



národní
úložiště
šedé
literatury

Probabilistic learning model PAC - lecture notes

Hakl, František
2015

Dostupný z <http://www.nusl.cz/ntk/nusl-364639>

Dílo je chráněno podle autorského zákona č. 121/2000 Sb.

Tento dokument byl stažen z Národního úložiště šedé literatury (NUŠL).

Datum stažení: 28.04.2024

Další dokumenty můžete najít prostřednictvím vyhledávacího rozhraní [nusl.cz](http://www.nusl.cz) .



Institute of Computer Science
The Czech Academy of Sciences

Probabilistic learning model PAC — lecture notes

František Hák

Technical report No. 1227

December 2015

Abstrakt:

This report summarises basic theoretical properties of Valiant's distribution-independent learnability PAC model which determines a probabilistic framework of supervised learning. The form of this text is self-contained and provides a detailed proves of the concerned theorems and lemmas. Essentials of VC-dimension of class sets are also included. The text presented is aimed as a study material for MSc. and PhD. students such as for computer science engineers.

Keywords:

Supervised learning, PAC model, Vapnik-Chervonenkis dimension, learning algorithms

**Probabilistic
learning
models**

František Hák

*lecture notes
Faculty of Nuclear Science and Physical Engineering
Czech Technical University, Prague
December 2015*

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

This research was supported by the Czech Science Foundation, Contract №15-18108S .

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

Contents

	Symbols, notation	6
1	PAC learning model	9
1.1	Concepts and concept classes	9
1.2	Basic Terms of Learning Algorithms	11
2	Vapnik-Chervonenkis dimension	17
2.1	VC-dimension of concept class	17
2.1.1	General properties of VC-dimension	19
2.1.2	VC-dimension of union and intersection	23
2.2	VC-dimension of linear concepts	28
2.2.1	Application of Cover’s lemma	31
2.3	VC-dimension of composed mapping	35
2.4	VC-dimension of symmetric difference	41
3	Sample Complexity and VC–dimension	45
3.1	Estimate of the Number of Samples	45
3.1.1	Delta rule learning algorithm	63
3.1.2	Lower bound for maximal steps of delta rule algorithm.	65
3.1.3	Linear separation and linear programming	69
4	Appendices	71
4.1	Source codes	71
4.1.1	Python code for $\omega(\epsilon)$ solver	71

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

List of Figures

1.1	Learning scenario.	9
1.2	Example of two hypotheses with the same error.	12
2.1	Corresponding graph for composed mapping.	36
3.1	Graph and values of the solution of $\omega(\epsilon)$	55
3.2	Nontrivial concept classes.	58
3.3	Construction of $2^k + 1$ dichotomies derived from a fixed dichotomy of $\{-1, +1\}^k$ (for $k = 2$).	66

Symbols and notation

The following fonts and symbols, everywhere in the text of this booklet have the meaning described below:

$\stackrel{\text{def}}{=}$	definition of newly introduced object, set or number
$\stackrel{2}{\Rightarrow}$	implication from equation, etc.
$\stackrel{4}{=}$	equality follows from 4
$\stackrel{5}{\neq}$	nonequality follows from 5
$\stackrel{\text{xx}}{\wedge}$	follows from xx
$\stackrel{\text{rs}}{\wedge}$	follows from rs
$\stackrel{\text{xx}}{\lessgtr}$	less than follows from xx
$\stackrel{\text{xx}}{\gtrless}$	great than follows from xx
$K = \{1, \dots, m\}$	set of natural numbers between 1 and m
B_α^N	ball in \mathfrak{R}^N space centered at initial point and with radius α
\mathfrak{R}^N	real N -dimensional space
Ω^m	m -tuple Cartesian product of set Ω
z	element in Cartesian (multiple) product of specified set
\bar{X}	general (sub)set of a given set
\mathcal{X}	system of subsets of a given set
\mathcal{X}	set of systems of sets over a given set
$2^{\bar{X}}$	potential set of \bar{X} , $2^{\bar{X}} \stackrel{\text{def}}{=} \{\bar{Z} \mid \bar{Z} \subset \bar{X}\}$
$[\Omega]_\lambda$	linear hull of set Ω of vectors from linear vector space (=set of all linear combination of finite set of vectors from Ω)
\bar{B}^\perp	orthogonal complement of the set \bar{B}
$[\bar{B}]_\kappa$	convex hull of set \bar{B}
\mathcal{H}_n^*	vektor space (with specified parameters)
$C_{\bar{A}}$	set of continuous functions defined on a given set \bar{A}
\mathbf{NP}_k^α	set of object with predefined properties, generally set of functions, mappings, etc., context dependent
B^*	general learning algorithm
$\bar{A} \Delta \bar{B}$	symetric difference of two sets \bar{A} and \bar{B}
$\bar{A} - \bar{B}$	difference of sets \bar{A} and \bar{B}
$o(\tilde{f})$	$\tilde{g} = o(\tilde{f})$ means that absolute value of function \tilde{g} is on asymptotic neighbourhood of zero majorized by function \tilde{f}
$O(\tilde{f})$	$\tilde{g} = O(\tilde{f})$ means that absolute value of function \tilde{g} is on asymptotic neighbourhood of infinity majorised by values of \tilde{f}
\mathbf{A}	matrix of real numbers, columns or rows will be denoted by notion LINES OF MATRIX
\vec{x}	vector from a given vector space (all vectors will be considered as columns)

\vec{x}_i	i-th indice of the vector)
$\mathbf{A}^{(n)}$	square matrix of order 2^n
$\vec{x}^{(n)}$	vector of dimension 2^n
$\langle \vec{x} \vec{y} \rangle$	scalar product of vectors, $\langle \vec{x} \vec{y} \rangle = \sum_{i=1}^n \vec{x}_i \vec{y}_i$
$\vec{x} \odot \vec{y}$	tensor product of vectors or matrices
$\overset{\wedge}{\text{AND}}$	binary function AND
$\overset{\vee}{\text{OR}}$	binary function OR
$\overset{\sqcup}{\text{XOR}}$	binary function XOR
$\overset{\neg}{\text{NOT}}$	unary function NOT
$[M]$	greatest whole number less than M (=whole bellow part)
$\lceil M \rceil$	smallest whole number greater than M (=whole upper part)
$vpm(j)$	$vpm(j)$ is vector from $\{-1, +1\}^n$, (resp. $\{0, 1\}^n$, by context), whose coordinates correspond with binary inscription of natural number j (coordinates of vector $vpm(j)$ correspond to pozitions of 1 in binary expansion of j , are equal to 1).
$vbini(j)$	$vbini(j)$ is vector from 0,1, whose coordinates correspond with binary inscription of natural number j (coordinates of vector $vbini(j)$ correspond to pozitions of 1 in binary expansion of j , are equal to 1).
$nint(\vec{v})$	whole number, derived from vector $\vec{v} \in \{-1, +1\}^n$, (resp. $\vec{v} \in \{0, 1\}^n$, by context), where ones components of vector $nint(\vec{v})$ correspond to 1 in binary expansion of number $nint(\vec{v})$, the rest components correspond to position of 0 in binary expansion of $nint(\vec{v})$
$\binom{a}{i}$	binomical coefficient
\tilde{f}	general mapping between two sets
$(\tilde{f} * \tilde{g})$	convolution of functions \tilde{f} a \tilde{g}
$\tilde{A}(b)$	value of function (mapping) \tilde{A} at point b
$\ln(x)$	natural logarithm of x
$\log_2(x)$	logarithm of radix 2
$x \equiv y \pmod{r}$	$x \equiv y \pmod{r}$ means that there exist whole number k such that $x = yk + r$, where x , y , and r are whole numbers
$Prob_{\tilde{\Pi}}(\bar{A})$	probability of set \bar{A} under probability distribution $\tilde{\Pi}$
$ \bar{V} $	cardinality of set \bar{V}
$ \delta $	absolute value of number δ
$\ \vec{z}\ _{max}$	maximum norm of vector
$\ \vec{z}\ _E$	Euclidean norm of vector
$\ \vec{z}\ _{\mathcal{L}_2^K}$	specified norm of vector
$\bar{A} _{\mathcal{P}}$	orthogonal projection of set \bar{A} into subspace \mathcal{P}

Newly mentioned object will be emphasized in the text by another font, like **NEWLY INTRODUCED NOTION** .

Preface

About this booklet. This work grew out of lecture notes for MSc. courses that I taught at the Mathematical department of FNSPE CTU in Prague since 1997. Consequently, this is a textbook covering enough material for two semester courses. The main focus of this textbook is on introduction to PAC learning model and its variants.

The audience. This text is intended for students in mathematics and other fields such as computer science and electrical engineering. Also a significant portion of the presented material is suitable for undergraduates.

Prerequisites. The necessary prerequisites is basic linear algebra and theory of probability. In some places knowledge of some basic analysis, functional analysis and topology is needed.

Chapter 1

PAC learning model

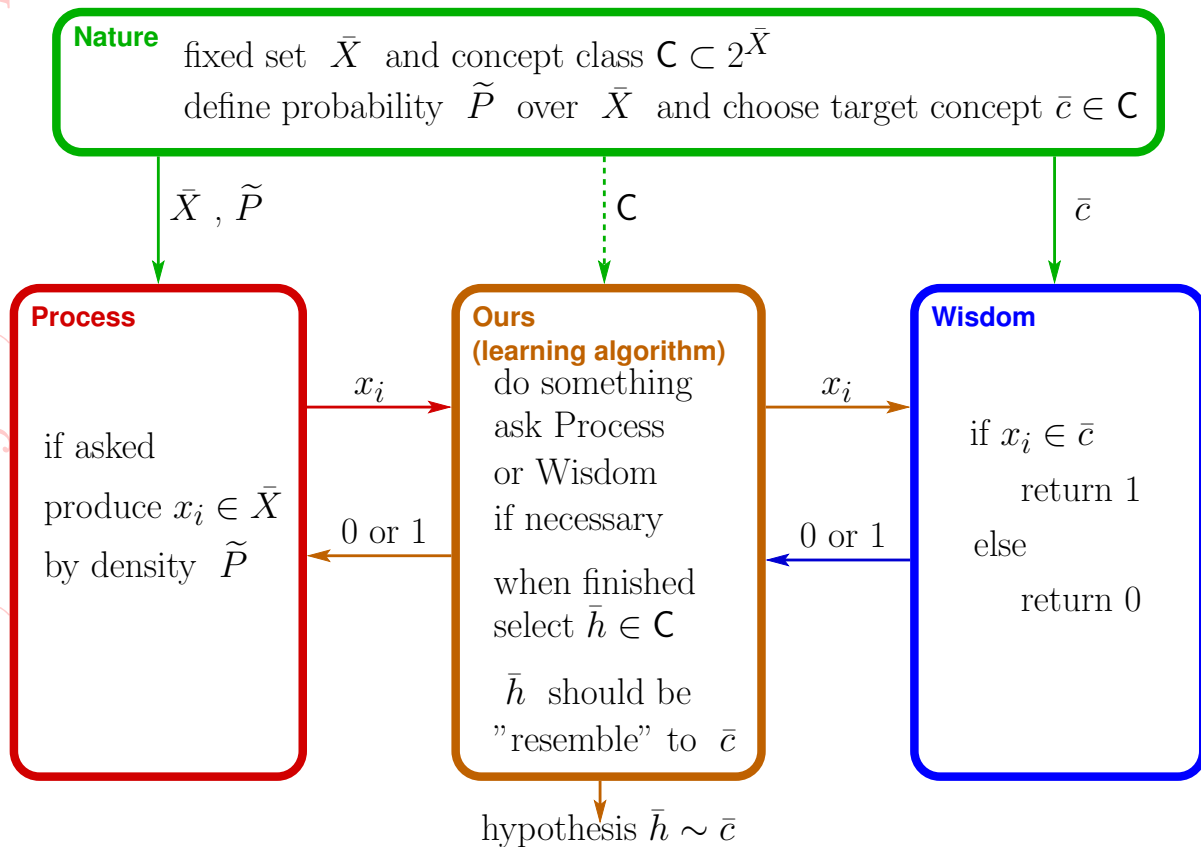


Figure 1.1: Learning scenario.

