

Learning Theory

Aarti Singh

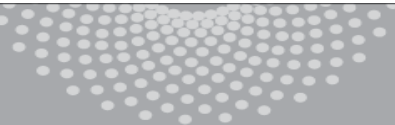
Machine Learning 10-701

Mar 29, 2023

Slides courtesy: Carlos Guestrin



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Learning Theory

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

PAC Learnability

Probably Approximately Correct (ϵ, δ)

- True function space, F
- Model space, H

F is **PAC Learnable** by a learner using H if

there exists a learning algorithm s.t. for all functions in F, for all distributions over inputs, for all $0 < \epsilon, \delta < 1$,

with ^{probably} probability $> 1 - \delta$, the algorithm outputs a model

$h \in H$ s.t. $\text{error}_{\text{true}}(h) \leq \epsilon$ *approximately correct*

in time and samples that are polynomial in $1/\epsilon, 1/\delta$. ←

↓
data
m

$$m \sim \frac{1}{\epsilon}$$

$$m \sim \frac{1}{\epsilon^2}$$

$$\sim \frac{1}{\sqrt{\epsilon}}$$

$$m \sim \frac{1}{2\epsilon}$$

A simple setting

- Classification
 - m i.i.d. data points
 - **Finite** number of possible classifiers in model class (e.g., dec. trees of depth d)
- Lets consider that a learner finds a classifier h that gets zero error in training
 - $\text{error}_{\text{train}}(h) = 0$
- What is the probability that h has more than ε true (= test) error?
 - $\text{error}_{\text{true}}(h) \geq \varepsilon$

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

How likely is a bad classifier to get m data points right?

- Consider a bad classifier h i.e. $\text{error}_{\text{true}}(h) \geq \varepsilon$
 $\text{error}_{\text{true}}(h) = P(h(x) \neq Y)$
- Probability that h gets one data point right
 $\leq 1 - \varepsilon$
- Probability that h gets m data points right
 $\leq \underline{\underline{(1 - \varepsilon)^m}}$

How likely is a learner to pick a bad classifier?

- Usually there are many (say k) bad classifiers in model class

$$h_1, h_2, \dots, h_k \quad \text{s.t. } \text{error}_{\text{true}}(h_i) \geq \varepsilon \quad i = 1, \dots, k$$

- Probability that learner picks a bad classifier = Probability that some bad classifier gets 0 training error

Prob(h_1 gets 0 training error OR
 h_2 gets 0 training error OR ... OR
 h_k gets 0 training error)

$$\leq \text{Prob}(h_1 \text{ gets 0 training error}) + \\ \text{Prob}(h_2 \text{ gets 0 training error}) + \dots + \\ \text{Prob}(h_k \text{ gets 0 training error})$$

$$\leq \underline{k} (1-\varepsilon)^m$$

$$\underline{P(A \cup B)} \leq P(A) + P(B)$$

Union
bound

Loose but
works

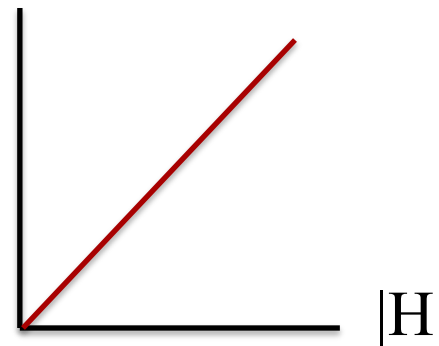
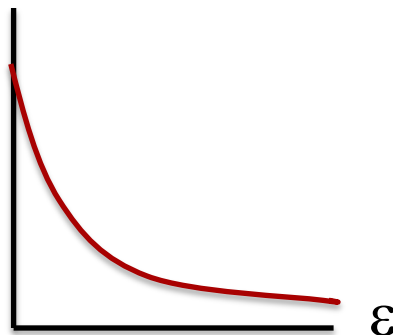
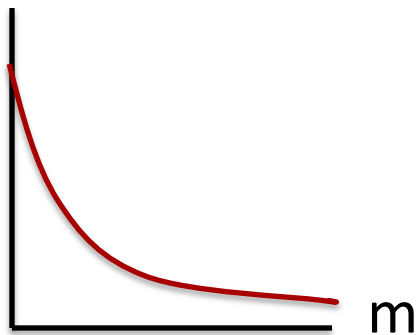
How likely is a learner to pick a bad classifier?

