10-301/601: Introduction to Machine Learning Lecture 15 – Learning Theory

Geoff Gordon

with thanks to Matt Gormley & Henry Chai

Statistical learning theory setup

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 $y^{(n)} = c^*(\mathbf{x}^{(n)}) \in \{-1, +1\}$ note: **binary** classification
- 3. The learning algorithm chooses the hypothesis (classifier) with lowest *training* error rate from a specified hypothesis set, **%**
- 4. Goal: return a hypothesis (or classifier) with low *true* error rate (measure on *test* set)

Types of Error

- True error rate
 - Actual quantity of interest for learning
 - How well your hypothesis will perform on new samples
- Training error rate
 - Used to choose $h \in \mathcal{H}$ (e.g., fit model parameters)
 - May be a very optimistic estimate of true error

Types of Error

- True error rate
 - Actual quantity of interest for learning
 - How well your hypothesis will perform on new samples
- Training error rate
 - Used to choose $h \in \mathcal{H}$ (e.g., fit model parameters)
 - May be a very optimistic estimate of true error
- Test error rate
 - Used to evaluate hypothesis performance
 - Good estimate of true error (w/ enough test data)
- Validation error rate
 - Used to help choose \(\mathcal{H}\) (e.g., set hyperparameters)
 - Somewhat optimistic estimate of true error

Error rate is also called risk

True error rate = (true) risk — unknown

$$R(h) = \mathbb{E}_{x \sim p^{\star}} \left[\mathbb{T} \left(c^{\star}(x) \neq h(x) \right) \right]$$

Training error rate = empirical risk — we can measure this

$$\hat{R}(h) = \mathbb{E}_{x \sim D} \left(\mathbb{I} \left(c^*(x) \neq h(x) \right) \right)$$

$$= \frac{1}{M} \sum_{i=1}^{M} \mathbb{I} \left[y^{(i)} \neq h(x^{(i)}) \right]$$

Three classifiers

1. The *true classifier*, c^* : best answer but may be unachievable

2. The (true) risk minimizer (best achievable answer):

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

3. The *empirical risk minimizer* (the only one of the three that we can actually know)

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

Three classifiers

1. The *true classifier*, c^* : best answer but may be unachievable

2. The (true) risk minimizer (best achievable answer):

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

3. The *empirical risk minimizer* (the only one of the three that we can actually know)

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

Overfitting

Recall: **overfitting** = difference between true error rate and training error rate = $\frac{1}{2}(\zeta) - \frac{1}{2}(\zeta)$

- Goal for today: predict and control overfitting for ERM
 - by finding (and proving) conditions that keep $\hat{R}(h)$ close to R(h)

Bound on overfitting

PAC = <u>Probably Approximately Correct criterion</u>

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

for some ϵ (difference between true and empirical risk) and δ (probability of "failure")

why the name?

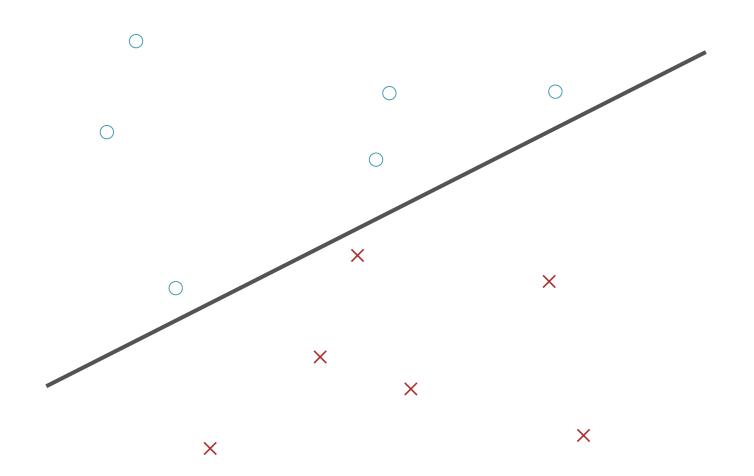
Sample Complexity

- ullet We will do ERM on some ${\mathscr H}$ with M training examples
- ullet We want to satisfy the PAC criterion with $small\ \epsilon$ and δ
- Chief levers: ${\mathcal H}$ and M
- •Sample complexity (of ERM on \mathcal{H}) = the M we need in order to satisfy the PAC criterion for a given ϵ and δ

