# Computational learning and discovery

CSI 873 / MATH 689

Instructor: I. Griva

Wednesday 7:20 - 10 pm

# 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

# Computational Learning Theory

**answers such questions as**

- How many training examples are needed to converge with high probability to a successful learner? **(Sample complexity).**

- How much computational effort are needed for a learner to converge to a successful hypothesis? **(Computational complexity).**

- How many training example will the learner misclassify before converging to a successful hypothesis? **(Mistake bounds).**

# 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: Candidate elimination algorithm

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 $VS$ 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

# General notations

$X$     is a set of instances over which the concept is defined

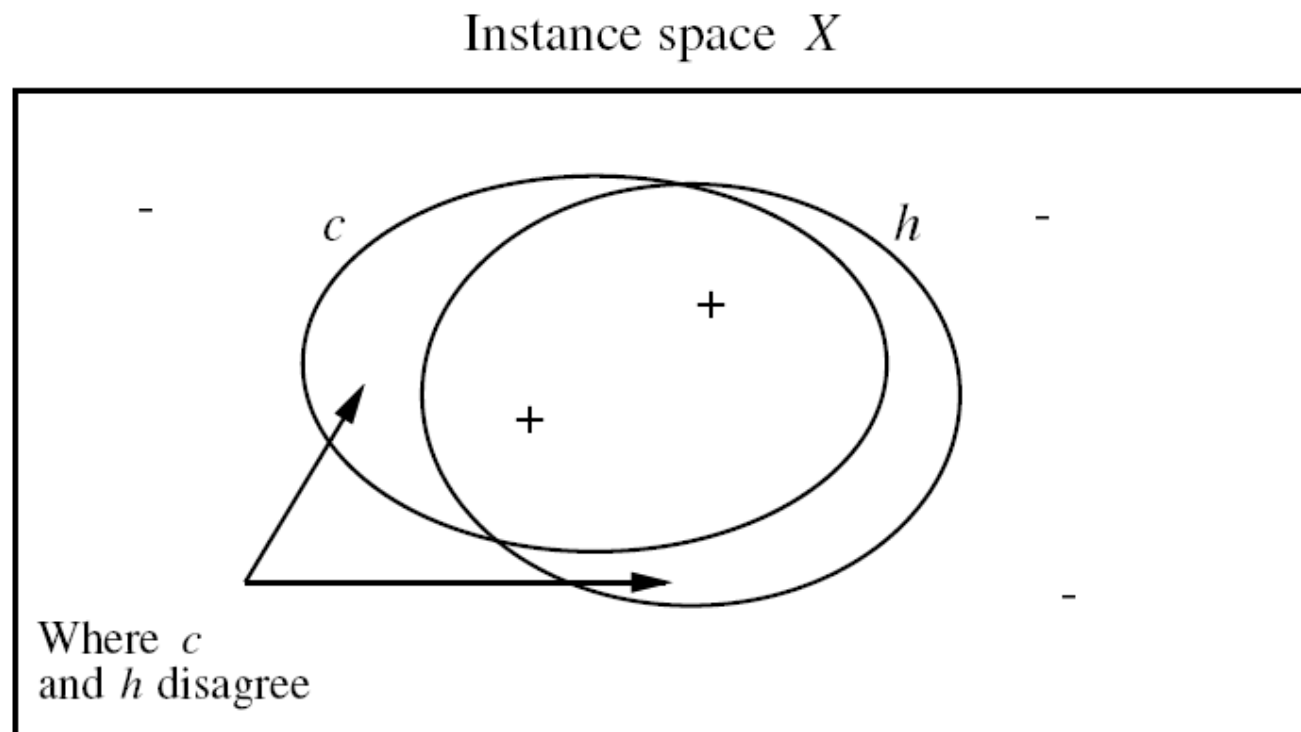$D$     is is the probability distribution that defines the probability of encountering each instance in $X$

$c$     is a target concept, $c\colon X \to \{0, 1\}$

$(x, c(x))$     is a training example

$H$     is the set of all possible hypotheses the learner considers to identify the target concept

The goal is to find a hypothesis $h$ in $H$ such that $h(x) = c(x)$ for all $x$ in $X$

# True error of the hypothesis

Instance space $X$



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 $\mathcal{D}$.

