

CS group meeting 02/10/2017

-1-

The Nature of Statistical Learning Theory V. Vapnik.

## Lecture 2

Two-sided empirical process

$$\xi^l = \sup_{\alpha \in \Lambda} \left| \int Q(z, \alpha) dF(z) - \frac{1}{l} \sum_{i=1}^l Q(z_i, \alpha) \right|$$

Uniform two-sided convergence  $\quad l = 1, 2, \dots$

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

$\forall \varepsilon > 0$

One-sided empirical process

$$\xi_+^l = \sup_{\alpha \in \Lambda} \left( \int Q(z, \alpha) dF(z) - \frac{1}{l} \sum_{i=1}^l Q(z_i, \alpha) \right)_+$$

Uniform one-sided convergence

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

$\forall \varepsilon > 0$

# "Law of Large Numbers"

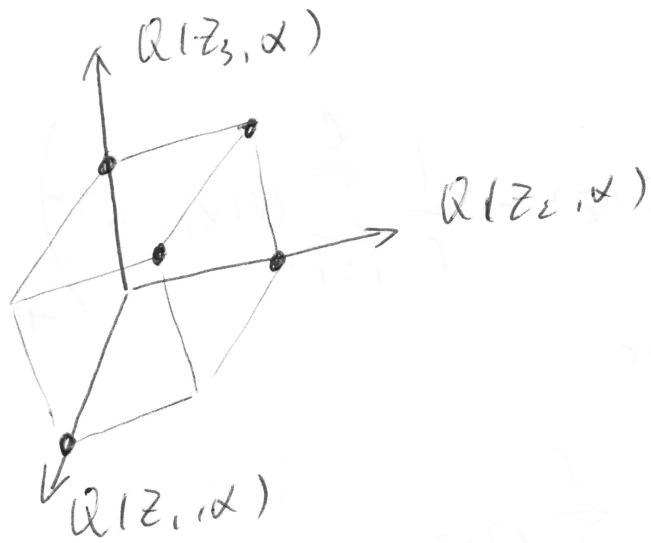
Entropy of the set of functions  $Q(z, \alpha)$ ,  $\alpha \in \Lambda$   
on a sample of size  $l$ .

- The case of a set of indicator functions

If  $Q(z, \alpha) \subset \Lambda$  is a set of indicator functions  
and  $z_1, \dots, z_l$  is a training sample

$$q(\alpha) = (Q(z_1, \alpha), \dots, Q(z_l, \alpha)) \in \Lambda$$

will be a binary vector that belongs to the  
vertices of an  $l$ -dimensional cube



$N^{\Lambda}(z_1, \dots, z_l)$  is the  
number of vertices that  
are different and can  
be obtained by the training  
sample  $z_1, \dots, z_l$  and  
the set of functions

$$Q(z, \alpha), \alpha \in \Lambda$$

Random entropy  $H^A(z_1, \dots, z_l) = \ln N^A(z_1, \dots, z_l)$

Entropy of the set of indicator functions  $\mathcal{Q}(z, \alpha)$   
 $\alpha \in \Lambda$  on samples of size  $l$

$$H^A(\ell) = \mathbb{E} \ln N^A(z_1, \dots, z_l).$$

- The general case and VC-entropy

Definition Let  $A \subseteq \mathcal{Q}(z, \alpha) \subseteq B$ ,  $\alpha \in \Lambda$  be  
 a set of bounded loss functions. Using this  
 set of functions and the training set  $z_1, \dots, z_l$   
 one can construct the following set of  $l$ -dimensional  
 vectors

$$q(\alpha) = (\mathcal{Q}(z_1, \alpha), \dots, \mathcal{Q}(z_l, \alpha)) \quad \alpha \in \Lambda$$

This set of vectors belongs to the  $l$ -dimensional  
 cube and has a finite minimal  $\varepsilon$ -net in the metric  
 $C^0$  (or  $L_p$ ). Let  $N = N^A(\varepsilon; z_1, \dots, z_l)$  be the  
 number of elements of the minimal  $\varepsilon$ -net of this  
 set of vectors  $q(\alpha)$ ,  $\alpha \in \Lambda$ .

## Random VC entropy

$$H^A(\varepsilon; z_1, \dots, z_e) = \ln N^A(\varepsilon; z_1, \dots, z_e)$$

## VC entropy

$$H^A(\varepsilon; \ell) = \bar{E} H^A(\varepsilon; z_1, \dots, z_e)$$

Conditions for uniform two-sided convergence

Theorem For uniform two-sided convergence it is necessary and sufficient that

$$\lim_{\ell \rightarrow \infty} \frac{H^A(\varepsilon; \ell)}{\ell} = 0 \quad \forall \varepsilon > 0$$

Conditions for uniform one-sided convergence

Theorem In order for uniform one-sided convergence of empirical means to their expectations to hold for the set of totally bounded functions  $Q(z, \alpha)$ , it is necessary and sufficient that for any positive  $S, \eta$  and  $\varepsilon$  there exists a set of functions  $Q^*(z, \alpha^*)$ ,  $\alpha^* \in \Lambda^*$  satisfying  $Q(z, \alpha) - Q^*(z, \alpha^*) \geq 0 \quad \forall z$

$$\int (Q(z, \alpha) - Q^*(z, \alpha^*)) dF(z) \leq S$$

such that

$$\lim_{l \rightarrow \infty} \frac{H^{\wedge^*}(\varepsilon, l)}{l} < \eta$$

holds for the  $\varepsilon$ -entropy of the set  $Q^*(z, \alpha)$ ,  $\alpha \in \Lambda^*$   
on samples of size  $l$ .

