



Computational Learning Theory

HTF: ?? (7.9)
B: Ch 7.1.5
RN, Chapter 18.5

R Greiner
Cmput 466 / 551

Thanks to A Blum



Computational Learning Theory

- Inductive Learning
 - Protocol
 - Error
- Probably Approximately Correct Learning
 - Consistency Filtering
 - Sample Complexity
 - Eg: Conjunction, Decision List
- Issues
 - Bound
 - Other Models



What General Laws constrain Inductive Learning?

- Sample Complexity
 - How many training examples are sufficient to learn target concept?
- Computational Complexity
 - Resources required to learn target concept?
- Want theory to relate:
 - Training examples
 - Quantity
 - Quality
 - How presented
 - Complexity of hypothesis/concept space
 - Accuracy of approx to target concept
 - Probability of successful learning

These results only useful wrt $O(\dots)$!

Protocol

- Given:

- space of examples X
- fixed (unknown) distribution D over X
- set of hypotheses H
- set of possible target concepts C

- Learner observes sample $S = \{ \langle x_i, c(x_i) \rangle \}$

- instances x_i drawn from distr. D
- Labeled $c(x)$ by target concept $c \in C$
(Learner does NOT know $c(\cdot), D$)

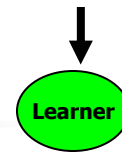
- Learner outputs $h \in H$ estimating c

- h is evaluated by performance on subsequent instances drawn from D

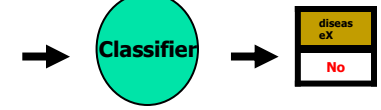
- For now:

- $C = H$ (so $c \in H$)
- Noise-free data

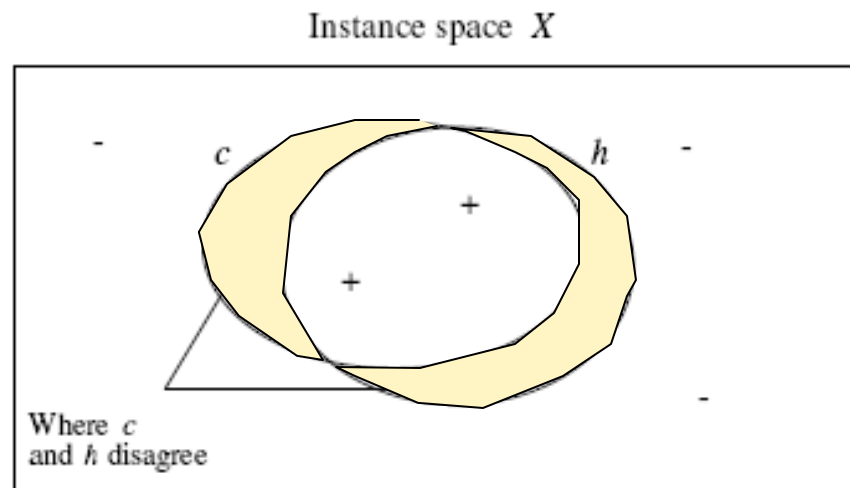
| tu m | Pr es | So | ... | Col our | dise aseX |
|---------|----------|----|-----|------------|--------------|
| 35 | 95 | Y | ... | Pale | No |
| 22 | 11 | N | ... | Cle | Yes |
| : | : | : | : | : | : |
| 10 | 87 | N | ... | Pale | No |



| - | - | - | - | - |
|---|---|---|-----|---|
| 3 | 9 | N | ... | a |
| 2 | 0 | | | l |



True Error of Hypothesis



Def'n: The true error of hypothesis h wrt

- target concept c
- distribution D

\equiv probability that h will misclassify instance drawn from D

$$\text{err}_D(h) = \Pr_{x \in D} [c(x) \neq h(x)]$$



Probably Approximately Correct

Goal:

PAC-Learner produces hypothesis \hat{h} that
is approximately correct,

$$\text{err}_D(\hat{h}) \approx 0$$

with high probability

$$P(\text{err}_D(\hat{h}) \approx 0) \approx 1$$

- Double "hedging"

- approximately

- probably

Need both!



PAC-Learning

- Learner L can draw labeled instance $\langle x, c(x) \rangle$ in unit time
 $x \in X$ drawn from distribution D labeled by target concept $c \in C$

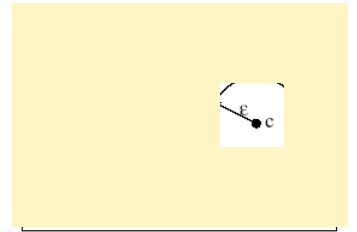
Def'n: Learner L PAC-learns class C (by H)

if

1. for any target concept $c \in C$,
any distribution D , any $\epsilon, \delta > 0$,
 L returns $h \in H$ s.t.
w/ prob. $\geq 1 - \delta$, $\text{err}_D(h) < \epsilon$
2. L 's run-time (and hence, sample complexity)
is $\text{poly}(\text{size}(x), \text{size}(c), 1/\epsilon, 1/\delta)$

- Sufficient:
 1. Only $\text{poly}(\dots)$ training instances – $|H| = 2^{\text{poly}(\dots)}$
 2. Only poly time / instance ...Often $C = H$

Simple Learning Algorithm: Consistency Filtering



- Draw $m_H(\epsilon, \delta)$ random (labeled) examples S_m
- Remove every hyp. that contradicts any $\langle x, y \rangle \in S_m$
- Return any remaining (consistent) hypothesis

Challenges:

- Q1: Sample size: $m_H(\epsilon, \delta)$
- Q2: Need to decide if $h \in H$ is consistent w/ all S_m
... efficiently ...

Boolean Functions (\equiv Concepts)

Eg: $h_{X_1 \vee \neg X_2}(X_1, X_2, X_3) = \begin{cases} 1 & \text{if } X_1 \vee \neg X_2 \\ 0 & \text{otherwise} \end{cases}$

| X_1 | X_2 | X_3 | $h_{X_1 \vee \neg X_2}(X_1, X_2, X_3)$ |
|-------|-------|-------|--|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |

$h_{X_1 \vee \neg X_2}(0, 1, 1) = 0$

$h_{X_1 \vee \neg X_2}(1, 1, 0) = 1$

Note: Hypothesis maps unlabeled-tuple to $\{0, 1\}$

Labeled-tuple is $\left\{ \begin{array}{l} \text{Consistent} \\ \text{Inconsistent} \end{array} \right\}$ w/ hyp.

So $\langle (0, 1, 1), 1 \rangle$ is

Inconsistent with $h_{X_1 \vee \neg X_2}$

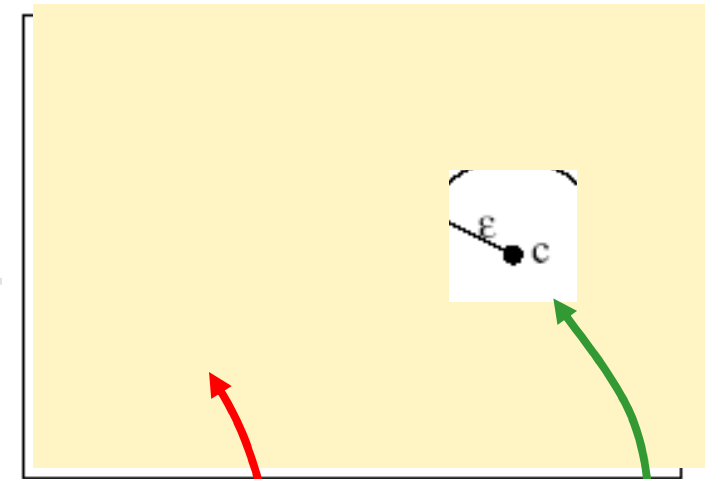
Consistent with $h_{X_2 \vee X_3}$

Bad Hypotheses

Idea: Find $m = m_H(\epsilon, \delta)$ s.t.
after seeing m examples,
every BAD hypothesis h ($\text{err}_{D,c}(h) > \epsilon$)
will be ELIMINATED
with high probability ($\approx 1 - \delta$)
leaving only good hypotheses

... then pick ANY of the remaining good ($\text{err}_{D,c}(h) < \epsilon$) hyp's

Find m large number that
very small chance that a "bad" hypothesis is consistent with m examples