1.1 Concepts and concept classes

Key subject relating to the term learning algorithm is any subset of firmly given set \bar{X} . The terms concept, hypothesis, concept class and hypothesis class are used to designate such systems in the literature relevant to our further explanation. That is also why this terminology will also be used in the following definitions and explanations.

Definition 1.1.1 Let \bar{X} be an arbitrary set. Then we call $\bar{c} \subset \bar{X}$ **CONCEPT** (over set \bar{X}). Nonempty collection of sets $C \subset 2^{\bar{X}}$ is **CONCEPT CLASS** (over set \bar{X}). Members of this class $\bar{c} \in C$ are concepts of the class C .

One of the methods for describing concept is to set up the concept as a sample of a given set in a representation. Formally, we then define the concept $C_{\tilde{f}}$ for Boolean function \tilde{f} .

Definition 1.1.2 Let \tilde{f} be a mapping from the set \bar{X} into two point set $\{-1, +1\}$. Then define concept $\bar{c}_{\tilde{f}}$ as

$$\bar{c}_{\tilde{f}} \stackrel{\text{def}}{=} \left\{ x \in \bar{X} \mid \tilde{f}(x) = 1 \right\}.$$

Obviously, the preceding definition can be used to define concept class corresponding to set of functions over \bar{X} .

Definition 1.1.3 Let us assume that \bar{F} is an arbitrary set of function which maps \bar{X} into $\{-1, +1\}$. Then we say that \bar{F} **REPRESENT CONCEPT CLASS**

$$C_{\bar{F}} \stackrel{\text{def}}{=} \left\{ \bar{A} \subset \bar{X} \mid (\exists \tilde{f} \in \bar{F})(\bar{A} = \bar{c}_{\tilde{f}}) \right\}.$$

Among others such classes we point out so called **HALFSPACE_n** which is most natural concept class in variety applications.

Definition 1.1.4 Let n be natural number and

$$\bar{F} \stackrel{\text{def}}{=} \left\{ \tilde{f} : \mathfrak{R}^n \rightarrow \{-1, +1\} \mid \tilde{f}(\vec{x}) = \text{sgn}(\langle \vec{x} \mid \vec{w} \rangle - t), \quad t \in \mathfrak{R}, \quad \vec{w}, \vec{x} \in \mathfrak{R}^n \right\}.$$

Then **HALFSPACE_n** $\stackrel{\text{def}}{=} C_{\bar{F}}$.

It is straightforward that **HALFSPACE_n** contains all halfspaces of n -dimensional Euclidean space.

Definition 1.1.5 Let n be natural number and

$$\bar{F} \stackrel{\text{def}}{=} \left\{ \tilde{f} : \mathfrak{R}^n \rightarrow \{-1, +1\} \mid \tilde{f}(\vec{x}) = \text{sgn}(\|\vec{x} - \vec{c}\|_E - r), \quad r \in \mathfrak{R}^+, \quad \vec{c}, \vec{x} \in \mathfrak{R}^n \right\}.$$

Then **BALL_n** $\stackrel{\text{def}}{=} C_{\bar{F}}$.

It is straightforward that **BALL_n** contains all balls in n -dimensional Euclidean space.

1.2 Basic Terms of Learning Algorithms

For better comprehension of the functionality of the Probably Approximately Correct learning model we will use a more formal example of the role of learning than the one presented at the beginning of this section. Let us consider a problem involving a specific given set $\bar{A} \in \mathfrak{R}^n$ knowing a priori that \bar{A} is a sphere in geometric terms. Let us assume to have a given sequence of points $\vec{x}_i \in \mathfrak{R}^n$, while each of them is known either to lie within the set \bar{A} or not. Our task is to decide how does the set \bar{A} look like, and - in this specific case - to find the centre and radius of the relevant sphere. It is evident that our decision is highly dependent on the properties of the given sequence \vec{x}_i . Let us assume that, based on an algorithm (deterministic or nondeterministic), we have opted for a sphere in \mathfrak{R}^n , claiming it is equal to the set \bar{A} . It is natural to expect that this particular set should contain all the points $\vec{x}_i \in \bar{A}$ and, on the other hand, should contain no \vec{x}_i that does not lie within \bar{A} . This requirement will be called consistency and a set which we assume to be equal to the sphere \bar{A} will be called a hypothesis (the equality of the sphere selected by us and \bar{A} is solely our hypothesis, based on the properties of the algorithm used and on the sequence of points \vec{x}_i).

Let us designate our chosen sphere \bar{K} . Then, it is evident that the condition of consistency must be met in case, when symmetric difference of the sets \bar{K} and \bar{A} will be an empty set, and it is natural to demand that this should hold for the generated hypothesis or that this should at least hold with a high rate of probability. The terms discussed above are treated in the following definitions.

Definition 1.2.1 Let C be a concept class and a set system H satisfies $C \subset H \subset 2^{\bar{X}}$. Then, H is HYPOTHESIS CLASS for concept class C .

Definition 1.2.2 Let $\tilde{x} \stackrel{\text{def}}{=} \{x_1, \dots, x_m\}$, $x_i \in \bar{X}$, $i \in \{1, \dots, m\}$, $\vec{z} \in \{-1, +1\}^m$ and let $\bar{c} \subset \bar{X}$. Then the ordered tuple

$$(\tilde{x}, \vec{z})$$

is a SAMPLE OF CONCEPT \bar{c} OF LENGTH m if and only if

$$(\forall i \in \{1, \dots, m\}) ((x_i \in \bar{c}) \Leftrightarrow (\vec{z}_i = 1)).$$

For concept class C define SAMPLE SPACE OF CONCEPT CLASS as

$$\bar{S}_C \stackrel{\text{def}}{=} \bigcup_{m \geq 1} \left\{ \bigcup_{\bar{c} \in C} \left\{ (\tilde{x}, \vec{z}) \mid (\tilde{x}, \vec{z}) \text{ is a sample of the length } m \text{ of concept } \bar{c} \right\} \right\}.$$

A set $\bar{b} \subset \bar{X}$ is CONSISTENT SET with sample (\tilde{x}, \vec{z}) if and only if for all $i \in \{1, \dots, m\}$, holds $(x_i \in \bar{b} \Leftrightarrow \vec{z}_i = 1)$.

A rate of difference between two sets has to be available for the purpose of further explanation. We define this rate of difference as follows:

Definition 1.2.3 Let \bar{X} be a probability space with probability distribution \tilde{P} . Furthermore, \mathbf{C} is concept class over \bar{X} and \mathbf{H} is hypothesis class for \mathbf{C} . Then, for each $\bar{c} \in \mathbf{C}$ and $\bar{h} \in \mathbf{H}$ the number

$$e_{\tilde{P}}(\bar{c}, \bar{h}) \stackrel{\text{def}}{=} \text{Prob}_{\tilde{P}}(\bar{h} \Delta \bar{c})$$

is ERROR OF HYPOTHESIS \bar{h} REGARDING A CONCEPT \bar{c} and probability \tilde{P} (symbol $\bar{h} \Delta \bar{c}$ denotes symmetric difference between sets \bar{h} and \bar{c} , e.g. $\bar{h} \Delta \bar{c} = (\bar{c} - \bar{h}) \cup (\bar{h} - \bar{c})$).

” \bar{h} should be resemble to \bar{c} ”

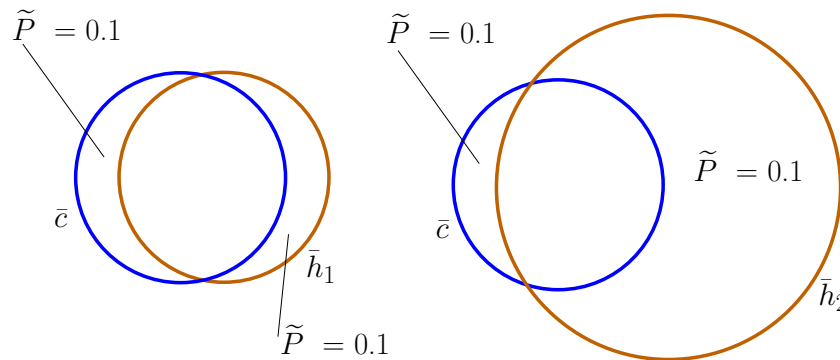


Figure 1.2: Example of two hypotheses with the same error.

Since the set \bar{A} has been described solely by the sequence of its elements (plus a priori information on its shape) our considered learning algorithm is essentially a mapping between the set of all the samples of the set \bar{A} into the set of all spheres in \mathbb{R}^n . We will not require a precise learning algorithm to provide an exact hypothesis for all the possible tasks of this type but we will only require the given learning algorithm to produce, in most cases, a “satisfactory” hypothesis (for instance a hypothesis consistent with the given sample). The following definition represents a quantification of this requirement.

Definition 1.2.4 Let us have a defined function $\tilde{m}(\epsilon, \delta)$ mapping a set $(0, 1) \times (0, 1)$ into a set of natural numbers, and let \mathbf{H} be a hypothesis class for the concept class \mathbf{C} , defined over the set \bar{X} . Then LEARNING ALGORITHM OF COMPLEXITY $\tilde{m}(\epsilon, \delta)$ of the concept class \mathbf{C} is each mapping $\tilde{A}^* : \bar{S}_{\mathbf{C}} \rightarrow \mathbf{H}$ such that for all $\bar{c} \in \mathbf{C}$, for all $0 < \epsilon < 1$, $0 < \delta < 1$ and for any probability \tilde{P} defined on \bar{X} is the probability of the set

$$\left\{ \bar{x} \in \bar{X}^m \mid \left(\bar{x}, \bar{z} \right) \text{ is sample } \bar{c} \text{ a } e_{\tilde{P}}(\bar{c}, \tilde{A}^*((\bar{x}, \bar{z}))) \geq \epsilon \right\}$$

smaller than the number δ . If such a learning algorithm exists we say that \mathbf{C} IS UNIFORMLY LEARNABLE according to the hypothesis class \mathbf{H} . We will call each such learning algorithm the (ϵ, δ) -LEARNING ALGORITHM.

