

Learning Theory

Machine Learning – CSE546

Carlos Guestrin

University of Washington

November 25, 2013

©Carlos Guestrin 2005-2013

1

What now...

- We have explored many ways of learning from data
- But...
 - How good is our classifier, really?
 - How much data do I need to make it “good enough”?

©Carlos Guestrin 2005-2013

2

A simple setting...

- Classification
 - N data points *iid*
 - **Finite** number of possible hypothesis (e.g., dec. trees of depth d)
- A learner finds a hypothesis h that is **consistent** with training data
 - Gets zero error in training $\leftarrow \text{error}_{\text{train}}(h) = 0$
- What is the probability that h has more than ϵ true error?
 - $\text{error}_{\text{true}}(h) \geq \epsilon$ *for any $\epsilon > 0$*

©Carlos Guestrin 2005-2013

3

How likely is a bad hypothesis to get N data points right?

- Hypothesis h that is **consistent** with training data \rightarrow got N i.i.d. points right
 - h "bad" if it gets all this data right, but has high true error *$\epsilon > 0$*
- Prob. h with $\text{error}_{\text{true}}(h) \geq \epsilon$ gets one data point right

less than $1 - \epsilon$
- Prob. h with $\text{error}_{\text{true}}(h) \geq \epsilon$ gets N data points right

less $(1 - \epsilon)^N$

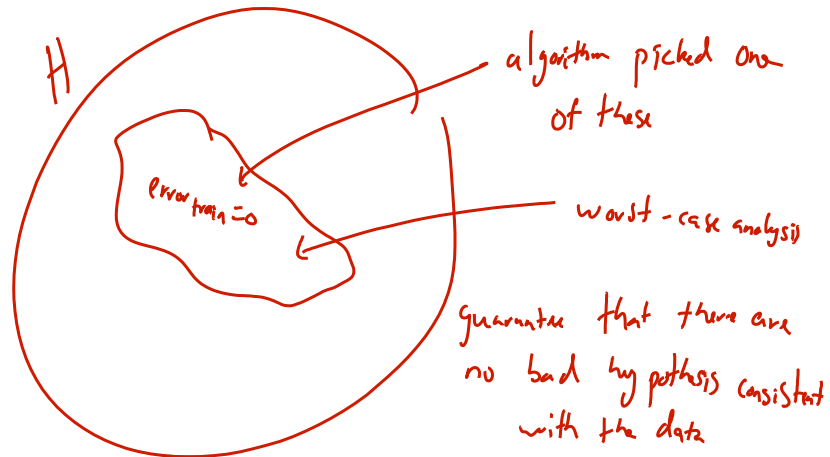


*if error true $\epsilon \geq 0.25$
75% prob h will get
1 point right*

©Carlos Guestrin 2005-2013

4

But there are many possible hypothesis that are consistent with training data



©Carlos Guestrin 2005-2013

5

How likely is learner to pick a bad hypothesis

- Prob. ^{a few to consider} h with $\text{error}_{\text{true}}(h) \geq \epsilon$ gets N data points right
 less $(1-\epsilon)^N$ — h_1, \dots, h_k
- There are k hypothesis consistent with data
 - How likely is learner to pick a bad one?

$$P(\exists h \text{ error}_{\text{train}} = 0 \ \& \ \text{error}_{\text{true}} \geq \epsilon)$$

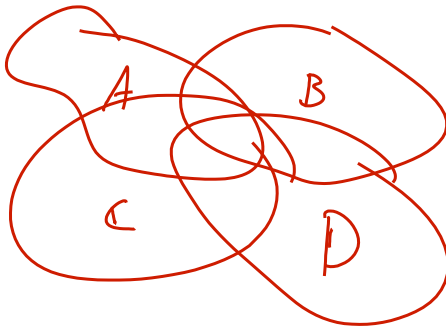
$$= P(\text{error}_{\text{true}}(h_1) \geq \epsilon \ \text{OR} \ \text{error}_{\text{true}}(h_2) \geq \epsilon \ \text{OR} \dots \ \text{OR} \ \text{error}_{\text{true}}(h_k) \geq \epsilon)$$