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

# True error vs training 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 $h$?
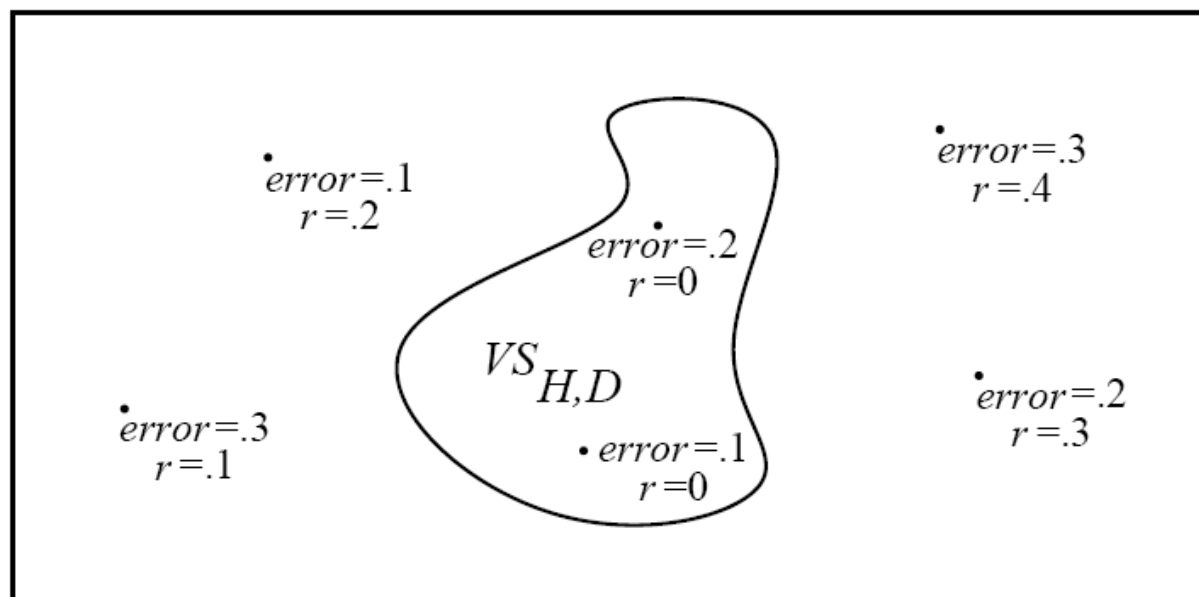
# PAC learning

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 $size(c)$.

# Exhausting the version space

Hypothesis space $H$



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

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

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

# Exhausting the version space

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

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

# Exhausting the version space

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 $\delta$

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

then

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

# Learning conjunction 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

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

# Learning *Enjoy Sport*

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

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

$$m \geq \frac{1}{.1}(\ln 973 + \ln(1/.05))$$
$$m \geq 10(\ln 973 + \ln 20)$$
$$m \geq 10(6.88 + 3.00)$$
$$m \geq 98.8$$

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:

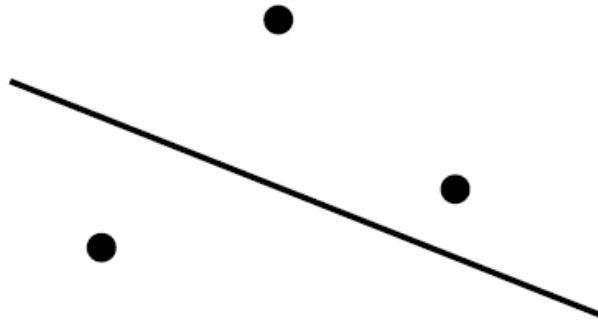$$Pr[error_D(h) > error_D(h) + \epsilon] \leq e^{-2m\epsilon^2}$$

# Shattering 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.
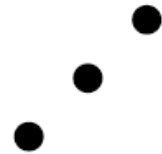
# Vapnik-Chervonenkis dimension

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

# VC dimension of linear decision surface



(a)                    (b)

# Sample complexity and VC dimension

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

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

**Sufficient condition**

# Sample complexity and VC dimension

**Theorem 7.3. Lower bound on sample complexity.** Consider any concept class $C$ such that $VC(C) \geq 2$, any learner $L$, and any $0 < \epsilon < \frac{1}{8}$, and $0 < \delta < \frac{1}{100}$. Then there exists a distribution $\mathcal{D}$ and 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 $\delta$, $L$ outputs a hypothesis $h$ having $error_{\mathcal{D}}(h) > \epsilon$.