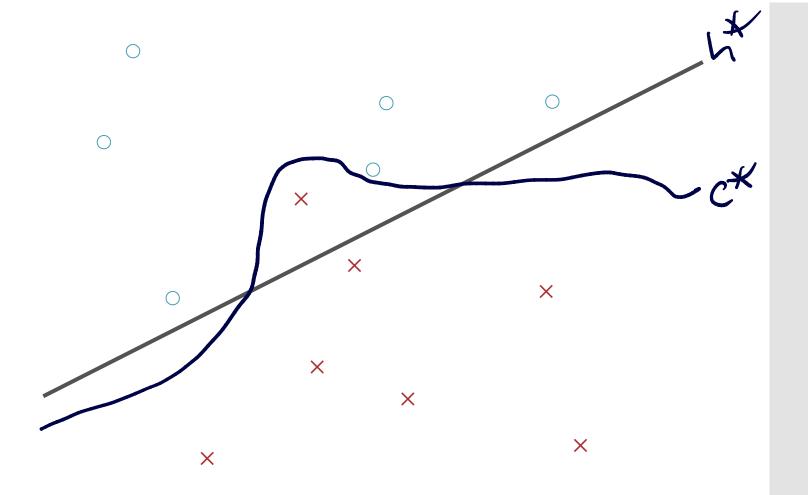
Four cases

- 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 $\left| \mathcal{H} \right| < \infty$
 - •Infinite \rightarrow $\left| \mathcal{H} \right| = \infty$

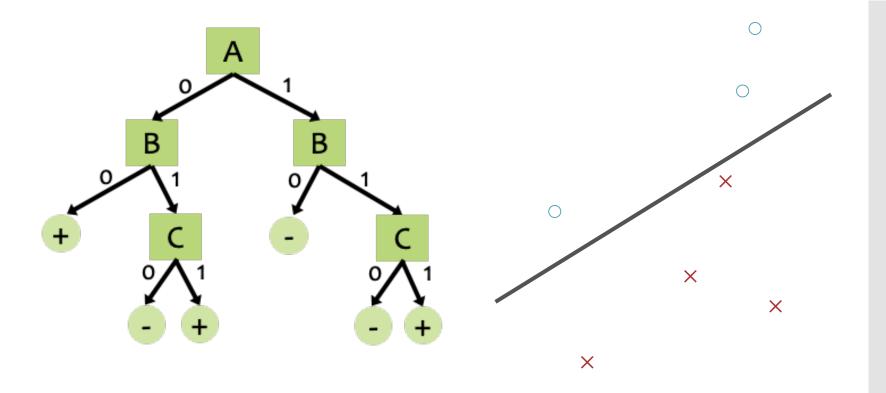
Realizable vs. agnostic



Realizable vs. agnostic



Finite vs. infinite $|\mathcal{H}|$



Decision trees of bounded depth on discrete attributes
 vs. linear separators in 2D

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**

ullet The true classifier, c^*

•The risk minimizer,

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

•The empirical risk minimizer,

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

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

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

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

$$\hat{R}(h) = 0$$
 have $R(h) \le \epsilon$

We will prove this over the next few slides

Theorem 1: Finite, Realizable Case

Bound is linear in $\frac{1}{\epsilon}$ thesis set \mathcal{H} s.t. $c^* \in \mathcal{H}$ and arbitrary finear in $\frac{1}{\epsilon}$ f the number of labelled training data points sfies

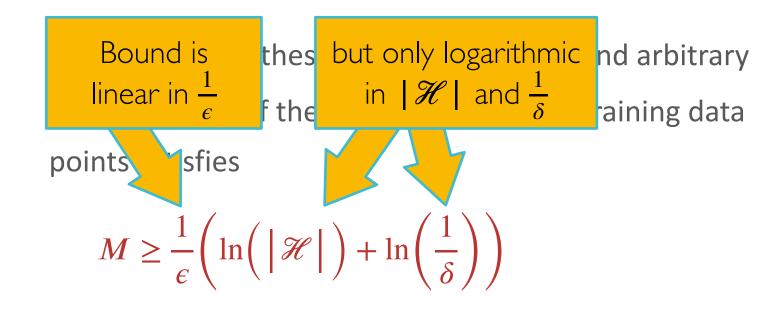
$$M \ge \frac{1}{\epsilon} \left(\ln \left(\left| \mathcal{H} \right| \right) + \ln \left(\frac{1}{\delta} \right) \right)$$