REMARKS:

In the literature, this particular model of computational complexity is designated as PAC (Probably Approximately Correct) learning.

It is evident that each set may be described very precisely by means of a great number of its members. Algorithms using very long samples may probably produce much more veritable hypotheses. If we are satisfied with the given accuracy of generating hypotheses, described by numbers ϵ and δ , it is natural to ask how many a sample must be used for the given concept class to guarantee this accuracy. That is why we will introduce the term sample complexity of learning algorithms as follows:

Definition 1.2.5 *Minimal value $\tilde{m}(\epsilon, \delta)$ for which A^* is a learning algorithm is a SAMPLE COMPLEXITY of learning algorithm A^* .*

This complexity is defined solely with a view to the length of the sample of the learning algorithm used; the issue of the inner computational complexity of the algorithm itself is not taken into consideration in this case! The roughest estimate of sample complexity is given by the subsequent statement, based solely on probability properties (providing a basic estimate of the number of samples that have to be used for the learning algorithm which produces a consistent hypothesis).

Theorem 1.2.1 *Let us assume that C is a concept class over finite set \bar{X} and $H = C$. Let learning algorithm A^* require at least*

$$\frac{1}{\epsilon} \ln \left(\frac{|C|}{\delta} \right) \quad (1.1)$$

queries and for any given concept $\bar{c} \in C$ and any probability density \tilde{P} defined on \bar{X} , produce a consistent hypothesis. Then

$$Prob_{\tilde{P}} \left(e_{\tilde{P}} \left(\bar{c}, \tilde{A}^* \left((\tilde{x}, \tilde{z}) \right) \right) \geq \epsilon \right) < \delta.$$

■ *Proof:*

Let us assume that \bar{h} is a hypothesis consistent with the sample $(\{x_1, \dots, x_m\}, \tilde{z})$, which is generated by a learning algorithm. Without loss of generality, it is possible to assume that for some $1 \leq k \leq m$ is $x_1, \dots, x_k \in \bar{c}$ and $x_{k+1}, \dots, x_m \notin \bar{c}$. Hence, it ensues from the consistency of \bar{h} that $x_1, \dots, x_k \in \bar{h}$ and $x_{k+1}, \dots, x_m \notin \bar{h}$. Therefore, the following estimate holds

$$e_{\tilde{P}}(\bar{c}, \bar{h}) = \sum_{x \in \bar{c} \Delta \bar{h}} Prob_{\tilde{P}}(x) \leq 1 - \sum_{i=1}^m Prob_{\tilde{P}}(x_i).$$

Hence, we obtain from the requirement $\epsilon \leq e_{\tilde{P}}(\bar{c}, \bar{h})$

$$\sum_{i=1}^m Prob_{\tilde{P}}(x_i) \leq 1 - \epsilon$$

and from the positiveness of the numbers $Prob_{\tilde{P}}(x_i)$ it further ensues that

$$(\forall i \in \{1, \dots, m\})(Prob_{\tilde{P}}(x_i) \leq 1 - \epsilon).$$

Hence, the probability of selection of m points $x_i \in \bar{X}$ such that they are consistent with hypothesis \bar{h} for which $e_{\bar{P}}(\bar{c}, \bar{h}) \geq \epsilon$ meets the inequality

$$Prob_{\bar{P}}(\{x_1, \dots, x_m \mid e_{\bar{P}}(\bar{c}, \bar{h}) \geq \epsilon \text{ and } \bar{h} \text{ is consistent}\}) = \prod_{i=1}^m Prob_{\bar{P}}(x_i) \leq (1 - \epsilon)^m$$

(selected queries may be repeated; this is an independent selection).

But we are interested to know the probability that such hypothesis really exists. The number of hypotheses is, however, finite and equals the number of concepts in the class \mathbf{C} . Therefore, the overall probability of the incidence of the hypothesis with the properties described above is lower than the number

$$|\mathbf{C}| \cdot (1 - \epsilon)^m.$$

We demand that probability should, at the most, equal the value of δ , hence that it should hold that $|\mathbf{C}| \cdot (1 - \epsilon)^m \leq \delta$. Let us choose

$$m \stackrel{\text{def}}{=} \left\lceil \frac{1}{\epsilon} \ln \left(\frac{|\mathbf{C}|}{\delta} \right) \right\rceil.$$

Then the inequality ($\delta < 1$) holds ¹.

$$|\mathbf{C}| \cdot (1 - \epsilon)^m \leq |\mathbf{C}| e^{-\epsilon m} \leq |\mathbf{C}| e^{\ln(\frac{\delta}{|\mathbf{C}|})} = \delta.$$

– q. e. d. –

Example 1.2.1 Let $\bar{X} \stackrel{\text{def}}{=} \{1, \dots, k\}$, $\mathbf{C} \stackrel{\text{def}}{=} \{\bar{a} \subset \bar{X} \mid (\exists i \in \bar{X}) (\bar{a} = \{j \in \bar{X} \mid i \leq j\})\}$. Further define learning algorithm A^* as

$$\tilde{A}^*((x_1, \dots, x_m, z_1, \dots, z_m)) = \bar{h} = \langle b, k \rangle \cap \bar{X}, \quad \text{where } b = \min_{i \in \{1, \dots, m\}} \{x_i \mid z_i = 1\}.$$

Obviously the algorithm A^* produces consistent hypotheses only. At the same time, $|\mathbf{C}| = k$. Let $0 < \epsilon, \delta < 1$ and

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

Then A^* is (ϵ, δ) -learning algorithm.

Example 1.2.2 Let $\bar{X} \stackrel{\text{def}}{=} \{1, \dots, k\}^n$, let

$$\mathbf{C} \stackrel{\text{def}}{=} \{\langle a_1, b_1 \rangle \times \langle a_2, b_2 \rangle \times \dots \times \langle a_n, b_n \rangle \mid (\forall j \in \{1, \dots, n\}) (a_j \leq b_j)\}.$$

Further define learning algorithm A^* as

$$\tilde{A}^*((\vec{x}_1, \dots, \vec{x}_m, z_1, \dots, z_m)) = \langle a_1, b_1 \rangle \times \langle a_2, b_2 \rangle \times \dots \times \langle a_n, b_n \rangle$$

¹Let $\tilde{f}(\epsilon) \stackrel{\text{def}}{=} e^{-\epsilon} - (1 - \epsilon)$. Then $\tilde{f}(0) = 0$ and $\tilde{f}'(\epsilon) = 1 - e^{-\epsilon} > 0$ on $(0, 1)$. Hence $e^{-\epsilon} > 1 - \epsilon$ on $(0, 1)$.

where

$$(\forall j \in \{1, \dots, n\}) \left(a_j = \min_{i \in \{1, \dots, m\}} \{(\vec{x}_i)_j | \vec{z}_i = 1\} \quad \text{and} \quad b_j = \max_{i \in \{1, \dots, m\}} \{(\vec{x}_i)_j | \vec{z}_i = 1\} \right).$$

Obviously the algorithm A^* produces consistent hypotheses only. At the same time, $|\mathcal{C}| = \binom{k}{2}^n$. Let $0 < \epsilon, \delta < 1$ and

$$m \geq \frac{1}{\epsilon} \tilde{\ln} \left(\frac{\binom{k}{2}^n}{\delta} \right).$$

Then A^* is (ϵ, δ) -learning algorithm.

Let us note that the estimate 1.1 is not dependent on the specific rate of probability \tilde{P} . This can be explained by the fact that the influence of different probabilities of the elements in \bar{X} is compensated by the probability of the selection of this element according to the density of probability \tilde{P} . In other words, elements from \bar{X} with considerable probability, which - in case of belonging to symmetric difference $\bar{c} \Delta \bar{h}$ greatly contribute to the actual size of error, being more frequent in an average sample, which - in view of the fact that algorithm always produces a consistent hypothesis - eliminates their impact. That is why the estimate 1.1 is independent of probability \tilde{P} .

In the subsequent explication, we will demonstrate that this estimate can be markedly improved, notably by applying the term VC-dimension.

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

Chapter 2

Vapnik-Chervonenkis dimension

We now begin to introduce the mathematical theory which is necessary to study PAC learning in cases of infinite concept and hypothesis classes. A main notion discussed in this chapter is Vapnik-Chervonenkis dimension which is a measure of possible set classification potential of set collection contained in concept class given.

2.1 VC-dimension of concept class

Definition 2.1.1 Let \bar{X} is a given set, $\bar{Y} \subset \bar{X}$, $\mathcal{C} \subset 2^{\bar{X}}$. Then \mathcal{C} IS SHATTERED BY \bar{Y} , iff:

$$(\forall \bar{z} \subset \bar{Y})(\exists \bar{c} \in \mathcal{C})(\bar{c} \cap \bar{Y} = \bar{z}).$$

Further define VAPNIK-CHERVONENKIS DIMENSION of the set system \mathcal{C} as

$$\text{VC}_{\dim}(\mathcal{C}) \stackrel{\text{def}}{=} \sup \{ |\bar{Y}| \mid \mathcal{C} \text{ is shattered by } \bar{Y} \}.$$

Note, that Vapnik-Chervonenkis dimension is equal to the size of maximal set $\bar{A} \subset \bar{X}$ satisfying

$$\{ \bar{B} \mid (\exists \bar{Z} \in \mathcal{C})(\bar{B} = \bar{A} \cap \bar{Z}) \} = 2^{\bar{A}},$$

or to $+\infty$, if an set of any size is shattered by \mathcal{C} .

The notion of VC-dimension is so important that we illustrate it using following examples (see [AB92]).

Example 2.1.1

1. Let $\bar{X} = \mathfrak{R}$ and \mathcal{C} are all intervals $(a, +\infty)$, where $a \in \mathfrak{R}$. In this case, obviously for an arbitrary $\{b, c \in \mathfrak{R} \mid b < c\}$ there does not exist an interval $(a, +\infty)$ which contains the point b and does not contain the point c . Hence $\text{VC}_{\dim}(\mathcal{C}) = 1$.
2. Let $\bar{X} = \mathfrak{R}$, s is a natural number, $s > 1$ and \mathcal{C} contains all union of s intervals. Let $\bar{S} = \{x_1, \dots, x_{2s} \mid x_i < x_{i+1}\}$. It is straightforward that \bar{S} is shattered by concept class \mathcal{C} . But for $\bar{S} = \{x_1, \dots, x_{2s+1} \mid x_i < x_{i+1}\}$ each concept containing points $x_1, x_3, \dots, x_{2s+1}$ also contains at least one point from x_2, x_4, \dots, x_{2s} . Clearly such \bar{S} is not shattered by the concept class \mathcal{C} . It follows $\text{VC}_{\dim}(\mathcal{C}) = 2s$.