©Carlos Guestrin 2005-2013

6

Union bound

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



©Carlos Guestrin 2005-2013

7

How likely is learner to pick a bad hypothesis

- Prob. a particular h with $\text{error}_{\text{true}}(h) \geq \epsilon$ gets N data points right *less than $(1-\epsilon)^N$*
- There are k hypothesis consistent with data
 - How likely is it that learner will pick a bad one out of these k choices?

$$\leq P(\exists h: \text{error}_{\text{train}}(h) = 0 \ \& \ \text{error}_{\text{true}}(h) \geq \epsilon) \stackrel{\text{union bound}}{\leq} k(1-\epsilon)^N$$

$$\leq |H| (1-\epsilon)^N$$

what's k ?
 $k \leq |H|$
 total # hypothesis
 (very loose)

©Carlos Guestrin 2005-2013

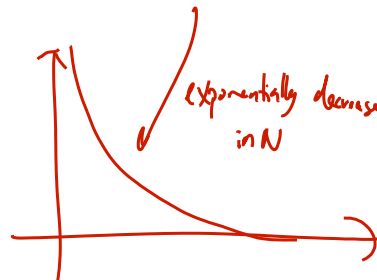
8

Generalization error in finite hypothesis spaces [Haussler '88]

- Theorem:** Hypothesis space H finite, dataset D with N i.i.d. samples, $0 < \epsilon < 1$: for any learned hypothesis h that is consistent on the training data: (error true(h) > epsilon)

$$P(\text{error}_{\text{true}}(h) \geq \epsilon) \leq |H|e^{-N\epsilon}$$

$0 \leq \epsilon \leq 1 \implies (1-\epsilon) \leq e^{-\epsilon}$
 $|H|(1-\epsilon)^N \leq |H|e^{-N\epsilon}$



©Carlos Guestrin 2005-2013

9

Using a PAC bound

- Typically, 2 use cases:
 - 1: Pick ϵ and δ , give you N
 - 2: Pick N and δ , give you ϵ

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-N\epsilon}$$

Prob(lose job) = prob(error true(h) > epsilon)
 $\leq |H|e^{-N\epsilon} \leq \delta$

$\ln |H| - N\epsilon \leq \ln \delta$
 (log dependence on |H|) (log dependence on tolerance)

$\implies N \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{\epsilon}$
 (collect) (epsilon linear dependence on accuracy)

$\epsilon \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{N}$

decrease at rate $O(\frac{1}{N})$
 great rate

more general setting, decrease as $O(\frac{1}{\sqrt{N}})$

©Carlos Guestrin 2005-2013

10

Summary: Generalization error in finite hypothesis spaces [Haussler '88]

- **Theorem:** Hypothesis space H finite, dataset D with N i.i.d. samples, $0 < \epsilon < 1$: for any learned hypothesis h that is consistent on the training data:

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-N\epsilon}$$

Even if h makes zero errors in training data, may make errors in test

©Carlos Guestrin 2005-2013

11

Limitations of Haussler '88 bound

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-N\epsilon}$$

- Consistent classifier

$\text{error}_{\text{train}}(h) = 0 \rightarrow$ highly unrealistic \rightarrow label noise \Rightarrow two data points with same x have different y
 \rightarrow overfit

- Size of hypothesis space

$\ln |H|$ can be bad \rightarrow finite $|H|$ if $|H|$ is very very large
 \rightarrow $|H|$ is infinite, e.g. in SVMs or LR

©Carlos Guestrin 2005-2013

12

What if our classifier does not have zero error on the training data?

- A learner with zero training errors may make mistakes in test set
- What about a learner with $error_{train}(h)$ in training set?

what happens $error_{train}(h) > 0$

$\Rightarrow error_{true}(h)?$

©Carlos Guestrin 2005-2013

13

Simpler question: What's the expected error of a hypothesis?

- The error of a hypothesis is like estimating the parameter of a coin!