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

$$\hat{R}(h) = 0$$
 have $R(h) \le \epsilon$

We will prove this over the next few slides

Theorem 1: Finite, Realizable Case

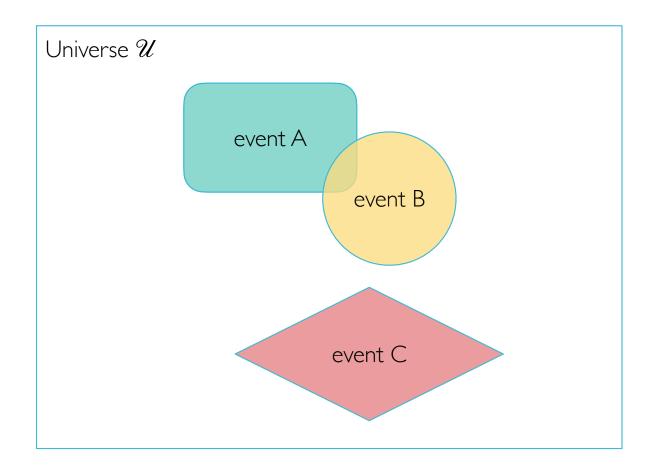


then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with

$$\hat{R}(h) = 0$$
 have $R(h) \le \epsilon$

We will prove this over the next few slides

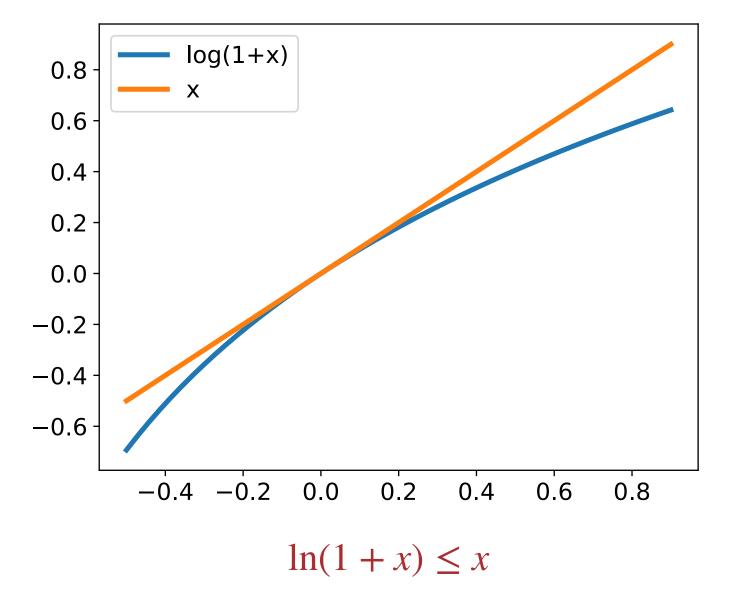
Union bound



$$P(A \text{ or } B \text{ or } C) \leq P(A) + P(B) + P(C)$$

 $P(\text{some event happens}) \leq \text{sum of probabilities}$

Bound on log



Notation for conclusion

Theorem said: if

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

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

$$\hat{R}(h) = 0$$
 have $R(h) \le \epsilon$

- Write E for event: $\exists h \in \mathcal{H}$ with $\hat{R}(h) = 0$, $R(h) > \epsilon$
- Theorem's conclusion is $P(E) < \delta$

Bound for one hypothesis

• Consider some h with $R(h) > \epsilon$. What's $P(\hat{R}(h) = 0)$?

Bound for *k* hypotheses

•Suppose there are k hypotheses with $R(h) > \epsilon$. What's

$$P(E)$$
, i.e., $P(\hat{R}(h_1) = 0 \text{ or } \hat{R}(h_2) = 0 \text{ or } \dots)$?
 $P(E) \leftarrow k (1 - E)^{M} \leftarrow |\mathcal{H}| (1 - E)^{M}$

Theorem will be true if $P(E) \leq \delta$

we have

Solve for M

$$P(E) < |\mathcal{H}| (1 - e)^{M}$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n P(E)| < |n |P(| + M |n (1 - e))$$

$$|n P(E)| < |n P(E)| < |n P(E)|$$

$$|n P(E)| < |n$$

Poll Question 2:

Recall

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

> 17(|230



 $x \in \{0,1\}^{2}$; examples of conjunctions are

$$h(\mathbf{x}) = x_1 (1 - x_2) x_4 x_{10}$$

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

 $h(x)=\left(1-x_3\right)\left(1-x_4\right)x_8$ • Assuming $c^*\in\mathcal{H}$, if M=10, $\epsilon=0.1$, and $\delta=0.01$, how many labelled examples do we need so that Theorem $1 \Rightarrow PAC$?