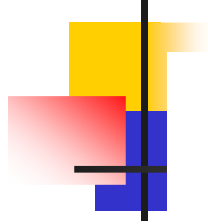


Eliminate ALL

$$H_{\text{bad}} = \{ h \in H \mid \text{err}_D(h) > \epsilon \}$$

Leave SOME

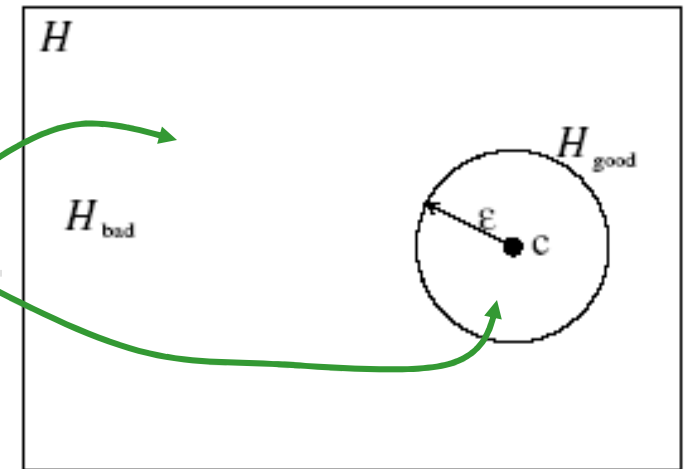
$$H_{\text{Good}} = \{ h \in H \mid \text{err}_D(h) \leq \epsilon \}$$



$$\mathcal{H}_{bad} = \{h \in \mathcal{H} \mid \text{err}_{\mathcal{D}}(h) > \epsilon\}$$

$$\mathcal{H}_{good} = \{h \in \mathcal{H} \mid \text{err}_{\mathcal{D}}(h) \leq \epsilon\}$$

Note $|\mathcal{H}_{Bad}| \leq |\mathcal{H}|$



| | x_1 | x_2 | \dots | x_m | Good/Bad | | |
|---------------|-------------|-------|---------|----------|----------|--------------|-----|
| \mathcal{H} | h_1 | ✓ | ✓ | \dots | ✓ | $< \epsilon$ | Ok |
| | h_2 | ✓ | - | \dots | ✓ | $< \epsilon$ | |
| | \vdots | | | \vdots | | \vdots | |
| | h_k | ✓ | ✓ | \dots | - | $< \epsilon$ | Bad |
| | $h_{bad,1}$ | - | - | \dots | ✓ | $> \epsilon$ | |
| | \vdots | | | \vdots | | \vdots | |
| | $h_{bad,r}$ | - | - | \dots | - | $> \epsilon$ | |

✓ : $h_i(x_j) = c(x_j)$

- : $h_i(x_j) \neq c(x_j)$

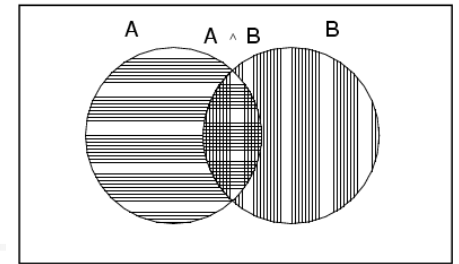


Sample Bounds – Derivation

- Let h_1 be ε -bad hypothesis ... $\text{err}(h_1) > \varepsilon$
 - $\Rightarrow h_1$ mis-labels example w/prob $P(h_1(x) \neq c(x)) > \varepsilon$
 - $\Rightarrow h_1$ **correctly** labels random example w/prob $\leq (1 - \varepsilon)$
- As examples drawn **INDEPENDENTLY**
 - $P(h_1 \text{ correctly labels } m \text{ examples}) \leq (1 - \varepsilon)^m$

Sample Bounds – Derivation II

True



- Let h_2 be another ε -bad hypothesis
- What is probability that *either* h_1 or h_2 survive m random examples?

$$\begin{aligned} P(h_1 \vee h_2 \text{ survives}) &= P(h_1 \text{ survives}) + P(h_2 \text{ survives}) \\ &\quad - P(h_1 \& h_2 \text{ survives}) \\ &\leq P(h_1 \text{ survives}) + P(h_2 \text{ survives}) \\ &\leq 2(1 - \varepsilon)^m \end{aligned}$$

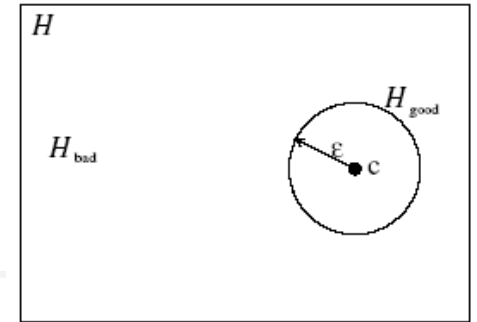
- If k ε -bad hypotheses $\{h_1, \dots, h_k\}$:
 $P(h_1 \vee \dots \vee h_k \text{ survives}) \leq k(1 - \varepsilon)^m$



Sample Bounds – Derivation

- Let h_1 be ε -bad hypothesis ... $\text{err}(h_1) > \varepsilon$
 - $\Rightarrow h_1$ mis-labels example w/prob $P(h_1(x) \neq c(x)) > \varepsilon$
 - $\Rightarrow h_1$ correctly labels random example w/prob $\leq (1 - \varepsilon)$
- As examples drawn INDEPENDENTLY
 - $P(h_1 \text{ correctly labels } m \text{ examples}) \leq (1 - \varepsilon)^m$
- Let h_2 be another ε -bad hypothesis
- What is probability that either h_1 or h_2 survive m random examples?
 - $P(h_1 \vee h_2 \text{ survives})$
 - $= P(h_1 \text{ survives}) + P(h_2 \text{ survives}) - P(h_1 \& h_2 \text{ survives})$
 - $\leq P(h_1 \text{ survives}) + P(h_2 \text{ survives})$
 - $\leq 2(1 - \varepsilon)^m$

Sample Bounds, con't



Let $H_{\text{bad}} = \{ h \in H \mid \text{err}(h) > \epsilon \}$

- Probability that any $h \in H_{\text{bad}}$ survives is

$P(\text{any } h_b \text{ in } H_{\text{bad}} \text{ is consistent with } m \text{ exs.})$

$$\leq |H_{\text{bad}}| (1 - \epsilon)^m \leq |H| (1 - \epsilon)^m$$

- This is $\leq \delta$ if $|H| (1 - \epsilon)^m \leq \delta$

$$\Rightarrow m_H(\epsilon, \delta) \geq \left(\log \frac{|H|}{\delta} \right) / -\log(1 - \epsilon) \geq \frac{1}{\epsilon} \left(\log \frac{|H|}{\delta} \right)$$

- $m_H(\epsilon, \delta)$ is "Sample Complexity" of hypothesis space H
- Fact: For $0 \leq \epsilon \leq 1$, $(1 - \epsilon) \leq e^{-\epsilon}$

Sample Complexity

- Hypothesis Space (expressiveness): H
- Error Rate of Resulting Hypthesis: ϵ
 - $\text{err}_{D,c}(h) = P(h(x) \neq c(x)) \leq \epsilon$
- Confidence of being ϵ -close: δ
 - $P(\text{err}_{D,c}(h) \leq \epsilon) > 1 - \delta$
- Sample size: $m_H(\epsilon, \delta)$

- Any hypothesis consistent with

$$m_H(\epsilon, \delta) = \frac{1}{\epsilon} \left(\log \frac{|H|}{\delta} \right)$$

examples,

has error of at most ϵ , with prob $\leq 1 - \delta$



Boolean Function... Conjunctions

- Boolean Instance: $\langle x_1, \dots, x_n \rangle$
 $\langle 1, 0, 1, 1 \rangle$ for $\langle x_1 = 1, x_2 = 0, x_3 = 1, x_4 = 1 \rangle$
- Boolean Function: $f(\langle x_1, \dots, x_n \rangle) \in \{0, 1\}$
- Conjunction (type of Boolean function)

$$f_{+-0-0+}(X) = x_1 \bar{x}_2 \bar{x}_4 x_6$$
$$= \begin{cases} 1 & \text{if } x_1(X) = t, x_2(X) = f, x_4(X) = f, \\ & \text{and } x_6(X) = t \\ 0 & \text{otherwise} \end{cases}$$

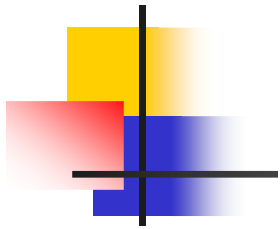
$$f_{+-0-0+}(\langle \underline{1}, \underline{0}, 1, \underline{0}, 0, \underline{1} \rangle) = 1$$

$$f_{+-0-0+}(\langle \underline{0}, \underline{0}, 1, \underline{0}, 0, \underline{1} \rangle) = 0$$

(Ie, must match each literal mentioned)

- Only 3^n possible conjunctions
out of 2^{2^n} boolean functions!

