# Machine Learning 10-701

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

March 1, 2011

Today:

- Computational Learning Theory
- VC dimension
- PAC results as quantitative model of overfitting

Recommended reading:

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

---

# 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):
$$error_{true}(h) \leq \frac{1}{m}(\ln|H| + \ln(1/\delta))$$

## 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 <u>time that is polynomial</u> 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 <u>time that is polynomial</u> 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

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

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

---

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

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 <u>guarantee</u> zero training error (regardless of the target function c)

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

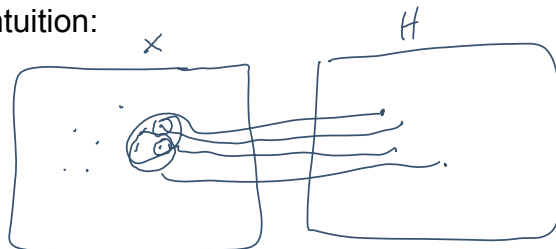Question: If H = {h | h: X → Y} is infinite, what measure of complexity should we use in place of |H| ?

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

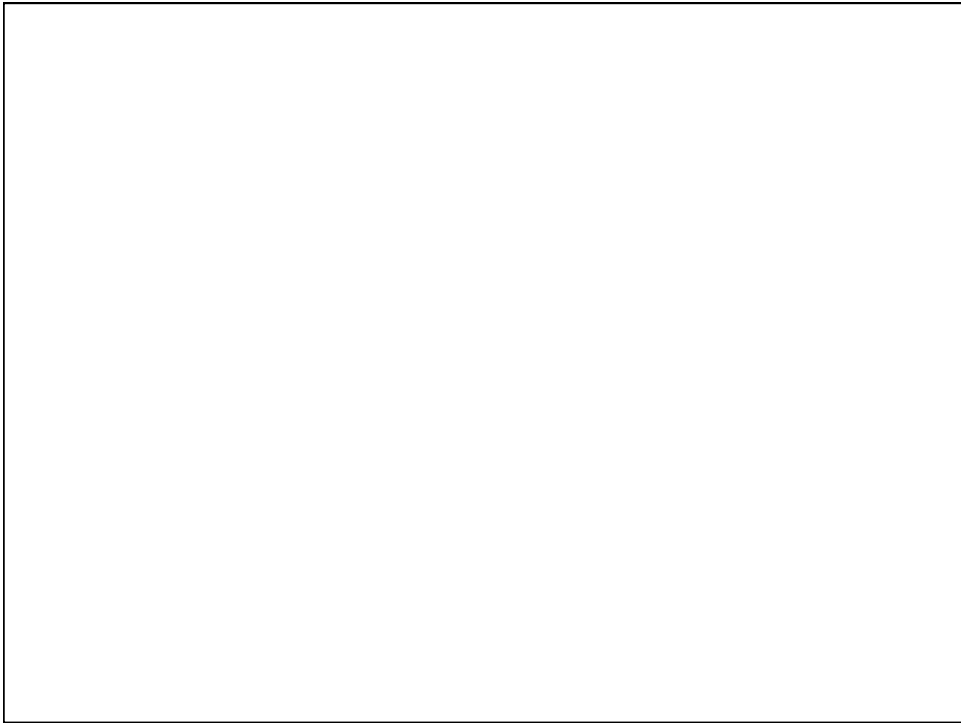**VC dimension of H is the size of this subset**

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)
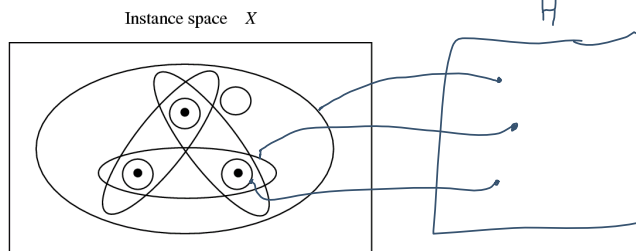
Informal intuition:

## Shattering a Set of Instances

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

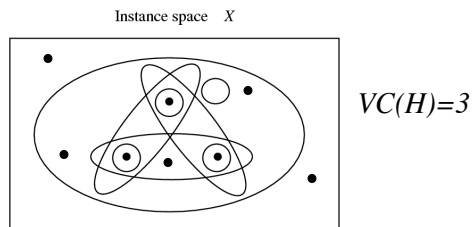a labeling of each member of S as positive or negative

*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.