- 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: 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. $\left|S\right| = M$, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have

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

with probability at least $1 - \delta$.

$$Plule = \frac{1}{M} \left(\ln \left| \frac{1}{f} \right| + \ln \frac{1}{\delta} \right)$$
Recall Theorem | said $M \ge \frac{1}{\epsilon} \left(\ln \left(\left| \frac{\mathcal{H}}{\delta} \right| \right) + \ln \left(\frac{1}{\delta} \right) \right)$

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

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

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

$$\left| R(h) - \hat{R}(h) \right| \leq \epsilon$$

•Bound is inversely *quadratic* in ϵ , e.g., halving ϵ means we need four times as many labelled training data points

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

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

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

$$\left| R(h) - \hat{R}(h) \right| \leq \epsilon$$

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

Theorem 2: corollary

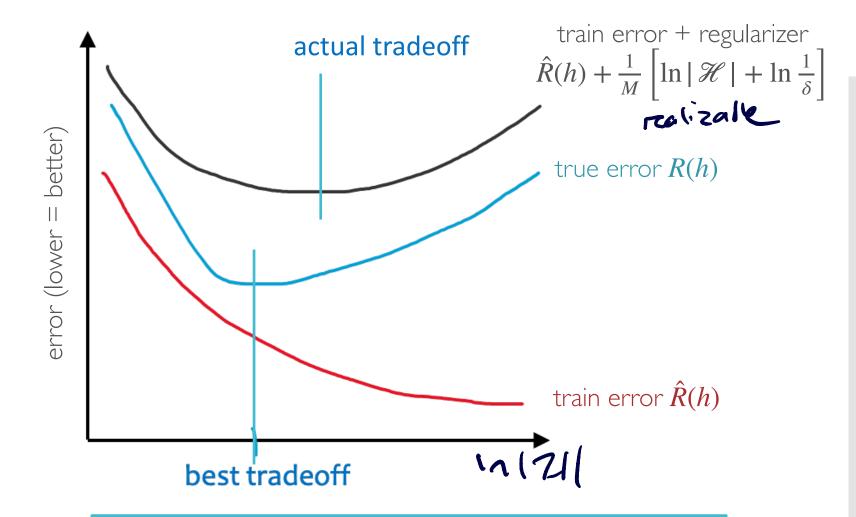
• For a finite hypothesis set $\mathcal H$ and arbitrary distribution p^* , given a training data set S s.t. $\left|S\right|=M$, all $h\in\mathcal H$ have

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

with probability at least $1 - \delta$.

Learning theory & model selection

Key point: we want to trade off low training error vs. keeping \mathscr{H} simple



Ex: $\mathscr{H}=$ conjunctions on d binary attributes: $\ln |\mathscr{H}|=d\ln 3$ Expert sorts attributes, most likely to be relevant first We allow conjunctions on first d of them: training error \downarrow as d increases regularizer \uparrow as d increases stop when PAC bound is smallest (best tradeoff)

What happens when

$$|\mathcal{H}| = \infty$$
?

Bounds:

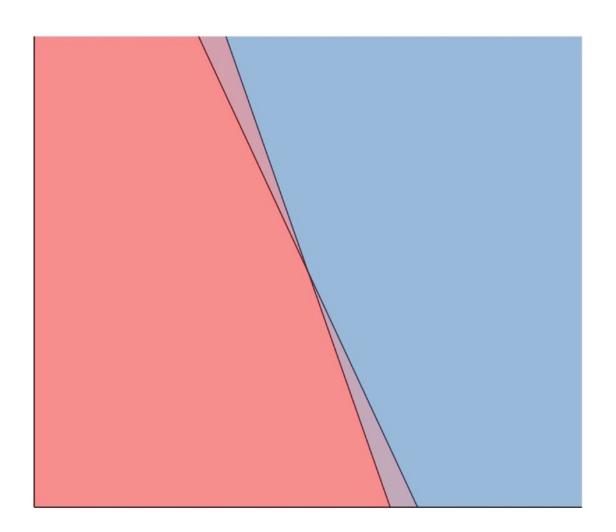
$$R(h) \le \frac{1}{M} \left(\ln \left(\left| \mathcal{H} \right| \right) + \ln \left(\frac{1}{\delta} \right) \right)$$

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

with probability at least $1 - \delta$.

