# Learning Theory

Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 27, 2013

1

---

# 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 – $error_{train}(h) = 0$
- What is the probability that $h$ has more than $\varepsilon$ true error?
  - $error_{true}(h) \geq \varepsilon$    For some  $\varepsilon > 0$

2

---

1

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

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

  *less than  1-ε*    | *if error ε = 0.25*
  *75% points are*
  *correct = 1-ε*

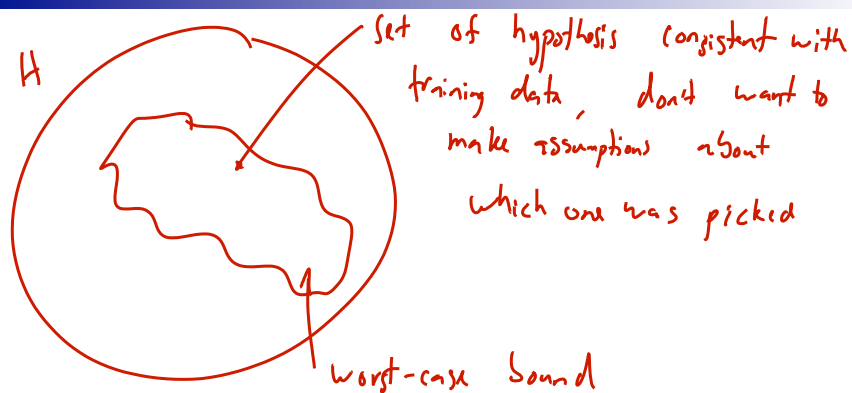- Prob. *h* with error$_{true}$(h) ≥ ε gets *N* data points right

  *less than* $(1-\varepsilon)^N$    *Prob bad h wins*    *decreases exponentially in N*

3

---

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

*H*

*Set of hypothesis consistent with training data, don't want to make assumptions about which one was picked*

*worst-case bound*

4

2

# How likely is learner to pick a bad hypothesis

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

  less than $(1-\varepsilon)^N$

  $\rightarrow h_1, \ldots h_k$

- There are $k$ hypothesis consistent with data
  - How likely is learner to pick a bad one?

  some bad, some good

  $\exists$ deal with worst case

$$P\left( \exists h \text{ consistent with data}^{\text{train}}, \text{error}_{\text{true}}(h) \geq \varepsilon \right)$$

$$= P\left( \text{error}_{\text{true}}(h_1) \geq \varepsilon \text{ OR } \text{error}_{\text{true}}(h_2) \geq \varepsilon \text{ OR} \ldots \text{ OR } \text{error}_{\text{true}}(h_k) \geq \varepsilon \right)$$
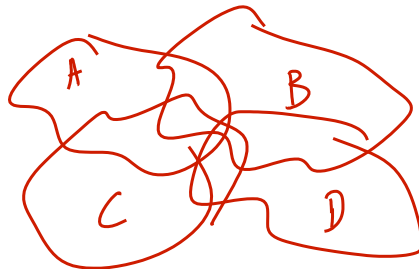
Bound?

5

# Union bound

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

6

3

# How likely is learner to pick a bad hypothesis

- Prob. a particular $h$ with error$_{true}$(h) $\geq \varepsilon$ gets $N$ data points right    *less than*    $(1-\varepsilon)^N$

- There are $k$ hypothesis consistent with data
  - How likely is it that learner will pick a bad one out of these $k$ choices?

$P(\exists h \text{ consistent with train data}, \text{error}_{true}(h) \geq \varepsilon) \leq k(1-\varepsilon)^N$

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

what's $k$?

$K \leq |H|$

total # hypothsis

(crazy loose)

# Generalization error in finite hypothesis spaces [Haussler '88]

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

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

prob. that you'll be fired

prob picking bad h

Decreases exponentially

$\leq |H| (1-\varepsilon)^N \leq |H| (e^{-\varepsilon})^N = |H| e^{-\varepsilon N}$

for $0 \leq \varepsilon \leq 1$

$1-\varepsilon \lesssim e^{-\varepsilon}$

# Using a PAC bound

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

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

*(handwritten annotations)*

$P(error_{true}(h) \geq \epsilon) \leq |H|e^{-N\epsilon} \leq \delta$

bad event ← upper bound ← acceptable prob.

