

Computational Learning Theory



Introduction

The PAC Learning Framework

Finite Hypothesis Spaces

Examples of PAC Learnable Concepts

Introduction

- Computational learning theory:
 - Provides a theoretical analysis of learning
 - Shows when a learning algorithm can be expected to succeed
 - Shows when learning may be impossible
 - ...

Introduction

- Some fundamental problems addressed by Computational Learning Theory:
 - **Sample Complexity:** How many examples we need to find a good hypothesis?
 - **Computational Complexity:** How much computational power we need to find a good hypothesis?
 - **Mistake Bound:** How many mistakes we will make before finding a good hypothesis?

Computational Learning Theory



Introduction

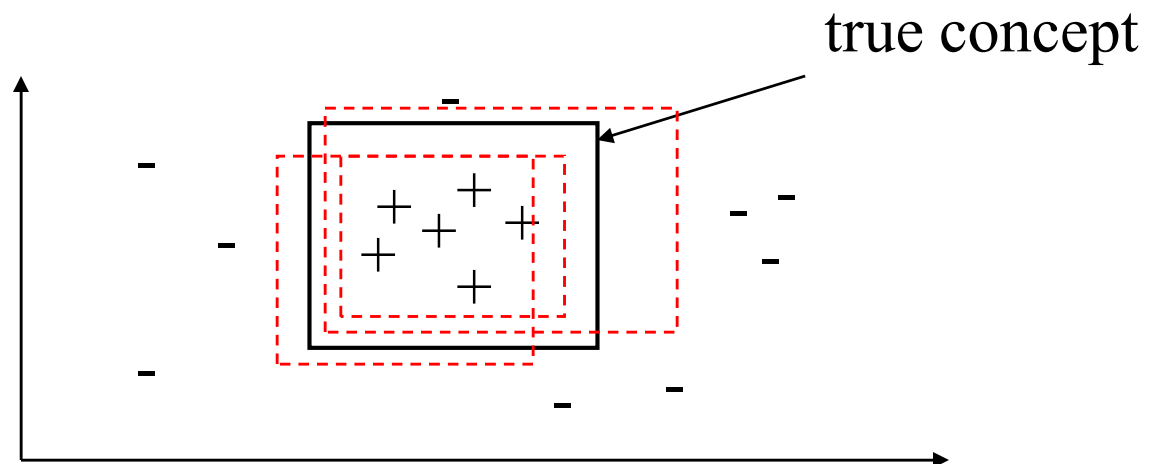
The PAC Learning Framework

Finite Hypothesis Spaces

Examples of PAC Learnable Concepts

The PAC Learning Framework

- Let's start with a simple problem: Assume a two dimensional space with positive and negative examples. Our goal is to find a rectangle that includes the positive examples but not the negatives (input space is \mathbb{R}^2):



Definitions

- *Class of Concepts **C***. Let C be a class of concepts that we wish to learn. In our example C is the family of all rectangles in \mathbb{R}^2 .
- *Distribution **D***. Assume instances are generated at random from a distribution D .
- *Class of Hypotheses **H***. The hypotheses our algorithm considers while learning the target concept.
- *True error of a hypothesis **h***
 - $\text{error}_D(h) = \Pr_D[c(x) \neq h(x)]$

True Error of A Hypothesis

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**

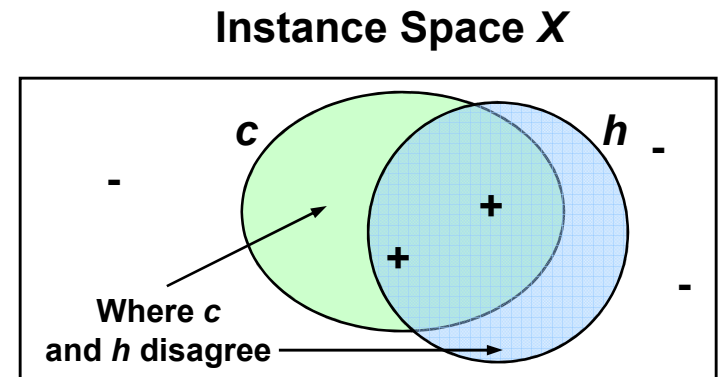
- True error of hypothesis h with respect to target concept c :

How often $h(x) \neq c(x)$ **over random instances drawn from distribution D**

Definition

- The true error (denoted $error_D(h)$) of hypothesis h with respect to target concept c and distribution D is the probability that h will misclassify an instance drawn at random according to D .

