

# Statistical Learning Theory

## Consistency and bounds on the rate of convergence for ERM methods

Miguel A. Vezanzones

Grupo Inteligencia Computacional  
Universidad del País Vasco

# Outline

- 1 Introduction
- 2 Consistency
  - Introduction
  - VC entropy
  - Necessary and sufficient conditions for uniform convergence
- 3 Theory of non-falsifiability
- 4 Bounds on the rate of convergence
  - Three milestones of learning theory

# Consistency of learning processes

- Consistency: convergence in probability to the best possible result.
- Consistency of learning processes:
  - To explain when a learning machine that minimizes empirical risk can achieve a small value of actual risk (to generalize) and when it can not.
  - Equivalently, to describe necessary and sufficient conditions for the consistency of learning processes that minimize the empirical risk.
- This guarantees that the constructed theory is general and cannot be improved from the conceptual point of view.

# Theory of non-falsifiability

- Kant's problem of demarcation (s. XVIII): is there a formal way to distinguish true theories from false theories?
  - One of the main questions of modern philosophy.
- Popper's theory of non-falsifiability (s. XX): criterion for demarcation between true and false theories.
- Strongly related to what happens if the ERM method is not consistent.

# Bounds on the rate of convergence

- It is required for any machine minimizing empirical risk to satisfy consistency conditions.
- But, consistency conditions say nothing about the rate of convergence of the obtained risk  $R(\alpha_l)$  to the minimal one  $R(\alpha_0)$ .
- It is possible to construct examples where the ERM principle is consistent, but where the risks have an arbitrary slow asymptotic rate of convergence.
- The theory of bounds on the rate of convergence tries to answer the following question:
  - Under what conditions is the asymptotic rate of convergence fast?





# Notation

- Let  $Q(\mathbf{z}, \alpha_l)$  be a function that minimizes the empirical risk functional

$$R_{emp} = \frac{1}{l} \sum_{i=1}^l Q(\mathbf{z}_i, \alpha)$$

for a given set of i.i.d. observations  $\mathbf{z}_1, \dots, \mathbf{z}_l$ .











# ERM consistency

- In order to create a theory of consistency of the ERM method depending only on the general properties (capacity) of the set of functions, a consistency definition excluding trivial consistency cases is needed.
- This is done by non-trivial consistency definition.

## Non-trivial consistency

- The ERM principle is nontrivially consistent for the set of functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , and the probability distribution function  $F(\mathbf{z})$  if for any nonempty subset  $\Lambda(c)$ ,  $c \in (-\infty, \infty)$  defined as

$$\Lambda(c) = \left\{ \alpha : \int Q(\mathbf{z}, \alpha) dF(\mathbf{z}) > c, \quad \alpha \in \Lambda \right\}$$

the convergence

$$\inf_{\alpha \in \Lambda(c)} R_{emp}(\alpha) \xrightarrow{P} \inf_{\alpha \in \Lambda(c)} R(\alpha) \quad (3)$$

is valid.

# Key theorem of learning theory

- Vapnik and Chervonenkins, 1989.

## Theorem

Let  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , be a set of functions that satisfy the condition

$$A \leq \int Q(\mathbf{z}, \alpha) dF(\mathbf{z}) \leq B \quad (A \leq R(\alpha) \leq B)$$

then for the ERM principle to be consistent, it is necessary and sufficient that the empirical risk  $R_{emp}(\alpha)$  converges uniformly to the actual risk  $R(\alpha)$  over the set  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , in the following sense:

$$\lim_{l \rightarrow \infty} P \left\{ \sup_{\alpha \in \Lambda} (R(\alpha) - R_{emp}(\alpha)) > \varepsilon \right\} = 0, \quad \forall \varepsilon > 0 \quad (4)$$

# Consistency of the ERM principle

- According to the key theorem, the uniform one-sided convergence (4) is a necessary and sufficient condition for (non-trivial) consistency of the ERM method.
- Conceptually, the conditions for consistency of the ERM principle are necessarily and sufficiently determined by the “worst” function of the set of functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ .

# Outline

- 1 Introduction
- 2 Consistency
  - Introduction
  - VC entropy
  - Necessary and sufficient conditions for uniform convergence
- 3 Theory of non-falsifiability
- 4 Bounds on the rate of convergence
  - Three milestones of learning theory

# Introduction

- The key theorem expresses that consistency of the ERM principle is equivalent to existence of uniform one-sided convergence.
- Conditions for uniform two-sided convergence play an important role in constructing conditions for uniform two-sided convergence.
- Necessary and sufficient conditions for both uniform one-sided and two-sided convergence are obtained on the basis of the VC entropy concept.

# Empirical process

- An empirical process is a stochastic process in the form of a sequence of random variables

$$\xi^l = \sup_{\alpha \in \Lambda} \left| \int Q(\mathbf{z}, \alpha) dF(\mathbf{z}) - \frac{1}{l} \sum_{i=1}^l Q(\mathbf{z}_i, \alpha) \right|, \quad l = 1, 2, \dots \quad (5)$$

that depend on both, the probability measure  $F(\mathbf{z})$  and the set of functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ .

- The problem is to describe conditions under which this empirical process converges in probability to zero.

# Consistency of an empirical process

- The necessary and sufficient conditions for an empirical process to converge in probability to zero imply that the equality

$$\lim_{l \rightarrow \infty} P \left\{ \sup_{\alpha \in \Lambda} \left| \int Q(\mathbf{z}, \alpha) dF(\mathbf{z}) - \frac{1}{l} \sum_{i=1}^l Q(\mathbf{z}_i, \alpha) \right| > \varepsilon \right\} = 0, \quad \forall \varepsilon > 0 \quad (6)$$

holds true.

# Law of large numbers and its generalization

- If the set of functions contains only one element, then the sequence of random variables  $\xi^l$  always converges in probability to zero: law of large numbers.
- Generalization of the law of large numbers for the case where a set of functions has a finite number of elements:

## Definition

The sequence of random variables  $\xi^l$  converges in probability to zero if the set of functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , contains a finite number  $N$  of elements.

# Law of large numbers and its generalization

- When  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , has an infinite number of elements, the sequence of random variables  $\xi^l$  does not necessarily converges in probability to zero.
- Problem of the existence of a law of large numbers in functional space (uniform two-sided convergence of the means to their probabilities): generalization of the classical law of large numbers.

# Entropy

- Necessary and sufficient conditions for both uniform one-sided convergence and uniform two-sided convergence are obtained on the basis of a concept called *the entropy of a set of functions*  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , for a sample of size  $l$ .

# Entropy of the set of indicator functions

## Diversity

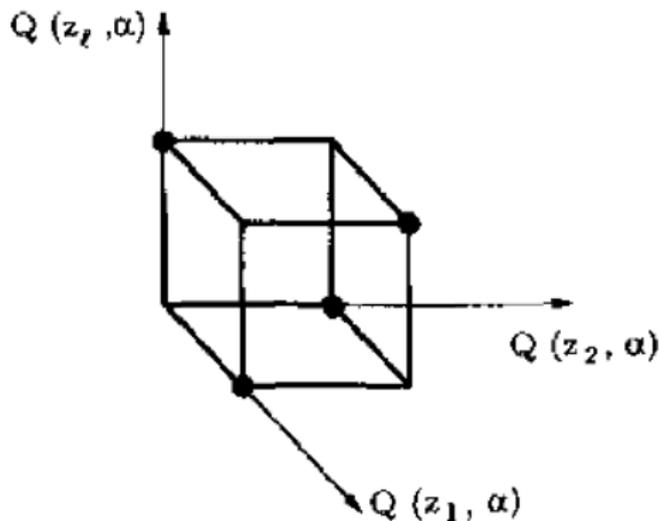
- Lets characterize the *diversity* of a set of indicator functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , on the given set of data by the quantity  $N^{\wedge}(\mathbf{z}_1, \dots, \mathbf{z}_l)$  that evaluates how many different separations of the given sample can be clone using functions from the set of indicator functions.
- Consider the set of  $l$ -dimensional binary vectors:

$$q(\alpha) = (Q(\mathbf{z}_1, \alpha), \dots, Q(\mathbf{z}_l, \alpha)), \quad \alpha \in \Lambda$$

Geometrically, the diversity is the number of different vertices of the  $l$ -dimensional cube that can be obtained on the basis of the sample  $\mathbf{z}_1, \dots, \mathbf{z}_l$  and the set of functions.

