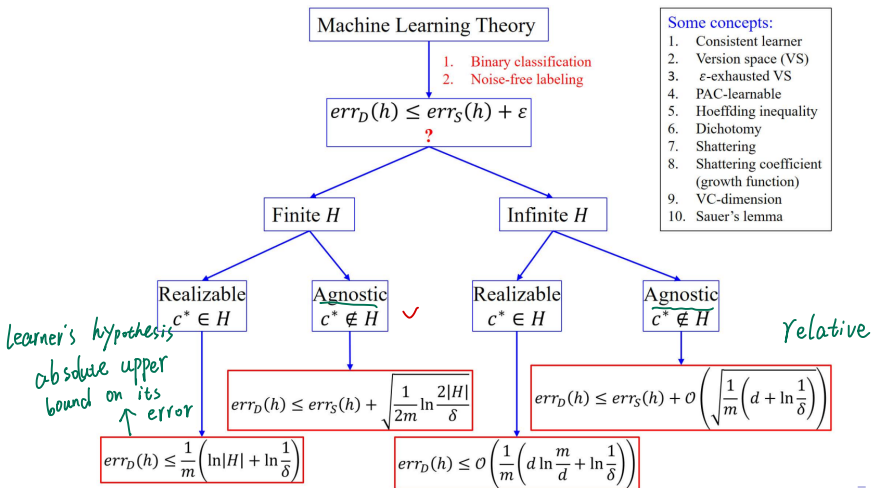


Discussions on Learning Theory

April 20, 2020

Big Picture of Today's Discussion

We recap some important concepts in the learning theory, which has been fully discussed in our classes.



What General Laws Constrain Inductive Learning

Generalizability of learning.

In machine learning it's really generalization error that we care about, but most learning algorithms fit their models to the training set.

- ✓ Sample Complexity
 - How many training examples are sufficient to learn target concept?
- ✓ Computational Complexity
 - Resources required to learn target concept? *Running time*
- We want theory to relate
 - ✓ Training examples.
 - Quantity *m*
 - Quality
 - How presented
 - Complexity of hypothesis space. *H*
 - Accuracy to which target function is approximated. *ϵ*
 - Probability of successful learning. *δ*

Problem Setting for Learning from Data

\mathcal{X} \mathcal{Y}

Given:

- Training data $S = ((x_1, y_1), \dots, (x_n, y_n)) \subset \mathcal{X} \times \mathcal{Y}$, where \mathcal{X} and \mathcal{Y} are termed input space and output space, respectively. *n pairs*
- Set of hypothesis $\mathcal{H} = \{h: \mathcal{X} \rightarrow \mathcal{Y}\}$. *$\{0, 1\}$*
- Set of possible target functions $\mathcal{C} = \{c: \mathcal{X} \rightarrow \mathcal{Y}\}$. *labeling function*
- Noise-free label $c(x)$.

Goal: Learner must output a hypothesis $h \in \mathcal{H}$ estimating c such that

$$h = \arg \min_{h \in \mathcal{H}} \text{error}_{\text{train}}(h).$$

- learner is trying to figure out the unknown labeling function.

Consistent Learner

With zero error on S
whenever \mathcal{H} contains such h .

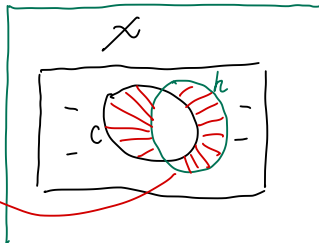
- Consistent Learner

- outputs hypothesis h that perfectly fits the training data S .

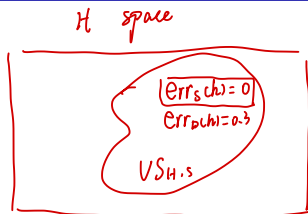
$$h(x) = c^*(x), \quad \forall x \in S.$$

- $\text{error}_{\text{train}}(h) = \Pr_{x \in S} [h(x) \neq c(x)]$

- $\text{error}_{\text{true}}(h) = \Pr_{x \in D} [h(x) \neq c(x)]$



Version Space



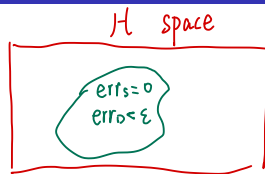
- Version Space (VS)

- set of all hypotheses $h \in H$ that correctly classify the training data S ,

$$VS_{H,S} = \{h \in H \mid \forall x \in S, \boxed{h(x)} = \underline{c^*}(x)\}.$$

- A consistent learner must produce a hypothesis in the $VS_{H,S}$ for H given S .