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



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 D

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

Set of 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}

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

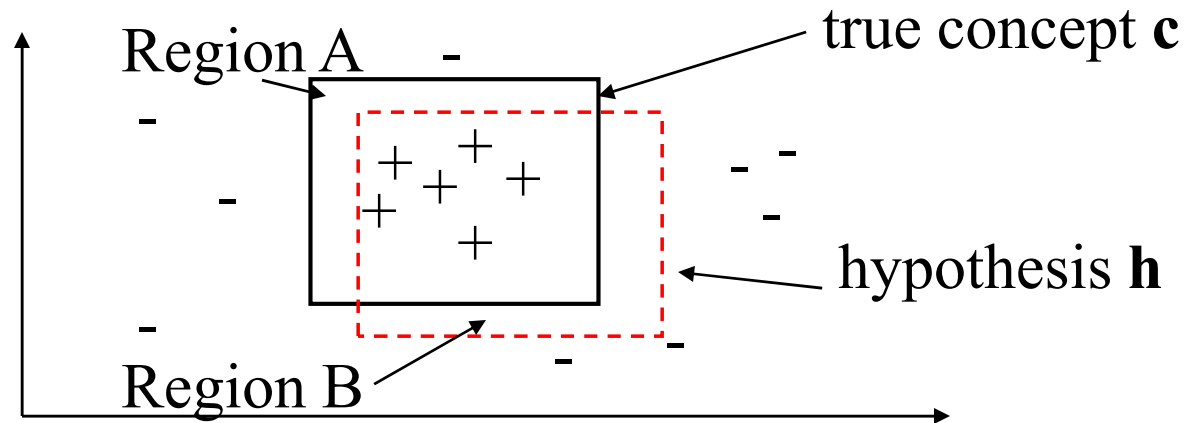
Probability distribution $P(x)$

True Error

Region A : false negatives

Region B : false positives

True error: probability in regions A and B.



Probably Approximately Correct learning

[Valiant84]

- *Concept class* \mathcal{C} of Boolean functions over domain X typically $\{0, 1\}^n$ or \mathbb{R}^n
- Unknown target concept $f \in \mathcal{C}$ to be learned from examples
- **Unknown** and **arbitrary** distribution \mathcal{D} over X

Learner has access to i.i.d. draws from \mathcal{D} , labeled according to f :

$(x^1, f(x^1))$

$(x^2, f(x^2))$

$(x^3, f(x^3))$

...

each x^1, x^2, \dots belongs to X ,
i.i.d. drawn from \mathcal{D}

PAC learning concept class \mathcal{C}

Learner's goal:

Efficiently come up with hypothesis that will have high accuracy on future examples.

- For any target function $f \in \mathcal{C}$,
- for any distribution \mathcal{D} over X ,
- with probability $1 - \delta$, learner outputs hypothesis $h : X \rightarrow \{0, 1\}$ that is ϵ -accurate w.r.t. \mathcal{D} :

$$\Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \epsilon.$$

Algorithm must be computationally efficient: should run in time

$$\text{poly}(n, \frac{1}{\epsilon}, \frac{1}{\delta}, \text{size}(f)).$$

Considerations

- We don't need a hypothesis with zero error. There might be some error as long as it is small (bounded by a constant ϵ).
- We don't need to always produce such a good enough hypothesis. The probability of failure should be bounded by a constant δ .
- **Goal:** With probability $1-\delta$, output a hypothesis h which satisfies $\text{error}_D(h) < \epsilon$.

$$\text{Probability}[\text{error}_D(h) > \epsilon] < \delta$$

-
- A learner finds a hypothesis **h** that is **consistent** with the training data

- $\text{Error}_{\text{Train}}(h) = 0$

- The probability that **h** has more than ε true error

- $\text{Error}_{\text{true}}(h) \geq \varepsilon$

□ Hypothesis h that is consistent with training data means it got m i.i.d points right

