

Computational Learning Theory

- Theoretical characterization of the **difficulties** and **capabilities** of learning algorithms.
- Questions:
 - Conditions for successful/unsuccessful learning
 - Conditions of success for particular algorithms
- Two frameworks:
 - Probably Approximately Correct (PAC) framework: classes of hypotheses that can be learned; complexity of hypothesis space and bound on training set size.
 - Mistake bound framework: number of training errors made before correct hypothesis is determined.

1

Specific Questions

- Sample complexity: How many training examples are needed for a learner to converge?
- Computational complexity: How much computational effort is needed for a learner to converge?
- Mistake bound: How many training examples will the learner misclassify before converging?

Issues: When to say it was successful? How are inputs acquired?

3

Computational Learning Theory

What general laws constrain inductive learning?

We seek theory to relate:

- Probability of successful learning
- Number of training examples
- Complexity of hypothesis space
- Accuracy to which target concept is approximated
- Manner in which training examples presented

2

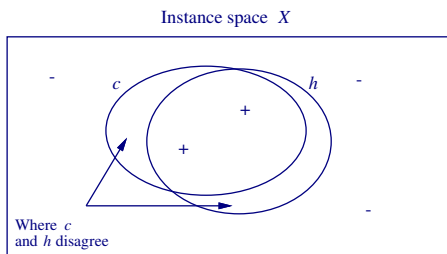
Sample Complexity

How many training examples are sufficient to learn the target concept?

1. If learner proposes instances, as queries to teacher
 - Learner proposes instance x , teacher provides $c(x)$
2. If teacher (who knows c) provides training examples
 - teacher provides sequence of examples of form $\langle x, c(x) \rangle$
3. If some random process (e.g., nature) proposes instances
 - instance x generated randomly, teacher provides $c(x)$

4

True Error of a Hypothesis

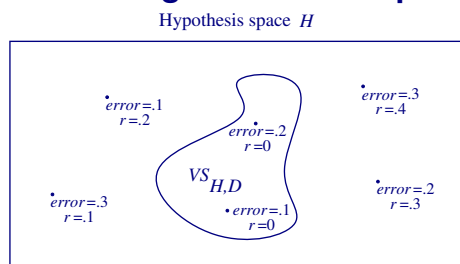


Definition: The **true error** (denoted $error_{\mathcal{D}}(h)$) of hypothesis h with respect to target concept c and distribution \mathcal{D} is the probability that h will misclassify an instance drawn at random according to \mathcal{D} .

$$error_{\mathcal{D}}(h) \equiv \Pr_{x \in \mathcal{D}} [c(x) \neq h(x)]$$

5

Exhausting the Version Space



(r = training error, $error$ = true error)

Definition: The version space $VS_{H,D}$ is said to be **ϵ -exhausted** with respect to c and \mathcal{D} , if every hypothesis h in $VS_{H,D}$ has error less than ϵ with respect to c and \mathcal{D} .

$$(\forall h \in VS_{H,D}) error_{\mathcal{D}}(h) < \epsilon$$

7

Two Notions of Error

Training error of hypothesis h with respect to target concept c

- How often $h(x) \neq c(x)$ over training instances

True error of hypothesis h with respect to c

- How often $h(x) \neq c(x)$ over future random instances

Our concern:

- Can we bound the true error of h given the training error of h ?
- First consider when training error of h is zero (i.e., $h \in VS_{H,D}$)

6

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

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

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

If we want this probability to be below δ

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

then

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

8

Proof of ϵ -Exhausting Theorem

Theorem: Prob. of $VS_{H,D}$ not being ϵ -exhausted is $\leq |H|e^{-\epsilon m}$.

Proof:

- Let $h_i \in H$ ($i = 1..k$) be those that have true error greater than ϵ w.r.t c ($k \leq |H|$).
- We fail to ϵ -exhaust the VS iff at least one h_i is consistent with all m sample training instances (note: they have true error greater than ϵ).
- Prob. of a single hypothesis with error $> \epsilon$ is consistent for one random sample is at most $(1 - \epsilon)$.
- Prob. of that hypothesis being consistent with m samples is $(1 - \epsilon)^m$.
- Prob. of at least one of k hypotheses with error $> \epsilon$ is consistent with m samples is $k(1 - \epsilon)^m$.
- Since $k \leq |H|$, and for $0 \leq \epsilon \leq 1$, $(1 - \epsilon) \leq e^{-\epsilon}$:

$$k(1 - \epsilon)^m \leq |H|(1 - \epsilon)^m \leq |H|e^{-\epsilon m}$$

9

Prototypical Concept Learning Task

- Given:**
 - Instances X : Possible days, each described by the attributes *Sky*, *AirTemp*, *Humidity*, *Wind*, *Water*, *Forecast*
 - Target function c : $EnjoySport : X \rightarrow \{0, 1\}$
 - Hypotheses H : Conjunctions of literals. E.g.