ϵ -exhausted Version Space



Definition

The version space $VS_{H,D}$ with respect to training data D is said to be ϵ -exhausted if every hypothesis h in $VS_{H,D}$ has true error less than ϵ .

$$(\forall h \in VS_{H,D}) \text{error}_{\text{true}}(h) < \epsilon.$$

ϵ -exhausted Version Space

- One can never be sure that $VS_{H,D}$ is ϵ -exhausted.
but one can bound the probability that it is not.

How many examples will ϵ -exhaust the VS?

Theorem: [Haussler, 1988].

If the hypothesis space H is finite and D is a sequence of $m \geq 1$ independent random examples of some target concept c , then for any $0 \leq \epsilon \leq 1$, the probability that the version space with respect to H and D is not ϵ -exhausted (with respect to c) is less than

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

Interesting! This bounds the probability that any consistent learner will output a hypothesis h with $\text{error}(h) \geq \epsilon$

Any(!) learner that outputs a hypothesis consistent with all training examples (i.e., an h contained in $VS_{H,D}$)

PAC-learnable (Probably Approximately Correct)

PAC - Framework.

Formal Def.

Consider a class C of possible target concepts defined over a set of instances X of length n , and a learner L using hypothesis space H .

Definition: C is **PAC-learnable** by L using H if for all $c \in C$, distributions \mathcal{D} over X , ϵ such that $0 < \epsilon < 1/2$, and δ such that $0 < \delta < 1/2$, learner L will with probability at least $(1 - \delta)$ output a hypothesis $h \in H$ such that $\text{error}_\mathcal{D}(h) \leq \epsilon$ in time that is polynomial in $1/\epsilon$, $1/\delta$, n and $\text{size}(c)$.

Sufficient condition:

Holds if learner L requires only a polynomial number of training examples, and processing per example is polynomial

a reasonable amount of computation!

- With high probability learns a close approximation to the C^* .
- ϵ, δ , at least $(1 - \delta)$, a learner can learn a concept with error at most ϵ .

Hoeffding's Inequality

Hoeffding bounds

HW3.

Consider coin of bias p flipped m times.
Let N be the observed # heads. Let $\epsilon \in [0,1]$.

Hoeffding bounds:

- $\Pr[N/m > p + \epsilon] \leq e^{-2m\epsilon^2}$, and
- $\Pr[N/m < p - \epsilon] \leq e^{-2m\epsilon^2}$.

$$\Pr \left[\left| \frac{N}{m} - p \right| > \epsilon \right] \leq 2e^{-2m\epsilon^2}.$$

Exponentially decreasing tails

Remark.

- **Tail inequality:** bound probability mass in tail of distribution (how concentrated is a random variable around its expectation).

if we use $\frac{N}{m}$, #heads
r.v.
Empirical average, as estimate of p , then
the probability of our being far from the true value is small with
sufficiently large m .

Dichotomy

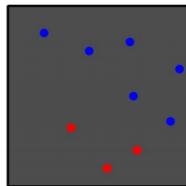
$$h: S \rightarrow \{+1, -1\}$$

Given a dataset $S = \{x_1, \dots, x_m\}$,

$(\underline{h(x_1)}, \dots, \underline{h(x_m)})$ tuple

A dichotomy of S

1. If H is diverse, we get many different dichotomies.
2. If H contains many similar functions, we only get a few dichotomies.



dichotomy

• $H(S)$ 是在 $S = \{x_1, \dots, x_m\}$ 所有 dichotomy 的一个集合.
Specified. set.

$$H(S) = \{ (h(x_1), \dots, h(x_m)) : h \in H \}$$

Shattering Coefficient (Growth Function)

- $H[m]$ - max number of ways to split m points using concepts in H

$$\underbrace{H[m]}_0 = \max_{|S|=m} |H[S]| \leq 2^m.$$

For arbitrary S . 在 m 个 points 上能表达的最大的
dichotomy 的数量.

H 能在 S 上实现所有的 dichotomy.

Shattering, VC-dimension

Definition: H shatters S if $|H[S]| = 2^{|S|}$.

A set of points S is shattered by H if there are hypotheses in H that split S in all of the $2^{|S|}$ possible ways, all possible ways of classifying points in S are achievable using concepts in H .

Definition: VC-dimension (Vapnik-Chervonenkis dimension)

The VC-dimension of a hypothesis space H is the cardinality of the largest set S that can be shattered by H .