$\theta \approx \hat{\theta} = \frac{3}{5}$

- Chernoff bound: for N i.i.d. coin flips, x^1, \dots, x^N , where $x^j \in \{0, 1\}$. For $0 < \epsilon < 1$:

$$P\left(\theta - \frac{1}{N} \sum_{j=1}^N x^j > \epsilon\right) \leq e^{-2N\epsilon^2}$$

\uparrow truth $\underbrace{\frac{1}{N} \sum_{j=1}^N x^j}_{\frac{3}{5} \text{ estimate of mean}}$ \uparrow ϵ \leftarrow decrease exponentially in N

©Carlos Guestrin 2005-2013

14

Using Chernoff bound to estimate error of a single hypothesis

$$P\left(\theta - \frac{1}{N} \sum_{j=1}^N x^j > \epsilon\right) \leq e^{-2N\epsilon^2}$$

$\theta = \text{error}_{\text{true}}(h)$
 $\hat{\theta} = \text{error}_{\text{train}}(h) \leftarrow \frac{1}{N} \sum_{j=1}^N \mathbb{1}(h(x^j) \neq y^j)$
 $= \int_{\mathcal{X}} P(x) \mathbb{1}(h(x) \neq f(x)) dx$
 $P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq e^{-2N\epsilon^2}$

in reality bad h (pointing to θ)
h looks good (pointing to $\hat{\theta}$)

©Carlos Guestrin 2005-2013

15

But we are comparing many hypothesis: **Union bound**

For each hypothesis h_i :

$$P(\text{error}_{\text{true}}(h_i) - \text{error}_{\text{train}}(h_i) > \epsilon) \leq e^{-2N\epsilon^2}$$

What if I am comparing two hypothesis, h_1 and h_2 ?

is there an h_2 that is truly better than my h_1

$$P(\text{error}_{\text{train}}(h_1) < \text{error}_{\text{train}}(h_2) \text{ but } \text{error}_{\text{true}}(h_1) > \text{error}_{\text{true}}(h_2))$$

instead worst case analysis

$$P([\text{error}_{\text{true}}(h_1) - \text{error}_{\text{train}}(h_1) > \epsilon] \text{ OR } [\text{error}_{\text{true}}(h_2) - \text{error}_{\text{train}}(h_2) > \epsilon])$$

$$\stackrel{\text{union bound}}{\leq} 2e^{-2N\epsilon^2}$$

©Carlos Guestrin 2005-2013

16

Generalization bound for $|H|$ hypothesis

- Theorem:** Hypothesis space H finite, dataset D with N i.i.d. samples, $0 < \epsilon < 1$: for any learned hypothesis h :

$$P(\text{error}_{\text{true}}(h_i) - \text{error}_{\text{train}}(h_i) > \epsilon) \leq e^{-2N\epsilon^2}$$

hold $\forall h$: $P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H| e^{-2N\epsilon^2}$

for $\delta > 0$ $\epsilon \geq \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}} \Rightarrow$ rate is only $O\left(\frac{1}{\sqrt{N}}\right)$
 much worse than $O\left(\frac{1}{N}\right)$
 but still works!!
 ;)

©Carlos Guestrin 2005-2013

17

PAC bound and Bias-Variance tradeoff

$$P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H| e^{-2N\epsilon^2}$$

or, after moving some terms around, ϵ from previous slide
 with probability at least $1-\delta$:

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}}$$

	"bias"	"variance"
"complex hypothesis space"	low	large $\Leftarrow H $ is large
"simple hypothesis space"	high	low $\Leftarrow H $ is small

- Important: PAC bound holds for all h , but doesn't guarantee that algorithm finds best h !!!**

©Carlos Guestrin 2005-2013

18

What about the size of the hypothesis space?

$$N \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{2\epsilon^2}$$

- How large is the hypothesis space?

$|H|$?