3. Let $\bar{X} = \mathfrak{R}^n$ and \mathcal{C} are all intervals in the space \mathfrak{R}^n . Obviously an set of $2n$ points located in the middle of faces of the unit cube is shattered by the concept class \mathcal{C} (we can move faces slightly in the direction of coordinates axes). But for arbitrary set \bar{A} of $2n + 1$ points, there exists minimal parallelepiped containing whole set \bar{A} . Straightforwardly there exists an point $\bar{\mathbf{x}} \in \bar{A}$ which is in the interior of minimal parallelepiped, or lies in the some of its face. Obviously the set $\bar{A} - \{\bar{\mathbf{x}}\}$ can not be separate off the set $\{\bar{\mathbf{x}}\}$. Hence $\mathbf{VC}_{\dim}(\mathcal{C}) = 2n$.

Example 2.1.2 For the completeness we show an example of the concept class with infinite VC-dimension.

1. Put

$$\mathcal{C} \stackrel{\text{def}}{=} \left\{ \bar{A}_\alpha \mid (\exists \alpha \in \mathfrak{R}^n) \left(\bar{A}_\alpha = \left\{ x \in \mathfrak{R} \mid \widetilde{\sin}(\alpha x) \geq 0 \right\} \right) \right\}.$$

Then for arbitrary natural l define sequence $\bar{Z}_l \stackrel{\text{def}}{=} \{z_i\}_1^l$, $z_i \stackrel{\text{def}}{=} \frac{1}{10^i}$. Further for any subset of the set $\bar{Q} \subset \bar{Z}_l$ the sequence $\delta_1, \dots, \delta_l$, $\delta_i \in \{0, 1\}$ satisfies $\delta_i = 1 \Leftrightarrow z_i \in \bar{Q}$. If we define

$$\alpha \stackrel{\text{def}}{=} \pi \left(\sum_{i=1}^l (1 - \delta_i) 10^i + 1 \right),$$

then equality

$$\alpha z_j = \alpha \frac{1}{10^j} = \pi \left(\sum_{i=1}^l (1 - \delta_i) 10^{i-j} + \frac{1}{10^j} \right).$$

holds. Hence

$$\alpha z_j = \alpha \frac{1}{10^j} = \pi \left(\sum_{i=1}^{j-1} \frac{1 - \delta_i}{10^{j-i}} + \frac{1}{10^j} + (1 - \delta_j) + \sum_{i=j+1}^l (1 - \delta_i) 10^{i-j} \right).$$

Obviously the expression $\sum_{i=j+1}^l (1 - \delta_i) 10^{i-j}$ is divisible by 2 while the expression $\sum_{i=1}^{j-1} \frac{1 - \delta_i}{10^{j-i}} + \frac{1}{10^j}$ is less than to 1.

It follows that $\widetilde{\sin}(\alpha z_j) < 0$ for $\delta_j = 0$ and $\widetilde{\sin}(\alpha z_j) > 0$ for $\delta_j = 1$. Thus the set \bar{A}_α separate each subset \bar{Z}_l defined via sequence $\delta_1, \dots, \delta_l$. Therefore for arbitrary l there exists the set of the size l which is shattered by the concept class \mathcal{C} , So $\mathbf{VC}_{\dim}(\mathcal{C}) = +\infty$.

2. Put

$$\mathcal{C} \stackrel{\text{def}}{=} \{[\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_m]_\kappa \mid m \in \mathbb{N} \text{ and } (\forall i \in \{1, \dots, m\}) (\bar{\mathbf{x}} \in \mathfrak{R}^n)\},$$

e.g. \mathcal{C} is the set of all convex hulls of all finite subsets of \mathfrak{R}^n . It is clear, that vertices of a given convex hull over m points can be shattered by \mathcal{C} , so $\mathbf{VC}_{\dim}(\mathcal{C}) = +\infty$.

2.1.1 General properties of VC-dimension

In the next text we analyze concept classes defined over an finite set \bar{X} . We start with the Sauer's lemma ([Sau72]) which provide us upper bound on number of concepts in a concept class with predefined VC-dimension. We feel necessary to point out here that this upper bound is genuine combinatorial property of finite sets.

Lemma 2.1.1 (Sauer) *Let \bar{X} be a finite set and $\mathcal{C} \subset 2^{\bar{X}}$. Then*

$$|\mathcal{C}| \leq \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C})} \binom{|\bar{X}|}{i}.$$

Further, there exists $\mathcal{C} \subset 2^{\bar{X}}$ such that equality holds.

■ *Proof:*

Let us define for arbitrary but fixed $y \in \bar{X}$ the following sets

$$\begin{aligned} \mathcal{C}(y) &\stackrel{\text{def}}{=} \{ \bar{A} \dot{-} \{y\} \mid \bar{A} \in \mathcal{C} \} \\ \mathcal{C}_y &\stackrel{\text{def}}{=} \{ \bar{A} \in \mathcal{C} \mid (\exists \bar{B} \in \mathcal{C})(\bar{A} \neq \bar{B} \text{ and } \bar{B} = \bar{A} \cup \{y\}) \} \\ \mathcal{C}^y &\stackrel{\text{def}}{=} \{ \bar{A} \in \mathcal{C} \mid (\exists \bar{B} \in \mathcal{C}_y)(\bar{A} = \bar{B} \cup \{y\}) \}. \end{aligned}$$

We prove the claim of lemma it tree steps.

■ add 1)

Firstly, we show that for arbitrary $y \in \bar{X}$ is

$$|\mathcal{C}| - |\mathcal{C}(y)| = |\mathcal{C}_y|.$$

Obviously

$$\mathcal{C}(y) = \{ \bar{A} \mid \bar{A} \in \mathcal{C}, y \notin \bar{A} \} \cup \{ \bar{A} \dot{-} \{y\} \mid \bar{A} \in \mathcal{C}, y \in \bar{A} \}.$$

Hence

$$\begin{aligned} |\mathcal{C}(y)| &= |\{ \bar{A} \mid \bar{A} \in \mathcal{C}, y \notin \bar{A} \}| + |\{ \bar{A} \dot{-} \{y\} \mid \bar{A} \in \mathcal{C}, y \in \bar{A} \}| - \\ &|\{ \bar{A} \mid \bar{A} \in \mathcal{C}, y \notin \bar{A} \} \cap \{ \bar{A} \dot{-} \{y\} \mid \bar{A} \in \mathcal{C}, y \in \bar{A} \}| = |\mathcal{C}| - |\mathcal{C}_y|. \end{aligned}$$

■ add 2)

Secondly, we prove

$$\text{VC}_{dim}(\mathcal{C}_y) = n - 1 \Rightarrow \text{VC}_{dim}(\mathcal{C}) \geq n. \quad (2.1)$$

It follows from the assumption $\text{VC}_{dim}(\mathcal{C}_y) = n - 1$ that there exists an set $\bar{A} \subset \bar{X} \dot{-} \{y\}$, $|\bar{A}| = n - 1$ which satisfy

$$\{ \bar{B} \mid (\exists \bar{Z} \in \mathcal{C}_y)(\bar{B} = \bar{A} \cap \bar{Z}) \} = 2^{\bar{A}}. \quad (2.2)$$

To verify it we prove the equality

$$\{ \bar{B} \mid (\exists \bar{Z} \in (\mathcal{C}_y \cup \mathcal{C}^y))(\bar{B} = (\bar{A} \cup \{y\}) \cap \bar{Z}) \} = 2^{\bar{A} \cup \{y\}}.$$

Inclusion \subset is obvious, so let us pay our attention to opposite inclusion \supset . Let $\bar{H} \in 2^{\bar{A} \cup \{y\}}$. If $y \notin \bar{H}$ then by 2.2 $\exists \bar{Z} \in \mathcal{C}_y$ such that $\bar{H} = \bar{A} \cap \bar{Z} = (\bar{A} \cup \{y\}) \cap \bar{Z}$. On the contrary, let $y \in \bar{H}$. If we recall 2.2 then for the set $\bar{H} - \{y\}$ there exists $\bar{M} \in \mathcal{C}_y$ such that $\bar{H} - \{y\} = \bar{A} \cap \bar{M}$. On account of $\bar{M} \in \mathcal{C}_y$ it must be $\bar{M} \cup \{y\} \in \mathcal{C}^y$ and therefore the set equality $\bar{H} = (\bar{A} \cup \{y\}) \cap (\bar{M} \cup \{y\})$ holds. So we get that $\text{VC}_{dim}(\mathcal{C}_y \cup \mathcal{C}^y) \geq n$ and because $(\mathcal{C}_y \cup \mathcal{C}^y) \subset \mathcal{C}$, the equality $\text{VC}_{dim}(\mathcal{C}) \geq n$ must be fulfilled.

■ add 3)

Now we are ready to conclude the proof of the lemma by induction on the set size $|\bar{X}|$.
 $|\bar{X}| = 1, 2$:

The claim of the lemma is straightforward.

$|\bar{X}| = k \mapsto |\bar{X}| = k + 1$:

Let us assume that induction assumption is valid for the case when if $|\bar{X}| = k$. The definition of the set system \mathcal{C}_y follows that the set contained in it does not contain y , so the whole set system \mathcal{C}_y is build over the set $\bar{X} - \{y\}$. If we apply induction assumption we obtain

$$|\mathcal{C}_y| \leq \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C}_y)} \binom{|\bar{X}| - 1}{i} \leq \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C}) - 1} \binom{|\bar{X}| - 1}{i}$$

(the second inequality holds owing to the fact $\text{VC}_{dim}(\mathcal{C}) - 1 \geq \text{VC}_{dim}(\mathcal{C}_y)$ which comes true due to 2.1).

At the same time, in view of induction assumption $\mathcal{C}(y)$ is defined over the set $\bar{X} - \{y\}$, the estimation

$$|\mathcal{C}(y)| \leq \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C}(y))} \binom{|\bar{X}| - 1}{i} \leq \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C})} \binom{|\bar{X}| - 1}{i},$$

holds, whereas the second inequality immediately follows from the definition of the VC-dimension. So we can write (remind $\binom{j}{i} = \binom{j-1}{i} + \binom{j-1}{i-1}$ and $\binom{j}{0} \stackrel{\text{def}}{=} 1$, $\binom{j}{-1} \stackrel{\text{def}}{=} 0$)

$$|\mathcal{C}| = |\mathcal{C}(y)| + |\mathcal{C}_y| \leq \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C})} \binom{|\bar{X}| - 1}{i} + \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C})} \binom{|\bar{X}| - 1}{i-1} = \sum_{i=0}^{\text{VC}_{dim}(\mathcal{C})} \binom{|\bar{X}|}{i}.$$

So we proved the desired inequality. Finally, we have to show that there exists a concept class for which the equality is true. Let us define a concept class

$$\mathcal{C} \stackrel{\text{def}}{=} \{ \bar{a} \subset \bar{X} \mid |\bar{a}| \leq k \}.$$

It is easy to verify that $|\mathcal{C}| = \sum_{i=0}^k \binom{|\bar{X}|}{i}$ and that $\text{VC}_{dim}(\mathcal{C}) = k$, which conclude the proof.