If arbitrarily large finite sets can be shattered by H , then $\text{VCdim}(H) = \infty$

$$\text{VCdim}(H) = \max \{ |S| : S \text{ shattered by } H \}$$

$$\text{VCdim}(h) = d$$

- ① Lower bound: show there exists a set of d points that can be shattered.
- ② upper bound: show no set of $d+1$ points that can be shattered

Sauer's Lemma

Sauer's Lemma:

Let $d = \text{VCdim}(H) < \infty$

- $m \leq d$, then $H[m] = 2^m$
- $m > d$, then $H[m] = O(m^d) \Rightarrow$

$$\max_{|S|=m} |H[S]|$$

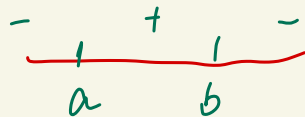
$$H[m] \leq \sum_{i=0}^d \binom{m}{i} = \mathbb{I}_d(m) = O(m^d)$$

- This Lemma states that $H[m]$ grows exponentially fast for $m \leq d$.
But then only grows like a polynomial for $m > d$.
- $\text{VCdim}(H) < \infty$, $m \rightarrow \infty$, implies "Learnability".

eg: The interval example.

It is the class of intervals on the \mathbb{R}
given $h = [a, b]$ with $a, b \in \mathbb{R}$, $a \leq b$

$$h(x) = \begin{cases} +1, & \text{if } a \leq x \leq b, \\ -1, & \text{otherwise} \end{cases}$$



$$H(m) = \binom{m}{2} + \binom{m}{1} + 1 = \binom{m}{2} + m + 1$$

↓ Sauer's Lemma is tight!

Proof:

- At least 1 +1 follows by at least

-1. $\binom{m}{2}$

- At least 1 +1 but no negative -1 followed the +1. $\binom{m}{1}$

- all the -1.

_____ []
- - - -



$$|H(m)| \leq \mathbb{I}_d(m)$$

- Base case (for two variable).

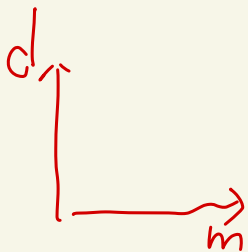
- $m=0$, for any d . $\mathbb{I}_d(m) = \sum_{i=0}^d \binom{0}{i}$

$$|H(0)| \leq 1. \text{ Since we can}$$

label 0 points at most 1 ways.

- $d=0$, for any m . $\mathbb{I}_d(m) = \sum_{i=0}^0 \binom{m}{i} = 1.$

$|H(m)| = 1$ since $\text{VC dim}(H) = 0$ implies we label everything with the same label.



- Inductive step:

$$\boxed{\binom{m}{k} = \binom{m-1}{k} + \binom{m-1}{k-1}} \quad \checkmark$$

included the first item.

• $\binom{m}{k} = 0$ if $k < 0$ or $k > m$

Assume lemma holds $\bar{m} + \bar{d} < m + d$.

Given $S = (x_1, \dots, x_m)$, we want to show $|\mathcal{H}(S)| \leq 2^d(m)$.

key ideas:

Construct two new Hypothesis: H_1 and H_2 .

H_1 和 H_2 在 $S' = \{x_1, \dots, x_{m-1}\}$ 忽略了 x_m

H_1 : 是 H 的一个最小子集且能 S' shattering.

H_2 : 如果有 2个 hypothesis label S' 以相同的方式. 一个假在 H_1 里, 一个假在 H_2 里.

e.g:



\mathcal{H}

	x_1	x_2	x_3	x_4	x_5
h_1	0	1	1	0	0
h_2	0	1	1	0	1
h_3	0	1	1	1	0
h_4	1	0	0	1	0
h_5	1	0	0	1	1
h_6	1	0	0	0	1

\mathcal{H}_1

x_1	x_2	x_3	x_4
0	1	1	0
0	1	1	1
1	0	0	1
1	0	0	0

\mathcal{H}_2

x_1	x_2	x_3	x_4
0	1	1	0
1	0	0	1

✓ if a set is shattered by \mathcal{H}_1 , then it is also shattered by \mathcal{H} . Because can generate \mathcal{H} by using the same as \mathcal{H}_1 is when we generate \mathcal{H}_1 . Thus,

$$VC \dim(\mathcal{H}_1) \leq VC \dim(\mathcal{H}) = d.$$

• ✓ if a set T is shattered by H_2 , then

$T \cup \{x_m\}$ is shattered by H .

Because, H 中会有 两个对应的 Hypothesis
和 H_2 有着相同的 element. 通过加上 x_{m-1}
and $x_m = 0$. Thus,

$$VC \dim(H) \geq VC \dim(H_2) + 1$$

Now, by inductive,

$$\begin{aligned} |H(s)| &= |H_1(s')| + |H_2(s')| \\ &\leq \sum_{\bar{v}=0}^d \binom{m-1}{\bar{v}} + \sum_{\bar{v}=0}^{d-1} \binom{m-1}{\bar{v}} \\ &= \sum_{\bar{v}=0}^d \binom{m-1}{\bar{v}} + \boxed{\sum_{\bar{v}=0}^d \binom{m-1}{\bar{v}-1}} \\ &= \sum_{\bar{v}=0}^d \binom{m}{\bar{v}} = \underline{\Phi_d(m)} \end{aligned}$$