**Necessary condition**

# VC dimension of neural network

**Theorem 7.4.   VC-dimension of directed acyclic layered networks.** (See Kearns and Vazirani 1994.) Let $G$ be a layered directed acyclic graph with $n$ input nodes and $s \geq 2$ internal nodes, each having at most $r$ inputs. Let $C$ be a concept class over $\Re^r$ of VC dimension $d$, corresponding to the set of functions that can be described by each of the $s$ internal nodes. Let $C_G$ be the $G$-composition of $C$, corresponding to the set of functions that can be represented by $G$. Then $VC(C_G) \leq 2ds \log(es)$, where $e$ is the base of the natural logarithm.

# Mistake bounds

So far: how many examples needed to learn?

What about: how many mistakes before convergence?

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?

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

FIND-S:

- Initialize $h$ to the most specific hypothesis $l_1 \wedge \neg l_1 \wedge l_2 \wedge \neg l_2 \ldots l_n \wedge \neg l_n$
- 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 $h$?

Consider the Halving Algorithm:

- Learn concept using version space CANDIDATE-ELIMINATION 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?

# 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 $M_A(C)$.

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

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

# Relative Mistake Bound for Weighted Majority

**Theorem 7.5.** **Relative mistake bound for WEIGHTED-MAJORITY.** Let $D$ be any sequence of training examples, let $A$ be any set of $n$ prediction algorithms, and let $k$ be the minimum number of mistakes made by any algorithm in $A$ for the training sequence $D$. Then the number of mistakes over $D$ made by the WEIGHTED-MAJORITY algorithm using $\beta = \frac{1}{2}$ is at most

$$2.4(k + \log_2 n)$$

# Summary

- The probably approximately correct (PAC) model considers algorithms that learn target concepts from some concept class $C$, using training examples drawn at random according to an unknown, but fixed, probability distribution. It requires that the learner probably (with probability at least $[1 - \delta]$) learn a hypothesis that is approximately (within error $\epsilon$) correct, given computational effort and training examples that grow only polynomially with $1/\epsilon$, $1/\delta$, the size of the instances, and the size of the target concept.

- Within the setting of the PAC learning model, any consistent learner using a finite hypothesis space $H$ where $C \subseteq H$ will, with probability $(1 - \delta)$, output a hypothesis within error $\epsilon$ of the target concept, after observing $m$ randomly drawn training examples, as long as

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

This gives a bound on the number of training examples sufficient for successful learning under the PAC model.

# Summary (continued)

- One constraining assumption of the PAC learning model is that the learner knows in advance some restricted concept class $C$ that contains the target concept to be learned. In contrast, the *agnostic learning* model considers the more general setting in which the learner makes no assumption about the class from which the target concept is drawn. Instead, the learner outputs the hypothesis from $H$ that has the least error (possibly nonzero) over the training data. Under this less restrictive agnostic learning model, the learner is assured with probability $(1-\delta)$ to output a hypothesis within error $\epsilon$ of the best possible hypothesis in $H$, after observing $m$ randomly drawn training examples, provided

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

- The number of training examples required for successful learning is strongly influenced by the complexity of the hypothesis space considered by the learner. One useful measure of the complexity of a hypothesis space $H$ is its Vapnik-Chervonenkis dimension, $VC(H)$. $VC(H)$ is the size of the largest subset of instances that can be shattered (split in all possible ways) by $H$.

# Summary (continued)

- An alternative upper bound on the number of training examples sufficient for successful learning under the PAC model, stated in terms of $VC(H)$ is

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

A lower bound is

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

- An alternative learning model, called the *mistake bound model,* is used to analyze the number of training examples a learner will misclassify before it exactly learns the target concept. For example, the HALVING algorithm will make at most $\lfloor\log_2|H|\rfloor$ mistakes before exactly learning any target concept drawn from $H$. For an arbitrary concept class $C$, the best worst-case algorithm will make $Opt(C)$ mistakes, where

$$VC(C) \leq Opt(C) \leq \log_2(|C|)$$

- The WEIGHTED-MAJORITY algorithm combines the weighted votes of multiple prediction algorithms to classify new instances. It learns weights for each of these prediction algorithms based on errors made over a sequence of examples. Interestingly, the number of mistakes made by WEIGHTED-MAJORITY can be bounded in terms of the number of mistakes made by the best prediction algorithm in the pool.