Instance space   $X$

H

## 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 $VC(H) \equiv \infty$.

Instance space $X$



$VC(H)=3$

---

## Sample Complexity based on VC dimension

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

ie., to guarantee that any hypothesis that perfectly fits the training data is probably $(1-\delta)$ approximately $(\varepsilon)$ correct

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

Compare to our earlier results based on |H|:

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

6

# VC dimension: examples

Consider X = $\mathbb{R}$, want to learn c:X→{0,1}

What is VC dimension of

- Open intervals:

  (H1:) if $x > a$ then $y = 1$ else $y = 0$   $VC = 1$

  (H2:) if $x > a$ then $y = 1$ else $y = 0$   $VC = 2$
  or, if $x > a$ then $y = 0$ else $y = 1$

  | 1 | 1 | 0 |
  |---|---|---|
  | 0 | 1 | 0 |
  | 1 | 0 | 1 |

- Closed intervals:

  (H3:) if $a < x < b$ then $y = 1$ else $y = 0$   $VC = 2$

  (H4:) if $a < x < b$ then $y = 1$ else $y = 0$
  or, if $a < x < b$ then $y = 0$ else $y = 1$

---

# VC dimension: examples

Consider X = <, want to learn c:X→{0,1}

What is VC dimension of

- Open intervals:

  H1: if $x > a$ then $y = 1$ else $y = 0$     VC(H1)=1

  H2: if $x > a$ then $y = 1$ else $y = 0$     VC(H2)=2
  or, if $x > a$ then $y = 0$ else $y = 1$

- Closed intervals:

  H3: if $a < x < b$ then $y = 1$ else $y = 0$     VC(H3)=2

  H4: if $a < x < b$ then $y = 1$ else $y = 0$     VC(H4)=3
  or, if $a < x < b$ then $y = 0$ else $y = 1$

# VC dimension: examples

$$X = \mathbb{R}^2$$

What is VC dimension of lines in a plane?

- $H_2 = \{ ((w_0 + w_1x_1 + w_2x_2) > 0 \rightarrow y=1) \}$   $VC \gtrsim 3$
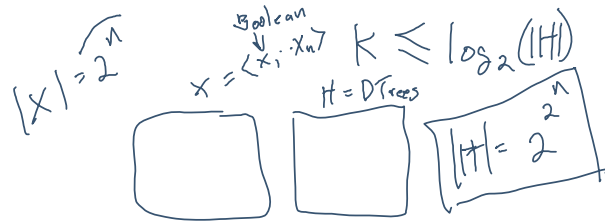


---

# VC dimension: examples

What is VC dimension of

- $H_2 = \{ ((w_0 + w_1x_1 + w_2x_2) > 0 \rightarrow y=1) \}$
  - $VC(H_2)=3$
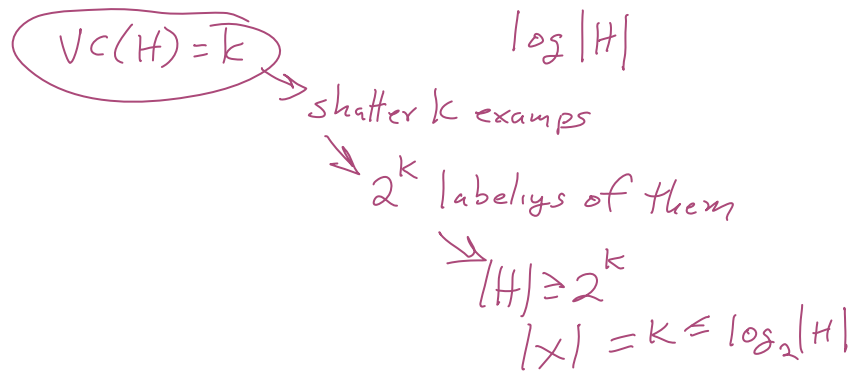- For $H_n$ = linear separating hyperplanes in n dimensions, $VC(H_n)=n+1$

**For any finite hypothesis space H, can you give an upper bound on VC(H) in terms of |H| ?**
**(hint: yes)**

$VC(H) = K$

$\Rightarrow$ $H$ can express $\geq 2^k$ fns.

