisfies

$$M \ge \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

then with probability at least  $1 - \delta$ , all  $h \in \mathcal{H}$  with  $\widehat{R}(h) = 0$  have  $R(h) \leq \epsilon$ 

• Making the bound tight and solving for  $\epsilon$  gives...

### Statistical Learning Theory Corollary

• For a finite hypothesis set  $\mathcal{H}$  s.t.  $c^* \in \mathcal{H}$  and arbitrary distribution  $p^*$ , given a training data set S s.t. |S| = M, all  $h \in \mathcal{H}$  with  $\hat{R}(h) = 0$  have

$$R(h) \le \frac{1}{M} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

with probability at least  $1 - \delta$ .

### Theorem 2: Finite, Agnostic Case

• For a finite hypothesis set  ${\mathcal H}$  and arbitrary distribution  $p^*$ , if the number of labelled training data points satisfies

$$M \ge \frac{1}{2\epsilon^2} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)$$

then with probability at least  $1-\delta$ , all  $h\in\mathcal{H}$  satisfy

$$\left| R(h) - \widehat{R}(h) \right| \le \epsilon$$

- Bound is inversely quadratic in  $\epsilon$ , e.g., halving  $\epsilon$  means we need four times as many labelled training data points
- Again, making the bound tight and solving for  $\epsilon$  gives...

3/10/25 **24** 

### Statistical Learning Theory Corollary

• For a finite hypothesis set  $\mathcal H$  and arbitrary distribution  $p^*$ , given a training data set S s.t. |S|=M, all  $h\in\mathcal H$  have

$$R(h) \le \hat{R}(h) + \sqrt{\frac{1}{2M}} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)$$

with probability at least  $1 - \delta$ .

### What happens when $|\mathcal{H}| = \infty$ ?

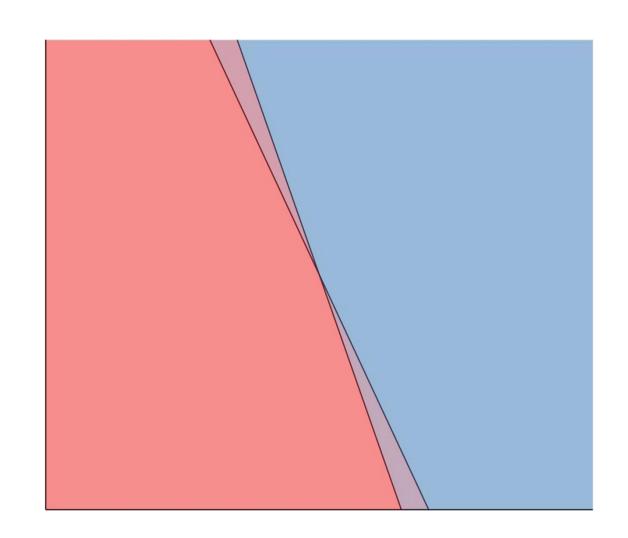
• For a finite hypothesis set  $\mathcal H$  and arbitrary distribution  $p^*$ , given a training data set S s.t. |S|=M, all  $h\in\mathcal H$  have

$$R(h) \le \hat{R}(h) + \sqrt{\frac{1}{2M}} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)$$

with probability at least  $1 - \delta$ .

#### Intuition

For most infinite hypothesis sets  $\mathcal{H}$ , many hypotheses in  $\mathcal{H}$  will behave very similarly

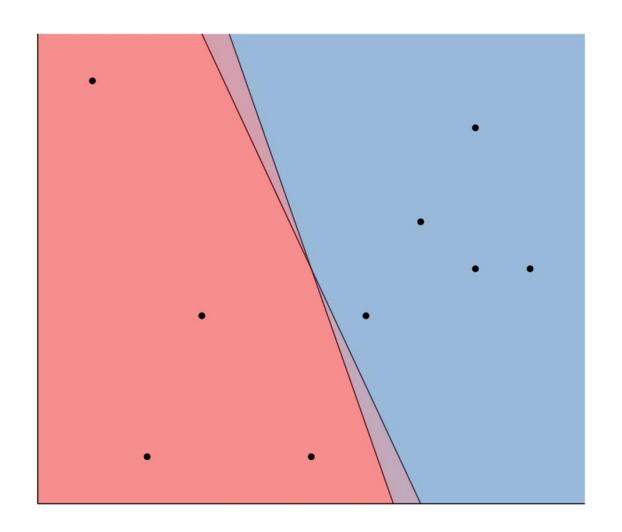


**27** 

#### Intuition

For most infinite hypothesis sets  $\mathcal{H}$ , many hypotheses in  $\mathcal{H}$  will behave very similarly

Relative to a given dataset, these two hypotheses are *identical*!



3/10/25 **2**