

10-301/601: Introduction to Machine Learning

Lecture 17 – Learning Theory (Finite Case)

Henry Chai

5/28/25

Front Matter

- Announcements:
 - HW4 released on 5/23, due 5/28 (today!) at 11:59 PM
 - Midterm on 5/30 at 9:30 AM in BH A36
 - Lectures 1 – 14 are in-scope; **this week's lectures will not be tested on the midterm**
 - Recitation on 5/29 will be a review of the practice problems

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

$$y^{(n)} = c^*(\mathbf{x}^{(n)})$$

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

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)

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

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

$$\begin{aligned}\hat{R}(h) &= P_{\vec{x} \sim D} (c^*(\vec{x}) \neq h(\vec{x})) \\ &= \frac{1}{N} \sum_{n=1}^N \underset{\substack{\uparrow \\ \text{indicator}}}{\mathbb{1}} (c^*(\vec{x}^{(n)}) \neq h(\vec{x}^{(n)}))\end{aligned}$$

where $D = \{(\vec{x}^{(n)}, y^{(n)})\}_{n=1}^N$ is a training dataset $\leadsto \vec{x} \sim D$ means uniform sampling

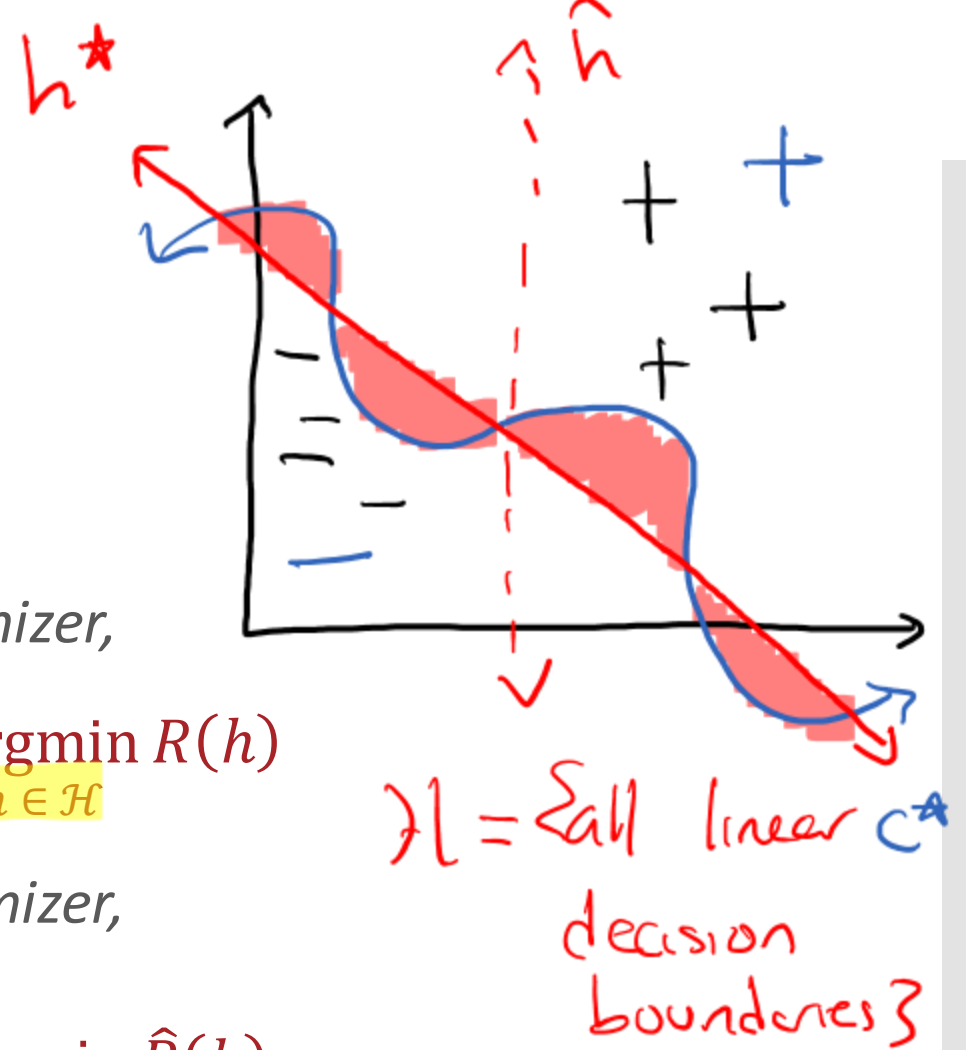
Three Hypotheses of Interest

1. The *true function*, c^*

$$h^* = \operatorname{argmin} R(h)$$

- ### 3. The *empirical risk minimizer*,

$$\hat{h} = \operatorname{argmin} \hat{R}(h)$$



0 surveys completed



0 surveys underway

Select all that apply: Which of the following statements is *always* true?

