

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

independent and
identically distributed

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

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

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

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)

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

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

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

$\vec{x} \sim D \Rightarrow \vec{x}$ is uniformly at random
chosen from $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

$$\hat{R}(h) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(h(\vec{x}^{(i)}) \neq c^*(\vec{x}^{(i)}))$$

indicator function training error rate

Three Hypotheses of Interest

1. The *true function*, c^*

2. The *expected risk minimizer*,

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

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. $c^* = h^*$ 38%

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

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

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

E. None of the above 42%

F. TOXIC

$\rightarrow R(c^*) = 0$

• The *true function*, c^*

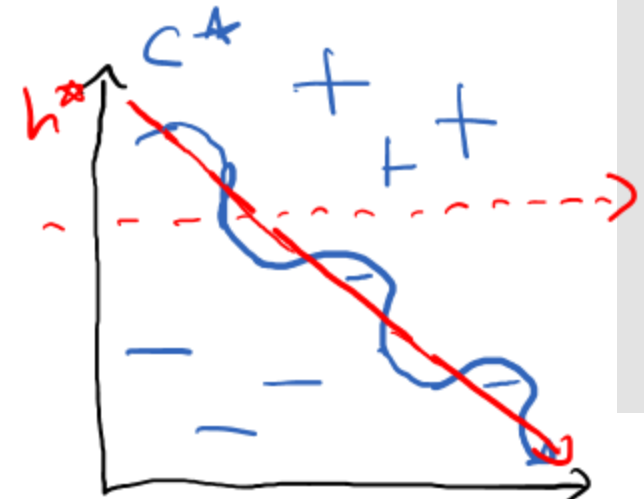
• The *expected risk minimizer*,

$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R(h) \stackrel{?}{=} 0$

• The *empirical risk minimizer*,

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

$\mathcal{H} = \{ \text{all linear decision boundaries in 2D} \}$



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:

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

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 → $c^* \in \mathcal{H}$

• Agnostic → c^* might or might not be in \mathcal{H}

→ • Finite vs. Infinite

• Finite → $|\mathcal{H}| < \infty$

• Infinite → $|\mathcal{H}| = \infty$

most realistic

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 $R(h) \leq \epsilon$

(sketch)

Proof of Theorem 1: Finite, Realizable Case

1. Assume the worst !!
Assume that every hypothesis in H is bad! ($R(h) > \epsilon$)
2. The probability that **a** bad hypothesis "tricks" me ($\hat{R}(h) = 0$) is tiny! And shrinks as M grows
3. The probability that any of the bad hypotheses "tricks" me is pretty small and also shrinks as $M \uparrow$

Proof of
Theorem 1:
Finite,
Realizable Case

do some math

⋮

$$P(\text{at least one "bad" hypothesis correctly classifies } M \text{ training data points}) \leq |H| (1 - \epsilon)^M \leq \delta$$

do some more math

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

Proof of
Theorem 1:
Finite,
Realizable Case

Given $M \geq \frac{1}{\epsilon} (\ln |H| + \ln(\frac{1}{\delta}))$ labelled
training data points sampled from p^* ,
the probability \exists some hypothesis $h \in H$
with $R(h) > \epsilon$ and $\hat{R}(h) = 0$ is $\leq \delta$



\Rightarrow Given $M \geq$

the probability that all bad hypotheses
 $h \in H$ with $R(h) > \epsilon$ have $\hat{R}(h) > 0$ is $\geq 1 - \delta$
sampled from p^* ,

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

Given $M \geq \frac{1}{\epsilon} (\ln(H) + \ln(\frac{1}{\delta}))$ labelled training data points sampled from \mathcal{P}^* , the probability that all $h \in H$ with $\underbrace{R(h) > \epsilon}_A \Rightarrow \underbrace{\hat{R}(h) > 0}_B$ is $\geq 1 - \delta$

\Leftrightarrow

$\underbrace{\neg B}_{\hat{R}(h) = 0} \Rightarrow \underbrace{\neg A}_{R(h) \leq \epsilon} \quad \square$

Poll Question 2:

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

• $h(\mathbf{x}) = x_1(1 - x_2)x_4x_{10}$ *

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

• 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. $10(10 \ln 3 + \ln 100) \approx \underline{156}$ I. $100(10 \ln 3 + \ln 10) \approx 1329$

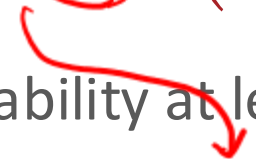
$\frac{1}{0.01} = 100$

$\ln 3^{10}$
 $= 10 \ln 3$

100%
13%
30%
38%

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 $R(h) \leq \epsilon$

- 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

$$\underbrace{|R(h) - \hat{R}(h)|}_{\leq \epsilon}$$

- 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

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

with probability at least $1 - \delta$.