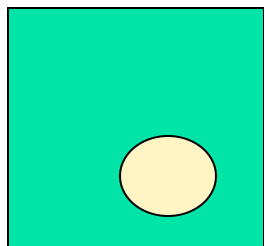
- \mathcal{H}_C = conjunctions of literals



- $|\mathcal{H}_C| = 3^n$:
 - Each variable can be
 - included positively " x_i ",
 - included negatively " \bar{x}_i ",
 - excluded
- $$\Rightarrow m_{\mathcal{H}_C}(\epsilon, \delta) = \frac{1}{\epsilon} \left[n \ln 3 + \ln \frac{1}{\delta} \right]$$

Alg: Collect $m_{\mathcal{H}_C}(\epsilon, \delta) = \frac{1}{\epsilon} \left[n \ln 3 + \ln \frac{1}{\delta} \right]$ labeled samples
 Let $h = x_1 \bar{x}_1 x_2 \bar{x}_2 \cdots x_n \bar{x}_n$
 For each +-example $y = \bigwedge_i \pm_i x_i$
 Remove from h any literal NOT included in y

| Data | Current Hyp | | | | | | |
|------|-------------|-------------|-------|-------------|-------|-------------|---------------------|
| | x_1 | \bar{x}_1 | x_2 | \bar{x}_2 | x_3 | \bar{x}_3 | |
| | | | | | | | Never true |
| | | | | | | | True only for "101" |
| | | | | | | | True only for "10*" |



- Just uses +-examples!
 - Finds "smallest" hypothesis (true for as few +examples as possible)
 - ... No mistakes on -examples
- As each step is efficient $O(n)$, only $\text{poly}(n, 1/\epsilon, 1/\delta)$ steps
 \Rightarrow algorithm is *efficient!*

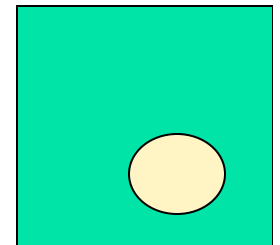
- \mathcal{H}_C = conjunctions of literals

- $|\mathcal{H}_C| = 3^n$:
 - Each variable can be
 - included positively " x_i ",
 - included negatively " \bar{x}_i ",
 - excluded
- $$\Rightarrow m_{\mathcal{H}_C}(\epsilon, \delta) = \frac{1}{\epsilon} \left[n \ln 3 + \ln \frac{1}{\delta} \right]$$

Alg: Collect $m_{\mathcal{H}_C}(\epsilon, \delta)$ samples
 Let $h =$
 For $i = 1, \dots, n$
 If literal $\pm x_i$ NOT included in y

Does NOT explicitly build all 3^n conjunctions, then throw some out...

| | Current Hyp | | | | | | |
|---|-------------|-------------|-------|-------------|-------|-------------|---------------------|
| | x_1 | \bar{x}_1 | x_2 | \bar{x}_2 | x_3 | \bar{x}_3 | |
| $\langle \langle 1 \ 0 \ 1 \rangle + \rangle$ | x_1 | | | \bar{x}_2 | x_3 | | Never true |
| $\langle \langle 0 \ 1 \ 1 \rangle - \rangle$ | x_1 | | | \bar{x}_2 | x_3 | | True only for "101" |
| $\langle \langle 1 \ 0 \ 0 \rangle + \rangle$ | x_1 | | | \bar{x}_2 | | | True only for "10*" |
| $\langle \langle 0 \ 0 \ 0 \rangle - \rangle$ | x_1 | | | \bar{x}_2 | | | |



- Just uses +-examples!
 - Finds "smallest" hypothesis (true for as few +examples as possible)
 - ... No mistakes on -examples
- As each step is efficient $O(n)$, only $\text{poly}(n, 1/\epsilon, 1/\delta)$ steps \Rightarrow algorithm is *efficient!*



PAC-Learning k-CNF

- $CNF \equiv$ Conjunctive Normal Form
 $(x_1 \vee \bar{x}_2 \vee x_7) \wedge (x_2 \vee x_4 \vee \bar{x}_9) \wedge \dots \wedge (x_7 \vee \bar{x}_8 \vee \bar{x}_9)$

- $k-CNF \equiv$ CNF where each clause has $\leq k$ literals
1-CNF \equiv Conjunctions

- As $\exists O\left(\binom{n}{k} 3^k\right)$ possible $\leq k$ -clauses,

$$|\mathcal{H}_{k-CNF}| = 2^{O\left(\binom{n}{k} 3^k\right)}$$

$$\binom{n}{k} = O(n^k)$$

$$\Rightarrow m_{\mathcal{H}_{k-CNF}} = O\left(\frac{1}{\epsilon} \left[(3n)^k + \ln \frac{1}{\delta}\right]\right)$$

Alg: Consistency Filtering:

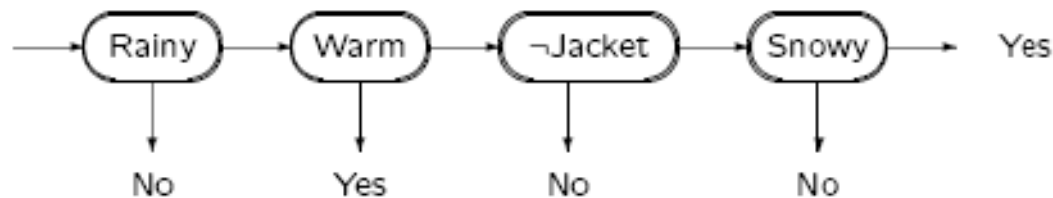
Let $T =$ all $O\left(\binom{n}{k} 3^k\right)$ possible k -clauses.
After each \pm -example y ,
Remove from T all clauses INCONSISTENT w/ y
Return $\bigwedge T$

- Similar for Disjunctions, k -DNF, ...

? What about CNF $\equiv n$ -CNF ?

Decision Lists

- When to go for walk?
 - Vars: *rainy*, *warm*, *jacket*, *snowy*
 - Don't go for walk if *rainy*.
Otherwise, go for walk if *warm* or
if \neg *jacket* and it is *snowy*.



Def'n: A *DL* \equiv list of "if-then rules"
where $\left\{ \begin{array}{l} \text{condition} \equiv \text{a literal} \\ \text{consequent is } + \text{ or } - \end{array} \right\}$
(\equiv decision tree with just one long path)

- How many DLs?
4n possible "rules", each of form " $\pm x_i \Rightarrow \pm$ "
 $\Rightarrow (4n)!$ orderings, so $|H_{DL}| \cdot (4n)!$
(Actually: $\leq n! 4^n$)

Example of Learning DL

Data:

| | x_1 | x_2 | x_3 | x_4 | x_5 | Label |
|-------|-------|-------|-------|-------|-------|-------|
| i_1 | 1 | 0 | 0 | 1 | 1 | A |
| i_2 | 0 | 1 | 1 | 0 | 0 | B |
| i_3 | 1 | 1 | 1 | 0 | 0 | A |
| i_4 | 0 | 0 | 0 | 1 | 0 | B |
| i_5 | 1 | 1 | 0 | 1 | 1 | A |
| i_6 | 1 | 0 | 0 | 0 | 1 | B |

1. When $x_1 = 0$, class is "B"

Form $h = \langle \neg x_1 \mapsto \mathbf{B} \rangle$

Eliminate i_2, i_4

2. When $x_2 = 1$, class is "A"

Form $h = \langle \neg x_1 \mapsto \mathbf{B}; x_2 \mapsto \mathbf{A} \rangle$

Eliminate i_3, i_5

3. When $x_4 = 1$, class is "A"

Form $h = \langle \neg x_1 \mapsto \mathbf{B}; x_2 \mapsto \mathbf{A}; x_4 \mapsto \mathbf{A} \rangle$

Eliminate i_1

4. Always have class "B"

Form $h = \langle \neg x_1 \mapsto \mathbf{B}; x_2 \mapsto \mathbf{A}; x_4 \mapsto \mathbf{A}; t \mapsto \mathbf{B} \rangle$

Eliminate rest (i_6)

PAC-Learning Decision Lists

Let: S = set of

$$m_{DL} = O\left(\frac{1}{\epsilon}[n \ln(n) + \ln \frac{1}{\delta}]\right)$$

training instances

h = empty list

R = all $4n$ possible rules

While $S \neq \{\}$ do

1. Find $r \in R$ s.t.
+ consistent w/ S
+ r applies to ≥ 1 $s \in S$

(If none, halt w/ "Failure")

2. $h := h \circ r$
(Put rule at BOTTOM of hypothesis)

3. $S := S - \{s \mid s \text{ classified by } h\}$
(Throw out examples classified by current hypothesis)

Covering Algorithm



Proof (PAC-Learn DL)

- **Correctness#1:** Enough data?

Yes. $\frac{1}{\epsilon} \ln \frac{|\mathcal{H}_{DL}|}{\delta}$

- **Correctness#2:** Consistency?

If \exists DL consistent w/data...

- $\exists \geq 1$ choice for step 1 (eg, first rule in L satisfied by ≥ 1 example)
 - \exists DL consistent w/ remaining data: original DL!
- **Efficiency:** Algorithm runs in poly time, since
 - each iteration requires poly time, and
 - each iteration removes ≥ 1 example (only poly examples)
 - **Generalization:** k-DL
 - ... whose nodes each contain CONJUNCTION of k literals
(So earlier DL \equiv 1-DL)

k-DL \supset k-CNF, k-DNF, k-depth DecTree, ...



Why Learning May Succeed

- Learner L produces classifier $h = L(S)$ that does well on training data S

Why?

1. If x appears a lot

- then x probably occurs in training data S

- As h does well on S , δ

$h(x)$ is probably correct on x

2. If example x appears rarely ϵ

($P(x) \approx 0$)

then h suffers only small penalty for being wrong.

