#### Announcements

#### Assignments

- HW10 (programming + "written")
  - Due Thu 4/30, 11:59 pm

# Introduction to Machine Learning

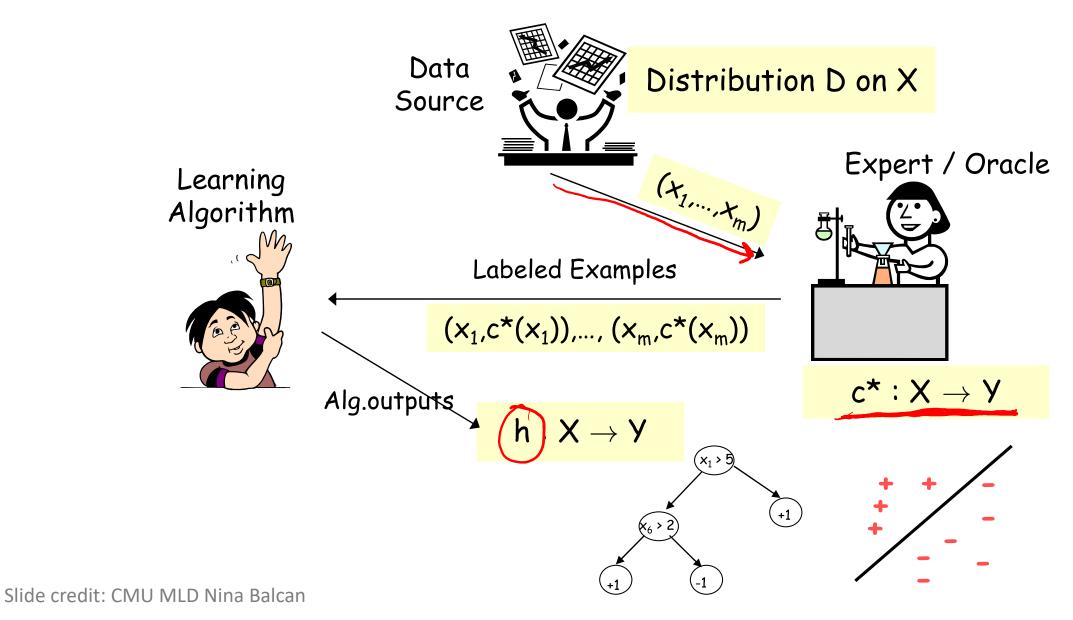
**Learning Theory** 

Instructor: Pat Virtue

# **Questions For Today**

- Given a classifier with zero training error, what can we say about true error (aka. generalization error)? (Sample Complexity, Realizable Case)
- Given a classifier with low training error, what can we say about true error (aka. generalization error)?
   (Sample Complexity, Agnostic Case)
- Is there a theoretical justification for regularization to avoid overfitting? (Structural Risk Minimization)

# Model for Supervised Learning



### Optimal Classification Function

Find the best  $h(x) \to \hat{y}$  by searching in the space of hypothesis functions  $h \in \mathcal{H}$ .

Optimal classifier:

$$h^*(x) = \underset{y}{\operatorname{argmax}} P(Y = y \mid X = x)$$

#### But why?

Goal: find a prediction function  $h^*: \mathcal{X} \to \mathcal{Y}$  that minimizes the expected loss for randomly drawn test data (X,Y)

$$h^*$$
 = argmin  $\mathbb{E}_{XY}[L(Y, h(X))] \leftarrow \mathbb{R}(h) = \mathbb{E}[Loss]$ 

 $L(y, \hat{y})$  is the loss or cost of predicting  $\hat{y}$  when the true value is y.

#### Loss Functions

$$h^* = \underset{h}{\operatorname{argmin}} \mathbb{E}_{XY} [\underline{L(Y, h(X))}]$$

#### Loss function:

$$L: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$$

#### Classification:

■ Two-class, 0,1 loss ←

$$y=0 \quad \hat{y}=0$$

$$y=0 \quad 0$$

$$y=1 \quad 0$$

Two-class, arbitrary loss

#### Regression:

### Optimal Classification Function

#### Expected loss is also called risk:

$$R(h) = \mathbb{E}_{XY}[L(Y, h(X))]$$

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

$$= \underset{h}{\operatorname{argmin}} \mathbb{E}_{XY}[L(Y, h(X))]$$

For 0,1 loss classification, risk is also error:

$$R(h) = E_{xy}[L(Y, h(X))] \leftarrow L(y, \hat{y}) = \begin{cases} 1 & \text{if } h(x) = y \\ 0 & \text{ow} \end{cases}$$

$$= \sum_{x,y} P(x,y) L(y, h(x))$$

$$R(h) = \sum_{x,y} \frac{1}{y} I(y^{(i)} \neq h(x^{(i)})) = \frac{1}{y} \sum_{x,y} \# error$$

