

# INTRODUCTION TO THE VC-DIMENSION

Wei-Ning Chen

December 28, 2018

- 1 RECAP
- 2 MOTIVATION
- 3 THE VC-DIMENSION
  - Definitions
  - Examples
- 4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY
- 5 ADVANCED TOPICS
  - Glivenko-Cantelli Theorem
  - VC-entropy and Growth Function
- 6 EXERCISES AND DISCUSSION

## DEFINITION (UNIFORM CONVERGENCE)

We say that a hypothesis class  $\mathcal{H}$  has the uniform convergence property (w.r.t. a domain  $Z$  and a loss function  $\ell$ ) if there exists a function  $m_{\mathcal{H}}^{UC}$  such that for every  $\epsilon, \delta \in (0, 1)$  and for every probability distribution  $\mathcal{D}$  over  $Z$ , if  $S$  is a sample of  $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$  examples drawn i.i.d. according to  $\mathcal{D}$ , then

$$\mathcal{P}(|L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon, \forall h \in \mathcal{H}) \geq 1 - \delta$$

Equivalently,

$$\lim_{m \rightarrow \infty} \mathcal{P}(\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon) = 0$$

*Remark:* Compare to the definition of PAC:

$$\mathcal{P}(L_{\mathcal{D}}(h_S) - \inf_{h \in \mathcal{H}} L_{\mathcal{D}}(h) \leq \epsilon) \geq 1 - \delta$$

## THEOREM (NO-FREE-LUNCH)

Let  $\mathcal{A}$  be any learning algorithm for the task of binary classification with respect to the 0 – 1 loss over a domain  $\mathcal{X}$ . Let  $m$  be any number smaller than  $|\mathcal{X}|/2$ , representing a training set size. Then, there exists a distribution  $\mathcal{D}$  over  $\mathcal{X} \times \{0, 1\}$  such that:

- There exists a function  $f : \mathcal{X} \rightarrow \{0, 1\}$  with  $L_{\mathcal{D}}(f) = 0$ .
- $\mathcal{P}(L_{\mathcal{D}}(\mathcal{A}(S)) \geq \frac{1}{8}) \geq \frac{1}{7}$

## 1 RECAP

## 2 MOTIVATION

## 3 THE VC-DIMENSION

- Definitions
- Examples

## 4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY

## 5 ADVANCED TOPICS

- Glivenko-Cantelli Theorem
- VC-entropy and Growth Function

## 6 EXERCISES AND DISCUSSION

In chapter 2, we see that every finite hypothesis class  $\mathcal{H}$  is learnable; moreover, the sample complexity is bounded by

$$m_{\mathcal{H}}(\delta, \epsilon) \leq \frac{\log(|\mathcal{H}|)/\delta}{\epsilon}$$

So, what if  $|\mathcal{H}| = \infty$ ?

## EXAMPLE: CONCENTRIC CIRCLE

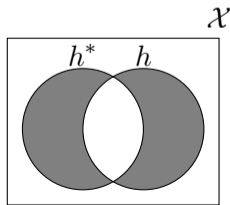
### EXAMPLE

Let  $\mathcal{X} = \mathbb{R}^2$ ,  $\mathcal{Y} = \{0, 1\}$ , and let  $\mathcal{H}$  be the class of concentric circles in the plane, that is,  $\mathcal{H} = \{h_r : r \in \mathbb{R}_+\}$ . Prove that  $\mathcal{H}$  is PAC learnable (assume realizability), and its sample complexity is bounded by

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \frac{\log(2\delta)}{\epsilon}.$$

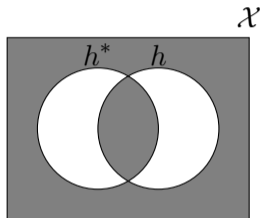
First, we specify  $\mathcal{H}_B$ .

By definition, if  $h \in \mathcal{H}_B$ , we have  $\mathcal{D}(h(x) \neq h^*(x)) \geq \epsilon$



