# Machine Learning 10-701

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

February 24, 2011

Today:

- Computational Learning Theory
- PAC learning theorem
- VC dimension

Recommended reading:

- Mitchell: Ch. 7
- suggested exercises: 7.1, 7.2, 7.7

---

## 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 function is approximated
- Manner in which training examples presented

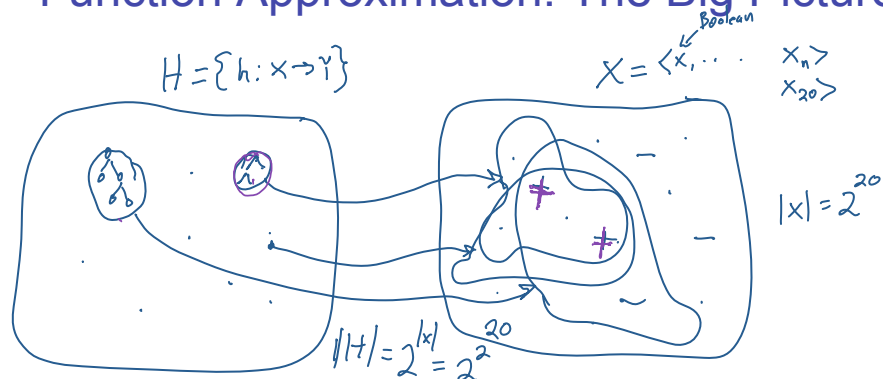* see Annual Conference on Learning Theory (COLT)

## 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 $c(x)$

---

# Function Approximation: The Big Picture

$$H = \{ h : X \to Y \}$$

$$X = \langle x_1, \dots \; x_n, \; x_{20} \rangle \quad \text{Boolean}$$

$$|x| = 2^{20}$$

$$\|H\| = 2^{|x|} = 2^{2^{20}}$$

How many labeled examples are needed in order to determine which of the $2^{2^{20}}$ hypotheses is the correct one?

All $2^{20}$ instances in X must be labeled!

There is no free lunch!

Inductive inference - generalizing beyond the training data is impossible unless we add more assumptions (eg. priors over H)

2

## Sample Complexity: 3

Given:

- set of instances $X$
- set of hypotheses $H$
- set of possible target concepts $C = \{c : X \to \{0,1\}\}$
- training instances generated by a fixed, unknown probability distribution $\mathcal{D}$ over $X$ $\leftarrow P(x)$

Learner observes a sequence $D$ of training examples of form $\langle x, c(x) \rangle$, for some target concept $c \in C$
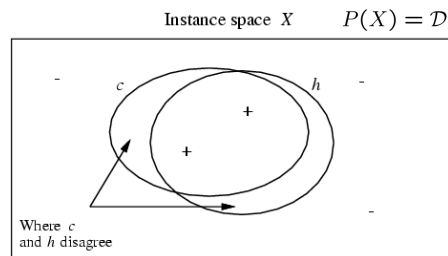
- instances $x$ are drawn from distribution $\mathcal{D}$
- teacher provides target value $c(x)$ for each

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



Instance space $X$      $P(X) = \mathcal{D}$

Where $c$ and $h$ disagree

**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 $\mathcal{D}$.

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

3

## 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 D

$$error_D(h) \equiv \Pr_{x \in D}[c(x) \neq h(x)] \equiv \frac{\sum_{x \in D} \delta(c(x) \neq h(x))}{|D|}$$

training examples

*True error* of hypothesis $h$ with respect to $c$

- How often $h(x) \neq c(x)$ over future instances drawn at random from $\mathcal{D}$

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

Probability distribution P(x)

---

## Two Notions of Error

Can we bound

$$error_{\mathcal{D}}(h)$$

in terms of

$$error_D(h)$$

??

*Training error* of hypothesis $h$ with respect to target concept $c$

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

$$error_D(h) \equiv \Pr_{x \in D}[c(x) \neq h(x)] \equiv \frac{\sum_{x \in D} \delta(c(x) \neq h(x))}{|D|}$$

training examples

*True error* of hypothesis $h$ with respect to $c$

- How often $h(x) \neq c(x)$ over future instances drawn at random from $\mathcal{D}$

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

Probability distribution P(x)

$$error_{\mathrm{D}}(h) \equiv \Pr_{x \in \mathrm{D}}[c(x) \neq h(x)] \equiv \frac{\sum_{x \in \mathrm{D}} \delta(c(x) \neq h(x))}{|\mathrm{D}|}$$

training examples

Can we bound $error_{\mathcal{D}}(h)$ _true_

in terms of $error_{\mathrm{D}}(h)$ _train_

??