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

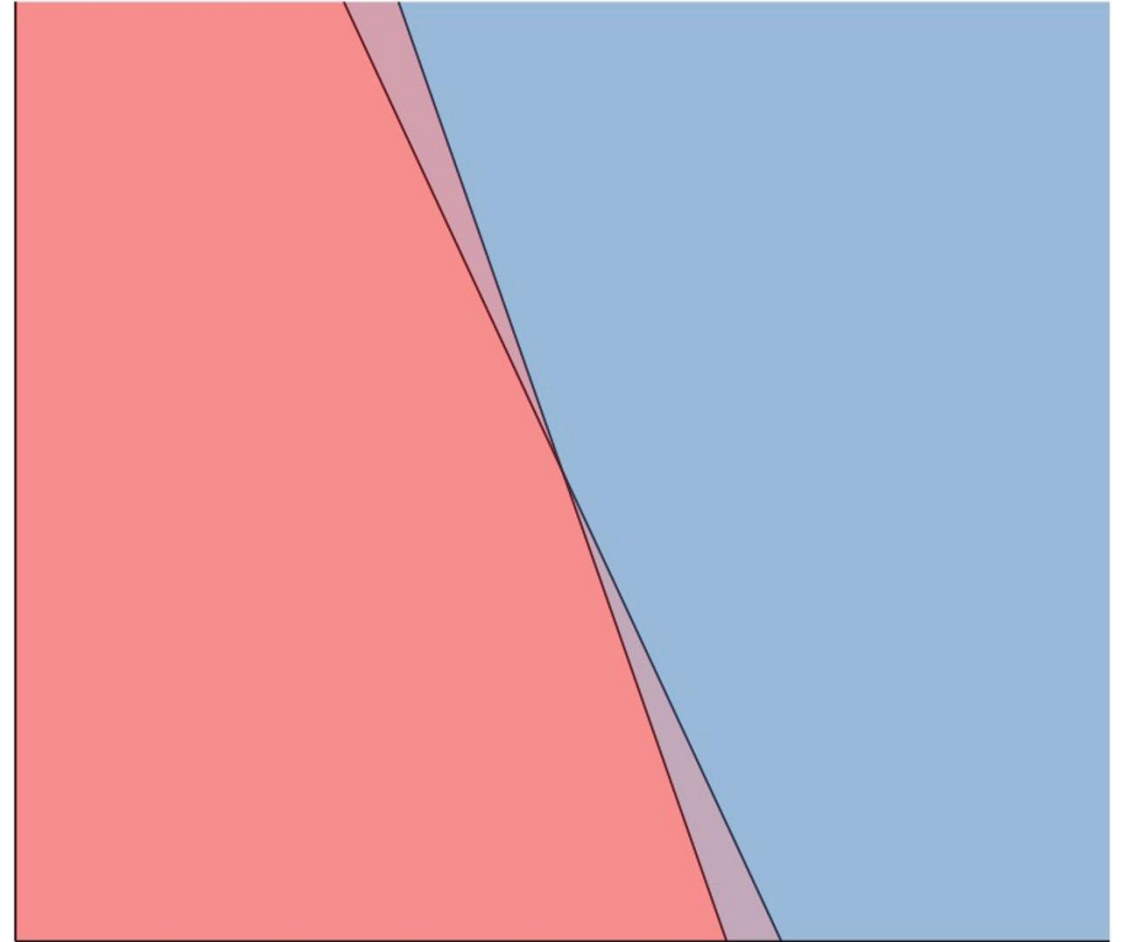
- For a finite hypothesis set \mathcal{H} and arbitrary distribution p^* , given a training data set S s.t. $|S| = M$, all $h \in \mathcal{H}$ have

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

with probability at least $1 - \delta$.

Intuition

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



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*!