## EXAMPLE: CONCENTRIC CIRCLE

Equivalently,  $\mathcal{D}(h(x) = h^*(x)) \leq 1 - \epsilon$



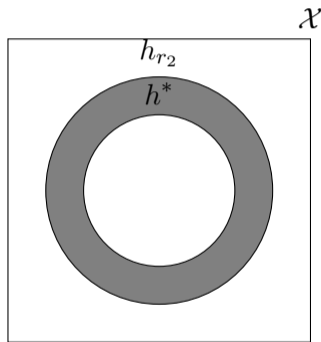
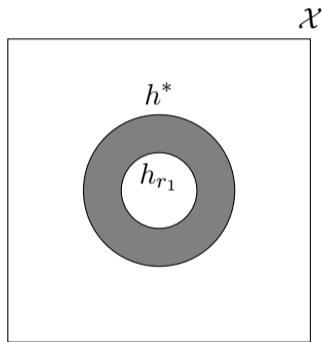
If now  $\mathcal{H}$  is finite, we can apply union bound:

$$\mathcal{D}^m\left(\bigcup_{h \in \mathcal{H}_B} \forall i = [m] | h(x_i) = h^*(x_i)\right) \leq |\mathcal{H}|(1 - \epsilon)^m \leq \delta$$

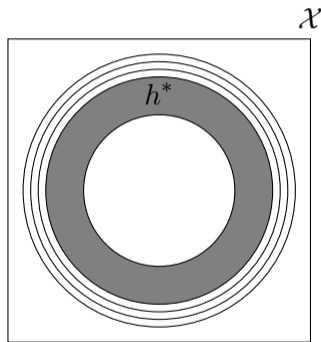
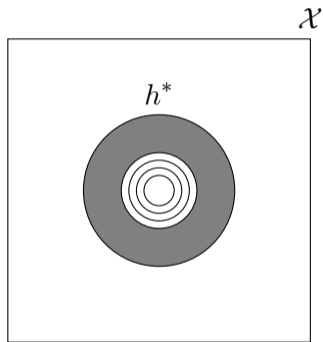
We can slightly modify the union bound for the case  $|\mathcal{H}| = \infty$ .



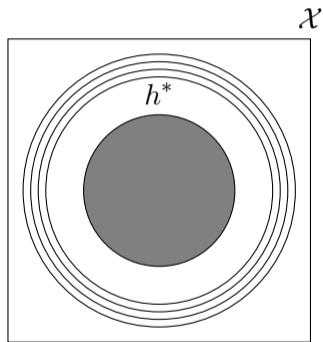
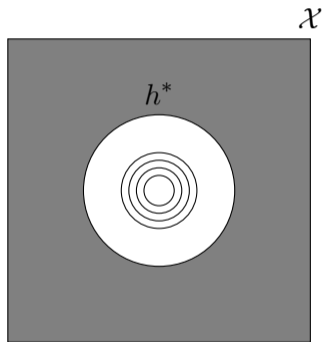
# EXAMPLE: CONCENTRIC CIRCLE



# EXAMPLE: CONCENTRIC CIRCLE



# EXAMPLE: CONCENTRIC CIRCLE



## EXAMPLE: CONCENTRIC CIRCLE

Let  $S \sim \mathcal{D}^m$ , and  $r_{\min} = \min_{x \in S} r_x$ ,  $r_{\max} = \max_{x \in S} r_x$ .

We have

$$\begin{aligned} \mathcal{P}_{S \sim \mathcal{D}^m}(L_{\mathcal{D}}(h_S) \geq \epsilon) &\leq \mathcal{P}_{S \sim \mathcal{D}^m}(r_{\min} \geq r_1 \cup r_{\max} \leq r_0) \\ &\leq \mathcal{P}_{S \sim \mathcal{D}^m}(r_{\min} \geq r_1) + \mathcal{P}_{S \sim \mathcal{D}^m}(r_{\max} \leq r_0) \\ &\leq 2(1 - \epsilon)^m \leq 2e^{-m\epsilon} \leq \delta \end{aligned}$$

Therefore, for all  $\epsilon$  and  $\delta$ , the sample complexity can be bounded by

$$m \leq \frac{\log(2/\delta)}{\epsilon} \quad \square$$

- In chapter 2, we see that every finite hypothesis class  $\mathcal{H}$  is learnable; moreover, the sample complexity is bounded by

$$m_{\mathcal{H}}(\delta, \epsilon) \leq \frac{\log(|\mathcal{H}|)/\delta}{\epsilon}$$

- Also, we see some examples that even the class is infinite-size, it may still be learnable.
- Therefore, we need a measure of  $\mathcal{H}$ 's complexity
- In this chapter, we will formally define the complexity of  $\mathcal{H}$  (VC dimension), and show that

$\mathcal{H}$  has uniform convergence property  $\iff \text{VCdim}(\mathcal{H}) < \infty$