- Usually there are many many (say k) bad classifiers in the class

$$h_1, h_2, \dots, h_k \quad \text{s.t. } \text{error}_{\text{true}}(h_i) \geq \varepsilon \quad i = 1, \dots, k$$

- Probability that learner picks a bad classifier

$$\leq k (1-\varepsilon)^m \leq \underbrace{|H|}_{\substack{\text{Size of} \\ \text{model class}}} (1-\varepsilon)^m \leq \underbrace{|H|}_{\text{Size of model class}} e^{-\varepsilon m}$$



PAC (Probably Approximately Correct) bound

- **Theorem [Haussler'88]:** Model class H finite, dataset D with m i.i.d. samples, $0 < \epsilon < 1$: for any learned classifier h that gets 0 training error:

$$P(\text{error}_{true}(h) \geq \epsilon) \leq \underbrace{|H|}_{\leq \frac{1}{\epsilon}} e^{-m\epsilon} \leq \delta$$

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

$$\text{error}_{true}(h) \leq \epsilon$$

Important: PAC bound holds for all h with 0 training error, but doesn't guarantee that algorithm finds best h !!!

Using a PAC bound

80% accuracy
es. 20% error

4. prob ≥ 0.9
 $\delta = 0.1$

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

- Given ϵ and δ , yields sample complexity

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

- Given m and δ , yields error bound

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

Poll

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

Assume m is the minimum number of training examples sufficient to guarantee that with probability $1 - \delta$ a consistent learner using model class H will output a classifier with true error at worst ϵ .

Then a second learner that uses model space H' will require $2m$ training examples (to make the same guarantee) if $|H'| = 2|H|$.

- A. True B. False

If we double the number of training examples to $2m$, the error bound ϵ will be halved.

- C. True D. False

Limitations of Haussler's bound

- Only consider classifiers with 0 training error

h such that zero error in training, $\text{error}_{\text{train}}(h) = 0$

- Dependence on size of model class |H|

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

what if |H| too big or H is continuous (e.g. linear classifiers)?

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 $error_{train}(h) \neq 0$ in training set?
- The error of a classifier is like estimating the parameter of a coin!

$$\begin{aligned} \underbrace{error_{true}(h)} &:= \underbrace{P(h(X) \neq Y)} && \equiv \underbrace{P(H=1)} =: \theta \\ \underbrace{error_{train}(h)} &:= \underbrace{\frac{1}{m} \sum_i \mathbf{1}_{h(X_i) \neq Y_i}} && \equiv \underbrace{\frac{1}{m} \sum_i Z_i} =: \hat{\theta} \end{aligned}$$

$E[\hat{\theta}] = \theta$

Hoeffding's bound for a single classifier

- Consider m i.i.d. flips x_1, \dots, x_m , where $x_i \in \{0, 1\}$ of a coin with parameter θ . For $0 < \epsilon < 1$:

$$P\left(\left|\theta - \frac{1}{m} \sum_i x_i\right| \geq \epsilon\right) \leq 2e^{-2m\epsilon^2}$$

Hoeffding
 $e^{-m\epsilon}$

- Central limit theorem: z_1, \dots, z_m iid $E[z_i] = \mu$, $\text{var}(z_i) = \sigma^2$
 $\sqrt{m} \left(\frac{1}{m} \sum_i z_i - \mu \right) \rightarrow N(0, \sigma^2)$ $\theta \leq \theta(1-\theta) \leq \frac{1}{4}$



$$e^{-\frac{\epsilon^2}{2\sigma^2}} \sim e^{-2m\epsilon^2} \leftarrow \frac{1}{m} \sum_i x_i - \theta \rightarrow N\left(0, \frac{\sigma^2}{m}\right) \rightarrow N\left(0, \frac{1}{4m}\right)$$