$$\langle ?, Cold, High, ?, ?, ? \rangle.$$

- Training examples D : Positive and negative examples of the target function

$$\langle x_1, c(x_1) \rangle, \dots \langle x_m, c(x_m) \rangle$$

- Determine:**
 - A hypothesis h in H such that $h(x) = c(x)$ for all x in D ?
 - A hypothesis h in H such that $h(x) = c(x)$ for all x in X ?

11

Learning Conjunctions of Boolean Literals

How many examples are sufficient to assure with probability at least $(1 - \delta)$ that

$$\text{every } h \text{ in } VS_{H,D} \text{ satisfies } error_{\mathcal{D}}(h) \leq \epsilon$$

Use our theorem:

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

Suppose H contains conjunctions of constraints on up to n boolean attributes (i.e., n boolean literals). Then $|H| = 3^n$, and

$$m \geq \frac{1}{\epsilon} (\ln 3^n + \ln(1/\delta))$$

or

$$m \geq \frac{1}{\epsilon} (n \ln 3 + \ln(1/\delta))$$

10

What should m be in *EnjoySport*?

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

If H is as given in *EnjoySport* then $|H| = 973$, and

$$m \geq \frac{1}{\epsilon} (\ln 973 + \ln(1/\delta))$$

... if we want to assure that with probability 95%, VS contains only hypotheses with $error_{\mathcal{D}}(h) \leq 0.1$, then it is sufficient to have m examples, where

$$m \geq \frac{1}{0.1} (\ln 973 + \ln(1/0.05))$$

$$m \geq 10(\ln 973 + \ln 20)$$

$$m \geq 10(6.88 + 3.00)$$

$$m \geq 98.8$$

12

PAC Learning

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 $error_{\mathcal{D}}(h) \leq \epsilon$, in time that is polynomial in $1/\epsilon$, $1/\delta$, n and $size(c)$.

13

PAC-Learnability of Conjunction of Boolean Literals

C : Conjunction of boolean variable or negation of those, e.g., $Old \wedge \neg Tall$.
Is it PAC-learnable?

- FIND-S learns the class in linear time of n (n literals), thus C is PAC-learnable.
- Proof:
 - Sample complexity is polynomial in n , $1/\delta$, and $1/\epsilon$, and independent of $size(c)$:

$$m \geq \frac{1}{\epsilon} (\ln |H| + \ln(1/\delta)) = \frac{1}{\epsilon} (n \ln 3 + \ln(1/\delta))$$

- To process each training instance, FIND-S algorithm requires effort linear in n , independent of $1/\delta$, $1/\epsilon$, and $size(c)$.

15

Agnostic Learning

So far, we assumed that $c \in H$. What if it is not the case?

Agnostic learning setting: don't assume $c \in H$

- What do we want then?
 - The hypothesis h that makes fewest errors on training data
- What is sample complexity in this case?

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

derived from Hoeffding bounds:

$$Pr[error_{\mathcal{D}}(h) > error_D(h) + \epsilon] \leq e^{-2m\epsilon^2}$$

14

PAC-Learnability of Unbiased Concept Class

- The concept class C is the set of all subsets of X :

$$|C| = 2^{|X|} = 2^{2^n}$$

- In that case H must equal C , for a learner to be successful.
- The sample complexity in that case has upper bound:

$$m \geq \frac{1}{\epsilon} (2^n \ln 2 + \ln(1/\delta)),$$

which is exponential in n . So, it is not PAC-learnable.

Other concept classes: k -DNF (polynomial sample complexity, but unsolvable in polynomial time), k -CNF (polynomial sample complexity, polynomial time complexity).

16

Shattering a Set of Instances

Definition: a **dichotomy** of a set S is a partition of S into two disjoint subsets.

Definition: a set of instances S is **shattered** by hypothesis space H if and only if for every dichotomy of S there exists some hypothesis in H consistent with this dichotomy.

17

The Vapnik-Chervonenkis Dimension

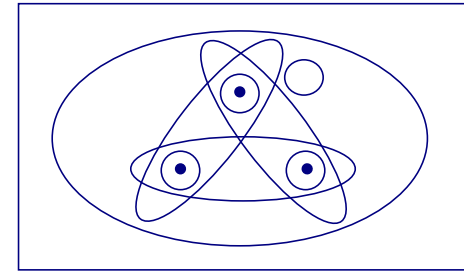
Definition: The **Vapnik-Chervonenkis dimension**, $VC(H)$, of hypothesis space H defined over instance space X is the size of the largest finite subset of X shattered by H . If arbitrarily large finite sets of X can be shattered by H , then $VC(H) \equiv \infty$.