1 RECAP

2 MOTIVATION

3 THE VC-DIMENSION

- Definitions
- Examples

4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY

5 ADVANCED TOPICS

- Glivenko-Cantelli Theorem
- VC-entropy and Growth Function

6 EXERCISES AND DISCUSSION

## DEFINITION (RESTRICTION $\mathcal{H}$ TO $C$ )

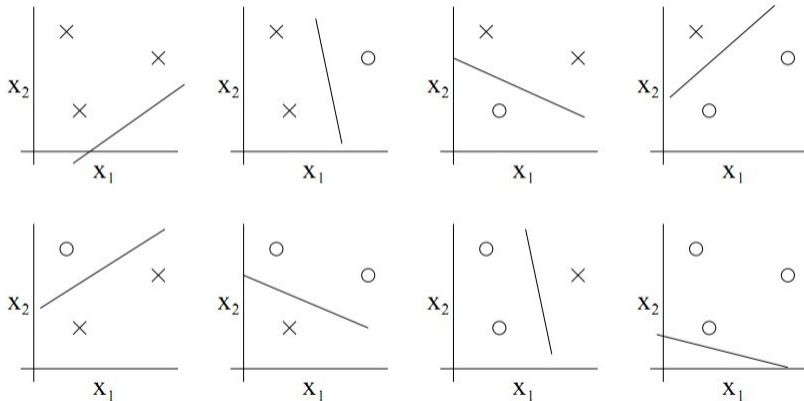
Let  $\mathcal{H}$  be a class of function from  $\mathcal{X}$  to  $\{0, 1\}$  and let  $C = \{c_1, \dots, c_m\} \subset \mathcal{X}$ . The restriction of  $\mathcal{H}$  to  $C$  is the set off all functions from  $C$  to  $\{0, 1\}$  that can be derived from  $\mathcal{H}$ . That is,

$$\mathcal{H}_C = \{h(c_1), \dots, h(c_m)\} : h \in \mathcal{H}$$

## DEFINITION (SHATTERING)

A hypothesis class  $\mathcal{H}$  shatters a finite set  $C$  if the restriction of  $\mathcal{H}$  to  $C$  is the set of all functions from  $C$  to  $\{0, 1\}$ . That is,  $|\mathcal{H}_C| = 2^{|C|}$ .

# THE VC-DIMENSION





## COROLLARY (COROLLARY 6.4)

Let  $\mathcal{H}$  be a hypothesis class of functions from  $\mathcal{X}$  to  $\{0, 1\}$ . Let  $m$  be a training set size. Assume that there exists a set  $C$  of size  $2m$  that is shattered by  $\mathcal{H}$ . Then, for any learning algorithm  $\mathcal{A}$

$$\mathcal{P}(L_{\mathcal{D}}(\mathcal{A}(S)) \geq \frac{1}{8}) \geq \frac{1}{7}$$

*Remark:* This is a direct result from NFL Theorem

*Remark2:* If for all  $m$ , there exists a set  $C$  of size  $2m$  that is shattered by  $\mathcal{H}$ , then  $\mathcal{H}$  is not PAC learnable

## DEFINITION (VC-DIMENSION)

The VC-dimension of a hypothesis class  $\mathcal{H}$ , denoted  $\text{VCdim}(\mathcal{H})$ , is the maximal size of a set  $C \subset \mathcal{X}$  that can be shattered by  $\mathcal{H}$ . If  $\mathcal{H}$  can shatter sets of arbitrarily large size, we say that  $\mathcal{H}$  has infinite VC-dimension.

Remark: If  $\text{VCdim}(\mathcal{H})=d$ , it means that

$\exists C \subset \mathcal{X}$  that can be shattered by  $\mathcal{H}$ ,

NOT

$\forall C \subset \mathcal{X}$  that can be shattered by  $\mathcal{H}$ ,

## EXAMPLE: THRESHOLD FUNCTIONS

- Let  $\mathcal{H}$  be the all threshold function on  $\mathbb{R}$ .
- For an arbitrary set  $C = \{c_1\}$ ,  $\mathcal{H}$  shatters  $C$ , therefore  $\text{VDdim}(\mathcal{H}) \geq 1$ .
- For an arbitrary set  $C = \{c_1, c_2\}$ ,  $\mathcal{H}$  does not shatter  $C$ . Therefore,  $\text{VDdim}(\mathcal{H}) < 2$ .