# Entropy of the set of indicator functions

Diversity (geometrics)



**Figure:** The set of  $l$ -dimensional binary vectors  $q(\alpha)$ ,  $\alpha \in \Lambda$ , is a subset of the set of vertices of the  $l$ -dimensional unit cube.

# Entropy of the set of indicator functions

## Random entropy and entropy

- The random entropy

$$H^{\wedge}(\mathbf{z}_1, \dots, \mathbf{z}_l) = \ln N^{\wedge}(\mathbf{z}_1, \dots, \mathbf{z}_l)$$

describes the diversity of the set of functions on the given data.

- The expectation of the random entropy over the joint distribution function  $F(\mathbf{z}_1, \dots, \mathbf{z}_l)$ :

$$H^{\wedge}(l) = E[\ln N^{\wedge}(\mathbf{z}_1, \dots, \mathbf{z}_l)] \quad (7)$$

is the *entropy* of the set or indicator functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , on samples of size  $l$ .

# Entropy of the set of real functions

## Diversity

- Let  $A \leq Q(\mathbf{z}, \alpha) \leq B$ ,  $\alpha \in \Lambda$ , a set of bounded loss functions.
- Considering this set of functions and the training set  $\mathbf{z}_1, \dots, \mathbf{z}_l$  one can construct the following set of  $l$ -dimensional vectors:

$$q(\alpha) = (Q(\mathbf{z}_1, \alpha), \dots, Q(\mathbf{z}_l, \alpha)), \quad \alpha \in \Lambda$$

- The diversity,  $N = N^\wedge(\varepsilon, \mathbf{z}_1, \dots, \mathbf{z}_l)$ , indicates the number of elements of the minimal  $\varepsilon$ -net of this set of vectors  $q(\alpha)$ ,  $\alpha \in \Lambda$ .

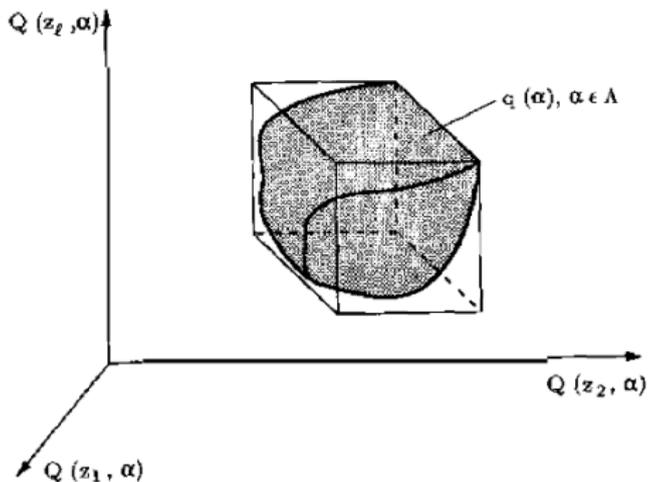
# Entropy of the set of real functions

## Minimal $\varepsilon$ -net

- The set of vectors  $q(\alpha)$ ,  $\alpha \in \Lambda$ , has a minimal  $\varepsilon$ -net  $q(\alpha_1), \dots, q(\alpha_N)$  if:
  - 1 There exist  $N = N^{\wedge}(\varepsilon, \mathbf{z}_1, \dots, \mathbf{z}_l)$  vectors  $q(\alpha_1), \dots, q(\alpha_N)$  such that for any vector  $q(\alpha^*)$ ,  $\alpha^* \in \Lambda$ , one can find among these  $N$  vectors one  $q(\alpha_r)$  that is  $\varepsilon$ -close to  $q(\alpha^*)$  in a given metric.
  - 2  $N$  is the minimum number of vectors that possesses this property.

# Entropy of the set of real functions

Diversity (geometrics)



**Figure:** The set of  $l$ -dimensional vectors  $q(\alpha)$ ,  $\alpha \in \Lambda$ , belongs to an  $l$ -dimensional cube.