- Assumption: Distribution is "stationary"
 - distr. for testing = distr. for training



Comments on Model

$$m_H(\varepsilon, \delta) = \frac{1}{\varepsilon} \left(\log \frac{|H|}{\delta} \right)$$

Simplify task:

- ~~1*. Assume $c \in H$, where H known~~
 - (Eg, lines, conjunctions, . . .)
- ~~2*. Noise free training data~~
- 3. Only require approximate correctness:
 - h is " ε -good": $P_x(h(x) \neq c(x)) < \varepsilon$
- 4. Allow learner to (rarely) be completely off
 - If examples NOT representative, cannot do well.
 - $P(h_L \text{ is } \varepsilon\text{-good}) \geq 1 - \delta$

Complicate task:

- 1. Learner must be computationally efficient
- 2. Over any instance distribution



Comments: Sample Complexity

$$m_H(\varepsilon, \delta) = \frac{1}{\varepsilon} \left(\log \frac{|H|}{\delta} \right)$$

- If k parameters, $\langle v_1, \dots, v_k \rangle$
 - $\Rightarrow |H_k| \approx B^k$
 - $\Rightarrow m_{H_k} \approx \log(B^k)/\varepsilon \approx k/\varepsilon$
- Too GENEROUS:
 - Based on pre-defined $C = \{c_1, \dots\} = H$
Where did this come from???
 - Assumes $c \in H$, noise-free
 - If $\text{err} \neq 0$, need $O(1/\varepsilon^2 \dots)$

Why is Bound so Lousy!

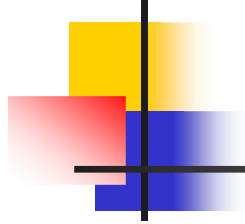
- Assumes error of all ε -bad hypotheses $\approx \varepsilon$
(Typically most bad hypotheses are really bad
 \Rightarrow get thrown out much sooner)
- Uses $P(A \text{ or } B) \leq P(A) + P(B)$.
(If hypotheses are correlated, then if one inconsistent, others probably inconsistent too)
- Assumes $|H_{\text{bad}}| = |H|$... see VCdimension
- WorstCase:
 - over all $c \in C$
 - over all distribution D over X
 - over all presentations of instances (drawn from D)
- Improvements
 - "Distribution Specific" learning
Known single dist (ε -cover)
Gaussian, . . .
 - Look at samples! \Rightarrow Sequential PAC Learning

If λ -bad, takes $\approx 1/\lambda$ to see evidence

Fundamental Tradeoff in Machine Learning

$$m_H(\epsilon, \delta) = \frac{1}{\epsilon} \left(\log \frac{|H|}{\delta} \right)$$

- Larger H is more likely to include
 - (approx to) target f
 - but it requires more examples to learn
- w/few examples,
cannot reliably find good hypothesis from large hyp. space
- To learn effectively (ϵ) from small # of samples (m),
only consider H where $|H| \approx e^{\epsilon m}$
- Restrict form of Boolean function
to reduce size of hypotheses space.
 - Eg, for $H_C =$ conjunctions of literals,
 $|H_C| = 3^n$, so only need poly number of examples!
 - Great if target concept is in H_C , but ...



Issues

- Computational Complexity
- Sampling Issues:

| | Finite | Countable | Uncountable |
|------------|---|--------------|-------------|
| Realizable | $\frac{1}{\epsilon} \ln \frac{ H }{\delta}$ | Nested Class | VC dim |
| Agnostic | $O\left(\frac{1}{\epsilon^2} \ln \frac{ H }{\delta}\right)$ | – | VC dim |



Learning = Estimation + Optimization

1. Acquire required relevant information by examining enough labeled samples
2. Find hypothesis $h \in H$ consistent with those samples
 - ... often “smallest” hypothesis
 - Spse H has 2^k hypotheses
Each hypothesis requires k bits
 - $\Rightarrow \log |H| \approx |h| = k$
 - \Rightarrow SAMPLE COMPLEXITY not problematic
 - But optimization often is... intractable!
 - Eg, consistency for 2term-DNF is NP-hard, ...
 - Perhaps find best hypothesis in $F \supset H$
 - 2-CNF \supset 2term-DNF
 - ... easier optimization problem!



Extensions to this Model

- Ockham Algorithm: Can PAC-learn H iff
 - can “compress” samples
 - have efficient consistency-finding algorithm

- Data Efficient Learner

Gathers samples sequentially, autonomously decides when to stop & return hypothesis

- Exploiting other information

- Prior background theory
- Relevance

- Degradation of Training/Testing Information

$$\left\{ \begin{array}{c} \text{Error} \\ \text{Omissions} \end{array} \right\} \text{ in } \left\{ \begin{array}{c} \text{Training} \\ \text{Testing} \end{array} \right\} \left\{ \begin{array}{c} \text{Attribute Value} \\ \text{ClassLabel} \end{array} \right\}$$



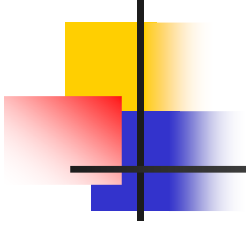
Other Learning Models

- Learning in the Limit [Recursion Theoretic]
 - Exact identification, no resource constraints
- On-Line learning
 - After seeing each unlabeled instance, learner returns (proposed) label
 - Learner then given correct label provided (penalized if wrong)
 - Q: Can learner converge, after making only k mistakes?
- Active Learners
 - Actively request useful information from environment
 - “Experiment”
- “Agnostic Learning”
 - What if target $\neg[f \in H]$?
 - Want to find CLOSEST hypotheses. . .
 - Typically NP-hard. . .
- Bayesian Approach: Model Averaging, . . .



Computational Learning Theory

- Inductive Learning is possible
 - With caveats: *error, confidence*
 - Depends on complexity of hypothesis space
- Probably Approximately Correct Learning
 - Consistency Filtering
 - Sample Complexity
 - Eg: Conjunctions, DecisionLists
- Many other meaningful models





Terminology

- **Labeled example:** Example of form $\langle x, f(x) \rangle$
- **Labeled sample:** Set of $\{ \langle x_i; f(x_i) \rangle \}$
- **Classifier:** Discrete-valued function.
Possible values $f(x) \in \{ 1, \dots, K \}$ called "classes";
"class labels"
- **Concept:** Boolean function.
 - x s.t. $f(x) = 1$ called "positive examples"
 - x s.t. $f(x) = 0$ called "negative examples"
- **Target function** (target concept): "True function" f generating the labels
- **Hypothesis:** Proposed function h believed to be similar to f .
- **Hypothesis Space:** Space of all hypotheses that can, in principle, be output by a learning algorithm

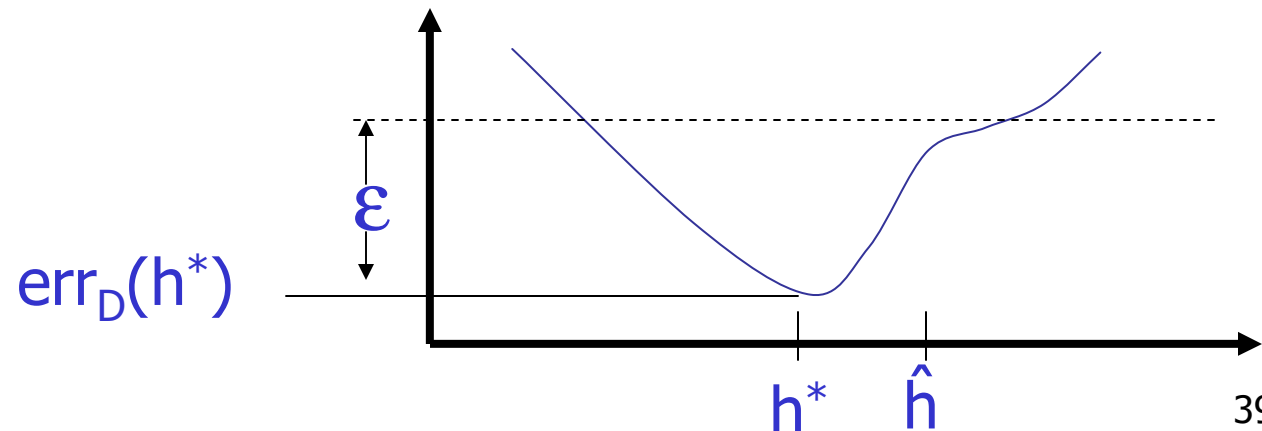


Computational Learning Theory

- Framework/Protocols
 1. Finite \mathcal{H} , Realizable case
 - 2. Finite \mathcal{H} , Unrealizable case
 3. Infinite \mathcal{H}
(Vapnik-Chervonenkis Dimension)
 4. Variable size Hypothesis Space
 5. Data-dependent Bounds
(Max Margin)
 6. Mistake Bound (Winnow)
- Topics:
 - Extensions to PAC
 - Other Learning Models
 - Occam Algorithms