## EXAMPLE: INTERVALS

- Let  $\mathcal{H}$  be the intervals over  $\mathbb{R}$ ; that is,  $\mathcal{H} = \{\mathbb{1}_{[a,b]}(x) \mid a, b \in \mathbb{R}\}$
- It is easy to show that  $\text{VCdim}(\mathcal{H}) = 2$

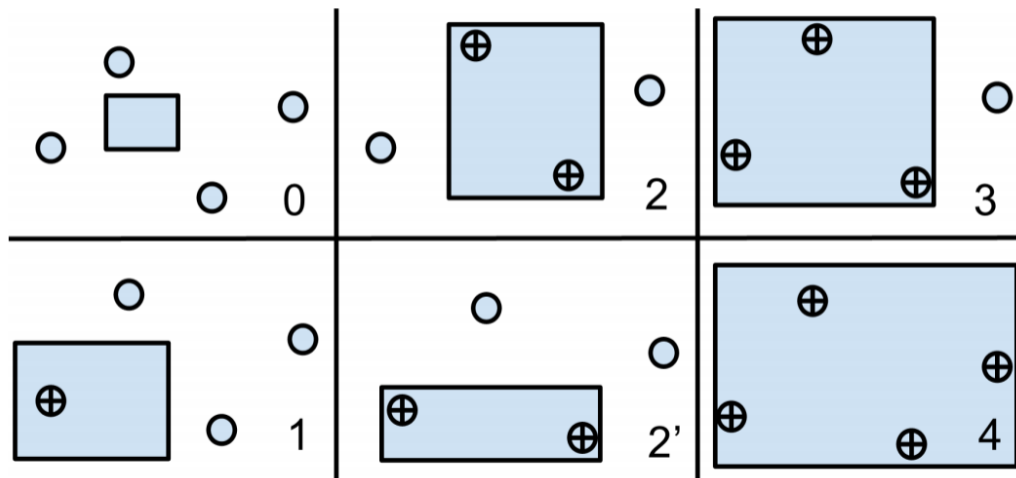
## EXAMPLE: AXIS ALIGNED RECTANGLES

- Let  $\mathcal{H}$  be the the class of axis aligned rectangles:

$$\mathcal{H} = \{\mathbb{1}_{[a,b] \times [c,d]} \mid a, b, c, d \in \mathbb{R}\}$$

- $\text{VDdim}(\mathcal{H}) = 4$

# EXAMPLE: AXIS ALIGNED RECTANGLES



1 RECAP

2 MOTIVATION

3 THE VC-DIMENSION

- Definitions
- Examples

4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY

5 ADVANCED TOPICS

- Glivenko-Cantelli Theorem
- VC-entropy and Growth Function

6 EXERCISES AND DISCUSSION

# THE FUNDAMENTAL THEOREM OF LEARNING THEORY

## THEOREM (THE FUNDAMENTAL THEOREM OF STATISTICAL LEARNING)

Let  $\mathcal{H}$  be a hypothesis class of functions from a domain  $\mathcal{X}$  to  $\{0, 1\}$  and let the loss function be the 0-1 loss. Then, the following are equivalent:

- 1  $\mathcal{H}$  has the uniform convergence property.
- 2 Any ERM rule is a successful agnostic PAC learner for  $\mathcal{H}$ .
- 3  $\mathcal{H}$  is agnostic PAC learnable.
- 4  $\mathcal{H}$  is PAC learnable.
- 5 Any ERM rule is a successful PAC learner for  $\mathcal{H}$ .
- 6  $\mathcal{H}$  has a finite VC-dimension.



# GROWTH FUNCTION AND SAUER'S LEMMA

The growth function measures the maximal "effective" size of  $\mathcal{H}$  on a set of  $m$  examples.

## DEFINITION (GROWTH FUNCTION)

Let  $\mathcal{H}$  be a hypothesis class. Then the growth function of  $\mathcal{H}$  is defined as

$$\tau_{\mathcal{H}}(m) = \sup_{C \subset \mathcal{X}: |C|=m} |\mathcal{H}_C|$$

In words,  $\tau_{\mathcal{H}}(m)$  is the number of different functions from a set  $C$  of size  $m$  to  $\{0, 1\}$  that can be obtained by restricting  $\mathcal{H}$  to  $C$ .

*Remark:* if  $\text{VCdim}(\mathcal{H}) = d$ , then for any  $m \leq d$  we have  $\tau_{\mathcal{H}}(m) = 2^m$ .  
However, what is interesting is the case  $m \geq d$ .

## THEOREM (SAUER'S LEMMA)

Let  $\mathcal{H}$  be a hypothesis with  $\text{VCdim}(\mathcal{H}) = d$ . Then for all  $m$ ,

$$\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i} \leq (em/d)^d$$

*Remark:* if  $\text{VCdim}(\mathcal{H})$  is finite, then the growth function is polynomial in  $m$ .