$\ln|H| - N\epsilon \leq \ln\delta$ ← log of |H| ← log depend on δ

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

amount data needed ← linear in $\frac{1}{\epsilon}$

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

decrease $O\left(\frac{1}{N}\right)$

$\Rightarrow$ very good rate

More general settings: $O\left(\frac{1}{\sqrt{N}}\right)$

9

---

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

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

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

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

10

5

# Limitations of Haussler '88 bound

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

- Consistent classifier

$error_{train}(h) = 0$ → highly unrealistic

label noise,
complex data, model
bias
fitting problems

→ Overfit

- Size of hypothesis space

$\ln|H|$ is bad

$|H|$ very very large

$|H|$ infinite $\left(\begin{matrix} e.g. \\ SVM \\ LR \end{matrix}\right)$

11

---

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

- A learner with <span style="color:red">zero</span> training errors may make mistakes in test set
- What about a learner with *error_train(h)* in training set?

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

:) $error_{true}(h)$ ?

12

6

# 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, $x^1,\ldots,x^N$, where $x^j \in \{0,1\}$. For $0<\epsilon<1$:

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

*true*   *mean of train data* $\hat{\theta}$   $\epsilon$   *something exp in $N$*

---

# Using Chernoff bound to estimate error of a single hypothesis

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

*true error$_{\text{true}}(h)$*   $\hat{\theta} = $ *error train*

$\theta = \int_{X} p(x)\, \mathbb{1}(h(x) \neq t(x))\, dx$   *Sample estimate of integral*

*if label noise:* $\theta = \int_{X}\int_{y_x} p(x)\, p(y_x|x)\, \mathbb{1}(h(x) \neq y_x)\, dx\, dy$

$\frac{1}{N}\sum_{6=1}^{N} \mathbb{1}(h(x^j) \neq y^j) = error_{train}(h)$

$t(x^j)$

$P\left(error_{true}(h) - error_{train}(h) \geq \epsilon\right) \le e^{-2N\epsilon^2}$

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

*(handwritten, top right)* over fit by more than ε

For each hypothesis $h_i$:
$$P(error_{true}(h_i) - error_{train}(h_i) > \epsilon) \le e^{-2N\epsilon^2}$$

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

*(handwritten)* is $h_1$ better than $h_2$?

*(handwritten)* Danger: $P\left(error_{train}(h_1) < error_{train}(h_2) \ , \ but \ error_{true}(h_1) > error_{true}(h_2)\right)$

*(handwritten)* But want $P\left(\left[error_{true}(h_1) - error_{train}(h_1) \ge \epsilon\right] \ OR \ \left[error_{true}(h_2) - error_{train}(h_2) \ge \epsilon\right]\right)$

*(handwritten)* $\le P(error_{true}(h_1) - error_{train}(h_1) \ge \epsilon) + P(error_{true}(h_2) - error_{train}(h_2) \ge \epsilon)$

*(handwritten)* $\le 2 \ e^{-2N\epsilon^2}$

15

---

# Generalization bound for |H| hypothesis

- **Theorem**: Hypothesis space *H* finite, dataset *D* with *N* i.i.d. samples, 0 < ε < 1 : for any learned hypothesis *h*:
$$P(error_{true}(h_i) - error_{train}(h_i) > \epsilon) \le e^{-2N\epsilon^2}$$

*(handwritten)* holds $\forall h$:
$$P(error_{true}(h) - error_{train}(h) > \epsilon) \le |H| \ e^{-2N\epsilon^2}$$

*(handwritten)* $\epsilon \ge \sqrt{\dfrac{\ln|H| + \ln \frac{1}{\delta}}{2N}}$ $\rightarrow O\left(\dfrac{1}{\sqrt{N}}\right)$ rate