# Two Types of Error

#### 1. True Error (aka. expected risk)

$$R(h) = P_{\mathbf{x}} \underbrace{p^*(\mathbf{x})} \underbrace{(c^*(\mathbf{x}) \neq h(\mathbf{x}))}_{\mathbf{x}}$$

#### 2. Train Error (aka. empirical risk)

$$\hat{R}(h) = P_{\mathbf{x}} \mathcal{S}(c^*(\mathbf{x}) \neq h(\mathbf{x}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(\underline{c^*(\mathbf{x}^{(i)})} \neq h(\mathbf{x}^{(i)}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(y^{(i)} \neq h(\mathbf{x}^{(i)}))$$

This quantity is always

We can measure this on the training data

where  $S = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}_{i=1}^N$  is the training data set, and  $\mathbf{x} \sim \mathcal{S}$  denotes that  $\mathbf{x}$  is sampled from the empirical distribution.

### PAC / SLT Model

1. Generate instances from unknown distribution  $p^*$ 

$$\mathbf{x}^{(i)} \sim p^*(\mathbf{x}), \, \forall i$$
 (1)

2. Oracle labels each instance with unknown function  $c^{*}$ 

$$y^{(i)} = c^*(\mathbf{x}^{(i)}), \,\forall i \tag{2}$$

3. Learning algorithm chooses hypothesis  $h \in \mathcal{H}$  with low(est) training error,  $\hat{R}(h)$  however  $\hat{R}(h)$ 

$$\hat{h} = \underset{h}{\operatorname{argmin}} \hat{R}(h) \tag{3}$$

4. Goal: Choose an h with low generalization error R(h)

# Three Hypotheses of Interest

The **true function**  $c^*$  is the one we are trying to learn and that labeled the training data:

$$y^{(i)} = \underline{c}^*(\mathbf{x}^{(i)}), \, \forall i \tag{1}$$

The **expected risk minimizer** has lowest true error:



$$h^* = \operatorname*{argmin}_{h \in \mathcal{H}} R(h)$$
 True or False: h\* and c\* are

**Question:** 

always equal.

The empirical risk minimizer has lowest training error:

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \, \hat{R}(h) \tag{3}$$

#### Piazza Poll 1

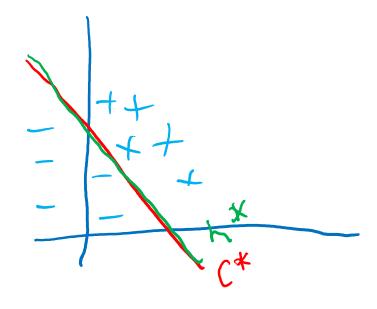
H linear sep

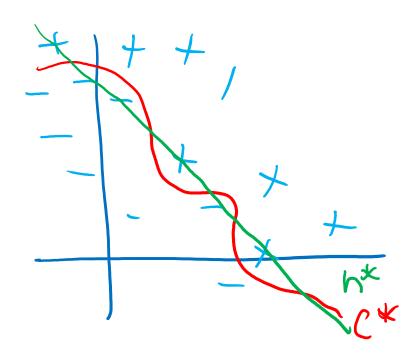
True or False: h\* and c\* are always equal.

A. True

B. False

C. Calanity





PAC Learning Know Can we bound R(h) in terms of  $\hat{R}(h)$ ? Yes! PAC: Probably Approximately Collect PAC Learner yield h EH which is approximately correct:  $R(h) \approx 0$ Definition: PAC Criterion: rall rallRandom? R(h) based Xtrain ~ P\*(x)

Definition: sample complexity is min num. of training evanpts N 5, t. the PAC criterion is satisfied for  $\epsilon$  and  $\delta$   $1-\delta=97%$   $\epsilon=0.1%$ 

Definition: consistent hypothesis

a h ∈ H is consistent with training data if R(h)=0

The **PAC criterion** is that our learner produces a high accuracy learner with high probability:

$$P(|R(h) - \hat{R}(h)| \le \epsilon) \ge 1 - \delta \tag{1}$$

Suppose we have a learner that produces a hypothesis  $h \in \mathcal{H}$  given a sample of N training examples. The algorithm is called **consistent** if for every  $\epsilon$  and  $\delta$ , there exists a positive number of training examples N such that for any distribution  $p^*$ , we have that:

$$P(|R(h) - \hat{R}(h)| > \epsilon) < \delta \tag{2}$$

The **sample complexity** is the minimum value of N for which this statement holds. If N is finite for some learning algorithm, then  $\mathcal H$  is said to be **learnable**. If N is a polynomial function of  $\frac{1}{\epsilon}$  and  $\frac{1}{\delta}$  for some learning algorithm, then  $\mathcal H$  is said to be **PAC learnable**.