# UNIFORM CONVERGENCE IN VC CLASS

## THEOREM (UNIFORM CONVERGES IN VC CLASS (THEOREM 6.11))

Let  $\mathcal{H}$  be a class and let  $\tau_{\mathcal{H}}(m)$  be its growth function. Then, for every  $\mathcal{D}$  and every  $\delta$

$$\mathcal{P}_{S \sim \mathcal{D}^m} (|L_{\mathcal{D}}(h) - L_S(h)| > \epsilon) \leq \delta$$

where  $\epsilon$  can be choose as  $\frac{4 + \sqrt{(\log(\tau_{\mathcal{H}}(2m)))}}{\delta\sqrt{2m}}$ . In other words, this theorem tells us that

$$VCdim(\mathcal{H}) < \infty \iff \lim_{\ell \rightarrow \infty} \frac{\log \tau_{\mathcal{H}}(\ell)}{\ell} = 0 \iff \text{uniform convergence property holds.}$$

# THE FUNDAMENTAL THEOREM OF LEARNING THEORY

## THEOREM (THE FUNDAMENTAL THEOREM-QUANTITATIVE VERSION)

Let  $\mathcal{H}$  be a hypothesis class from a domain  $\mathcal{X}$  to  $\{0, 1\}$  and let the loss function be the 0-1 loss. Then, there are absolute constants  $C_1, C_2$  such that:

- 1  $\mathcal{H}$  has the uniform convergence property with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}^{UC} \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

- 2  $\mathcal{H}$  is agnostic PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}} \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

- 3  $\mathcal{H}$  is PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}} \leq C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$$

1 RECAP

2 MOTIVATION

3 THE VC-DIMENSION

- Definitions
- Examples

4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY

5 ADVANCED TOPICS

- Glivenko-Cantelli Theorem
- VC-entropy and Growth Function

6 EXERCISES AND DISCUSSION

# GLIVENKO-CANTELLI THEOREM

## DEFINITION (EMPIRICAL DISTRIBUTION)

Let  $X_1, \dots, X_n$  be i.i.d. random variables in  $\mathbb{R}$  with common cdf  $F(x)$ . The empirical distribution function for  $X_1, \dots, X_n$  is given by

$$F_n(x) = \frac{1}{n} \sum_i \mathbb{1}_{(-\infty, x]}(X_i)$$

## THEOREM (GLIVENKO-CANTELLI THEOREM)

$$\|F_n - F\|_\infty = \sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0 \text{ almost surely}$$

Remark:  $\forall x, F_n(x) \rightarrow F(x)$  trivially by LLN

# GLIVENKO-CANTELLI THEOREM

- More generally, consider a space  $\mathcal{X}$  and  $\sigma$ -field  $\mathcal{F}$  generated by borel set with probability measure  $P$  and empirical measure  $P_n$
- Then  $F(x) = P((-\infty, x])$ , and  $F_n(x) = P_n((-\infty, x])$
- We can rewrite  $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$  as  $\sup_{c \in \mathcal{C}} |P_n(C) - P(C)|$ ,  
where  $\mathcal{C} = \{(-\infty, x] | x \in \mathbb{R}\}$
- What happens for the general  $\mathcal{C}$  ?

# GLIVENKO-CANTELLI CLASS

## DEFINITION (GC-CLASS)

Let  $\mathcal{C} \subset \{C \mid C \text{ measurable in } \mathcal{X}\}$ . Then if

$$\|P_n - P\|_{\mathcal{C}} = \sup_{C \in \mathcal{C}} |P_n(C) - P(C)|$$

we say class  $\mathcal{C}$  is a Glivenko-Cantelli class.