Note that $|H|$ can be infinite, while $VC(H)$ finite!

19

Three Instances Shattered

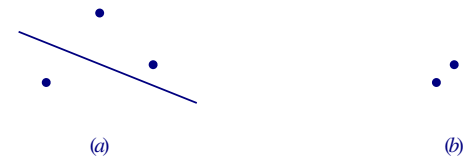
Instance space X



Each closed contour indicates one dichotomy. What kind of hypothesis space H can shatter the instances?

18

VC Dim. of Linear Decision Surfaces



- When H is a set of lines, and S a set of points, $VC(H) = 3$.
- (a) can be shattered, but (b) cannot be. However, if at least one subset of size 3 can be shattered, that's fine.
- Set of size 4 cannot be shattered, for any combination of points (think about an XOR-like situation).

20

VC Dimension: Another Example

$S = \{3.1, 5.7\}$, and hypothesis space includes intervals $a < x < b$.

- Dichotomies: both, none, 3.1, or 5.7.
- Are there intervals that cover all the above dichotomies?

What about $S = x_0, x_1, x_2$ for an arbitrary x_i ? (cf. collinear points).

21

Mistake Bounds

So far: how many examples needed to learn?

What about: how many mistakes before convergence?

- This is an interesting question because some learning systems may need to start operating while still learning.

Let's consider similar setting to PAC learning:

- Instances drawn at random from X according to distribution \mathcal{D} .
- Learner must classify each instance before receiving correct classification from teacher.
- Can we bound the number of mistakes learner makes before converging?

23

Sample Complexity from VC Dimension

How many randomly drawn examples suffice to ϵ -exhaust $VS_{H,D}$ with probability at least $(1 - \delta)$?

$$m \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 8VC(H) \log_2(13/\epsilon))$$

$VC(H)$ is directly related to the sample complexity:

- More expressive H needs more samples.
- More samples needed for H with more tunable parameters.

22

Mistake Bounds: Find-S

Consider Find-S when $H =$ conjunction of boolean literals

Find-S:

- Initialize h to the most specific hypothesis $l_1 \wedge \neg l_1 \wedge l_2 \wedge \neg l_2 \dots l_n \wedge \neg l_n$
- For each positive training instance x
 - Remove from h any literal that is not satisfied by x
- Output hypothesis h .

How many mistakes before converging to correct h ?

24

Mistake Bound of FIND-S

FIND-S starts with $\langle \emptyset, \emptyset, \dots, \emptyset \rangle$, the most specific hypothesis:

Conjunction of $2n$ literals.

- After the first sample is observed, half of the $2n$ gets eliminated.
- For the remaining positive examples that is misclassified, at least one or more of the remaining n literals must be eliminated.
- Thus, the total number of mistakes can be at most $n + 1$ in the worst case.

25

Mistake Bound of Halving Algorithm

- Start with version space = H .
- Mistake is made when more than half of the $h \in H$ misclassified.
- In that case, at most half of $h \in VS$ will be eliminated.
- That is, each **mistake** reduces the VS by half.
- Initially $|VS| = |H|$, and each mistake halves the VS , so it takes $\log_2 |H|$ mistakes to reduce $|VS|$ to 1.
- Actual worst-case bound is $\lfloor \log_2 |H| \rfloor$.

27

Mistake Bounds: Halving Algorithm

Consider the Halving Algorithm:

- Learn concept using version space *Candidate-Elimination* or *List-Then-Eliminate* algorithm.
- Classify new instances by majority vote of version space members.

How many mistakes before converging to correct h ?

- ... in worst case?
- ... in best case?

26

Optimal Mistake Bounds

Let $M_A(C)$ be the max number of mistakes made by algorithm A to learn concepts in C . (maximum over all possible $c \in C$, and all possible training sequences)

$$M_A(C) \equiv \max_{c \in C} M_A(c)$$

Definition: Let C be an arbitrary non-empty concept class. The **optimal mistake bound** for C , denoted $Opt(C)$, is the minimum over all possible learning algorithms A of $M_A(C)$.

$$Opt(C) \equiv \min_{A \in \text{learning algorithms}} M_A(C)$$

28

Mistake Bounds and VC Dimension

Littlestone (1987) showed:

$$VC(C) \leq Opt(C) \leq M_{Halving}(C) \leq \log_2(|C|)$$