$|H|$ is really large

$\Rightarrow \ln |H| = \text{only large}$

$\Rightarrow \text{OK}$

$|H| = \text{really really large}$

$\Rightarrow \ln |H| = \text{really large}$

$\Rightarrow \text{lots of data needed}$

Boolean formulas with m binary features

x_1, \dots, x_m $x_1 \wedge \neg x_2 \vee x_3 \wedge \neg x_2 \dots$

$$N \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{2\epsilon^2}$$

H : any boolean formula $|H|$?

what does one h represent:

$x_1 \dots x_m$	y	
0 0 0 0 0 1	1	} 2^m rows
0 0 0 0 0 0	1	
:	0	} 2 possibilities for y in each 2^m row
1 1 1 1 1 1	0	
	1	} $ H = 2^{2^m}$
	1	

$\ln |H| = 2^m \ln 2$
 \propto exponentially large

H : conjunctions only:

$x_1 \wedge \neg x_2 \wedge x_2$

Each feature \rightarrow positive } 3 possibilities
 \rightarrow negated }
 \rightarrow absent }

$|H| = 3^m \leftarrow \text{only really large}$

$\ln |H| = m \ln 3$

\uparrow
 linear in # of features

Number of decision trees of depth k

$$N \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{2\epsilon^2}$$

Recursive solution

Given m attributes

H_k = Number of decision trees of depth k

$H_0 = 2$

$H_{k+1} = (\text{\#choices of root attribute}) * (\text{\# possible left subtrees}) * (\text{\# possible right subtrees})$

$$= m * H_k * H_k$$

Write $L_k = \log_2 H_k$

$L_0 = 1$

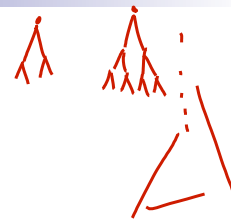
$L_{k+1} = \log_2 m + 2L_k$

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

Simplify to
 $\ln |H| \leq 2^k \log m$
 really really large in depth
 really nice in terms of # of features

PAC bound for decision trees of depth k

$$N \geq \frac{2^k \log m + \ln \frac{1}{\delta}}{\epsilon^2}$$



■ Bad!!!

- Number of points is exponential in depth!

■ But, for N data points, decision tree can't get too big...

no reason to have more than N leaves

Number of leaves never more than number data points

Number of Decision Trees with k Leaves

- Number of decision trees of depth k is really really big:

- $\ln |H|$ is about $2^k \log m$

- Decision trees with up to k leaves:

- $|H|$ is about $m^k k^{2k}$ ← only really large

- A very loose bound

$$\ln |H| \leq k \ln m + 2k \ln k$$

much better

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

$$\ln |H_{\text{DTs } k \text{ leaves}}| \leq 2k(\ln m + \ln k)$$

$$error_{true}(h) \leq error_{train}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}}$$

$\ln |H|$

$$error_{true}(h) \leq error_{train}(h) + \sqrt{\frac{2k(\ln m + \ln k) + \ln \frac{1}{\delta}}{2N}}$$

	"bias"	"variance"
$k \approx N$	goes to ϕ	LARGE bound $> ?$
$k \ll N$	potentially larger	much smaller

What did we learn from decision trees?

- Bias-Variance tradeoff formalized

$$error_{true}(h) \leq error_{train}(h) + \sqrt{\frac{2k(\ln m + \ln k) + \ln \frac{1}{\delta}}{2N}}$$

- Moral of the story:

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

- Complexity N – no bias, lots of variance
- Lower than N – some bias, less variance

What about continuous hypothesis spaces?

$$error_{true}(h) \leq error_{train}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}}$$

- Continuous hypothesis space:

- $|H| = \infty$ ← SVMs
- Infinite variance???

- **As with decision trees, only care about the maximum number of points that can be classified exactly!**

- Called VC dimension... see readings for details

What you need to know

- Finite hypothesis space
 - Derive results
 - Counting number of hypothesis
 - Mistakes on Training data
- Complexity of the classifier depends on number of points that can be classified exactly
 - Finite case – decision trees
 - Infinite case – VC dimension
- Bias-Variance tradeoff in learning theory
- Remember: will your algorithm find best classifier?