– q. e. d. –

The following simple lemma highlights usefulness of the Sauer's lemma.

Lemma 2.1.2 *Let \mathcal{C} be a finite concept class. Then $\text{VC}_{dim}(\mathcal{C}) \leq \log_2(|\mathcal{C}|)$.*

■ *Proof:*

To shatter any set of the size d we need at least 2^d different concepts, hence the set of the size greater than $\log_2(|C|)$ can not be shattered by the set C . Therefore $\text{VC}_{dim}(C) \leq \log_2(|C|)$.

– q. e. d. –

The notion introduces in the next definition has very narrow relationship to VC-dimension and also to Sauer's lemma.

Definition 2.1.2 Let C be a nonempty concept class over the set \bar{X} and $\bar{S} \subset \bar{X}$. Then we define

$$\Pi_C(\bar{S}) \stackrel{\text{def}}{=} \{\bar{S} \cap \bar{c} \mid \bar{c} \in C\}.$$

Further for a fixed $m \geq 0$ let us define

$$\Pi_C(m) \stackrel{\text{def}}{=} \max \{|\Pi_C(\bar{S})| \mid \bar{S} \subset \bar{X}, |\bar{S}| = m\},$$

(maximum is over all sets $\bar{S} \subset \bar{X}$ of the size m).

Apparently $\Pi_C(\bar{S})$ is the system of all subset of the set \bar{S} whose are possible to separate from their complement in \bar{S} using concept from concept class C . The number $\Pi_C(m)$ expresses the cardinality of maximal such system $\Pi_C(\bar{S})$ under condition that the set $\bar{S} \subset \bar{X}$ is finite set of the size m . There exists evident relationship between the number discussed above and the VC-dimension

$$\text{VC}_{dim}(C) = \sup \{m \mid \Pi_C(m) = 2^m\}.$$

Definition 2.1.3 For all $d \geq 0$ a $m \geq 0$ put $\Phi_{d,m} \stackrel{\text{def}}{=} \sum_{i=0}^d \binom{m}{i}$, if $m \geq d$ and $\Phi_{d,m} \stackrel{\text{def}}{=} 2^m$ if $m < d$.

Estimations contained in the successive lemma are usefull in deriving upper bounds on necessary sample size which quarrants that learning algorithm produce an hypothesis with acceptable small error (in the sense of probability of the symmetric difference between the original concept \bar{c} and produced hypothesis \bar{h}).

Lemma 2.1.3 If C is an arbitrary concept class over \bar{X} and $\text{VC}_{dim}(C) = d < +\infty$, then

1. $\Pi_C(m) \leq \Phi_{d,m}$ for all $d, m \geq 0$.
2. $\Phi_{d,m} \leq m^d + 1$ for $d, m \geq 0$ and $\Phi_{d,m} \leq m^d$ for $d \geq 0$ and $m \geq 2$.
3. $\Phi_{d,m} \leq 2 \frac{m^d}{d!} \leq \left(\frac{em}{d}\right)^d$ for $m \geq d \geq 1$.

■ *Proof:*

• add 1)

Let m be arbitrary and let $\bar{S} \subset \bar{X}$ be arbitrary subset such that $|\bar{S}| = m$ and $|\Pi_C(\bar{S})| = \Pi_C(m)$. Put $\bar{X}^* \stackrel{\text{def}}{=} \bar{S}$ and $C^* \stackrel{\text{def}}{=} \{\bar{c} \cap \bar{S} \mid \bar{c} \in C\}$. As $\Pi_{C^*}(\bar{S}) = \{\bar{S} \cap \bar{c} \mid \bar{c} \in C^*\}$, we have $|\Pi_{C^*}(\bar{S})| \leq |C^*|$. It implies that $\Pi_{C^*}(m) \leq |C^*|$. At the same time $\Pi_{C^*}(\bar{S}) =$

$\{\bar{S} \cap \bar{c} | \bar{c} \in \mathbf{C}^*\} = \{\bar{S} \cap \bar{c} | \bar{c} \in \mathbf{C}\} = \Pi_{\mathbf{C}}(\bar{S})$ which follows $\Pi_{\mathbf{C}^*}(m) = \Pi_{\mathbf{C}}(m)$ (note that we assume $|\Pi_{\mathbf{C}}(\bar{S})| = \Pi_{\mathbf{C}}(m)$ and obviously $\Pi_{\mathbf{C}^*}(m) \leq \Pi_{\mathbf{C}}(m)$). In addition $\mathbf{VC}_{dim}(\mathbf{C}^*) \leq \mathbf{VC}_{dim}(\mathbf{C})$. So owing to Sauer's lemma we get the estimation

$$\Pi_{\mathbf{C}}(m) = \Pi_{\mathbf{C}^*}(m) \leq |\mathbf{C}^*| \leq \sum_{i=0}^{\mathbf{VC}_{dim}(\mathbf{C}^*)} \binom{|\bar{X}^*|}{i} \leq \sum_{i=0}^{\mathbf{VC}_{dim}(\mathbf{C})} \binom{|\bar{X}^*|}{i} = \Phi_{\mathbf{VC}_{dim}(\mathbf{C}), m}.$$

Hence the first inequality in the lemma is clearly fulfilled for any m .

• add 2)

This part will be proved by induction on m and d using the relationship $\binom{m}{i} = \binom{m-1}{i} + \binom{m-1}{i-1}$ (recall that $\binom{j}{0} \stackrel{\text{def}}{=} 1$, $\binom{j}{-1} \stackrel{\text{def}}{=} 0$). So we can write

$$\begin{aligned} \sum_{i=0}^d \binom{m}{i} &= \sum_{i=0}^d \binom{m-1}{i} + \sum_{i=0}^d \binom{m-1}{i-1} = \sum_{i=0}^d \binom{m-1}{i} + \sum_{i=1}^{d-1} \binom{m-1}{i} \stackrel{\text{ind.}}{\leq} \\ &\stackrel{\text{ind.}}{\leq} (m-1)^d + (m-1)^{(d-1)} = m^d \left(\frac{(m-1)^d}{m^d} + \frac{(m-1)^{(d-1)}}{m^d} \right) = \\ &= m^d \left[\left(1 - \frac{1}{m}\right)^d + \frac{1}{m} \left(1 - \frac{1}{m}\right)^{(d-1)} \right] = m^d \left(1 - \frac{1}{m}\right)^{d-1} \leq m^d. \end{aligned}$$

• add 3)

First, we prove the inequality $\Phi_{d,m} \leq 2 \frac{m^d}{d!}$. To do it we apply induction on m and d .

If $d = 1$ then $\Phi_{d,m} = m + 1 \leq 2m$ and inequality holds.

Further, for $m = d > 1$ is $\Phi_{d,m} = 2^d$. If we recall binomical formulae we can observe that for $d > 1$ the expression $2 \leq \left(1 + \frac{1}{d-1}\right)^{d-1}$ is true. If we apply induction on d we derive

$$2^d \leq \left(\frac{d}{d-1}\right)^{d-1} 2^{d-1} \stackrel{\text{induction}}{\leq} \stackrel{\text{assumption}}{\leq} 2 \left(\frac{d}{d-1}\right)^{d-1} \frac{(d-1)^{d-1}}{(d-1)!} = 2 \frac{d^d}{d!},$$

which concludes the claim for the case $m = d > 1$.

Now let us assume that $m > d > 1$. Because of $\Phi_{d,m} = \Phi_{d-1,m-1} + \Phi_{d,m-1}$ it is sufficient to verify that

$$2 \frac{(m-1)^{d-1}}{(d-1)!} + 2 \frac{(m-1)^d}{d!} \leq 2 \frac{m^d}{d!}.$$

After multiplying by the number $d!$, the previous expression is equivalent to the expressions bellow:

$$\begin{aligned} d(m-1)^{d-1} + (m-1)^d &\leq m^d \Leftrightarrow \\ (d+m-1)(m-1)^{d-1} &\leq m^d \Leftrightarrow \\ \frac{d+m-1}{m-1} &\leq \frac{m^d}{(m-1)^d} \Leftrightarrow \\ 1 + \frac{d}{m-1} &\leq \left(1 + \frac{1}{m-1}\right)^d, \end{aligned}$$

• where the last estimation is based on binomical formula again.

The second estimation $2 \frac{m^d}{d!} \leq \left(\frac{em}{d}\right)^d$ is evident for $d = 1$. For $d \geq 2$ we use Stirling's formulae

$$n! = n^n \sqrt{2\pi n} \cdot e^{-n + \frac{\delta(n)}{4}}, \quad 0 < \tilde{\delta}(n) < 1.$$

whose applying yields

$$2 \frac{m^d}{d!} < \frac{2m^d}{d^d e^{-d} \sqrt{2\pi d}} = \sqrt{\frac{2}{\pi d}} \left(\frac{em}{d}\right)^d < \left(\frac{em}{d}\right)^d.$$

– q. e. d. –

Corollary 2.1.4 Let \bar{X} be a finite set, $C \subset 2^{\bar{X}}$ and $\text{VC}_{\dim}(C) > 0$. Then

$$\text{VC}_{\dim}(C) > \frac{\ln(|\bar{C}|)}{1 + \ln(|\bar{X}|)}.$$

■ *Proof:*

Sauer lemma asserts that $|C| \leq \sum_{i=0}^{\text{VC}_{\dim}(C)} \binom{|\bar{X}|}{i}$. Applying last part of 2.1.3, we conclude

$$|C| \leq \left(\frac{e|\bar{X}|}{\text{VC}_{\dim}(C)} \right)^{\text{VC}_{\dim}(C)}.$$

In other words,

$$\ln(|C|) \leq \text{VC}_{\dim}(C) (1 + \ln(|\bar{X}|) - \ln(\text{VC}_{\dim}(C))) < \text{VC}_{\dim}(C) (1 + \ln(|\bar{X}|)).$$

– q. e. d. –

2.1.2 VC-dimension of union and intersection

Lemma 2.1.5 Let Q concept class over \bar{X} and P is defined as

$$P \stackrel{\text{def}}{=} \left\{ \bar{a} \subset \bar{X} \mid (\exists \bar{b} \in Q) (\bar{a} = \bar{X} - \bar{b}) \right\}.$$

Then $\text{VC}_{\dim}(Q) = \text{VC}_{\dim}(P)$.