□ h is consistent but "bad": it gets all training data right but has high true error

□ Prob. h with $\text{Error}_{\text{true}}(h) \geq \epsilon$ gets one data point right:

■ $P(h \text{ gets one point right}) \leq 1 - \epsilon$

■ $P(h \text{ gets } m \text{ iid points right}) \leq (1 - \epsilon)^m$

■ We want this to be less than δ . So lets set:

$$(1 - \epsilon)^m \leq \delta$$

■ Since $(1-x) \leq e^{-x}$ we have that

$e^{-\epsilon m} \leq \delta$ or equivalently (taking \ln of each side)

$$m \geq \frac{1}{\epsilon} \ln \left(\frac{1}{\delta} \right)$$

■ The result grows linearly in $1/\epsilon$ and logarithmically $1/\delta$

Valiant's "PAC" model [Val84]

"PAC" = Probably Approximately Correct:

- learning problem is identified with a "concept class" \mathcal{C} , which is a set of functions ("concepts") $f : \{0,1\}^n \rightarrow \{0,1\}$
- nature/adversary chooses one particular $f \in \mathcal{C}$ and a *probability distribution* on inputs \mathcal{D}
- the *learning algorithm* now takes as inputs ϵ and δ , and also gets random examples $\langle x, f(x) \rangle$, x drawn from \mathcal{D}
- goal: with probability $1-\delta$, output a hypothesis h which satisfies $\Pr_{x \leftarrow \mathcal{D}} [h(x) \neq f(x)] < \epsilon$
- efficiency: running time of algorithm, counting time 1 for each example; hopefully $\text{poly}(n, 1/\epsilon, 1/\delta)$

Example – learning conjunctions

As an example, we present an algorithm ([Val84]) for learning the concept class of **conjunctions** – i.e., **C** is the set of all AND functions.

- start with the hypothesis $h = x_1 \wedge x_2 \wedge \cdots \wedge x_n$
- draw $O((n / \epsilon) \log(1/\delta))$ examples:
 - whenever you see a *positive* example; e.g.,
 $\langle 11010110, 1 \rangle$,
you know that the zero coordinates (in this case, x_3, x_5, x_8) can't be in the target AND;
delete them from the hypothesis

It takes a little reasoning to show this works,
but it does.

Learning DNF formulas

Probably the most important concept class we would like to learn is **DNF formulas**: e.g., the set of all functions like

$$f = (x_1 \wedge \overline{x_2} \wedge \overline{x_6}) \vee (\overline{x_1} \wedge x_3) \vee (x_4 \wedge x_5 \wedge \overline{x_7} \wedge x_8).$$

(We actually mean poly-sized DNF: the number of *terms* should be $n^{O(1)}$, where n is the number of variables.)

Why so important?

- natural form of knowledge representation for people
- historical reasons: considered by Valiant, who called the problem “tantalizing” and “apparently [simple]”
- yet has proved a great challenge over the last 20 years

ADVANCED The original PAC model

The trouble with this model is that, despite Valiant's initial optimism, PAC-learning DNF formulas appears to be very hard.

The fastest known algorithm is due to Klivans and Servedio [KS01], and runs in time $\exp(n^{1/3} \log^2 n)$.

Technique: They show that for any DNF formula, there is a polynomial in x_1, \dots, x_n of degree at most $n^{1/3} \log n$ which is *positive* whenever the DNF is true and *negative* whenever the DNF is false. *Linear programming* can be used to find a hypothesis consistent with every example in time $\exp(n^{1/3} \log^2 n)$.

Note: Consider the model, more difficult than PAC, in which the learner is forced to output a hypothesis which itself is a DNF. In this case, the problem is *NP-hard*.



REST IS ADVANCED

Computational Learning Theory



Introduction

The PAC Learning Framework

Finite Hypothesis Spaces

Examples of PAC Learnable Concepts

Sample Complexity for Finite Hypothesis Spaces

- Definition: A *consistent learner* outputs the hypothesis \mathbf{h} in \mathbf{H} that perfectly fits the training examples (if possible).
- How many examples do we need to be approximately correct in finding a hypothesis output by a consistent learner that has low error?