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

Probability distribution P(x)

_test new_

if D was a set of examples drawn from $\mathcal{D}$ and **_independent_** of h, then we could use standard statistical confidence intervals to determine that with 95% probability, $error_{\mathcal{D}}(h)$ lies in the interval:

$$error_{\mathrm{D}}(h) \pm 1.96 \sqrt{\frac{error_{\mathrm{D}}(h)(1 - error_{\mathrm{D}}(h))}{n}}$$

but D is the **_training data_** for h ….
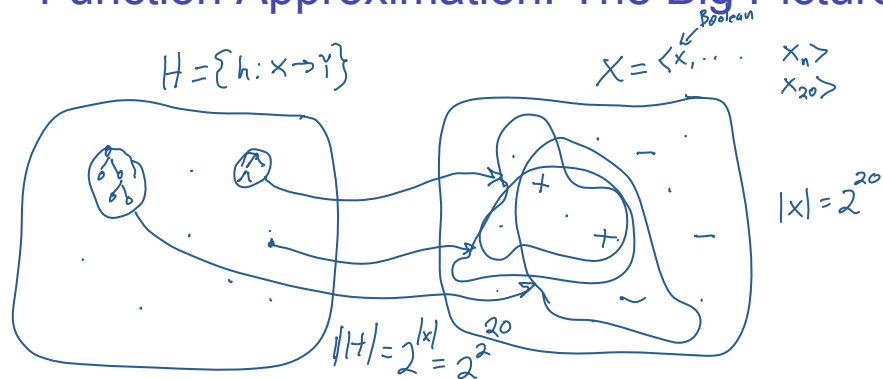
## Version Spaces

c: X → {0,1}

A hypothesis $h$ is **consistent** with a set of training examples $D$ of target concept $c$ if and only if $h(x) = c(x)$ for each training example $\langle x, c(x) \rangle$ in $D$.

$$Consistent(h, D) \equiv (\forall \langle x, c(x) \rangle \in D)\, h(x) = c(x)$$

The **version space**, $VS_{H,D}$, with respect to hypothesis space $H$ and training examples $D$, is the subset of hypotheses from $H$ consistent with all training examples in $D$.

$$VS_{H,D} \equiv \{h \in H | Consistent(h, D)\}$$
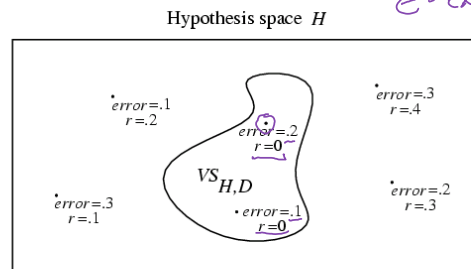
## Function Approximation: The Big Picture

$H = \{h: X \to Y\}$

Boolean

$X = \langle x_1 \cdots \quad x_n \rangle$ $x_{20} \rangle$

$|x| = 2^{20}$

$\|H\| = 2^{|x|} = 2^{2^{20}}$

How many labeled examples are needed in order to determine which of the $2^{2^{20}}$ hypotheses is the correct one?

All $2^{20}$ instances in X must be labeled!

There is no free lunch!

Inductive inference - generalizing beyond the training data is impossible unless we add more assumptions (eg. priors over H)

---

## Exhausting the Version Space

ε - exhausted for ε > .3

Hypothesis space $H$

- error=.1, r=.2
- error=.2, r=0
- error=.3, r=.4
- $VS_{H,D}$
- error=.3, r=.1
- error=.1, r=0
- error=.2, r=.3

$(r = \text{training error}, error = \text{true error})$

**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 true error less than ε with respect to $c$ and $\mathcal{D}$.

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

How many examples will $\epsilon$-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 $\epsilon$-exhausted (with respect to $c$) is less than

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

---

Interesting! This bounds the probability that <u>any consistent learner</u> will output a hypothesis $h$ with $error(h) \geq \epsilon$

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

7

hyp space H          let $h_1, \ldots h_k$ be the hyps hett with true error $\geq \epsilon$
instances X
fn  $c: X \rightarrow \{0, 1\}$
m labeled examps
error tolerance $\epsilon$

Prob that $h_1$ will be consistent with first training example
$$\leq (1-\epsilon)$$
"    $h_1$  will be cons. w/ m indep drawn examps?
$$\leq (1-\epsilon)^m$$
" that at least of $h_1 \ldots h_k$ will be consist w/m  "  ?