■ *Proof:*

Let $\bar{z} \subset \bar{X}$ be an arbitrary set shattered by concept class Q , an \bar{z}_1, \bar{z}_2 be arbitrary dichotomy of \bar{z} . Hence, there exists sets $\bar{a} \in Q$ and $\bar{b} \in Q$ such that $\bar{z}_1 \subset \bar{a}$, $\bar{z}_1 \cap \bar{b} = \emptyset$ and $\bar{z}_2 \subset \bar{b}$, $\bar{z}_2 \cap \bar{a} = \emptyset$. In other words $\bar{z}_1 \cap (\bar{X} - \bar{a}) = \emptyset$, $\bar{z}_1 \subset (\bar{X} - \bar{b})$ and $\bar{z}_2 \cap (\bar{X} - \bar{b}) = \emptyset$, $\bar{z}_2 \subset (\bar{X} - \bar{a})$. It follows that each subset $\bar{z} \subset \bar{X}$ shattered by concept class Q is also shattered by concept class P and vice versa. Therefore $\text{VC}_{\dim}(Q) = \text{VC}_{\dim}(P)$.

– q. e. d. –

Definition 2.1.4 Let $C_i \subset 2^{\bar{X}}$, $i \in \{1, \dots, k\}$ are concept classes. Then define the following sets

$$U_{C_1, \dots, C_k} \stackrel{\text{def}}{=} \left\{ \bigcup_{i=1}^k \bar{c}_i \mid (\forall i \in \{1, \dots, k\}) (\bar{c}_i \in C_i) \right\}$$

and

$$I_{C_1, \dots, C_k} \stackrel{\text{def}}{=} \left\{ \bigcap_{i=1}^k \bar{c}_i \mid (\forall i \in \{1, \dots, k\}) (\bar{c}_i \in C_i) \right\}.$$

Finally denote U_C as the set of all finite unions of subsets C and I_C as the set of all finite intersections of subsets C . Further, for each natural number k define the sets

$$U_{k,C} \stackrel{\text{def}}{=} \left\{ \bigcup_{i=1}^k \bar{c}_i \mid (\forall i \in \{1, \dots, k\}) (\bar{c}_i \in C) \right\}$$

and

$$I_{k,C} \stackrel{\text{def}}{=} \left\{ \bigcap_{i=1}^k \bar{c}_i \mid (\forall i \in \{1, \dots, k\}) (\bar{c}_i \in C) \right\}.$$

Lemma 2.1.6 Let C be a concept class over \bar{X} and $(\forall \bar{a} \in C) (\bar{X} \dot{-} \bar{a} \in C)$. Then the equality

$$\text{VC}_{\text{dim}}(U_{k,C}) = \text{VC}_{\text{dim}}(I_{k,C})$$

holds.

■ *Proof:*

Remember that for all sets the following de Morgan's formulas holds

$$\bar{X} \dot{-} \bigcup_{i=1}^k \bar{a}_i = \bigcap_{i=1}^k (\bar{X} \dot{-} \bar{a}_i) \quad \text{and} \quad \bar{X} \dot{-} \bigcap_{i=1}^k \bar{a}_i = \bigcup_{i=1}^k (\bar{X} \dot{-} \bar{a}_i).$$

Using these formulas it is easy to prove that $(\forall \bar{a} \in U_{k,C}) (\exists \bar{b} \in I_{k,C}) (\bar{a} = \bar{X} \dot{-} \bar{b})$.

Obviously the claim of lemma follows from lemma 2.1.5.

– q. e. d. –

Main properties of VC-dimension of concept classes union provides next theorems.

Theorem 2.1.7 Let \bar{X} be arbitrary set and P, Q are concept classes over \bar{X} . Then

$$\text{VC}_{\text{dim}}(U_{P,Q}) \leq \text{VC}_{\text{dim}}(P) + \text{VC}_{\text{dim}}(Q) + 1.$$

■ *Proof:*

We show this theorem by contradiction. Let $\text{VC}_{\text{dim}}(P) = m$, $\text{VC}_{\text{dim}}(Q) = n$, and $\text{VC}_{\text{dim}}(U_{P,Q}) = n+m+2$. It follows that there exists a points set $\{x_1, \dots, x_{m+n+2}\} \subset \bar{X}$ which is shattered by concept class $U_{P,Q}$, where $x_i \neq x_j$, $i, j \in \{1, \dots, m+n+2\}$. Because of $\text{VC}_{\text{dim}}(P) = m$

if $\text{VC}_{\dim}(\mathbf{Q}) = n$, there exist sets $\bar{a} \subset \{x_1, \dots, x_{m+1}\}$ and $\bar{b} \subset \{x_{m+2}, \dots, x_{m+n+2}\}$ such that

$$(\forall \bar{p} \in \mathbf{P}) (\bar{a} \neq \{x_1, \dots, x_{m+1}\} \cap \bar{p}) \text{ and } (\forall \bar{q} \in \mathbf{Q}) (\bar{b} \neq \{x_{m+2}, \dots, x_{m+n+2}\} \cap \bar{q}) .$$

Because the set $\{x_1, \dots, x_{m+n+2}\}$ is shattered by $\mathbf{U}_{\mathbf{P}, \mathbf{Q}}$ there exist sets $\bar{z} \in \mathbf{U}_{\mathbf{P}, \mathbf{Q}}$, $\bar{z}_{\mathbf{P}} \in \mathbf{P}$, $\bar{z}_{\mathbf{Q}} \in \mathbf{Q}$, such that

$$\bar{z} = \bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}} \quad \text{and} \quad (\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \{x_1, \dots, x_{m+n+2}\} = \bar{a} \cup \bar{b} .$$

For the sake of clarity, let us define sets $\bar{A} \stackrel{\text{def}}{=} \{x_1, \dots, x_{m+1}\}$ and $\bar{B} \stackrel{\text{def}}{=} \{x_{m+2}, \dots, x_{m+n+2}\}$. As we mentioned before we can write $\bar{A} \cap \bar{B} = \emptyset$, $\bar{a} \subset \bar{A}$, and $\bar{b} \subset \bar{B}$. So

$$\begin{aligned} (\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \{x_1, \dots, x_{m+n+2}\} &= (\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap (\bar{A} \cup \bar{B}) = \\ &= ((\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \bar{A}) \cup ((\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \bar{B}) = \bar{a} \cup \bar{b} . \end{aligned}$$

Owing to that $((\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \bar{B}) \dot{-} \bar{B} = \emptyset$ and $(\bar{a} \cup \bar{b}) \dot{-} \bar{B} = \bar{a}$ we get

$$((\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \bar{A}) = \bar{a} \quad \text{and using the same argumentation,} \quad ((\bar{z}_{\mathbf{P}} \cup \bar{z}_{\mathbf{Q}}) \cap \bar{B}) = \bar{b} ,$$

or equivalently

$$(\bar{z}_{\mathbf{P}} \cap \bar{A}) \cup (\bar{z}_{\mathbf{Q}} \cap \bar{A}) = \bar{a} \quad \text{and} \quad (\bar{z}_{\mathbf{P}} \cap \bar{B}) \cup (\bar{z}_{\mathbf{Q}} \cap \bar{B}) = \bar{b} .$$

– q. e. d. –

Theorem 2.1.8 Let \mathbf{P} , \mathbf{Q} be concept classes over an set \bar{X} and let $(\forall \bar{a} \in \mathbf{P}) (\bar{X} \dot{-} \bar{a} \in \mathbf{P})$.

Then

$$\text{VC}_{\dim}(\mathbf{U}_{\mathbf{P}, \mathbf{Q}}) \leq \text{VC}_{\dim}(\mathbf{P}) + \text{VC}_{\dim}(\mathbf{Q}) .$$

■ *Proof:*

Assume that \mathbf{P} fulfills assumptions of the theorem, e.g. $\bar{p} \in \mathbf{P} \Leftrightarrow \bar{X} \dot{-} \bar{p} \in \mathbf{P}$. Further let $\text{VC}_{\dim}(\mathbf{P}) = n$ and $\text{VC}_{\dim}(\mathbf{Q}) = m$.

To prove the theorem by contradiction let us assume that there exist mutually different points $\{x_1, \dots, x_{m+n+1}\}$ that are shattered by concept class $\mathbf{U}_{\mathbf{P}, \mathbf{Q}}$. The fact that $\text{VC}_{\dim}(\mathbf{Q}) = m$ implies that there exists an set $\bar{a} \subset \{x_1, \dots, x_{m+1}\}$ which can not be separated from $(\{x_1, \dots, x_{m+1}\} \dot{-} \bar{a})$ by any concept in \mathbf{Q} . Consequently for any set $\bar{Y} \subset \{x_{m+2}, \dots, x_{m+n+1}\}$ there exists an concept $\bar{p} \in \mathbf{P}$, which separates the set $\bar{a} \cup \bar{Y}$ from the set $\{x_1, \dots, x_{m+n+1}\} \dot{-} (\bar{a} \cup \bar{Y})$. Now let us take into mind two different cases:

■ add a)

Let $\bar{a} \neq \emptyset$. Jelikož výše uvedená množina je neprázdná, můžeme z ní vybrat libovolný bod x_0 . Potom ovšem pro libovolnou množinu $\bar{Y} \subset \{x_{m+1}, \dots, x_{m+n+1}\}$ existuje koncept $\bar{p} \in \mathbf{P}$ oddělující $x_0 \cup \bar{Y}$ od $\{x_1, \dots, x_{m+n+1}\} \dot{-} (x_0 \cup \bar{Y})$. Jelikož \mathbf{P} obsahuje s každým svým konceptem i jeho doplněk, zajistíme tímto konceptem i opačné oddělení. Jelikož toto platí pro libovolnou množinu $\bar{Y} \subset \{x_{m+2}, \dots, x_{m+n+1}\}$, je množina $\{x_0, x_{m+1}, \dots, x_{m+n+1}\}$ (jejíž mohutnost je $n + 1$) třídou konceptů \mathbf{P} rozdělena, což je v rozporu s tím, že $\text{VC}_{\dim}(\mathbf{P}) = n$

■ add b)

Let $\bar{a} = \emptyset$. Každá množina $\bar{Y} \subset \{x_{m+2}, \dots, x_{m+n+1}\}$ je odseparovatelná od množiny $\{x_1, \dots, x_{m+n+1}\} - \bar{Y}$ nějakým konceptem z \mathbf{P} . Vezmeme libovolný bod $x_0 \in \{x_1, \dots, x_{m+n+1}\}$. Potom pro libovolnou množinu $\bar{Y} \subset \{x_{m+2}, \dots, x_{m+n+1}\}$ existuje koncept $\bar{p} \in \mathbf{P}$, který tuto množinu odseparuje od množiny $x_0 \cup \{x_{m+1}, \dots, x_{m+n+1}\} - \bar{Y}$. Dále koncept $\bar{X} - \bar{p}$ zajistí opačnou separaci. So the set $\{x_0, x_{m+2}, \dots, x_{m+n+1}\}$ is shattered by concept class \mathbf{P} , which leads to contradiction.

