

Learning Theory

Aarti Singh

Machine Learning 10-315
Nov 30, 2020

Slides courtesy: Carlos Guestrin



MACHINE LEARNING DEPARTMENT

The Carnegie Mellon logo, consisting of a grid of small white dots that form a larger, faint shape of the university's name.

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”?

*error rate /
accuracy*

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 \hat{h} that gets zero error in training
 - $\text{error}_{\text{train}}(\hat{h}) = 0$ $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\hat{h}(x_i) \neq y_i}$
- What is the probability that \hat{h} has more than ε true (= test) error?
 - $\text{error}_{\text{true}}(\hat{h}) \geq \varepsilon$ $\mathbb{P}(\hat{h}(x) \neq y) = \mathbb{E}[\mathbb{1}_{\hat{h}(x) \neq y}]$

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. $\text{error}_{\text{true}}(h) \geq \varepsilon$
- Probability that h gets one data point right $\leq 1 - \varepsilon$
- Probability that h gets m data points right

$$\begin{aligned} &P(\text{1st data pt right} \\ &\quad \wedge \text{2nd} \quad \dots \\ &\quad \wedge \dots \quad \dots) \\ &= \prod_{i=1}^m P(x_i \text{ right}) \\ &\leq (1 - \varepsilon)^m \end{aligned}$$

$$\begin{aligned} &P(h \text{ correctly classifies} \\ &\quad \text{all data pts}) \\ &= P(h \text{ has 0 training error}) \end{aligned}$$

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

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

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

$$P(A \cup B \dots) \leq P(A) + P(B) \dots$$

Union bound

Loose but works

$$\leq k (1-\varepsilon)^m$$

$$\leq k (1-\varepsilon)^m$$

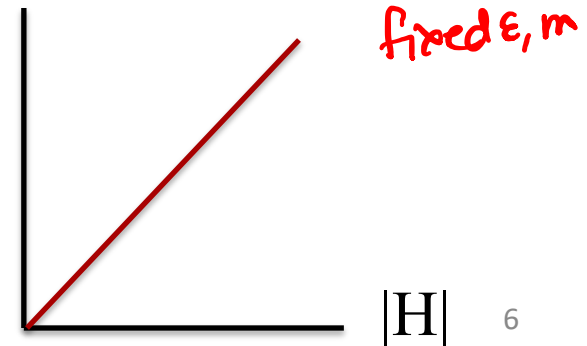
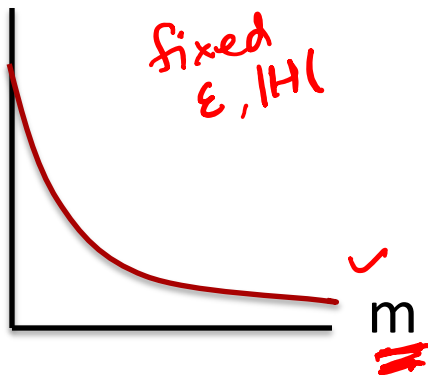
How likely is a learner to pick a bad classifier?

- 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

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



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 \hat{h} that gets 0 training error:

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

- Equivalently, with probability $\geq 1 - \delta$
 $\text{error}_{true}(\hat{h}) \leq \epsilon$

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

Using a PAC bound $= \delta$

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

- Given ϵ and δ , yields sample complexity

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

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

- Given m and δ , yields error bound

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

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)?

PAC bounds for finite model classes

H - Finite model class

e.g. decision trees of depth k

histogram classifiers with binwidth h

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

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! $E[\mathbb{1}_{h(X) \neq Y}]$

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

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 $0 < \epsilon < 1$:

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

- 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 :

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

- 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 $h \in H$:

$$\forall h \in H \quad P(A \cup B) \leq P(A) + P(B)$$

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

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

Summary of PAC bounds for finite model classes

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

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

With probability $\geq 1-\delta$,

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

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

Haussler's bound

2) For all $h \in H$

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

Hoeffding's bound

ϵ
 \rightarrow model complexity
 $|H|$

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 $\geq 1 - \delta$

$\forall h \in H$

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

- Fixed m

Model class	\downarrow	\downarrow
complex	small	large
simple	large	small

What about the size of the model class?

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

- Sample complexity

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

- How large is the model class? ←

1341

Number of decision trees of depth k

Recursive solution:

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

$$H_0 = 2$$

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

$$* (\# \text{ possible left subtrees})$$

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

Write $L_k = \log_2 H_k$

$$L_0 = 1$$

$$L_k = \log_2 n + 2L_{k-1} = \log_2 n + 2(\log_2 n + 2L_{k-2})$$

$$= \log_2 n + 2\log_2 n + 2^2\log_2 n + \dots + 2^{k-1}(\log_2 n + 2L_0)$$

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

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

PAC bound for decision trees of depth k

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

- 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

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

H_k = Number of binary decision trees with k leaves

$$H_1 = 2$$

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

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

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

Loose bound (using Sterling's approximation):

$$\log_2 H_k \leq n^{k-1} 2^{2k-1}$$

depth k
 $\log_2 H_k \sim 2^k$

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

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

With prob $\geq 1-\delta$ $\text{error}_{true}(h) \leq \text{error}_{train}(h) + \underbrace{\sqrt{\frac{\ln |H| + \ln \frac{2}{\delta}}{2m}}}_{\epsilon}$

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

$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{\underbrace{(k-1)}_{\text{Bias}} \ln n + (2k-1) \ln 2 + \ln \frac{2}{\delta}}{2m}}$$

	↓	↓	
$k = m$	0	-	large ($\sim > \frac{1}{2}$)
$k < m$	> 0	+	small ($\sim < \frac{1}{2}$)

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

Summary of PAC bounds for finite model class

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