Case 2: Finite \mathcal{H} , Unrealizable

- What if perfect classifier \notin hyp. space \mathcal{H} ?
 - either none exists (data inconsistent) or
 - hypothesis space is restricted
- Let: $h^* = \operatorname{argmin}_{h \in \mathcal{H}} \{ \operatorname{err}_D(h) \}$ be optimal $h \in \mathcal{H}$
- Want: \hat{h} s.t. $\operatorname{err}_D(\hat{h}) \leq \operatorname{err}_D(h^*) + \epsilon$





Case 2: Finite \mathcal{H} , Unrealizable

- What if perfect classifier \notin hyp. space \mathcal{H} ?
 - either none exists (data inconsistent) or
 - hypothesis space is restricted
- Let: $h^* = \operatorname{argmin}_{h \in \mathcal{H}} \{ \operatorname{err}_D(h) \}$ be optimal $h \in \mathcal{H}$
- Want: \hat{h} s.t. $\operatorname{err}_D(\hat{h}) \leq \operatorname{err}_D(h^*) + \varepsilon$

- Alg:

Draw $m = m(\varepsilon, \delta)$ instances S

Return $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \{ \underline{\operatorname{err}}_S(h) \}$

with optimal empirical score, over S

($\underline{\operatorname{err}}_S(h) = 1/m \sum_{x \in S} \operatorname{err}(h, x)$ is EMPIRICAL score)

- Issues:
 - How many instances?
 - Computational cost of $\operatorname{argmin}_{h \in \mathcal{H}} \{ \underline{\operatorname{err}}_S(h) \}$



Sample Complexity

Goal: Want enough instances that, w/prob $\geq 1 - \delta$

$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \{ \underline{\operatorname{err}}_S(h) \}$ is within ε of $h^* = \operatorname{argmin}_{h \in \mathcal{H}} \{ \operatorname{err}_D(h) \}$

- **Step1:** Sufficient to estimate ALL h 's to within $\varepsilon/2$.

$$| \operatorname{err}_D(h) - \underline{\operatorname{err}}_S(h) | \leq \varepsilon/2$$

If so, then

$$\begin{aligned} e_D(\hat{h}) - e_D(h^*) &= e_D(\hat{h}) - \underline{e}_S(\hat{h}) + \underline{e}_S(\hat{h}) - \underline{e}_S(h^*) + \underline{e}_S(h^*) - e_D(h^*) \\ &\leq \varepsilon/2 + 0 + \varepsilon/2 = \varepsilon \end{aligned}$$



Sample Complexity, con't

Goal: Want enough instances that, w/prob $\geq 1 - \delta$

$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \{ \underline{\operatorname{err}}_S(h) \}$ is within ε of $h^* = \operatorname{argmin}_{h \in \mathcal{H}} \{ \operatorname{err}_D(h) \}$

- **Step2:** Sufficient to estimate EACH h 's to within $\varepsilon/2$ with prob $\geq 1 - \delta / |\mathcal{H}|$

If so, then

$$\begin{aligned} & P(\exists h \in \mathcal{H} \mid \operatorname{err}_D(h) - \underline{\operatorname{err}}_S(h) \leq \varepsilon/2) \\ & \leq \sum_{h \in \mathcal{H}} P(\operatorname{err}_D(h) - \underline{\operatorname{err}}_S(h) \leq \varepsilon/2) \\ & \leq |\mathcal{H}| \delta / |\mathcal{H}| = \delta \end{aligned}$$

- **Step3:** How many instances s.t.
 $P(|\operatorname{err}_D(h) - \underline{\operatorname{err}}_S(h)| \leq \varepsilon/2) \leq \delta / |\mathcal{H}| ?$

Complexity of “Agnostic Learning”

- **Sample Complexity:** Good news!

- Hoeffding Inequality \Rightarrow Need only $m(\epsilon, \delta) = \frac{2}{\epsilon^2} \ln \frac{2|H|}{\delta}$

instances to estimate EACH h 's to within $\epsilon/2$
with prob $\geq 1 - \delta / |\mathcal{H}|$

$$P(\text{err}_D(h) - \underline{\text{err}}_S(h) \leq \epsilon/2)$$

$$\leq 2 \exp(-2 m (\epsilon/2)^2) \leq \delta / |\mathcal{H}|$$

- **Computational Complexity:** Bad news!

NP-hard to find

CONJUNCTION $h \in \mathcal{H}$ that is BEST FIT to DNF $c \in \mathcal{C}$

(target space = DNF; hypothesis space = Conjunctions)

- Note: Sample size typically poly;
Hardness tends to be Consistency/Optimization

Case 3: ∞ Hypothesis Spaces

\Rightarrow VC Dim

- Learning an initial subinterval.
“Factory is ok *iff* Temperature $\leq a$ ”
for some (unknown) $a \in [0, 100]$
 \Rightarrow target concept is some initial interval
 $C = H = \{ [0, a] \mid a \in [0, 100] \}$



Observe M instances

Return $[0, b]$,

where b is largest positive example seen.

- Clearly poly time per example.
How many examples?

Sample Complexity of Learning Initial Segment

- Approach#1: Use $m_H(\epsilon, \delta) = \frac{1}{\epsilon} \left(\log \frac{|H|}{\delta} \right)$ instances ?

But \mathcal{H} is UNCOUNTABLE!

- Approach#2:

- Let a_ϵ be real value $< a$ s.t. $[a_\epsilon, a]$ has probability ϵ

$$P([a_\epsilon, a]) = \epsilon$$



- Alg succeeds *iff* it sees example in $[a_\epsilon, a]$

$$P(\text{failure}) = P(\text{none of } M \text{ examples in } [a_\epsilon, a]) = (1 - \epsilon)^M$$

So for $P(\text{failure}) \leq \delta$, need

$$M \geq \frac{1}{\epsilon} \ln \frac{1}{\delta}$$



Uniform Convergence

- Simultaneously estimating all $\{ [a_\epsilon, a] \mid a \in [0, 100] \}$!

Q: Why possible?

A: Only one “degree of freedom”

⇒ each sample provides LOTS of information about many hypothesis

Q: How much is a degree of freedom worth?

Are they all worth the same?

A: Look at “effective number” of concepts, as fn of number of data points seen.

Only grows linearly....

- Number of “effective degrees of freedom”: called “VC-dimension”

Shattering a Set of Instances

- Hypothesis class \mathcal{H} trivially fit

$$\mathbf{X} = \{x_1, \dots, x_k\}$$

if

\forall labeling of examples in \mathbf{X} ,

$\exists h \in \mathcal{H}$ matching labeling

| | h_1 | h_2 | \dots | h_{2^k-1} | h_{2^k} |
|----------|----------|-------|---------|-------------|-----------|
| x_1 | 0 | 0 | \dots | 1 | 1 |
| x_2 | 0 | 0 | \dots | 1 | 1 |
| \vdots | \vdots | | | \vdots | \vdots |
| x_k | 0 | 1 | \dots | 0 | 1 |

- k instances; $|\mathcal{H}| \geq 2^k$

Any subset of size $k - 1$ is unconstrained!

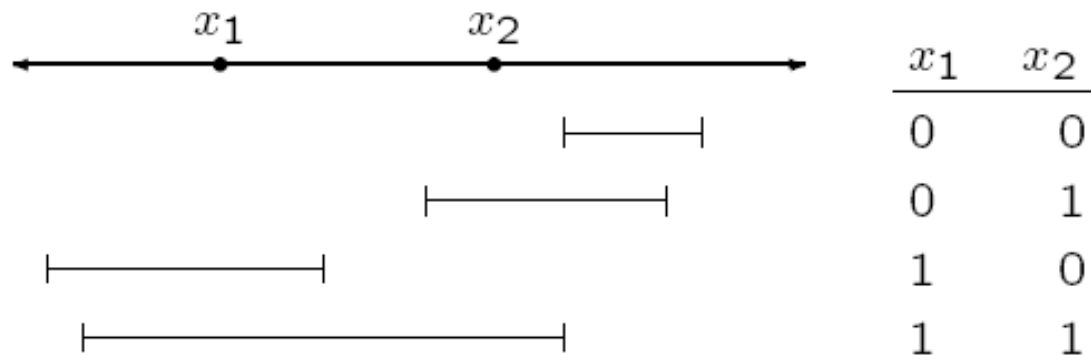
- Defn: Set of points $\mathbf{X} = \{x_i\}$ is shattered by hypothesis class \mathcal{H} if

$\forall S \subset X, \exists h_S \in \mathcal{H}$ s.t.