$|X| = 2^n$

$X = \langle x_1 : x_n \rangle$ (Boolean)

$K \leq \log_2(|H|)$

$H = DTrees$

$|H| = 2^{2^n}$

---

**Can you give an upper bound on VC(H) in terms of |H|, for any hypothesis space H?**
**(hint: yes)**

$\boxed{VC(H) = k}$

$\log |H|$

$\rightarrow$ shatter $k$ examps

$\rightarrow$ $2^k$ labelings of them

$\rightarrow$ $|H| \geq 2^k$

$|X| = k \leq \log_2 |H|$

## More VC Dimension Examples to Think About

- Logistic regression over n continuous features
  - Over n boolean features?

- Linear SVM over n continuous features

- Decision trees defined over n boolean features
  
  F: $<X_1, \dots X_n> \rightarrow Y$

- Decision trees of depth 2 defined over n features

- How about 1-nearest neighbor?

## Tightness of Bounds on Sample Complexity

How many examples *m* suffice to assure that any hypothesis that fits the training data perfectly is probably (1-δ) approximately (ε) correct?

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

How tight is this bound?

## Tightness of Bounds on Sample Complexity

How many examples $m$ suffice to assure that any hypothesis that fits the training data perfectly is probably (1-δ) approximately (ε) correct?

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

How tight is this bound?

**Lower bound on sample complexity** (Ehrenfeucht et al., 1989):

Consider any class C of concepts such that VC(C) > 1, any learner L, any 0 < ε < 1/8, and any 0 < δ < 0.01. Then there exists a distribution $\mathcal{D}$ and a target concept in C, such that if L observes fewer examples than

$$\max\left[\frac{1}{\epsilon}\log(1/\delta), \frac{VC(C)-1}{32\epsilon}\right]$$

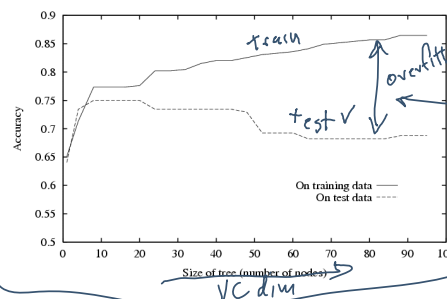Then with probability at least δ, L outputs a hypothesis with $error_{\mathcal{D}}(h) > \epsilon$

---

## Agnostic Learning: VC Bounds ✓

[Schölkopf and Smola, 2002]
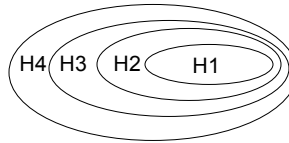
With probability at least (1-δ) every $h \in H$ satisfies

$$error_{true}(h) \leq error_{train}(h) + \sqrt{\frac{VC(H)(\ln\frac{2m}{VC(H)} + 1) + \ln\frac{4}{\delta}}{m}}$$



train

overfitting

test ✓

On training data ——
On test data -----

Size of tree (number of nodes)

VC dim

11

## Structural Risk Minimization [Vapnik]

Which hypothesis space should we choose?

• Bias / variance tradeoff



SRM: choose H to minimize bound on expected true error!

$$error_{true}(h) < error_{train}(h) + \sqrt{\frac{VC(H)(\ln\frac{2m}{VC(H)} + 1) + \ln\frac{4}{\delta}}{m}}$$

\* unfortunately a somewhat loose bound...

---

## PAC Learning: What You Should Know

• PAC learning: Probably (1-δ) Approximately (error ε) Correct
• Problem setting

• Finite H, perfectly consistent learner result ✓
• If target function is not in H, *agnostic learning* ✓
• If |H| = ∞ , use VC dimension to characterize H ✓

• Most important:
   – Sample complexity grows with complexity of H
   – Quantitative characterization of overfitting

• Much more: see Prof. Blum's course on Computational Learning Theory