*(handwritten)* with probability at least $1-\delta$ : $error_{true}(h) - error_{train}(h) \le \epsilon$

16

# PAC bound and Bias-Variance tradeoff

$$P\left(error_{true}(h) - error_{train}(h) > \epsilon\right) \le e^{-2N\epsilon^2}$$

bound on $\epsilon$

**or, after moving some terms around, with probability at least 1-δ:**

If want small

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

"bias"   "variance"

| | "bias" | "variance" |
|---|---|---|
| "Complex hypothesis space" | low | large ⟸ |H| is large |
| "simple H" | large | low ⟸ |H| is small |

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

©Carlos Guestrin 2005-2013          17

---

# What about the size of the hypothesis space?

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

- How large is the hypothesis space?

|H|?

|H|= really large                    but          |H|= really really large

⟹ log |H| = only large                          ⟹ ln |H| = really large

⟹ ok                                               ⟹ lots data

# Boolean formulas with *m* binary features

$$X_1 \wedge \neg X_2 \vee X_7 \wedge X_2 \cdots$$

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

H: any boolean formula : |H| ?

$X_1, X_2 \cdots X_m$

```
     Y
6 0 0 6 0     0 or 1
0 0 0 0 1     0 or 1
```

$2^m$ rows

each row 2 possibilities

$|H| = 2^{2^m}$

≡ really really big

$\ln |H| = 2^m \ln 2$

↑ exp. many database

H: all conjunctions with negation

$X_1 \wedge \neg X_3 \wedge X_7$

each feature : 3 possibilities → positive → negated → absent

$|H| = 3^m$ ⇒ really large

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

↗ linear in # of features

©Carlos Guestrin 2005-2013          19

---

# 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
$H_0 = 2$
$H_{k+1}$ = (#choices of root attribute) *
          (# possible left subtrees) *
          (# possible right subtrees)
     $= m * H_k * H_k$

Write $L_k = \log_2 H_k$
$L_0 = 1$
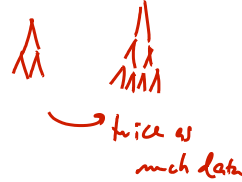$L_{k+1} = \log_2 m + 2L_k$
So $L_k = (2^k - 1)(1 + \log_2 m) + 1$

Simplify

$\ln |H| \leq 2^k \log m$

↗ really really big in terms of depth

↑ very nice in terms of num techs

©Carlos Guestrin 2005-2013          20

10

# PAC bound for decision trees of depth k

*exponential in depth*

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

*twice as much data*

- 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

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

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

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

| max number of leaves k | "bias" | "Variance" |
|---|---|---|
| $K \approx N$ | goes to zero | LARGE greater than 1 |
| $K << N$ | potentially larger | potentially small |

©Carlos Guestrin 2005-2013

23

---

# 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 $K \approx N$
  - Lower than $N$ – some bias, less variance $K << N$
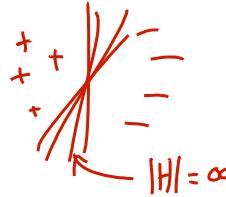
©Carlos Guestrin 2005-2013

24

12

# What about continuous hypothesis spaces?

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

- Continuous hypothesis space:
  - $|H| = \infty$
  - 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 ← # of leaves
  - Infinite case – VC dimension
- Bias-Variance tradeoff in learning theory
- Remember: will your algorithm find best classifier?

# Markov Decision Processes (MDPs)

Machine Learning – CSE546

Carlos Guestrin

University of Washington

December 2, 2013

27

---

e.g. Classification : supervised learning

$$X \rightarrow Y$$

(GPA, grade) $\rightarrow$ {hire, not hire}

unsupervised case : e.g. clustering

just X (GPA, grade) $\Rightarrow$ groups of people with similar Xs

# Reinforcement Learning

## training by feedback

weak feedback   $X_1$ ⊢ is good

$X_2$ ⊢ is bad

28

14

# Learning to act

- Reinforcement learning
- An agent
  - Makes sensor observations
  - Must select action
  - Receives rewards
    - positive for "good" states
    - negative for "bad" states