Intuition

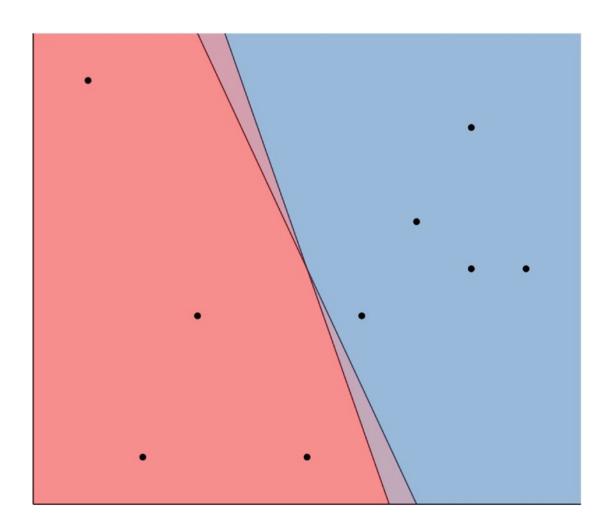
For "nice" infinite hypothesis sets \mathcal{H} , many hypotheses in \mathcal{H} will behave similarly



Intuition

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

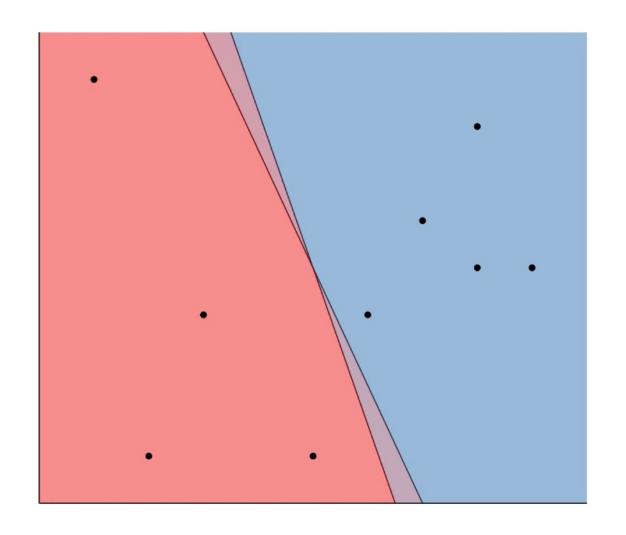
Relative to this dataset, these two hypotheses are *identical*!



Intuition

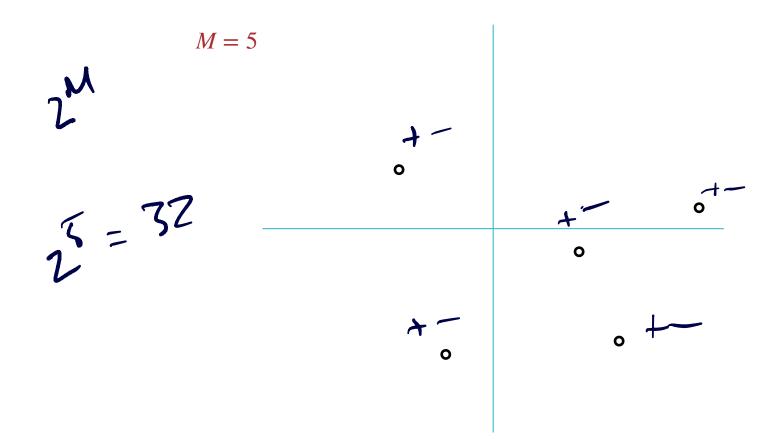
For "nice" infinite hypothesis sets \mathcal{H} , many hypotheses in \mathcal{H} will behave similarly

Relative to this dataset, these two hypotheses are *identical*!



Idea: instead of using full size of \mathcal{H} , count how many **actually distinct** hypotheses there are

How many distinct $h \in \mathcal{H}$ can there be?



• What's the largest possible number of distinct hypotheses on M points?

What does our bound tell us if $|\mathcal{H}| = 2^m$?