Hoeffding's bound for a single classifier

- Consider m i.i.d. flips x_1, \dots, x_m , where $x_i \in \{0, 1\}$ of a coin with parameter θ . For $0 < \epsilon < 1$:

$$P \left(\left| \theta - \frac{1}{m} \sum_i x_i \right| \geq \epsilon \right) \leq 2e^{-2m\epsilon^2}$$

- For a single classifier h

$$P \left(\left| \text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) \right| \geq \epsilon \right) \leq 2e^{-2m\epsilon^2}$$

Hoeffding's bound for $|H|$ classifiers

- For each classifier h_i :

$$P(|\text{error}_{true}(h_i) - \text{error}_{train}(h_i)| \geq \epsilon) \leq 2e^{-2m\epsilon^2}$$

- What if we are comparing $|H|$ classifiers?

Union bound

- **Theorem:** Model class H finite, dataset D with m i.i.d. samples, $0 < \epsilon < 1$: for any learned classifier $h \in H$:

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

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

Summary of PAC bounds for finite model classes

With probability $\geq 1-\delta$,

1) For all $h \in H$ s.t. $\text{error}_{\text{train}}(h) = 0$,

$$\text{error}_{\text{true}}(h) \leq \varepsilon = \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$$

Haussler's bound

2) For all $h \in H$

$$|\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)| \leq \varepsilon = \sqrt{\frac{\ln |H| + \ln \frac{2}{\delta}}{2m}}$$

Hoeffding's bound

PAC bound and Bias-Variance tradeoff

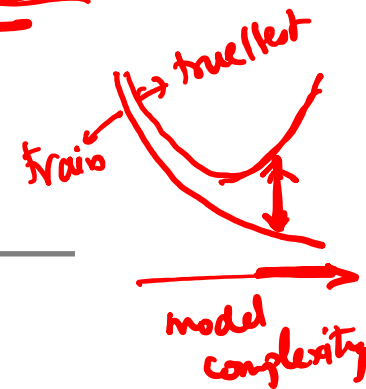
$$P(|\text{error}_{true}(h) - \text{error}_{train}(h)| \geq \epsilon) \leq 2|H|e^{-2m\epsilon^2} \leq \delta$$

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

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

- Fixed m

Model class		
complex	small	large
simple	large	small



What about the size of the model class?

$$2|H|e^{-2m\epsilon^2} \leq \delta$$

- Sample complexity

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

- How to measure the complexity of a model class?
 - E.g. decision trees:
 - trees with depth k
 - trees with k leaves

Number of decision trees of depth k

Recursive solution: *features*
 Given n **binary** attributes

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

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) = $n * H_{k-1} * H_{k-1}$

$$\hookrightarrow \log_2 \rightarrow \log_2 n + 2 \log_2 H_{k-1}$$

L_{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 $L_k = (2^k - 1)(1 + \log_2 n) + 1 \Rightarrow H_k \sim 2^{2^k}, m \sim 2^k$

PAC bound for decision trees of depth k

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

- 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, so we are over-counting a lot!

Number of decision trees with k leaves

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

H_k = Number of binary decision trees with k leaves

$$H_1 = 2$$

$$H_k = (\text{\#choices of root attribute})^*$$

$$\begin{aligned} & [(\text{\# left subtrees wth 1 leaf})^*(\text{\# right subtrees wth k-1 leaves}) \\ & + (\text{\# left subtrees wth 2 leaves})^*(\text{\# right subtrees wth k-2 leaves}) \\ & + \dots \\ & + (\text{\# left subtrees wth k-1 leaves})^*(\text{\# right subtrees wth 1 leaf})] \end{aligned}$$

$$H_k = n \sum_{i=1}^{k-1} H_i H_{k-i} = n^{k-1} C_{k-1} \quad (C_{k-1} : \text{Catalan Number})$$

Loose bound (using Sterling's approximation):

$$H_k \leq n^{k-1} \underline{\underline{2^{2k-1}}} \quad \Rightarrow m \sim 2k$$