– q. e. d. –

Classes of VC-dimension one

We turn now attention to most simplest, but still applicable, concept classes of the $\mathbf{VC}_{dim}(\mathbf{C}) = 1$. To sake of completeness let us characterize concept classes with $\mathbf{VC}_{dim}(\mathbf{C}) \leq 1$.

Theorem 2.1.9 *Let \mathbf{C} be nonempty concept class over \bar{X} . Then*

1. $\mathbf{VC}_{dim}(\mathbf{C}) = 0$ if and only if \mathbf{C} contains exactly one set.
2. Let one of the following conditions is true:
 - (a) \mathbf{C} is linearly ordered by inclusion, or
 - (b) any two sets in \mathbf{C} are disjoint.

Then $\mathbf{VC}_{dim}(\mathbf{C}) = 1$

■ *Proof:*

■ add 1: \Rightarrow)

Let \mathbf{C} contains a least two different sets \bar{a}, \bar{b} . Without loss of generality we can assume that there exists an point $z \in \bar{a}$ such that $z \notin \bar{b}$. In this case $\bar{a} \cap \{z\} = z$ and $\bar{b} \cap \{z\} = \emptyset$. So we have $\mathbf{VC}_{dim}(\mathbf{C}) = 1$ which is contradiction.

■ add 1: \Leftarrow)

In opposite, let be $\mathbf{C} \stackrel{\text{def}}{=} \{\{\bar{a}\}\}$. Let there exists $y \in \bar{X}$ which is shattered by \mathbf{C} . If $y \in \bar{a}$, then $y \cap \bar{a} \neq \emptyset$, and $\{y\}$ is not shattered by \mathbf{C} . If $y \notin \bar{a}$, then $y \cap \bar{a} = \emptyset$, and resemblely $\{y\}$ is not shattered by \mathbf{C} . Hence no one-point subset of \bar{X} can be shattered by \mathbf{C} .

■ add 2: a)

Let $\{a, b\} \subset \bar{X}$, $a \neq b$, be shattered by \mathbf{C} . So there exist sets $\bar{c}_1, \bar{c}_2 \in \mathbf{C}$ such that $\bar{c}_1 \cap \{a, b\} = \{b\}$ and $\bar{c}_2 \cap \{a, b\} = \{a\}$. But \mathbf{C} is linearly ordered by inclusion, hence without loss of generality, $\bar{c}_1 \subset \bar{c}_2$. It follows $\bar{c}_2 \cap \{a, b\} = \{a, b\}$ and we have a contradiction.

■ add 2: b)

Let $\{a, b\} \subset \bar{X}$, $a \neq b$, be shattered by \mathbf{C} . So there exist sets $\bar{c}_1, \bar{c}_2 \in \mathbf{C}$ such that $\bar{c}_1 \cap \{a, b\} = \{b\}$ and $\bar{c}_2 \cap \{a, b\} = \{a, b\}$. Hence $\bar{c}_1 \cap \bar{c}_2 = \{b\}$ which contradicts assumptions b) of the second part of the theorem.

– q. e. d. –

Theorem 2.1.10 Let \bar{X} be an arbitrary set, $C \subset 2^X$, and $\text{VC}_{\dim}(C) = 1$. Then $\text{VC}_{\dim}(U_{n,C}) \leq n$.

■ *Proof:*

Let $\text{VC}_{\dim}(U_{n,C}) = n + 1$. Hence there exists a set $\bar{A} \subset \bar{X}$, $|\bar{A}| = n + 1$, which is shattered by the concept class $U_{n,C}$. So there exists an $\bar{c} \in U_{n,C}$ in such a way that $\bar{c} \cap \bar{A} = \bar{A}$. Because of the set \bar{c} is union of n sets, $\bar{c} = \bar{c}_1 \cup \dots \cup \bar{c}_n$, and the size of the set \bar{A} is $n + 1$, there exists an index j , such that $|\bar{c}_j \cap \bar{A}| \geq 2$. Therefore there exists $a, b \in \bar{c}_j$ such that $a \neq b$, $a, b \in \bar{A}$. In accordance with the fact that \bar{A} is shattered by $U_{n,C}$, the set $\{a, b\}$ is shattered by concept class $U_{n,C}$ too. Because of $\text{VC}_{\dim}(C) = 1$ the set $\{a, b\}$ is not shattered by concept class C , so

$$(\exists \bar{g} \in 2^{\{a,b\}}) (\forall \bar{c} \in C) (\bar{c} \cap \{a, b\} \neq \bar{g}) . \quad (2.3)$$

In addition $\bar{g} \neq \{a, b\}$ due to $a, b \in \bar{c}_j$. Hence $\bar{g} \in \{\emptyset, \{a\}, \{b\}\}$. Let $\bar{c}_1 \cup \dots \cup \bar{c}_n \in U_{n,C}$ be arbitrary. So we have (accordingly to 2.3)

$$(\bar{c}_1 \cup \dots \cup \bar{c}_n) \cap \{a, b\} = (\bar{c}_1 \cap \{a, b\}) \cup \dots \cup (\bar{c}_n \cap \{a, b\}) \neq \bar{g}$$

because \bar{g} is empty set or one-point set. This contradicts that $\{a, b\}$ is shattered by $U_{n,C}$.

– q. e. d. –

Theorem 2.1.11 Let \bar{X} be an arbitrary set, $C_i \subset 2^X$, $i \in \{1, \dots, n\}$, and set systems C_i are linearly ordered by inclusion. Then $\text{VC}_{\dim}(I_{C_1, \dots, C_n}) \leq n$.

■ *Proof:*

Let $\text{VC}_{\dim}(I_{C_1, \dots, C_n}) = n + 1$. Hence there exists a set $\bar{A} \subset \bar{X}$, $|\bar{A}| = n + 1$, which is shattered by the concept class I_{C_1, \dots, C_n} . Further, let us define the following sets:

$$H = \left\{ \bar{h} \mid (\exists a \in \bar{A}) (\bar{h} = \bar{A} - \{a\}) \right\} , \quad G_i \stackrel{\text{def}}{=} \left\{ \bar{A} \cap \bar{c}_j \mid \bar{c}_j \in C_i, j \geq 1 \right\} ,$$

$$B \stackrel{\text{def}}{=} \left\{ \bar{b} \mid (\exists \bar{c} \in I_{C_1, \dots, C_n}) (\bar{b} = \bar{A} \cap \bar{c} \text{ and } |\bar{b}| = n) \right\} .$$

Straightforwardly, we have

1. The set H contains all subsets \bar{h} of the set \bar{A} with $|\bar{h}| = n$ and $|H| = n + 1$.
2. Each G_i , $i \in \{1, \dots, n\}$, is a finite set system linearly ordered by inclusion.
3. For any $i, j \in \{1, \dots, n\}$ there exist only one set $\bar{e} \in G_i$ such that $|e| = j$.
4. Because \bar{A} is shattered by I_{C_1, \dots, C_n} the equality of collection sets $H = B$ must be true. Hence $|B| = n + 1$.

Now let $\bar{b} \in \mathbf{B}$ be arbitrary. It follows, that $\bar{b} = \bar{A} \cap \bar{c}_1 \cap \dots \cap \bar{c}_n$, where $\bar{c}_i \in \mathbf{C}_i$, $i \in \{1, \dots, n\}$. Obviously $(\forall i \in \{1, \dots, n\}) (\bar{b} \subset \bar{c}_i)$. Let us assume, for a while, that $(\forall i \in \{1, \dots, n\}) (\bar{A} - \bar{b} \in \bar{c}_i)$ (note, that $\bar{A} - \bar{b}$ is one point set). But this imply that $\bar{A} = \bar{A} \cap \bar{c}_1 \cap \dots \cap \bar{c}_n = \bar{b}$, where $\bar{c}_i \in \mathbf{C}_i$, which is contradiction ($\bar{b} \neq \bar{A}$). So there must exists $j \in \{1, \dots, n\}$ such that $\bar{A} \cap \bar{c}_j = \bar{b}$, where $\bar{c}_j \in \mathbf{C}_j$, e.g. $\bar{b} \in \mathbf{G}_j$. As we mentioned in the item 3) on the property list above, such \bar{b} with cardinality n is unique in set system \mathbf{G}_j . Because we have only n set systems \mathbf{G}_j available, we get that $|\mathbf{B}| = n$, which contradicts item 4) of list mentioned. Confessedly the set \bar{B} can not be shattered by the system $\mathbf{l}_{\mathbf{C}_1, \dots, \mathbf{C}_n}$ and $\mathbf{VC}_{dim}(\mathbf{l}_{\mathbf{C}_1, \dots, \mathbf{C}_n}) \leq n$.

– q. e. d. –

Lemma 2.1.12 *Let $\mathbf{C} \subset 2^{\bar{X}}$ be a concept class, let $\mathbf{VC}_{dim}(\mathbf{C}) = d \geq 1$ be finite, and $k \geq 1$. Then $\mathbf{VC}_{dim}(\mathbf{U}_{k, \mathbf{C}}) \leq 2dk \log_2(3k)$ and $\mathbf{VC}_{dim}(\mathbf{l}_{k, \mathbf{C}}) \leq 2dk \log_2(3k)$.*

■ *Proof:*

We prove the statement of the lemma for union only, the case of intersection is resemble. The proof for the case $k = 1$ is trivial so let us assume that $k \geq 2$ and let $\bar{T} \subset \mathbf{X}$ be a finite set, $|\bar{T}| = m \geq 1$. Accordance with lemma 2.1.3, $\Pi_{\mathbf{C}}(|\bar{T}|) \leq \Phi_{d, m}$. Further, every set in the system $\Pi_{\mathbf{U}_{k, \mathbf{C}}}(\bar{T})$ is of the form $\bigcup_{i=1}^k \bar{a}_i$, where $\bar{a}_i \in \Pi_{\mathbf{C}}(\bar{T})$, $1 \leq i \leq k$. It follows that

$$|\Pi_{\mathbf{U}_{k, \mathbf{C}}}(\bar{T})| \leq |\Pi_{\mathbf{C}}(\bar{T})|^k \leq (\Phi_{d, m})^k.$$

If $(\Phi_{d, m})^k < 2^m$, then \bar{T} can not be shattered by $\mathbf{U}_{k, \mathbf{C}}$ and $\mathbf{VC}_{dim}(\mathbf{U}_{k, \mathbf{C}})$ is less than m . So, using statement of the 2.1.3 it suffices to show that $(\frac{em}{d})^{dk} < 2^m$ for $m = 2dk \log_2(3k)$, which is equivalent to the expression $\log_2(3k) < \frac{9k}{2e}$. This inequality is clearly satisfied for the value $k = 2$ and it is obvious that is satisfied for any k greater.

– q. e. d. –

2.2 VC-dimension of linear concepts

The very useful tool which can be used to construct an upper bound estimation of VC-dimension of the concept of halfspaces in Euclidean space \mathfrak{R}^n is the Radon's lemma (see. [Lej85]) which proof is based on the notion of affine independent¹ vector set.