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

$$1 \wedge 2^{M}$$

$$= M 1 \wedge 2$$

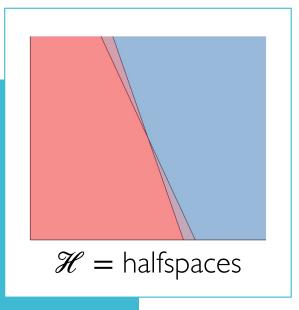
$$R(h) \leq 1 \wedge 2 + \frac{1}{M} 1 \wedge \frac{1}{\delta}$$

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

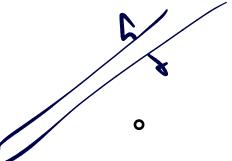
What does our bound tell us if $|\mathcal{H}| = 2^m$?

Vacuous! Need a tighter count of $|\mathcal{H}|$

Not surprising since $|\mathcal{H}| = 2^m$ is a kind of memorization learner

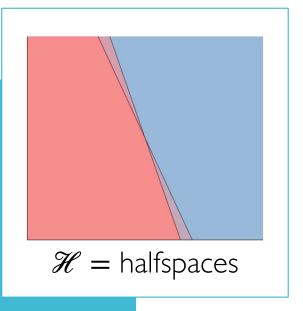


$$M = 1$$

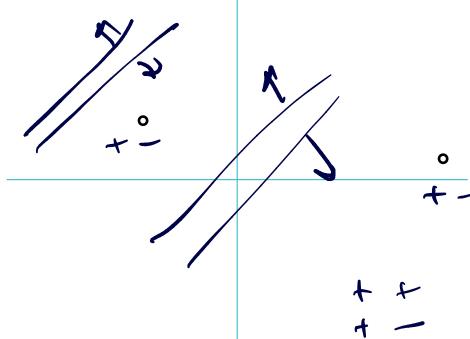


$$\mathscr{D} =$$

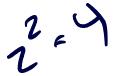
- ullet Fix ${\mathcal H}$
- Consider datasets \mathcal{D} of size M = 1, 2, ...
- For each dataset, count how many *actually distinct* hypotheses $h \in \mathcal{H}$ there are: $|\mathcal{H}(\mathcal{D})|$



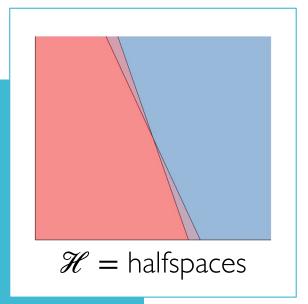
M = 2

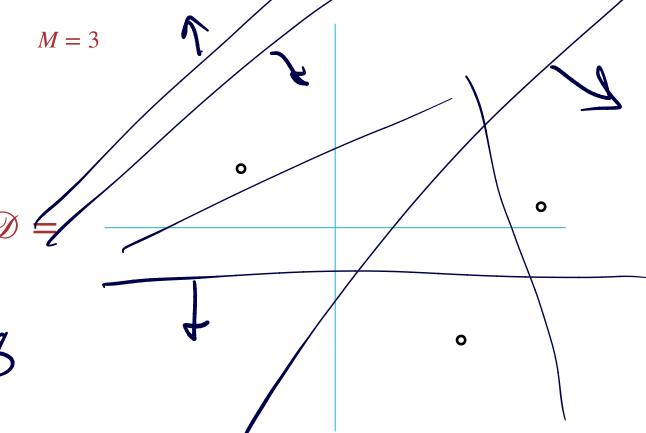


Counting *h*



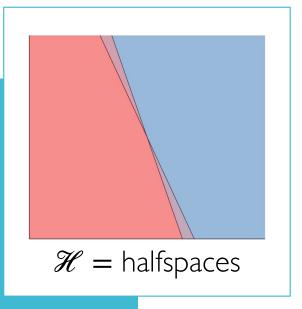
- ullet Fix ${\mathcal H}$
- Consider datasets \mathcal{D} of size M = 1, 2, ...
- For each dataset, count how many *actually distinct* hypotheses $h \in \mathcal{H}$ there are: $|\mathcal{H}(\mathcal{D})|$





Counting *h*

- ullet Fix ${\mathcal H}$
- Consider datasets \mathcal{D} of size M = 1, 2, ...
- For each dataset, count how many *actually distinct* hypotheses $h \in \mathcal{H}$ there are: $|\mathcal{H}(\mathcal{D})|$



$$M = 4$$

$$\mathcal{D} =$$

Counting h

- ullet Fix ${\mathcal H}$
- Consider datasets \mathcal{D} of size M = 1, 2, ...
- For each dataset, count how many *actually distinct* hypotheses $h \in \mathcal{H}$ there are: $|\mathcal{H}(\mathcal{D})|$