Question What if  $\lim_{l \rightarrow \infty} \frac{H^{\wedge}(\varepsilon_0, l)}{l} \neq 0$  for some  $\varepsilon_0$ ?

ERM does not hold  $\Rightarrow$  non-consistent

But Why non-consistent?

Theory of Non-falsifiability in Philosophy

Kant's problem of Demarcation

(i). Deductive : General  $\rightarrow$  Particular

(ii). Inductive : Particular  $\rightarrow$  General

Kant's problem: What is the difference between the cases of a justified inductive step, and those for which the inductive step is not justified?

All Natural Science is a result of inductive inference.  
So Kant's problem = Is there a formal way to distinguish true theories and false theories?

## Popper's solution

A necessary condition for justifiability of a theory is the feasibility of its falsification.

falsification = existence of a collection of particular assertions which cannot be explained by the given theory although they fall into its domain

If a given theory can be falsified  $\rightarrow$  scientific

If a given theory cannot be falsified  $\rightarrow$  non-scientific

### ▷ Complete Nonfalsifiability

We know for indicator functions

$$H^A(l) = \mathbb{E} \ln N^A(z_1, \dots, z_l), \quad N^A(z_1, \dots, z_l) \leq 2^l$$

If  $\lim_{l \rightarrow \infty} \frac{H^A(l)}{l} = \ln z$  (maximal entropy)

Then for almost all samples  $z_1, \dots, z_l$  we have

$$N^A(z_1, \dots, z_l) = 2^l$$

That is, almost all sample  $z_1, \dots, z_l$  of arbitrary size  $l$  can be separated in all possible ways by functions

of this set, so minimum of empirical risk for this machine is 0

Such a machine can give a general explanation for almost any data.

Non-falsifiable machine  $\rightarrow$  Not scicific  
 Minimal value of the empirical risk is 0 independent  
 of the expected risk.

That is why we need variance and bias!

▷ Partial Non-falsifiability

Theorem For the set of indicator functions

$a(z, \alpha)$ ,  $\alpha \in \Lambda$  let the convergence

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

be valid. Then there exists a subset  $Z^*$  of the set  $Z$  for which  $P(Z^*) = a(c) \neq 0$

and for almost any training set  $z_1, \dots, z_e$

we have  $z_1^*, \dots, z_k^* = \{z_1, \dots, z_e\} \cap Z^*$

and any sequence of binary values  $s_1, \dots, s_k$   $s_i \in \{0, 1\}$

- there exists a function  $Q(z, \alpha^*)$  for which  $\hat{Q}^8$

$$s_i = Q(z_i^*, \alpha^*) \quad i=1, 2, \dots, k$$

holds true.

So if the conditions of uniform two-sided convergence fail, then there exists some subspace of the input space where the learning machine is non-falsifiable.

▷ There is also a more sophisticated "potentially non falsifiable" related to general VC entropy that I will skip here.

Three milestones in learning theory.

Consider the case of indicator functions  $Q(z, \alpha)$

VC entropy  $H^A(l) = \overline{\mathbb{E}} \ln N^A(z_1, \dots, z_l)$

Annealed VC entropy  $H_{\text{ann}}^A(l) = \ln \overline{\mathbb{E}} N^A(z_1, \dots, z_l)$

Growth function  $G^A(l) = \ln \sup_{z_1, \dots, z_l} N^A(z_1, \dots, z_l)$

We have  $H^A(l) \leq H_{\text{ann}}^A(l) \leq G^A(l)$

First milestone in learning theory

$$\lim_{l \rightarrow \infty} \frac{H^{\wedge}(l)}{l} = 0 \Rightarrow \text{Consistency of the ERM principle}$$

Second milestone in learning theory

$$\lim_{l \rightarrow \infty} \frac{H_{\text{ann}}^{\wedge}(l)}{l} = 0 \Rightarrow \text{fast convergence of } R(\alpha_l) \text{ towards } R(\alpha_0)$$

$$\mathbb{P}(R(\alpha_l) - R(\alpha_0) > \varepsilon) < e^{-c\varepsilon^2 l} \quad (c > 0)$$

Third milestone in learning theory

$$\lim_{l \rightarrow \infty} \frac{G^{\wedge}(l)}{l} = 0 \Rightarrow \begin{aligned} &\text{fast convergence of } \\ &R(\alpha_l) \text{ towards } R(\alpha_0) \\ &\text{independent of the choice of} \\ &\text{the probability measure.} \end{aligned}$$