$$c^* = h^*$$

$$c^* = \hat{h}$$

$$h^* = \hat{h}$$

All of the above ($c^* = h^* = \hat{h}$)

None of the above

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:

$$P(|R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta \quad \forall h \in \mathcal{H}$$

for some ϵ (difference between expected and empirical risk) and δ (probability of “failure”)

- We want the PAC criterion to be satisfied for \mathcal{H} with small values of ϵ and δ

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 δ and ϵ
 - 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

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

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

Proof of Theorem 1: Finite, Realizable Case

1. Assume there are K "bad" hypotheses in $H: \{h_1, h_2, \dots, h_K\}$ with $R(h_i) > \epsilon$
2. Think about one bad hypothesis, h_i
 - a. $P(h_i \text{ correctly classifies the first training data point}) < 1 - \epsilon$
 - b. $P(h_i \text{ correctly classifies all } M \text{ training data points}) < (1 - \epsilon)^M$
3. $P(\text{that at least one bad hypothesis correctly classifies all } M \text{ training data points}) =$
 $P(\hat{R}(h_1) = 0 \cup \hat{R}(h_2) = 0 \cup \dots \cup \hat{R}(h_K) = 0)$

Proof of Theorem 1: Finite, Realizable Case

4. Use the union bound

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \\ \leq P(A) + P(B)$$



$$P(\hat{R}(h_1)=0 \cup \hat{R}(h_2)=0 \cup \dots \cup \hat{R}(h_K)=0) \\ \leq \sum_{k=1}^K P(\hat{R}(h_k)=0) < \sum_{k=1}^K (1-\epsilon)^M$$

$$= K(1-\epsilon)^M \\ P(\text{at least one bad hypothesis tricks us}) < K(1-\epsilon)^M$$

Proof of
Theorem 1:
Finite,
Realizable Case

$$5. \quad K(1-\epsilon)^M < |H|(1-\epsilon)^M$$

6. Use the fact that $1-x \leq \exp(-x)$
 $\forall x$

$$\Rightarrow (1-\epsilon)^M \leq \exp(-\epsilon M)$$

$$|H|(1-\epsilon)^M \leq |H| \exp(-\epsilon M) \leq \delta$$

$$7. \Rightarrow \exp(-\epsilon M) \leq \frac{\delta}{|H|}$$

$$\Rightarrow -\epsilon M \leq \ln(\delta/|H|)$$

$$\Rightarrow M \geq \frac{1}{\epsilon} \ln(\delta/|H|)$$

$$\Rightarrow M \geq \frac{1}{\epsilon} \ln(|H|/\delta) = \frac{1}{\epsilon} (\ln(|H|) + \ln(1/\delta))$$

Proof of
Theorem 1:
Finite,
Realizable Case

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



Given $M \geq \frac{1}{\epsilon} (\ln(|H|) + \ln(1/\delta))$ t.d.p., the probability that all bad hypotheses with $R(h_i) > \epsilon$ have $\hat{R}(h_i) > 0$ 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 \Rightarrow Henry brings an umbrella”
is the same as saying
“Henry didn’t bring an umbrella \Rightarrow it’s not raining”

Proof of Theorem 1: Finite, Realizable Case

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 $\geq 1 - \delta$

A \Rightarrow B

Proof of Theorem 1: Finite, Realizable Case

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 $\geq 1 - \delta$

$\neg B \Rightarrow \neg A$
(proof by contrapositive)



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

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

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

- Making the bound tight and solving for ϵ 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

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

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

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

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ satisfy $|R(h) - \hat{R}(h)| \leq \epsilon$

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

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

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

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

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

with probability at least $1 - \delta$.

Key Takeaways

- Statistical learning theory model
- Expected vs. empirical risk of a hypothesis
- Four possible cases of interest
 - realizable vs. agnostic
 - finite vs. infinite
- Sample complexity bounds and statistical learning theory corollaries for finite hypothesis sets