Growth function

- Def'n: the *growth function* $S_{\mathcal{H}}(M)$ is the maximum number of *distinct* $h \in \mathcal{H}$ for a dataset of size M
 - for halfspaces in 2D,

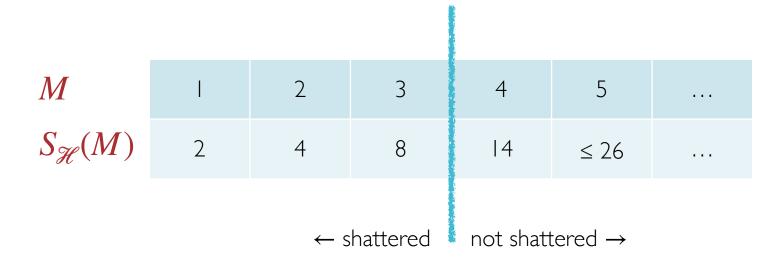
M	I	2	3	4	5	
$S_{\mathcal{H}}(M)$	2	4	8	14	≤ 26	

• for larger M, it turns out $S_{\mathcal{H}}(M) = O(M^3)$

not obvious!

Growth function

Def'n: **shattering**



**** shatters** a set of points if it can classify them all possible ways

Two kinds of behavior

ullet For many hypothesis classes ${\mathscr H}$, similar behavior:

$$S_{\mathcal{H}}(M) = \begin{cases} 2^M & M \le d \text{ shattered} \\ \ll 2^M & M > d \end{cases}$$
 not shattered

- e.g., intervals (or rectangles or hyperrectangles)
- e.g., bounded-depth decision trees
- e.g., fixed-architecture neural networks
- ullet For many other classes ${\mathcal H}$, instead $S_{{\mathcal H}}(M)=2^M$ for all M
 - e.g., unbounded-depth decision trees

can shatter a set of each size

e.g., unbounded-size neural networks

Two kinds of behavior

Learnable (can't memorize more than d points)

ullet For many hypothesis classes ${\mathscr H}$, similar behavior:

$$S_{\mathcal{H}}(M) = \begin{cases} 2^{M} & M \le d \text{ shattered} \\ \ll 2^{M} & M > d \text{ not shattered} \end{cases}$$

- e.g., intervals (or rectangles or hyperrectangles)
- e.g., bounded-depth decision trees
- e.g., fixed-architecture neural networks
- ullet For many other classes ${\mathscr H}$, instead $S_{\mathscr H}(M)=2^M$ for all M
 - e.g., unbounded-depth decision trees

can shatter a set of each size

e.g., unbounded-size neural networks

Not learnable (can memorize at any $|\mathcal{D}|$)

Sauer's lemma

- •Suppose $S_{\mathcal{H}}(M) = 2^M$ for $M \le d$, but $S_{\mathcal{H}}(d+1) < 2^{d+1}$
 - \rightarrow Then $S_{\mathcal{H}}(M) = O(M^d)$

"Suppose we grow exponentially (i.e., shatter) only up to M=d. Then for M>d we grow polynomially, with degree d."

related results derived multiple times: Sauer, Shelah, Perles, Vapnik/Chervonenkis

Sauer's lemma

d is called the **VC-dimension** of \mathcal{H}



$$ightharpoonup$$
 Then $S_{\mathcal{H}}(M) = O(M^d)$

"Suppose we grow exponentially (i.e., shatter) only up to M=d. Then for M>d we grow polynomially, with degree d."

related results derived multiple times: Sauer, Shelah, Perles, Vapnik/Chervonenkis

What do our bounds tell us w/ Sauer's lemma?

12 CM= 12 C + dhM

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

• Infinite realizable case:

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

$$\frac{1}{M} \left(\ln C + d \ln M + \ln \frac{1}{\delta} \right)$$

What do our bounds tell us w/ Sauer's lemma?