## DEFINITION (UNIFORMLY GC)

A class is called uniformly Glivenko-Cantelli if the convergence occurs uniformly over all probability measures  $\mathcal{P}$  on  $(\mathcal{X}, \mathcal{F})$ :

$$\sup_{P \in \mathcal{P}(S, \mathcal{A})} \mathbb{E} \|P_n - P\|_{\mathcal{C}} \rightarrow 0$$



## DEFINITION (VC CLASS)

A class with finite VC dimension is called a Vapnik-Chervonenkis class or VC class

## THEOREM (VAPNIK AND CHERVONENKIS, 1968)

*A class of sets  $\mathcal{C}$  is uniformly GC if and only if it is a Vapnik-Chervonenkis class*

## 1 RECAP

## 2 MOTIVATION

## 3 THE VC-DIMENSION

- Definitions
- Examples

## 4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY

## 5 ADVANCED TOPICS

- Glivenko-Cantelli Theorem
- VC-entropy and Growth Function

## 6 EXERCISES AND DISCUSSION

# VC ENTROPY AND GROWTH FUNCTION

In previous lecture, we define the growth function as

$$\tau_{\mathcal{H}}(m) = \sup_{C \subset \mathcal{X}: |C|=m} |\mathcal{H}_C|$$

Let's rewrite the growth function as another form:

## DEFINITION (VAPNIK)

Let  $x_1, \dots, x_m$  be  $m$  samples from  $\mathcal{X}$ . Then define the number

$$N^{\mathcal{H}}(x_1, \dots, x_m) = |\{h(x_1), \dots, h(x_m) \mid h \in \mathcal{H}\}| = |\mathcal{H}_{\{x_1, \dots, x_m\}}|$$

Obviously

$$\tau_{\mathcal{H}}(m) = \sup_{C \subset \mathcal{X}: |C|=m} |\mathcal{H}_C| = \sup_{\{x_1, \dots, x_m\} \subset \mathcal{X}} N^{\mathcal{H}}(x_1, \dots, x_m)$$

## VC ENTROPY AND GROWTH FUNCTION

The supremum is taken so that the bound works even in the worst distribution. In general case, we can replace supremum by expectation w.r.t. a specific distribution, which gives another value:

### DEFINITION (VC-ENTROPY, ANNEALED VC-ENTROPY, GROWTH FUNCTION)

Let  $N^{\mathcal{H}}(x_1, \dots, x_m)$  be defined as previous. Then we defined VC-entropy as

$$H^{\mathcal{H}}(m) = \mathbb{E} \log N^{\mathcal{H}}(x_1, \dots, x_m)$$

The annealed VC-entropy as

$$H_{ann}^{\mathcal{H}}(m) = \log \mathbb{E} N^{\mathcal{H}}(x_1, \dots, x_m)$$

And the growth function (with logarithm) as

$$G^{\mathcal{H}}(m) = \log \sup_{x_1, \dots, x_m} N^{\mathcal{H}}(x_1, \dots, x_m) (= \log \tau_{\mathcal{H}}(m))$$

## COROLLARY

$$H^{\mathcal{H}}(m) \leq H_{ann}^{\mathcal{H}}(m) \leq G^{\mathcal{H}}(m)$$

*Rmark:* The fundamental theorem of learning theory tells us

$$\lim_{m \rightarrow \infty} \frac{G^{\mathcal{H}}(m)}{m} = \lim_{m \rightarrow \infty} \frac{\log \tau_{\mathcal{H}}(m)}{m} = 0 \iff \lim_{m \rightarrow \infty} \mathcal{P}(\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon) = 0$$

The convergence is uniform for all distribution  $P$  in  $(\mathcal{X}, \mathcal{F})$

# THREE MILESTONES OF LEARNING THEORY

## THEOREM (VAPNIK)

1

$$\lim_{m \rightarrow \infty} \mathcal{D}(\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon) = 0$$

*is sufficient and necessary that*

$$\lim_{m \rightarrow \infty} \frac{H^{\mathcal{H}}(m)}{m} = 0$$

2

$$\lim_{m \rightarrow \infty} \mathcal{D}(\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon) \leq e^{-c\epsilon^2 m} \text{ (fast decay)}$$

*is sufficient if*

$$\lim_{m \rightarrow \infty} \frac{H_{ann}^{\mathcal{H}}(m)}{m} = 0$$

## THEOREM (VAPNIK)

3

$$\lim_{m \rightarrow \infty} P(\sup_{h \in \mathcal{H}} |L_S(h) - L_D(h)| > \epsilon) = 0, \text{ for all } P \in (\mathcal{X}, \mathcal{F})$$

*is sufficient and necessary that*

$$\lim_{m \rightarrow \infty} \frac{G^{\mathcal{H}}(m)}{m} = 0$$

## 1 RECAP

## 2 MOTIVATION

## 3 THE VC-DIMENSION

- Definitions
- Examples

## 4 THE FUNDAMENTAL THEOREM OF LEARNING THEORY

## 5 ADVANCED TOPICS

- Glivenko-Cantelli Theorem
- VC-entropy and Growth Function

## 6 EXERCISES AND DISCUSSION