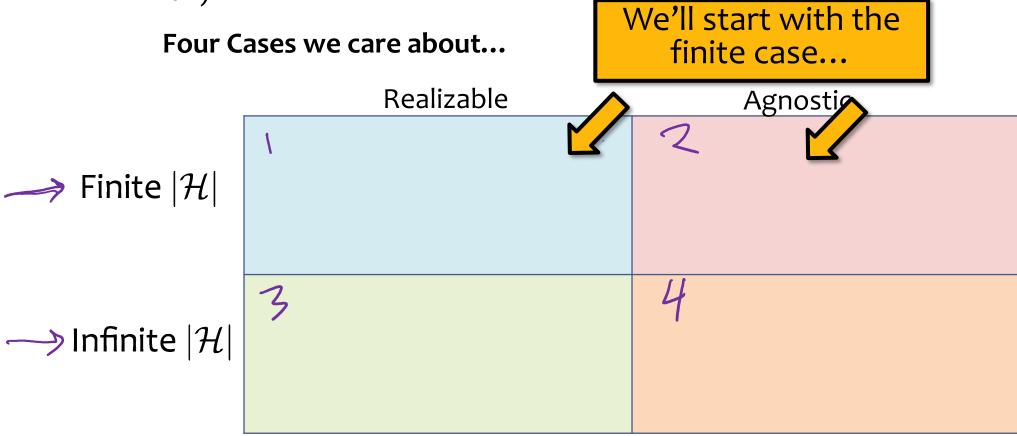
#### Four types of problems

Two cases 
$$c^*$$
• Realizable:  $c^* \in H$ 
• Agnostic:  $c^* \notin H$ 

Two cases  $|H|$ 
• Finite  $|H| = \infty$ 

# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).



Theorem 1: Sample Complexity (Realizable, Finite  $|\mathcal{H}|$ )

Theorem 1. Sample Complexity (Realizable, Finite 157)

$$N \ge \frac{1}{E} \left( \log |H| + \log / s \right)$$
 Tabelled examples

are sufficient to excuse that with prob 1-8

all  $h \in H$  with  $\hat{R}(h) = 0$  have  $R(h) \le E$ 

Proof of Theorem 1

{To the whiteboard...}

# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	<b>Thm.</b> 1 $N \geq \frac{1}{\epsilon} \left[ \log( \mathcal{H} ) + \log(\frac{1}{\delta}) \right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$ .	
Infinite $ \mathcal{H} $		

#### Piazza Poll 2

#### **Question:**

Suppose  $H = class of conjunctions over x in {0,1}^M$ 

Example hypotheses:

If M = 10,  $\varepsilon = 0.1$ ,  $\delta = 0.01$ , how many examples suffice according to Theorem 1?

#### **Answer:**

- A.  $10*(2*ln(10)+ln(100)) \approx 92$
- B.  $10*(3*ln(10)+ln(100)) \approx 116$
- C.  $10*(10*ln(2)+ln(100)) \approx 116$
- D. 10\*(10\*ln(3)+ln(100)) ≈ 156 ←
- E.  $0100*(2*ln(10)+ln(10)) \approx 691$
- F.  $100*(3*ln(10)+ln(10)) \approx 922$
- G.  $100*(10*ln(2)+ln(10)) \approx 924$
- H.  $100*(10*ln(3)+ln(10)) \approx 1329$