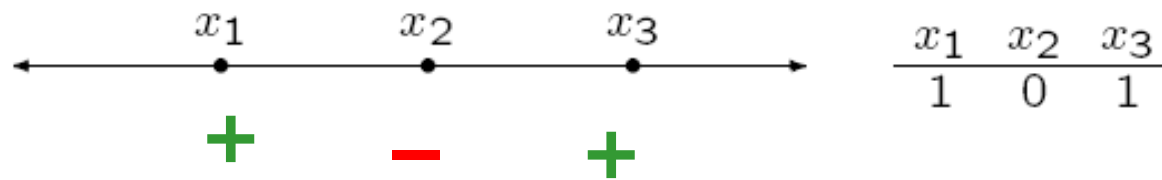
- $h_S(x) = 1 \quad \forall x \in S$
- $h_S(x) = 0 \quad \forall x \notin S$

Example of Shattering

- $\mathcal{H} = \{ [a, b] \mid a < b \} =$ intervals on real line
- Can shatter (any!) 2 points:



- \exists 3 points that can NOT be shattered:



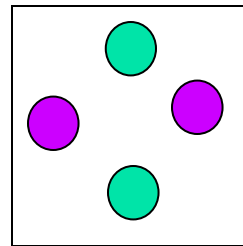
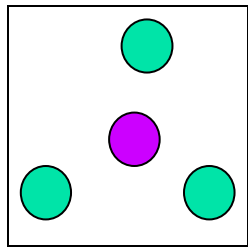
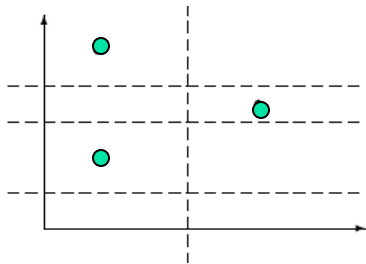


Vapnik-Chervonenkis Dimension

- *Def'n*: VCdim of concept class \mathcal{H}
≡ largest # of points shattered by \mathcal{H}
 - If arbitrarily large finite sets of \mathbf{X} shattered by \mathcal{H} ,
then $\text{VCdim}(\mathcal{H}) = \infty$
 - $\text{VCdim}(\mathcal{H}) = d \iff$
 - \exists set of d points that can be shattered,
but no set of $d+1$ points can be shattered
- Note: $\text{VCdim}(\mathcal{H}) \leq \log_2 |\mathcal{H}|$
- $\text{VCdim}(\mathcal{H})$ measures complexity of \mathcal{H}
... how many distinctions can its elements exhibit

VC-dimension: Linear Separator

- $\mathcal{H}_{\mathcal{L}S_2} = \{ [w_0, w_1, w_2] \in \mathcal{R}^3 \}$
= linear separators in 2-D
- Trivial to fit (any non-linear!) 3 points



- But cannot shatter ANY set of 4 points
 - If one point inside convex hull of others, can not make inside “-” and outside “+”
 - Otherwise, label alternatingly in cycle
- $\Rightarrow VC(\mathcal{H}_{\mathcal{L}S_2}) = VC(\text{LinearSeparator in 2Dim}) = 3$



Some VC Dims

- $\text{VCdim}(\text{LinearSeparator in } k\text{-Dim}) = k + 1$
- Multi-layer perceptron network over n inputs of depth s :
 $d \leq 2(n+1)s(1+\ln s)$
- Exact value for sigmoid units is ?unknown?
... probably slightly larger...
- Typically $\text{VCdim}(\text{model}) \approx$
of non-redundant tunable parameters



VCdim of . . .

- $H_{\text{int}} = \{ \text{intervals of real line} \}$
 - 2

- $H_{\text{box}} = \{ \text{axis-parallel boxes in 2-D} \}$
 - 4

Consider 5 points. Draw smallest enclosing axis-parallel box.
For each side of box, pick one point.. colored red.
Must be at least one pt left – blue.
Can't have Red=+, Blue = —

- $H_{\text{md}} = \{ \text{monotone disjunctions (} n \text{ features)} \}$
 - n

Clearly $\geq n$ as $\{100, 010, 001\}$.
Can not be $>n$ as only 2^n monotone disjunctions

- $H_{\text{all}} = \{ \text{all boolean functions on } n \text{ features} \}$
 - 2^n



How does VCdim measure Complexity?

- Def'n: $H[m]$ = maximum number of ways to split m points using concepts in H
- For $m \leq \text{VCdim}(H)$, $H[m] = 2^m$
For $m \geq \text{VCdim}(H)$, ...
- Theorem: $H[m] = O(m^{\text{VCdim}(H)})$
 - Ie, only $C[m]$ "different" concepts in H wrt any set of m examples.

\Rightarrow ? Replace $\ln(|H|)$ by $\ln(H[m]) \approx O(\text{VCdim}(H))$ in PAC bounds
YES (kinda)! . . . but NOT OBVIOUS,
since different data \Rightarrow different concepts

Upper/Lower Bounds using VCdim

- **Theorem 1:** Given class C ,
for any distribution D , target concept in C ,
given a sample size:

$$\frac{1}{\epsilon} \left(4 \log_2 \left(\frac{2}{\delta} \right) + 8 \text{VCdim}(C) \log_2 \left(\frac{13}{\epsilon} \right) \right)$$

then with prob $\geq 1 - \delta$,
any consistent $h \in C$ has error $\leq \epsilon$.

- **Theorem 2:** If $|C| \leq 2$, then
for any learning alg A ,
 \exists distribution D over X , distribution over C s.t.
expected error of A is $> \epsilon$
if A sees sample of size under

$$\frac{\text{VCdim}(C) - 1}{32\epsilon}$$



Comments on VC Dimension

- VCdim provides good measure of complexity of class:

Upper/Lower (worst case) bounds:

$$\tilde{\Theta}(VC \dim(C))$$

- Does this mean. . .
 - ... can't learn classes of infinite VCdimension?
A: No: just use poly dependence on $\text{size}(c)$
 - ... complicated hypotheses are bad?
A: No. Just need a lot of data to learn complicated concept classes...



Proof of Theorem#2 (Sketch)

- Theorem 2: ... need at least $m = \frac{\text{VCdim}(\mathcal{C}) - 1}{8\epsilon}$

(#examples needed for uniform convergence
... for all bad $h \in \mathcal{C}$ to look bad ...)

Proof: Consider $d = \text{VCdim}(\mathcal{C})$ points $\{x_1, x_2, \dots, x_d\}$
that can be shattered by target concepts $\{c_i\}_{i=1}^{2^k}$

- Define distribution D :
 - $1 - 4\epsilon$ on x_1
 - $4\epsilon / (d - 1)$ on each other
- Given m instances, expect to see only $1/2$ of $\{x_2, \dots, x_d\}$
so $E[\#\text{notSeen}] \geq (d - 1) / 2$
- As can only do 50/50 on instances NOT seen,
expected error is $\#\text{notSeen} \cdot 1/2 \cdot 4\epsilon / (d - 1) = \epsilon$



Summary of Training vs Test Error

- ϵ = “true” error of hyp h
 ϵ^* = minimum true error of any member of \mathcal{H}
 ϵ_T = “training set” error of hyp h
- After m examples, w/ probability $\geq 1 - \delta$, ...

– Finite Hypothesis Class; “Realizable”

$$\epsilon \leq \frac{1}{m} \left[\ln |\mathcal{H}| + \ln \frac{1}{\delta} \right]$$

– Finite Hypothesis Class; “UnRealizable”

$$\epsilon \leq \epsilon^* + \sqrt{\frac{1}{2m} \left[\ln |\mathcal{H}| + \ln \frac{1}{\delta} \right]}$$

– $d = \text{VCdim}(\mathcal{H})$

$$\epsilon \leq 2\epsilon_T + \frac{4}{m} \left[d \log \frac{2em}{d} + \ln \frac{4}{\delta} \right]$$

Case 4:

Why SINGLE Hypothesis Space?

- Large H is likely to include (approx to) target c but . . .
 - w/few examples, cannot reliably find good hypothesis from large hypothesis space
 - That is...
 - **Underfitting**: Every $h \in H$ has high ϵ_T
 \Rightarrow consider larger hypothesis space $H' \supset H$
 - **Overfitting**: Many $h \in H$ have $\epsilon_T \approx 0$
 \Rightarrow consider smaller $H'' \subset H$ to get lower d
- \Rightarrow To learn effectively ($> 1 - \epsilon$) from m instances, only consider H s.t. $|H| \approx e^{\epsilon m}$



How Learning Algorithms Manage This Tradeoff

S1: Start with small hypothesis space \mathcal{H}_1

S2: Grow hypothesis space

$$\mathcal{H}_1 \subset \mathcal{H}_2 \subset \mathcal{H}_3 \subset \dots$$

until finding a good (nearly consistent) hypothesis

Eg1 \mathcal{H}_1 = "leaf", then

\mathcal{H}_2 = "one DecTree node", then

\mathcal{H}_3 = "two DecTree nodes", then ...

Eg2 \mathcal{H}_1 = "constants", then

\mathcal{H}_2 = "linear functions", then

\mathcal{H}_3 = "quadratic functions", then ...