$k \leq |H|$          $\leq k (1-\epsilon)^m$
$$\leq |H| (1-\epsilon)^m$$          if $0 \leq \epsilon \leq 1$
$$\boxed{\leq |H| e^{-\epsilon m}}$$          then $(1-\epsilon) \leq e^{-\epsilon}$

---

## What it means

[Haussler, 1988]: probability that the version space is not ε-exhausted after *m* training examples is at most $|H| e^{-\epsilon m}$

$$\Pr[(\exists h \in H) s.t. (error_{train}(h) = 0) \wedge (error_{true}(h) > \epsilon)] \leq |H| e^{-\epsilon m}$$

↑

Suppose we want this probability to be at most δ

1. How many training examples suffice?
$$m \geq \frac{1}{\epsilon}(\ln |H| + \ln(1/\delta))$$

2. If $error_{train}(h) = 0$ then with probability at least (1-δ):
$$error_{true}(h) \leq \frac{1}{m}(\ln |H| + \ln(1/\delta))$$

## Example: H is Conjunction of Boolean Literals

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

Consider classification problem f:X→Y:

- instances: $X = <X_1\ X_2\ X_3\ X_4>$ where each $X_i$ is boolean
- learned hypotheses are rules of the form:
  - IF $<X_1\ X_2\ X_3\ X_4> = <0,?,1,?>$ ,  THEN Y=1, ELSE Y=0
  - i.e., rules constrain any subset of the $X_i$

How many training examples *m* suffice to assure that with probability at least 0.99, *any* consistent learner will output a hypothesis with true error at most 0.05?

$$m \geq \frac{1}{.05}\left(\ln |H| + \ln\left(\frac{1}{.01}\right)\right)$$

$$N=4 \rightarrow m \geq 180 \leftarrow 3^4 \quad 3^n \rightarrow \ln|H| = n \ln 3$$
$$\qquad\qquad\qquad n=4$$
$$N=10 \qquad \geq 312$$
$$n=100 \qquad \geq 2290$$

## Example: H is Decision Tree with depth=2

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

Consider classification problem f:X→Y:

- instances: $X = <X_1 \dots X_N>$ where each $X_i$ is boolean
- learned hypotheses are decision trees of depth 2, using only two variables

$$\binom{N}{2} \cdot 16 = \frac{N \cdot (N-1)}{2} \cdot 168$$

How many training examples *m* suffice to assure that with probability at least 0.99, *any* consistent learner will output a hypothesis with true error at most 0.05?

$$|H| = 8N^2 - 8N$$

$$m \geq \frac{1}{0.05}\left(\ln\left(8N^2 - 8N\right) + \ln\left(\frac{1}{0.01}\right)\right)$$

$$N=4 \qquad m \geq 184$$
$$N=10 \qquad m \geq 224$$
$$N=100 \qquad m \geq 318$$

## 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 $X$, $\epsilon$ such that $0 < \epsilon < 1/2$, and $\delta$ 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 $size(c)$.

## 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 $X$, $\epsilon$ such that $0 < \epsilon < 1/2$, and $\delta$ 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 $size(c)$.

Sufficient condition:

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

## 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?

note ε here is the difference between the training error and true error

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

derived from Hoeffding bounds:

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

true error    training error    degree of overfitting

---

## Additive Hoeffding Bounds – Agnostic Learning

- Given *m* independent coin flips of coin with true Pr(heads) = θ
  bound the error in the maximum likelihood estimate $\widehat{\theta}$

$$\Pr[\theta > \widehat{\theta} + \epsilon] \leq e^{-2m\epsilon^2}$$

- Relevance to agnostic learning: for any *single* hypothesis *h*

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

- But we must consider all hypotheses in H

$$\Pr[(\exists h \in H) error_{true}(h) > error_{train}(h) + \epsilon] \leq |H|e^{-2m\epsilon^2}$$

- So, with probability at least (1-δ) every h satisfies

$$error_{true}(h) \leq error_{train}(h) + \sqrt{\frac{\ln|H| + \ln\frac{1}{\delta}}{2m}}$$

## General Hoeffding Bounds

- When estimating parameter $\theta$ inside [a,b] from *m* examples

$$P(|\widehat{\theta} - E[\widehat{\theta}]| > \epsilon) \leq 2e^{\frac{-2m\epsilon^2}{(b-a)^2}}$$

- When estimating a probability $\theta$ is inside [0,1], so

$$P(|\widehat{\theta} - E[\widehat{\theta}]| > \epsilon) \leq 2e^{-2m\epsilon^2}$$

- And if we're interested in only one-sided error, then

$$P((E[\widehat{\theta}] - \widehat{\theta}) > \epsilon) \leq e^{-2m\epsilon^2}$$

---

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

Question: If H = {h | h: X $\rightarrow$ Y} is infinite, what measure of complexity should we use in place of |H| ?

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

Question: If H = {h | h: X → Y} is infinite,
what measure of complexity should we
use in place of |H| ?

Answer: The largest subset of X for which H can guarantee
zero training error (regardless of the target function c)