

Understanding Machine Learning – A theory Perspective (part 3)

Shai Ben-David

University of Waterloo

MLSS at MPI Tübingen, 2017

The fundamental theorem (qualitative)

Theorem: Given a class H of binary valued functions the following statements are equivalent:

- a. H has the Uniform Convergence Property
- b. ERM is an agnostic PAC learner for H
- c. H is agnostic PAC learnable
- d. H is PAC learnable
- e. $VCdim(H)$ is finite

Main tool for (e) implies (a)

The Shatter function

For a class H define a function $\Pi_H: \mathbb{N} \rightarrow \mathbb{N}$
as $\Pi_H(m) = \max_{\{A: |A|=m\}} |\{h|_A : h \text{ in } H\}|$

Some basic properties of the shatter function:

1. For every $m \leq \text{VCdim}(H)$, $\Pi_H(m) = 2^m$
2. For every $m > \text{VCdim}(H)$, $\Pi_H(m) < 2^m$

The Sauer (Shelah, Perles) lemma

For every class H of finite VC-dimension, d ,

For every m ,

$$\Pi_H(m) \leq \sum_{i=0}^d \binom{m}{i} \leq m^d$$

Quantitative version of the Fundamental Theorem

For some constants C_1, C_2 , for every d and every class H of binary valued functions such that $\text{VCdim}(H)=d$,

1. H has **Uniform Convergence** property with

$$C_1(d+\log(1/\delta))/\varepsilon^2 < m^{\text{uc}}_H(\varepsilon, \delta) < C_1(d+\log(1/\delta))/\varepsilon^2$$

2. H is **agnostic PAC learnable** with

$$C_1(d+\log(1/\delta))/\varepsilon^2 < m_H(\varepsilon, \delta) < C_1(d+\log(1/\delta))/\varepsilon^2$$

3. H is **PAC learnable** with

$$C_1(d+\log(1/\delta))/\varepsilon < m_H(\varepsilon, \delta) < C_1(d+\log(1/\delta))/\varepsilon$$

How to compute VC dimension

As a rule of thumb, the VC dimension of a class is often equal to the number of parameters need to be set to specify a specific h in H .

(think of H_{init} , $H_{\text{intervals}}$, $H_{\text{rectangles}}$, HS^n)

Is the story complete now?

- Issue 1 – finite VC classes may be too restricted.
- Issue 2 – computational complexity

Non-Uniform Learn - Definition

A class H is *NonUniformly learnable* if

There is a function $m_H: H \times (0,1)^2 \rightarrow \mathbb{N}$

and a learning algorithm A ,

s.t. for every distribution P over $X \times Y$

and every $\varepsilon, \delta > 0$, for every h in H

for samples S of size $m > m_H(h, \varepsilon, \delta)$

generated i.i.d. by P ,

$$\Pr[L_P(A(S)) > L_P(h) + \varepsilon] < \delta$$

Non-Uni characterization - Statement

Theorem: A class H is NonUniformly learnable if and only if there exist classes

$\{H_n : n \text{ in } \mathbb{N}\}$ such that:

1. Each H_n has the uniform convergence property.

And,

2. $H = \bigcup_{n \text{ in } \mathbb{N}} H_n$

Some NonUni Learnable classes

- The class of all polynomials epi-sets
 $H = \{h_p : p \text{ a polynomial in } x\}$
where, $h_p(x) = 1$ if and only if $p(x) > 0$.
- The class of all (characteristic functions of) finite subsets of (any) X .
- The class of all finite unions of rectangles.

Some classes are not NUL

If H shatters an infinite set, then H is not (even) NonUni learnable.

(in particular, the class of ALL functions over any infinite domain).

Proof of easy direction

Assume H is NonUni learnable,

Define, for every n ,

$$H_n = \{h \text{ in } H : m_H(h, 1/7, 1/8) < n\}$$

Note that each of these classes must have finite VCdim, and therefore has Uniform Convergence, and their union covers H .

Hard direction

Step 1: Weight functions.

We define a *weight function* to be any function $w : \mathbb{N} \rightarrow [0,1]$ such that

$$\sum_n w(n) \leq 1.$$

Examples: $w(n) = 1/2n^2$

or $w(n) = 1/2^n$

Rewriting the m function

Given a class H and a representation of H as a union of H_n 's, each enjoying uniform convergence, define for any n

$$\varepsilon_n(m, \delta) = \min\{\varepsilon: m_{H_n}(\varepsilon, \delta) < m\}$$

(namely, the minimal error that an m -size sample can guarantee)

Hard direction

The NonUniform generalization (loss) bound

For every weight function w , every prob. Dist P
every δ and every m , with probability

$> (1-\delta)$,

For all h in H

$$L_P(h) \leq L_S(h) + \min_{\{n: h \text{ in } H_n\}} \epsilon_n(m, w(n) \delta)$$

Hard direction

The bound minimization algorithm –

Structural Risk Minimization (SRM):

Given H , a decomposition of H to H_n 's of finite VCdim each and a weight function w ,

On a labeled training sample S of size m ,

Find h in H that minimizes the above error bound:

$$L_S(h) + \min_{\{n: h \in H_n\}} \epsilon_n(m, w(n) \delta)$$

Hard direction

The resulting sample complexity function:

$$m(h, \varepsilon, \delta) = m^{\text{uc}}_{H_{n(h)}}(\varepsilon/2, w(n(h)) \delta)$$

Applications of SRM

SRM has many applications, usually referred to as “ERM with regularization”:

Adding to the empirical error a “penalty” on complex (or otherwise, undesirable) h 's.

Example include

1. Norm of a linear classifier
2. Description length
3. Small margins
4. Low prior likelihood

Description length - definition

A description language for a class H is a function

$G: H \rightarrow$ Finite binary strings

Such that the range of G is prefix-free.

Kraft inequality

Any collection T of binary strings that is prefix-free, satisfies

$$\sum_{\sigma \in T} 2^{-|\sigma|} \leq 1$$

Corollary: For $H = \{h_1, h_2, \dots, h_n, \dots\}$,

We can use any description language for H to define weights $w(n) = 2^{-|G(h_n)|}$

Description length bound and Ocam's Razor

The resulting SRM algorithm is :
pick h that minimizes

$$L_S(h) + \text{sqrt}\{(|G(h)| + \ln(1/\delta))/2m\}$$