Approaches

1. Easy: $\bigcup_i \mathcal{H}_i$ countable, and realizable
2. General: Structural Risk Minimization
3. "Occam Algorithms"



#4a: Dealing w/ ∞ Set of Hypotheses

- Incremental algorithms:
 $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots \subset \mathcal{H}_n \subset \dots$
 $1 - DNF \subset 2 - DNF \subset 3 - DNF \subset \dots$

Assume: $m(\mathcal{H}_i, \epsilon, \delta)$ instances sufficient to PAC(ϵ, δ)-learn \mathcal{H}_i

Alg? Assume target in \mathcal{H}_1

Draw $m(\mathcal{H}_1, \epsilon, \delta)$ instances

Stop if find good $h_1 \in \mathcal{H}_1$

Otherwise...

Assume target in \mathcal{H}_2

Draw $m(\mathcal{H}_2, \epsilon, \delta)$ more instances

Stop if find good $h_2 \in \mathcal{H}_2$

Otherwise...

...

Assume target in \mathcal{H}_i

Draw $m(\mathcal{H}_i, \epsilon, \delta)$ more instances

Stop if find good $h_i \in \mathcal{H}_i$

Otherwise...

...



Correct Algorithm?

- Q: Suppose find “good” h_k at iteration k .

What is prob of making mistake?

- A: $P(\text{mistake}) = \sum_{i=1..k} P(\text{mistake @ iteration } i)$
 $\leq \sum_{i=1..k} \delta \leq k \delta$

\Rightarrow Need to use δ_i s.t. $\sum_{i=1..k} \delta_i \leq \delta$ for any k

- Eg: $\delta_i = \delta/2^i$

- Note: $P(\text{mistake}) \leq \sum_{i=1..k} \delta_i = \delta \sum_{i=1..k} 1/2^i = \delta$

- Takes k bits to identify member of 2^k -size hypothesis space

- . . . takes k bits just to express such a hypothesis

\Rightarrow reasonable to allow learning alg'm time poly in

$1/\epsilon$, $1/\delta$ and SIZE OF HYPOTHESIS



#4b: Structural Risk Minimization

- Consider
 - nested series: $H_1 \subset H_2 \subset \dots \subset H_k \subset \dots$
 - with VCdim: $d_1 \leq d_2 \leq \dots \leq d_k \leq \dots$
 - training errors: $\epsilon_1 \geq \epsilon_2 \geq \dots \geq \epsilon_k \geq \dots$

- Choose $h_k \in H_k$ that minimizes

$$\epsilon \leq 2\epsilon^k + \frac{4}{m} \left[d_k \log \frac{2em}{d_k} + \ln \frac{4}{\delta} \right]$$



Structural Risk Minimization

For $h \in \mathcal{H}$

$L(h)$ Probability of miss-classification

$\hat{L}_n(h)$ Empirical fraction of miss-classifications

Vapnik and Chervonenkis 1971: For **any** distribution with prob. $1 - \delta$, $\forall h \in \mathcal{H}$,

$$L(h) < \underbrace{\hat{L}_n(h)}_{\text{emp. error}} + c \underbrace{\sqrt{\frac{\text{VCdim}(\mathcal{H}) \log n + \log \frac{1}{\delta}}{n}}}_{\text{complexity penalty}}$$



An Improved VC Bound II

Canonical hyper-plane:

$$\min_{1 \leq i \leq n} |\mathbf{w}^\top \mathbf{x}_i + b| = 1$$

(No loss of generality)

Improved VC Bound (Vapnik 95) VC dimension of set of canonical hyper-planes such that

$$\|\mathbf{w}\| \leq A$$

$$\mathbf{x}_i \in \text{Ball of radius } L$$

is

$$\text{VCdim} \leq \min(A^2 L^2, d) + 1$$

Observe: Constraints reduce VC-dim bound
Canonical hyper-planes with **minimal norm** yields best bound

Suggestion: Use hyper-plane with minimal norm



Case 5: Data Dependent Bounds

- So far, bounds depend only on
 - ϵ_T
 - quantities computed prior to seeing S
(eg, size of H)
 \Rightarrow "worst case"
as must work for all but δ of possible training sets
- Data dependent bounds consider how h fits data
 - If S is not worst case training set
 \Rightarrow tighter error bound!

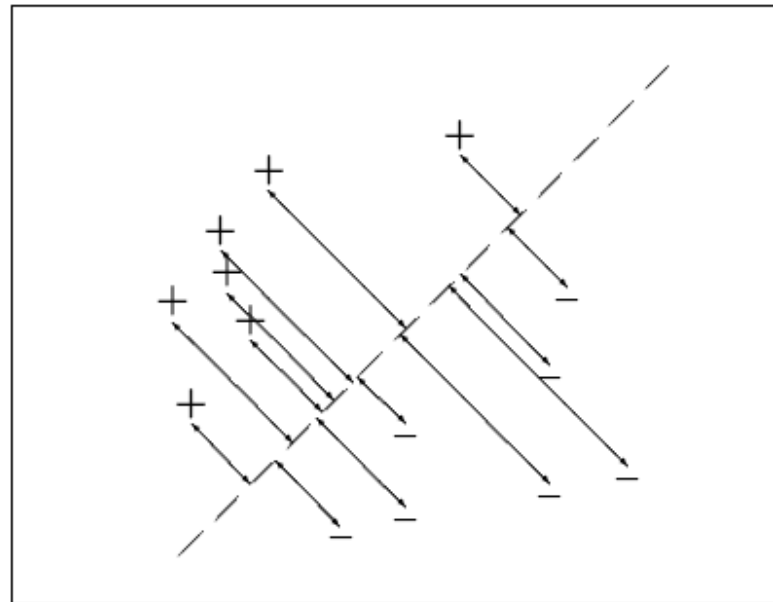
Margin Bounds

- $g(x)$ is real-valued function
“thresholded at 0” to produce $h(x)$:

$$\begin{aligned}g(x) > 0 &\Rightarrow h(x) = +1 \\g(x) < 0 &\Rightarrow h(x) = -1\end{aligned}$$

- **Margin** of $h(x)$ wrt S is

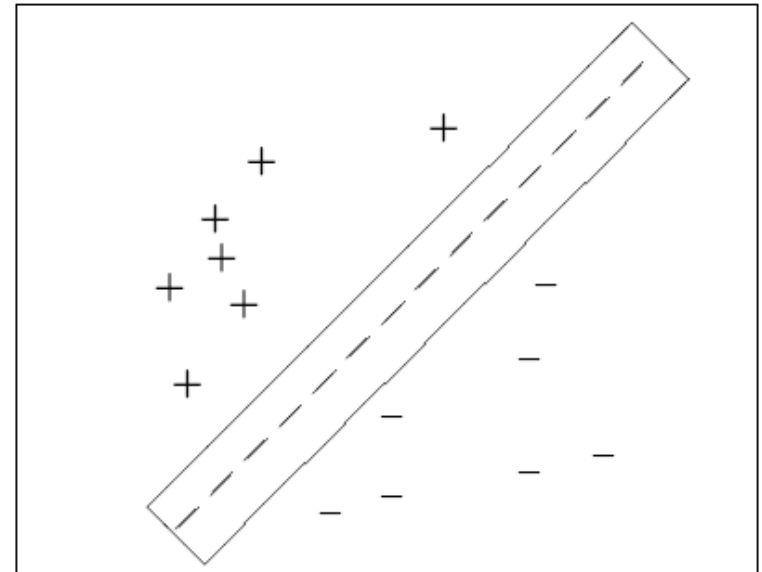
$$\gamma(g, S) = \min_i \{y_i g(x_i)\}$$



Margin Bounds: Key Intuition

Let $G = \{g(x)\}$ = set of real-valued functions that can be thresholded at 0 to give $h(x)$.

- Consider “thickening” each $g \in G$
... must correctly classify every point w/ margin $\geq \gamma$



- **fat shattering dimension:** $\text{fat}_\gamma(G)$
 \equiv VCdim of these “fat” separators

Note $\text{fat}_\gamma(G) \leq \text{VCdim}(G)$



Noise Free Margin Bound

- Spse find $g \in G$
with margin $\gamma = \gamma(g, S)$
for a training set of size m

- Then, with probability $1 - \delta$

$$\epsilon \leq \frac{2}{m} \left[d \log \frac{2em}{d\gamma} \log \frac{32m}{\gamma^2} + \log \frac{4}{\delta} \right]$$

$d = \text{fat}_{\gamma/8}(G)$ with margin $\gamma/8$

- Note $\text{fat}(\cdot)$ kinda-like $\text{VCdim}(\cdot)$!



Soft Margin Classification (2)

- Error rate of linear separator with unit weight vector and margin γ on training data lying in a sphere of radius R is, with probability $\geq 1 - \delta$,

$$\epsilon \leq \frac{C}{m} \left[\frac{R^2 + \|\xi\|^2}{\gamma^2} \log^2 m + \log \frac{1}{\delta} \right]$$

(constant C)

\Rightarrow we should

- maximize margin γ
- minimize slack $\|\xi\|^2$

... see support vector machines!