Number of decision trees

- With k leaves $m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{2}{\delta} \right)$

$$\log_2 H_k \leq (k - 1) \log_2 n + 2k - 1 \quad \text{linear in } k$$

number of points m is linear in #leaves

- With depth k

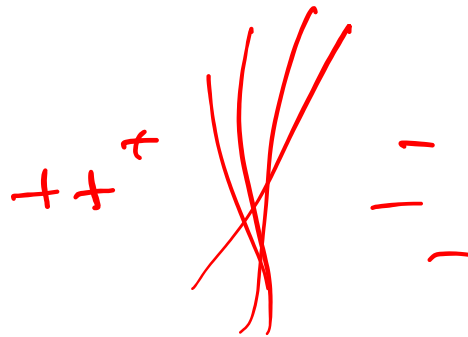
$$\log_2 H_k = (2^k - 1)(1 + \log_2 n) + 1 \quad \text{exponential in } k$$

number of points m is exponential in depth

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 can be classified using a classifier from this model space



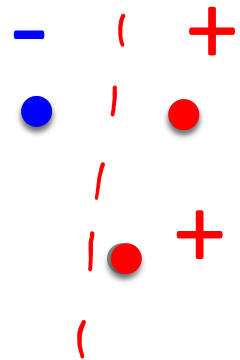
Rademacher Complexity

size of model space

- Instead of ~~all possible labelings~~, measure complexity by how accurately a model space can match a random labeling of the data.

For each data point i , draw random label

$$\sigma_i \quad \text{s.t.} \quad P(\sigma_i = +1) = \frac{1}{2} = P(\sigma_i = -1)$$



Then empirical Rademacher complexity of H is

$$\hat{R}_m(H) = \mathbb{E}_\sigma \left[\sup_{h \in H} \left(\frac{1}{m} \sum_{i=1}^m \sigma_i h(X_i) \right) \right]$$

Max correlation possible with random labels

Rademacher Bounds

- With probability $\geq 1-\delta$,

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \hat{R}_m(H) + 3\sqrt{\frac{\log(2/\delta)}{m}}$$

Handwritten notes: A red arrow points from the term $\frac{\ln(H) + \ln(1/\delta)}{m}$ (circled in red) to the $\hat{R}_m(H)$ term in the equation. Another red arrow points from the same term to the $\log(2/\delta)$ term in the equation.

where empirical Rademacher complexity of H

$$\hat{R}_m(H) = \mathbb{E}_\sigma \left[\sup_{h \in H} \left(\frac{1}{m} \sum_{i=1}^m \sigma_i h(X_i) \right) \right] \checkmark$$

is purely data-dependent.

Finite model class

- Rademacher complexity can be upper bounded in terms of model class size $|H|$:

$$\hat{R}_m(H) \leq \sqrt{\frac{2 \ln |H|}{m}}$$

- Often Rademacher bounds are significantly better, e.g. ...

Linear models with bounded norm

- Consider $h(X_i) = \langle \underline{w}, X_i \rangle$ with fixed $\|w\|, \|X_i\| \leq R$

$$\hat{R}_m(H) = \mathbb{E}_\sigma \left[\sup_{h \in H} \left(\frac{1}{m} \sum_{i=1}^m \sigma_i h(X_i) \right) \right]$$

⋮

$$\leq \frac{\|w\| R}{\sqrt{m}}$$

vs. $\sqrt{\frac{\ln |H|}{m}} = \infty$

Complexity increases with number of parameters d and norm of weights

Summary of PAC bounds

With probability $\geq 1-\delta$,

1) for all $h \in H$ s.t. $\text{error}_{\text{train}}(h) = 0$,

$$\text{error}_{\text{true}}(h) \leq \varepsilon = \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$$

Finite hypothesis space

2) for all $h \in H$,

$$|\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)| \leq \varepsilon = \sqrt{\frac{\ln |H| + \ln \frac{2}{\delta}}{2m}}$$

3) For all $h \in H$,

Infinite hypothesis space

$$|\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)| \leq \varepsilon = \hat{R}_m(H) + 3\sqrt{\frac{\log(2/\delta)}{m}}$$