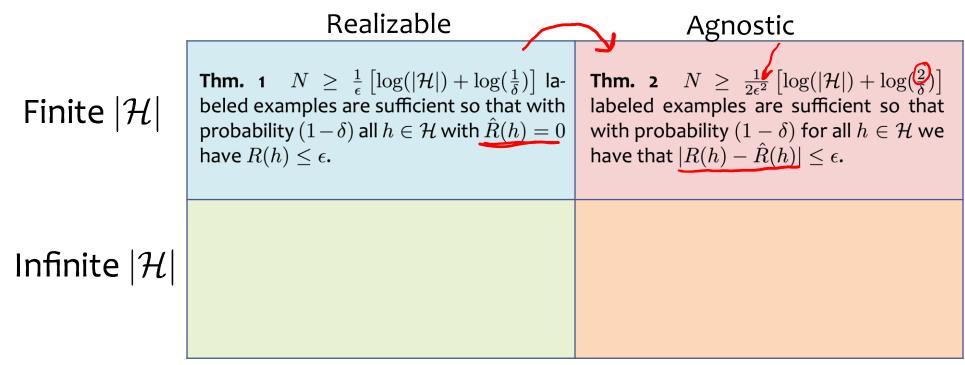


Thm. 1  $N \geq \left(\frac{1}{\epsilon}\right) \left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$  labeled examples are sufficient so that with probability  $(1-\delta)$  all  $h \in \mathcal{H}$  with  $\hat{R}(h) = 0$  have  $R(h) \leq \epsilon$ .

# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...



- Bound is inversely linear in epsilon (e.g. halving the error requires double the examples)
- 2. Bound is **only logarithmic in**|H| (e.g. quadrupling the hypothesis space only requires double the examples)
- Bound is **inversely quadratic in epsilon** (e.g. halving the error requires 4x the examples)
- Bound is only logarithmic in |H| (i.e. same as Realizable case)

10



Realizable

**Thm.** 1  $N \geq \frac{1}{\epsilon} \left[ \log(|\mathcal{H}|) + \log(\frac{1}{\delta}) \right]$  labeled examples are sufficient so that with probability  $(1-\delta)$  all  $h \in \mathcal{H}$  with  $\hat{R}(h) = 0$  have  $R(h) \leq \epsilon$ .

Agnostic

Infinite  $|\mathcal{H}|$ 

**Thm. 2**  $N \geq \frac{1}{2\epsilon^2} \left[ \log(|\mathcal{H}|) + \log(\frac{2}{\delta}) \right]$  labeled examples are sufficient so that with probability  $(1-\delta)$  for all  $h \in \mathcal{H}$  we have that  $|R(h) - \hat{R}(h)| \leq \epsilon$ .

# Using a PAC bound

$$|H|e^{-m\epsilon} \leq \delta$$

• Given  $\varepsilon$  and  $\delta$ , yields sample complexity

#training data, 
$$m \geq \frac{\ln|H| + \ln\frac{1}{\delta}}{\epsilon}$$

• Given m and  $\delta$ , yields error bound

error, 
$$\epsilon \geq \frac{\ln|H| + \ln\frac{1}{\delta}}{m}$$

### Summary of PAC bounds for finite model classes

With probability  $\geq 1-\delta$ ,

1) For all 
$$h \in H$$
 s.t.  $error_{train}(h) = 0$ ,  $error_{true}(h) \le \varepsilon = \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$ 

Haussler's bound

2) For all 
$$h \in H$$
 
$$|error_{true}(h) - error_{train}(h)| \le \varepsilon = \sqrt{\frac{\ln|H| + \ln\frac{2}{\delta}}{2m}}$$

Hoeffding's bound

# PAC bound and Bias-Variance tradeoff

$$P\left(|\operatorname{error}_{true}(h) - \operatorname{error}_{train}(h)| \ge \epsilon\right) \le 2|H|e^{-2m\epsilon^2} \le \delta$$

• Equivalently, with probability  $\geq 1 - \delta$ 

$$\frac{\widehat{\mathfrak{g}}(h)}{\operatorname{error}_{true}(h)} \leq \frac{\widehat{\mathfrak{g}}(h)}{\operatorname{error}_{train}(h)} + \sqrt{\frac{\ln|H| + \ln\frac{2}{\delta}}{2m}}$$

bias

Fixed |H|

Model class Train Size

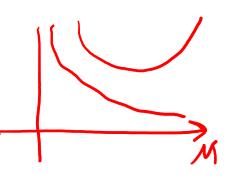
m small

small

m large

small large large small

variance



#### PAC bound and Bias-Variance tradeoff

$$P\left(|\operatorname{error}_{true}(h) - \operatorname{error}_{train}(h)| \ge \epsilon\right) \le 2|H|e^{-2m\epsilon^2} \le \delta$$

• Equivalently, with probability  $\geq 1 - \delta$ 

$$error_{true}(h) \le error_{train}(h) + \sqrt{\frac{\ln|H| + \ln\frac{2}{\delta}}{2m}}$$