Fat Shattering for Linear Separators: Noise-Free

Spse support for $P(\mathbf{x})$ within sphere of radius R

$$\|\mathbf{x}\| \leq R$$

$$G = \{g \mid g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} \ \& \ \|\mathbf{w}\| = 1\}$$

$$\text{Then } \text{fat}_\gamma(G) = \left(\frac{R}{\gamma}\right)^2$$

$$\Rightarrow \epsilon \leq \frac{2}{m} \left[\frac{64R^2}{\gamma^2} \log \frac{em\gamma}{8R^2} \log \frac{32m}{\gamma^2} + \log \frac{4}{\delta} \right]$$
$$\in \tilde{O}\left(\frac{R^2}{m\gamma^2}\right)$$

\Rightarrow For fixed R, m :

seek g that maximizes γ !

maximum margin classifier

- Even with kernel $K(\cdot, \cdot)$... where $\|\mathbf{x}\| = \sqrt{K(\mathbf{x}, \mathbf{x})}$

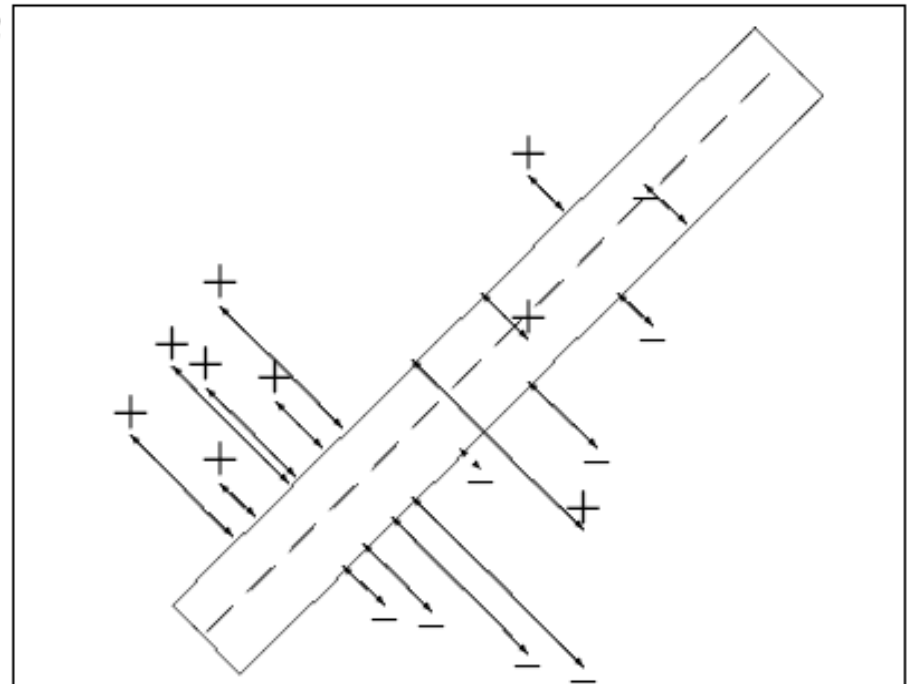
Soft Margin Classification

- Extension of margin analysis:
When data is not linearly separable:

- $\xi_i = \max\{0, \gamma - y_i, g(\mathbf{x}_i)\}$
"margin slack variable" for $\langle \mathbf{x}_i, y_i \rangle$

Note: $\xi_i > \gamma \Rightarrow \mathbf{x}_i$ misclassified by h

- $\xi = \langle \xi_1, \dots, \xi_m \rangle$
"margin slack vector for h on S "





Irrelevant Features

- Consider learning $CD(n)$ = disjunction of n features
 - “List-then-Eliminate” makes $O(n)$ mistakes
 - PAC-learning: $O(n/\epsilon \log(1/\delta))$
- Spse n is HUGE
 - Words in text
 - Boolean combination of “atomic” features
 - Features extracted in 480x560 image
 - ... but only $r \ll n$ features “relevant”
 - Eg: concept $x_4 \vee \neg x_{91} \vee \neg x_{203} \vee x_{907}$
- \exists learning alg that makes $O(r \ln n)$ mistakes!
“Winnow”



Winnow Algorithm

- Initialize weights w_1, \dots, w_n to 1
- Do until bored:
 - Given example $\mathbf{x} = [x_1, \dots, x_n]$,
If $w_1x_1 + w_2x_2 + \dots + w_nx_n \geq n$
output 1 otherwise 0
 - If mistake:
 - (a) If predicts 0 on 1-example, then
for each $x_i = 1$, set $w_i := w_i \times 2$
 - (b) If predicts 1 on 0-example, then
for each $x_i = 1$, set $w_i := w_i / 2$



Winnow's Effectiveness

Theorem Winnow MB-learns $CD(n)$, making at most $2+3r(1+\lg n)$ mistakes when target concept is disjunction of r var's.

Proof: 1. Any mistake made on 1-example must double params

- ≥ 1 weights in target function (the relevant weights),
 - & mistake on 0-example will not halve these weights.
 - Each "relevant" weight can be doubled $\leq 1+\lg n$ times, since only weights $\leq n$ can be doubled.
(Never double any weight $w_i > n$ as that weight alone \Rightarrow class is 1)
- \Rightarrow Winnow makes $\leq r(1+\lg n)$ mistakes on 1-examples

2. Negative examples?

- Let sw_t be sum of weights $\sum w_i = n$, at time t .
Initially $sw_0 = n$.
Each mistake on 1-example increases sw by $\leq n$
(. . . before doubling, we know $w_1x_1 + w_2x_2 + \dots + w_nx_n < n$)
Each mistake on 0-example decreases sw by $\geq n/2$
(. . . before halving, we know $w_1x_1 + w_2x_2 + \dots + w_nx_n \geq n$)
- As $sw \geq 0$, number of mistakes made on 0-examples $\leq 2 + 2$ number of mistakes made on 1-examples.

- So total # of mistakes is $r(1+\ln n) + [2+2r(1+\lg n)]$

Incorporating Winnow Into PAC Model

- Given a **MB(M)-learner**, can **PAC(ϵ, δ)-learn**
 - Return any h_i that makes $\frac{1}{\epsilon} \log(\frac{M}{\delta})$ correct predictions
 - Requires $m = \frac{M}{\epsilon} \log(\frac{M}{\delta}) = \frac{r \log(n)}{\epsilon} \log(\frac{r \log(n)}{\delta})$ instances
- Better PAC-learner: $O(\frac{1}{\epsilon}[r \log(n) + \log(\frac{1}{\delta})])$

1. Draw $m_1 = 4/\epsilon \max \{ M, 2 \ln(2/\delta) \}$ instances, S_1
2. Run Winnow (a MB-learner) on S_1 ,
generating $\leq M$ hypotheses $H = \{ h_1, \dots, h_M \}$
3. Draw $m_2 = O(8/\epsilon \log(2M/\delta))$ more instances S_2
4. Use S_2 to find best hypothesis, h^* in H
5. Return h^*

- Why: Most ϵ -bad hypotheses have error $\gg \epsilon$
 \Rightarrow reveal "badness" in $< \frac{1}{\epsilon} \log(\frac{M}{\delta})$ instances



Proof

- m_1 guarantees that ≥ 1 of H is good
 m_2 distinguishes good h^* from bad members of H .
- After m_1 instances, ≥ 1 of H has error $\leq \epsilon/2$
PROOF: Spse first $k - 1$ hyp's all have error $> \epsilon/2$,
and h_k had error $\leq \epsilon/2$
What is prob that h_k occurs after m_1 instances?

Worst if $k = M$ and each $\text{err}_D(h_i) = \epsilon/2$

Chernoff bounds $\Rightarrow \delta/2$:

- Consider flipping (sequence of M) $\epsilon/2$ weighted coins
- (each "head" \equiv error)
- After m_1 flips, expect $m_1 \times \epsilon/2 \leq 2M$ "heads"
- Prob of getting under M ($\leq 1/2$ exp. number) heads
 $\leq P(Y_M \leq (1 - 1/2) \epsilon/2) \leq \exp(-M \epsilon/2 \cdot 1/2)/2$
 $\leq \exp(-M \epsilon/8) \leq \delta$



Proof (II)

Use m_2 , select h^* w/ $\text{err}_S(h^*) \leq 3/4 \epsilon$

With prob $\geq 1 - \delta/2$ $\text{err}_D(h^*) \leq \epsilon$

PROOF: Need to show $\text{err}_S(h_i)$

[average # mistakes made by h_i over m_2 samples]

is within 3/4 of $\mu_i = \text{err}_D(h_i)$

- $P(\text{err}_S(h_i) < \text{err}_D(h_i) \times (1 - 1/4)) \leq \exp(- (m_2 \epsilon^2 / 4) / 2) \leq \delta / (2M)$
- So prob ANY $h_i \in H$ is off by $< 3/4$ is under $\delta / 2$
- m_1 is leading term
 $\Rightarrow O(1/\epsilon [r \log(n) + \log(1/\delta)])$
- Best known bound for learning r of n disjuncts!
- Note: Might NOT find 0 error r -disjunction. . .