Version Space ϵ -exhausted

- This is the same as asking how many examples we need to make the Version Space contain no hypothesis with error greater than ϵ .
- When a Version Space **VS** is such that no hypothesis has error greater than ϵ , we say the version space is ϵ -exhausted.
- How many examples do we need to make a version space VS be ϵ -exhausted?

Probability of Version Space being ϵ -exhausted

- The probability that the version space is not ϵ -exhausted after seeing m examples is the same as asking the probability that no hypothesis in **VS** has error greater than ϵ .
- Since the size of the **VS** is less than the size of the whole hypothesis space **H**, then that probability is clearly less than

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

- If we make this less than δ , then we have that

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

Haussler, 1988

- **Theorem:** Hypothesis space H is finite, dataset D with m i.i.d samples, $0 < \epsilon < 1$: for any learned hypothesis h that is consistent on the training data:

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$$

- **Limitations of Haussler '88:**
 - Consistent classifier
 - Size of hypothesis space

Agnostic Learning

- What happens if our hypothesis space \mathbf{H} does not contain the target concept \mathbf{c} ?
- Then clearly we can never find a hypothesis \mathbf{h} with zero error.
- Here we want an algorithm that simply outputs the hypothesis with minimum training error.

Using PAC bound

- Pick ε and δ , gives you m
- Pick m and δ , gives you ε

Computational Learning Theory



Introduction

The PAC Learning Framework

Finite Hypothesis Spaces

Examples of PAC Learnable Concepts

Examples of PAC-Learnable Concepts

- Can we say that concepts described by conjunctions of Boolean literals are PAC learnable?
- First, how large is the hypothesis space when we have n Boolean attributes?
 - Answer: $|H| = 3^n$

Examples of PAC-Learnable Concepts

- If we substitute this in our analysis of **sample complexity** for finite hypothesis spaces we have that:

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

- Thus the set of conjunctions of Boolean literals is
- PAC learnable.

K-Term DNF not PAC Learnable

- Consider now the class of functions of **k-term DNF** expressions. These are expressions of the form

$$T_1 \vee T_2 \vee \dots \vee T_k$$

- where \vee stands for disjunction each of the **k** terms and
- and T_i is a conjunction of Boolean attributes.

- E.g. A **3-term** DNF:

$$(x_1 \wedge \neg x_2) \vee (x_6 \wedge x_7) \vee (x_9)$$

K-Term DNF not PAC Learnable

- The size of $|H|$ is $k3^n$
- Using the equation for the sample complexity of finite hypothesis spaces:

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

- Although the sample complexity is polynomial in the main parameters, this problem is known to be NP-complete.

K-Term CNF is PAC Learnable

- But it is interesting to see that a larger family of functions, the class of **k-CNF** expressions is PAC learnable.
 - **k-CNF**: Conjunction of disjunctions where each disjunct has $\leq k$ literals
- This is interesting because the class of k-CNF expressions is strictly larger than the class of k-term DNF expressions.
 - Can convert k-term DNF into k-CNF by distributivity laws.

k -CNF Expressions

Definition: expressions $T_1 \wedge \dots \wedge T_j$ of arbitrary length j with each term T_i a disjunction of at most k boolean attributes.

Algorithm: reduce problem to that of learning conjunctions of boolean literals. New variables:

$$a_i(X_1) \vee \dots \vee a_i(X_n) \rightarrow Y_{a_i(X_1), \dots, a_i(X_n)}.$$

- the transformation is a bijection;

k -Term DNF Terms and k -CNF Expressions

- **Observation:** any k -term DNF formula can be written as a k -CNF expression. By associativity,

$$\bigvee_{i=1}^k a_i(X_1) \wedge \cdots \wedge a_i(X_n) = \bigwedge_{i_1, \dots, i_k=1}^n a_1(X_{i_1}) \vee \cdots \vee a_k(X_{i_k}).$$

- **Example:** $(u_1 \wedge u_2 \wedge u_3) \vee (v_1 \wedge v_2 \wedge v_3) = \bigwedge_{i,j=1}^3 (u_i \vee v_j).$
- But, the number of new variables is exponential in k : $O(n^k).$