Fixed m

Model class

|H| large (complex)

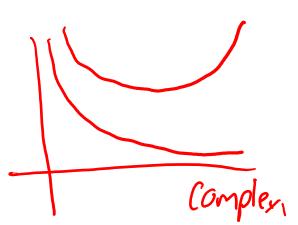
|H| small (simple)

small

large

large

small



## Number of decision trees of depth k

Recursive solution:

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

Given *n* binary attributes

 $H_k$  = Number of **binary** decision trees of depth k

$$H_0 = 2$$

 $H_k = (\#choices of root attribute)$ 

\*(# possible left subtrees)

\*(# possible right subtrees) =  $n * H_{k-1} * H_{k-1}$ 

Write 
$$L_k = log_2 H_k$$

$$L_0 = 1$$

$$L_{k} = \log_{2} n + 2L_{k-1} = \log_{2} n + 2(\log_{2} n + 2L_{k-2})$$

$$= \log_{2} n + 2\log_{2} n + 2^{2}\log_{2} n + \dots + 2^{k-1}(\log_{2} n + 2L_{0})$$

So 
$$L_k = (2^k-1)(1+\log_2 n) +1$$

### PAC bound for decision trees of depth k

$$m \ge \frac{\ln 2}{2\epsilon^2} \left( (2^k - 1)(1 + \log_2 n) + 1 + \log_2 \frac{2}{\delta} \right)$$

- Bad!!!
  - Number of points is exponential in depth k!

• But, for m data points, decision tree can't get too big...

Number of leaves never more than number data points

#### Number of decision trees with k leaves

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

 $H_k$  = Number of binary decision trees with k leaves

$$H_1 = 2$$

 $H_k = (\# choices of root attribute) *$ 

[(# left subtrees wth 1 leaf)\*(# right subtrees wth k-1 leaves)

+ (# left subtrees wth 2 leaves)\*(# right subtrees wth k-2 leaves)

+ ...

+ (# left subtrees wth k-1 leaves)\*(# right subtrees wth 1 leaf)]

$$H_k = n \sum_{i=1}^{k-1} H_i H_{k-i} = n^{k-1} C_{k-1}$$
 (C<sub>k-1</sub>: Catalan Number)

#### Loose bound (using Sterling's approximation):

$$H_k \le n^{k-1} 2^{2k-1}$$

#### Number of decision trees

With k leaves

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

$$\log_2 H_k \leq (k-1)\log_2 n + 2k - 1$$
 linear in k number of points m is linear in #leaves

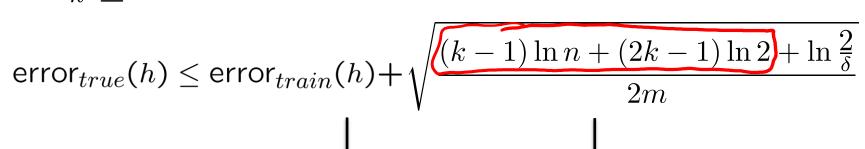
With depth k

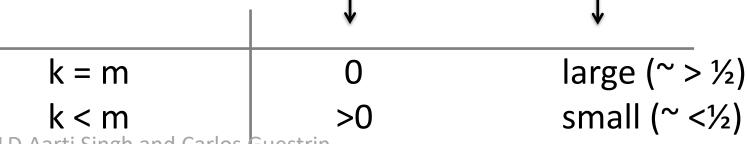
$$log_2 H_k = (2^{\bigcirc}1)(1+log_2 n) +1$$
 exponential in k number of points m is exponential in depth

# PAC bound for decision trees with k leaves – Bias-Variance revisited

With prob 
$$\geq 1-\delta$$
 error<sub>true</sub> $(h) \leq \operatorname{error}_{train}(h) + \sqrt{\frac{\ln(H) + \ln \frac{2}{\delta}}{2m}}$ 

With 
$$H_k \le n^{k-1} 2^{2k-1}$$
, we get





#### What did we learn from decision trees?

Moral of the story:

Complexity of learning not measured in terms of size of model space, but in maximum *number of points* that allows consistent classification

# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	<b>Thm.</b> 1 $N \geq \frac{1}{\epsilon} \left[ \log( \mathcal{H} ) + \log(\frac{1}{\delta}) \right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$ .	<b>Thm.</b> 2 $N \geq \frac{1}{2\epsilon^2} \left[ \log( \mathcal{H} ) + \log(\frac{2}{\delta}) \right]$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  \leq \epsilon$ .
Infinite $ \mathcal{H} $		