[Ng et al. '05]

©Carlos Guestrin 2005-2013                                         29

# Markov Decision Process (MDP) Representation

*position of peasant*

- State space:
  - Joint state **x** of entire system

  $x = (x_1, \ldots, x_n)$
  *gold*

- Action space:
  - Joint action **a**= {$a_1$,..., $a_n$} for all agents

  *build castle ...*

- Reward function:
  - Total reward R(**x**,**a**)
    - sometimes reward can depend on action

  *positive reward when X is win game*

- Transition model:
  - Dynamics of the entire system P(**x**'|**x**,**a**)

  *have gold no castle → build castle*
  *have castle but no gold*

  *want a policy*
  $\pi(x) \Rightarrow a$
  *what action at each state*

©Carlos Guestrin 2005-2013                                         30

15

# Discount Factors

$\gamma \in [0,1)$

People in economics and probabilistic decision-making do this all the time.

The "Discounted sum of future rewards" using discount factor $\gamma$" is

(reward now) +

$\gamma$ (reward in 1 time step) +

$\gamma^2$ (reward in 2 time steps) +

$\gamma^3$ (reward in 3 time steps) +

:

:     (infinite sum)

*maximize sum of discounted rewards*

*exp. less valuable in future*

*future*

---

# The Academic Life

**Assume Discount Factor $\gamma$ = 0.9**



0.6   0.2   0.6   0.2   0.7

A. Assistant Prof R(A) 20

B. Assoc. Prof R(A) 60

T. Tenured Prof R(A) 400   *in cents*

0.2   0.2

S. On the Street 10

D. Dead 0

0.7   0.3   0.3

Define:

$V_A$ = Expected discounted future rewards starting in state A   $= 20 + \gamma(0.6\, V_A + 0.2\, V_B + 0.2\, V_S)$

$V_B$ = Expected discounted future rewards starting in state B   $= 60 + \gamma(0.6\, V_B + 0.2\, V_S + 0.2\, V_T)$

$V_T$ = " " " " " " T

$V_S$ = " " " " " " S

$V_D$ = " " " " " " D

How do we compute $V_A$, $V_B$, $V_T$, $V_S$, $V_D$ ?

*n states.*
*n unknowns*
*n equations*
*linear*
*⇒ e.g. matrix inversion*

# Policy

Policy: $\pi(\mathbf{x}) = \mathbf{a}$ ➡ At state **x**, action **a** for all agents

$\pi(\mathbf{x}_0)$ = both peasants get wood

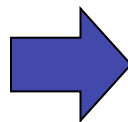$\pi(\mathbf{x}_1)$ = one peasant builds barrack, other gets gold

$\pi(\mathbf{x}_2)$ = peasants get gold, footmen attack

©Carlos Guestrin 2005-2013

33

---

# Value of Policy

Value: $V_\pi(\mathbf{x})$ ➡ Expected long-term reward starting from **x**

formal view of recursion

and act according to π

$$V_\pi(\mathbf{x}_0) = \mathbf{E}_\pi[R(\mathbf{x}_0) + \gamma R(\mathbf{x}_1) + \gamma^2 R(\mathbf{x}_2) + \gamma^3 R(\mathbf{x}_3) + \gamma^4 R(\mathbf{x}_4) + \ldots]$$

Future rewards discounted by $\gamma$ in [0,1)

Start from $\mathbf{x}_0$ a:

$\pi(\mathbf{x}_0)$

$R(\mathbf{x}_0)$

getting wood

$R(\mathbf{x}_1)$

$\pi(\mathbf{x}_1) = a'$

bad luck opponent

$\pi(\mathbf{x}_1')$

$R(\mathbf{x}_1')$

$\pi(\mathbf{x}_1'')$

$R(\mathbf{x}_1'')$

over time

$R(\mathbf{x}_2)$

$\pi(\mathbf{x}_2)$

$R(\mathbf{x}_3)$

$\pi(\mathbf{x}_3)$

$R(\mathbf{x}_4)$

17