

Intro. to machine learning (CSI 5325)

Lecture 15: Learning theory

Greg Hamerly

Spring 2008

Some content from Tom Mitchell.

- 1 Computational learning theory
 - Setting 1: learner poses queries to teacher
 - Setting 2: teacher chooses examples
 - Setting 3: randomly generated instances, labeled by teacher
- 2 Hypothesis error
- 3 Probably approximately correct (PAC) learning

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

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 ?

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 $c(x)$

Sample Complexity: 1

Learner proposes instance x , teacher provides $c(x)$
(assume c is in learner's hypothesis space H)

Optimal query strategy: play 20 questions

- pick instance x such that half of hypotheses in V_S classify x positive, half classify x negative
- When this is possible, need $\lceil \log_2 |H| \rceil$ queries to learn c
- when not possible, need even more

Sample Complexity: 2

Teacher (who knows c) provides training examples
(assume c is in learner's hypothesis space H)

Optimal teaching strategy: depends on H used by learner

Consider the case $H =$ conjunctions of up to n boolean literals and their negations

e.g., $(AirTemp = Warm) \cap (Wind = Strong)$, where
 $AirTemp, Wind, \dots$ each have 2 possible values.

- if n possible boolean attributes in H , $n + 1$ examples suffice
- why?

Sample Complexity: 3

Given:

- set of instances X
- set of hypotheses H
- set of possible target concepts C
- training instances generated by a fixed, unknown probability distribution \mathcal{D} over X

Learner observes a sequence D of training examples of form $\langle x, c(x) \rangle$, for some target concept $c \in C$

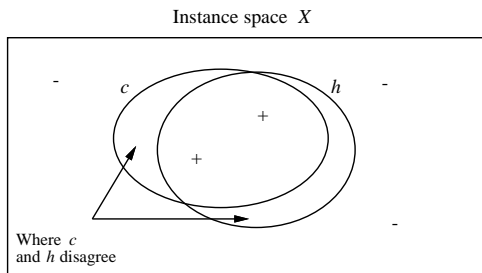
- instances x are drawn from distribution \mathcal{D}
- teacher provides target value $c(x)$ for each

Learner must output a hypothesis h estimating c

- h is evaluated by its performance on subsequent instances drawn according to \mathcal{D}

Note: randomly drawn instances, noise-free classifications

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} .

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

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

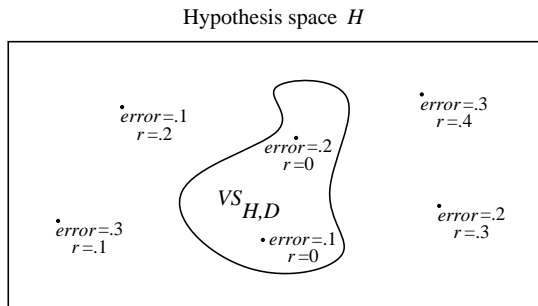
True error of hypothesis h with respect to c

- 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., $h \in VS_{H,D}$)

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} .

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

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 $error(h) \geq \epsilon$

If we want to this probability to be below δ

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

then

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

Proving Haussler's result

- Assume that h_1, h_2, \dots, h_k all have $error_{\mathcal{D}}(h_i) \geq \epsilon$
- Training has failed to ϵ -exhaust the version space iff one or more of h_i is consistent with all m training examples

$$Pr[\text{a single } h_i \text{ consistent w/one random example}] \leq (1 - \epsilon)$$

$$Pr[\text{a single } h_i \text{ consistent w/all } m \text{ examples}] \leq (1 - \epsilon)^m$$

$$Pr[\text{any } h_i \text{ consistent w/all } m \text{ examples}] \leq k(1 - \epsilon)^m$$

Further, $k \leq |H|$, so we get

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

Also, if $0 \leq \epsilon \leq 1$, then $(1 - \epsilon) \leq e^{-\epsilon}$, so

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

Learning Conjunctions of Boolean Literals

How many examples are sufficient to assure with probability at least $(1 - \delta)$ that every h in $VS_{H,D}$ 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

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

How About *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 want to assure that with probability 95%, VS contains only hypotheses with $error_{\mathcal{D}}(h) \leq .1$, then it is sufficient to have m examples, where

$$\begin{aligned} m &\geq \frac{1}{.1} (\ln 973 + \ln(1/.05)) \\ &= 10 (\ln 973 + \ln 20) \\ &\approx 10 (6.88 + 3.00) \\ &\approx 98.8 \end{aligned}$$

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

Agnostic Learning

So far, assumed $c \in H$

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:

$$Pr[\text{error}_D(h) > \text{error}_D(h) + \epsilon] \leq e^{-2m\epsilon^2}$$

Comparing consistent and agnostic sample bounds

To ϵ -exhaust the version space:

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

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

For a difference of only $\frac{1}{2\epsilon}$ times more examples.

- where does this difference come from?