# Entropy of the set of real functions

## Random entropy and entropy

- The random VC entropy of the set of functions  $A \leq Q(\mathbf{z}, \alpha) \leq B$ ,  $\alpha \in \Lambda$ , on the sample  $\mathbf{z}_1, \dots, \mathbf{z}_l$  is given by:

$$H^{\wedge}(\varepsilon; \mathbf{z}_1, \dots, \mathbf{z}_l) = \ln N^{\wedge}(\varepsilon; \mathbf{z}_1, \dots, \mathbf{z}_l)$$

- The expectation of the random VC entropy over the joint distribution function  $F(\mathbf{z}_1, \dots, \mathbf{z}_l)$ :

$$H^{\wedge}(\varepsilon; l) = E[\ln N^{\wedge}(\varepsilon; \mathbf{z}_1, \dots, \mathbf{z}_l)]$$

is the *VC entropy* of the set of real functions  $A \leq Q(\mathbf{z}, \alpha) \leq B$ ,  $\alpha \in \Lambda$ , on samples of size  $l$ .

# Outline

- 1 Introduction
- 2 Consistency
  - Introduction
  - VC entropy
  - Necessary and sufficient conditions for uniform convergence
- 3 Theory of non-falsifiability
- 4 Bounds on the rate of convergence
  - Three milestones of learning theory

# Conditions for uniform two-sided convergence

## Theorem

*Under some conditions of measurability on the set of real bounded functions  $A \leq Q(\mathbf{z}, \alpha) \leq B$ ,  $\alpha \in \Lambda$ , for uniform two-sided convergence it is necessary and sufficient that the equality*

$$\lim_{l \rightarrow \infty} \frac{H^{\wedge}(\varepsilon; l)}{l} = 0, \quad \forall \varepsilon > 0 \quad (8)$$

*be valid.*

# Conditions for uniform two-sided convergence

## Corollary

### Corollary

*Under some conditions of measurability on the set of indicator functions  $Q(\mathbf{z}, \alpha)$ ,  $\alpha \in \Lambda$ , for uniform two-sided convergence it is necessary and sufficient that*

$$\lim_{l \rightarrow \infty} \frac{H^{\wedge}(l)}{l} = 0$$

*which is a particular case of (8).*

# Uniform one-sided convergence

- Uniform two-sided convergence can be described as

$$\lim_{l \rightarrow \infty} P \left\{ \left[ \sup_{\alpha} (R(\alpha) - R_{emp}(\alpha)) \right] \vee \left[ \sup_{\alpha} (R_{emp}(\alpha) - R(\alpha)) \right] \right\} = 0 \quad (9)$$

which includes uniform one-sided convergence, and it's sufficient condition for ERM consistency.

- But for consistency of ERM principle, left-hand side of (9) can be violated.





# Conditions for uniform one-sided convergence

## Theorem

*Under some conditions of measurability on the set of real bounded functions  $A \leq Q(\mathbf{z}, \alpha) \leq B$ ,  $\alpha \in \Lambda$ , for uniform one-sided convergence it is necessary and sufficient that for any positive  $\delta$ ,  $\eta$  and  $\varepsilon$  there exist a set of functions  $Q^*(\mathbf{z}, \alpha^*)$ ,  $\alpha^* \in \Lambda^*$ , satisfying (10) such that the following holds:*

$$\lim_{l \rightarrow \infty} \frac{H^{\wedge}(\varepsilon; l)}{l} < \eta \quad (11)$$



# Outline

- 1 Introduction
- 2 Consistency
  - Introduction
  - VC entropy
  - Necessary and sufficient conditions for uniform convergence
- 3 Theory of non-falsifiability
- 4 Bounds on the rate of convergence
  - Three milestones of learning theory

# For Further Reading



The Nature of Statistical Learning Theory. Vladimir N. Vapnik. ISBN: 0-387-98780-0. 1995.



Statistical Learning Theory. Vladimir N. Vapnik. ISBN: 0-471-03003-1. 1998.

# Questions?

*Thank you very much for your attention.*

- Contact:
  - Miguel Angel Veganzones
  - Grupo Inteligencia Computacional
  - Universidad del País Vasco - UPV/EHU (Spain)
  - E-mail: miguelangel.veganzones@ehu.es
  - Web page: <http://www.ehu.es/computationalintelligence>