• Finite agnostic case:

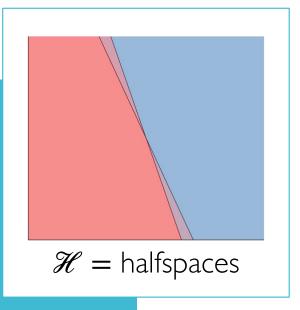
$$R(h) \le \hat{R}(h) + \sqrt{\frac{1}{2M}} \left(\ln\left(\left|\mathcal{H}\right|\right) + \ln\left(\frac{2}{\delta}\right) \right)$$

• Infinite agnostic case:

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

Finding the VC-dimension

- We defined <u>VC-dimension</u> of \mathcal{H} , $VC(\mathcal{H})$, as the size of the largest set S that \mathcal{H} can shatter
 - If $\mathcal H$ can shatter arbitrarily large sets, $VC(\mathcal H)=\infty$
- To prove that $VC(\mathcal{H}) = d$, need to show
 - 1. \exists some set of d data points that \mathcal{H} can shatter and
 - 2. \nexists a set of d+1 data points that \mathscr{H} can shatter



$$M = 3$$

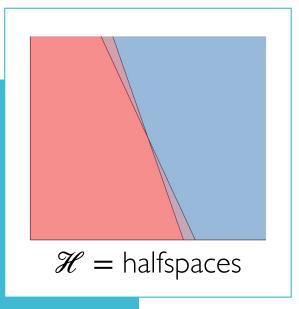
$$\mathcal{D} =$$

VC-dimension example

0

0

Before, we looked at this dataset of size 3



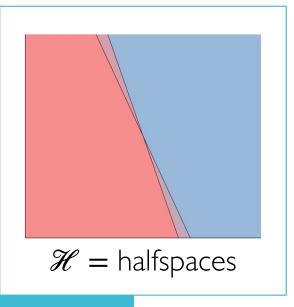
$$M = 3$$

$$\mathcal{D} =$$

0

VC-dimension example

• But what if we had looked at this one?



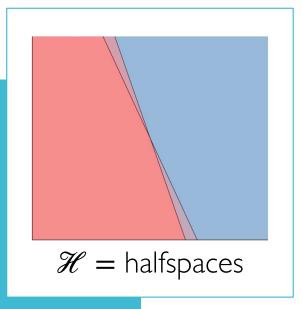
$$M = 3$$

$$\mathcal{D} =$$

0

VC-dimension example

- •Only 6 distinct hypotheses $< 2^3$
- Tempting to say $VC(\mathcal{H}) < 3$ but would be **wrong**

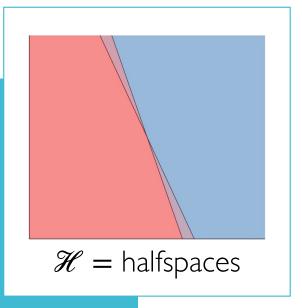


$$M = 4$$

$$\mathscr{D} =$$

VC-dimension example

Similarly, looked at this dataset of size 4

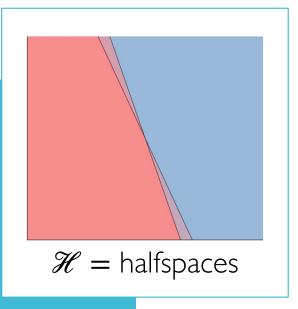


$$M = 4$$

$$\mathcal{D} =$$

VC-dimension example

But really should have checked this one as well

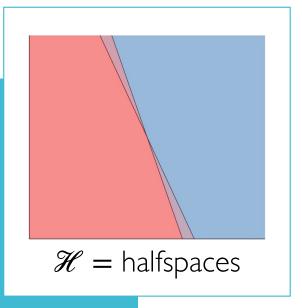


$$M = 4$$

$$\mathcal{D} =$$

VC-dimension example

And this one



$$M = 4$$

0

$$\mathscr{D} =$$

0

VC-dimension example

0

And this one

Halfspaces (linear separators)

- Just argued that halfspaces in 2D have VC = 3
- In general, halfspaces in d dimensions: VC = d+1

More VCdimension examples

- Try this at home: what is $VC(\mathcal{H})$ for
 - $\bullet \mathcal{H}$ = half-lines where positive class is on right
 - • \mathcal{H} = real intervals, positive when $x \in (a, b)$
 - \mathcal{H} = axis-parallel rectangles in 2D (+ on interior)

Learning objectives

- You should be able to...
 - Identify properties of a learning setting, assumptions needed to ensure low generalization error
 - Distinguish true error, train error (and test, validation errors)
 - Define PAC: what is approximately correct and what occurs with high probability
 - Apply sample complexity bounds to real-world learning examples
 - Theoretically motivate regularization