Lemma 2.2.2 (Radon) *Let $\bar{S} \stackrel{\text{def}}{=} \{\bar{x}_1, \dots, \bar{x}_k\} \subset \mathfrak{R}^n$, $k \geq n + 2$, \bar{x}_i are mutually different. Then there exists sets \bar{S}_1 and \bar{S}_2 such that $\bar{S}_1 \cup \bar{S}_2 = \bar{S}$, $\bar{S}_1 \cap \bar{S}_2 = \emptyset$ and*

$$[\bar{S}_1]_{\kappa} \cap [\bar{S}_2]_{\kappa} \neq \emptyset,$$

¹Let \bar{Y} be a vector space of the dimension d . Then a linear combination $\sum_{i=1}^n \alpha_i \bar{x}_i$ of arbitrary finite number of vectors $\bar{x}_1, \dots, \bar{x}_n \in \bar{Y}$ is called affine iff $\sum_{i=1}^n \alpha_i = 1$. The set \bar{A} of all affine combinations of the vectors $\bar{x}_1, \dots, \bar{x}_n$ forms affine hull of vectors $\bar{x}_1, \dots, \bar{x}_n$. Affinne hull is in fact affine subspace of \bar{Y} . Let $\bar{v} \in \bar{A}$. Then the set $\bar{P} \stackrel{\text{def}}{=} \{\bar{x} \in \bar{Y} | (\exists \bar{a} \in \bar{A}) (\bar{x} = \bar{a} - \bar{v})\}$ is a linear subspace of the space \bar{Y} . The dimension of this subspace is called affine rank of the system $\bar{x}_1, \dots, \bar{x}_n$. Obviously, the affine rank of any vector set in d -dimensional vector space can be at most d . The system $\bar{x}_1, \dots, \bar{x}_n$ is called affine independent iff its affine rank is equal to the number $n - 1$. So, any $(d + 2)$ -tuple of vectors in d dimensional space is affine dependent. The following theorem holds:

(the symbol $[\bar{S}]_\kappa$ denotes a convex hull of the set \bar{S}).

■ *Proof:*

The maximal number of affine independent vectors in linear space of dimension n is equal to $n+1$. Hence, any set of $n+2$ vectors is affine dependent. Therefore there exist numbers $\alpha_1, \dots, \alpha_{n+2}$ such that

$$\sum_{i=1}^{n+2} \alpha_i \cdot \vec{x}_i = \vec{0}, \quad \text{and} \quad \sum_{i=1}^{n+2} \alpha_i = 0.$$

In addition, there exists $j \in \{1, \dots, n+2\}$ satisfying $\alpha_j \neq 0$. Without loss of generality we can assume that $\alpha_j > 0$. Further let us define a dichotomy of the set \bar{S} as

$$\bar{S}_1 \stackrel{\text{def}}{=} \{\vec{x}_i \in \bar{S} \mid \alpha_i > 0\} \quad \text{and} \quad \bar{S}_2 \stackrel{\text{def}}{=} \{\vec{x}_i \in \bar{S} \mid \alpha_i \leq 0\},$$

($\sum_{i=1}^{n+2} \alpha_i = 0$ and $\alpha_j > 0$ implies that both of the sets \bar{S}_1, \bar{S}_2 are nonempty). Now we define

$$\omega \stackrel{\text{def}}{=} \sum_{\{i \mid \vec{x}_i \in \bar{S}_1\}} \alpha_i, \quad \vec{z}_1 \stackrel{\text{def}}{=} \sum_{\{i \mid \vec{x}_i \in \bar{S}_1\}} \frac{\alpha_i}{\omega} \cdot \vec{x}_i, \quad \vec{z}_2 \stackrel{\text{def}}{=} \sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} \frac{-\alpha_i}{\omega} \cdot \vec{x}_i.$$

The fact that

$$0 = \sum_{i=1}^{n+2} \alpha_i = \sum_{\{i \mid \vec{x}_i \in \bar{S}_1\}} \alpha_i + \sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} \alpha_i = \omega + \sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} \alpha_i$$

follows $\omega = -\sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} \alpha_i$. Hence

$$1 = \sum_{\{i \mid \vec{x}_i \in \bar{S}_1\}} \frac{\alpha_i}{\omega} = \sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} \frac{-\alpha_i}{\omega}.$$

So, it imply that $\vec{z}_1 \in [\bar{S}_1]_\kappa$ and $\vec{z}_2 \in [\bar{S}_2]_\kappa$. Further

$$\sum_{i=1}^{n+2} \alpha_i \cdot \vec{x}_i = \vec{0} \Rightarrow \sum_{\{i \mid \vec{x}_i \in \bar{S}_1\}} \alpha_i \cdot \vec{x}_i = \sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} -\alpha_i \cdot \vec{x}_i \Rightarrow \sum_{\{i \mid \vec{x}_i \in \bar{S}_1\}} \frac{\alpha_i}{\omega} \cdot \vec{x}_i = \sum_{\{i \mid \vec{x}_i \in \bar{S}_2\}} \frac{-\alpha_i}{\omega} \cdot \vec{x}_i.$$

So finally $\vec{z}_1 = \vec{z}_2$ and $[\bar{S}_1]_\kappa \cap [\bar{S}_2]_\kappa \neq \emptyset$.

– q. e. d. –

We are going to use Radon lemma to prove the following important fact.

Theorem 2.2.3 $\text{VC}_{\dim}(\text{HALFSPACE}_n) = \text{VC}_{\dim}(\text{BALL}_n) = n + 1$.

• **Theorem 2.2.1** Vectors $\vec{x}_1, \dots, \vec{x}_n$ are affine independent iff

$$(\forall \alpha_1, \dots, \alpha_n) \left[\left(\sum_{i=1}^n \alpha_i \vec{x}_i = \vec{0} \quad \text{a} \quad \sum_{i=1}^n \alpha_i = 0 \right) \quad \text{if and only if} \quad \alpha_1 = \alpha_2 = \dots = \alpha_n = 0 \right].$$

■ *Proof:*

■ add HALFSPACE)

First of all let $\{\vec{x}_1, \dots, \vec{x}_k\} \subset \mathfrak{R}^n$, $\vec{a} \in \mathfrak{R}^n$, $t \in \mathfrak{R}$. Let a vector \vec{y} is defined as $\vec{y} \stackrel{\text{def}}{=} \sum_{i=1}^k \alpha_i \vec{x}_i$, where $\sum_{i=1}^k \alpha_i = 1$ and $(\forall i \in \{1, \dots, k\}) (\alpha_i \geq 1)$. Let in addition $\langle \vec{x}_i | \vec{a} \rangle < t$, $i \in \{1, \dots, k\}$. Then

$$\langle \vec{y} | \vec{a} \rangle = \sum_{i=1}^k \alpha_i \langle \vec{x}_i | \vec{a} \rangle < \sum_{i=1}^k \alpha_i \cdot t = t. \quad (2.4)$$

Now we show that $\text{VC}_{dim}(\mathbf{HALFSPACE}_n) \leq n + 1$. Assume contradiction, so put $\text{VC}_{dim}(\mathbf{HALFSPACE}_n) > n + 1$. In other words, $\mathbf{HALFSPACE}_n$ shatters an set \bar{S} containing at least $n + 2$ points. Claim of Radon lemma implies that there exists two disjoint sets \bar{S}_1 a \bar{S}_2 such that $\bar{S}_1 \cup \bar{S}_2 = \bar{S}$, a $[\bar{S}_1]_\kappa \cap [\bar{S}_2]_\kappa \neq \emptyset$. At the same time, \bar{S} is shattered by $\mathbf{HALFSPACE}_n$, hence there exists a hyperplane separating sets \bar{S}_1 and \bar{S}_2 . But such hyperplane must separate whole convex hulls of \bar{S}_1 and \bar{S}_2 (see 2.4) which contradicts Radons lemma.

We finish the proof by construction of a set of $n + 1$ points, which is shattered by $\mathbf{HALFSPACE}_n$. Let $\bar{S} \stackrel{\text{def}}{=} \{\vec{0}, \vec{e}_1, \dots, \vec{e}_n\}$, where $\vec{0}$ is zero vector and \vec{e}_i je i -th vector of standard basis. Let \bar{S}_1 denote arbitrary subset of the set \bar{S} . Now we can define vector $\vec{\alpha} \in \mathfrak{R}^n$ and number t as:

$$\vec{\alpha}_i \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \vec{x}_i \in \bar{S}_1 \\ -1 & \text{if } \vec{x}_i \notin \bar{S}_1 \end{cases} \quad \text{a} \quad t \stackrel{\text{def}}{=} \begin{cases} \frac{1}{2} & \text{if } \vec{0} \in \bar{S}_1 \\ -\frac{1}{2} & \text{if } \vec{0} \notin \bar{S}_1. \end{cases}$$

Apparently a hyperplane $\{\vec{x} | \langle \vec{x} | \vec{\alpha} \rangle - t = 0\}$ in the space \mathfrak{R}^n separate sets $\bar{S} - \bar{S}_1$ and \bar{S}_1 . Since the set \bar{S}_1 was chosen arbitrarily, we get

$$\text{VC}_{dim}(\mathbf{HALFSPACE}_n) = n + 1.$$

■ add BALL)

Let us assume a contradiction, so put $\text{VC}_{dim}(\mathbf{BALL}_n) > n + 1$. In this case there exists $\bar{S} = \{\vec{x}_1, \dots, \vec{x}_{n+2}\} \subset \mathfrak{R}^n$ which is shattered by the concept class \mathbf{BALL}_n . As in the previous case, the Radon lemma implies that there exists two disjoint sets \bar{S}_1 a \bar{S}_2 such that $\bar{S}_1 \cup \bar{S}_2 = \bar{S}$, a $[\bar{S}_1]_\kappa \cap [\bar{S}_2]_\kappa \neq \emptyset$. Because \bar{S} is shattered by \mathbf{BALL}_n , $(\exists \bar{B}_1, \bar{B}_2 \in \mathbf{BALL}_n) (\bar{S}_1 = \bar{B}_1 \cap \bar{S} \text{ and } \bar{S}_2 = \bar{B}_2 \cap \bar{S})$. Let us assume that there exists $j \in \{1, \dots, n + 2\}$ such that $\vec{x}_j \in \bar{B}_1 \cap \bar{B}_2$. Than $\vec{x}_j \in \bar{S}_1$ and at the same time $\vec{x}_j \in \bar{S}_2$, which contradicts the fact that $\bar{S}_1 \cap \bar{S}_2 = \emptyset$. So $\bar{S} \cap (\bar{B}_1 \cap \bar{B}_2) = \emptyset$. Let

$$\bar{B}_1 = \{\vec{x} \in \mathfrak{R}^n \mid \|\vec{x} - \vec{c}_1\|_E \leq r_1\} \quad \text{and} \quad \bar{B}_2 = \{\vec{x} \in \mathfrak{R}^n \mid \|\vec{x} - \vec{c}_2\|_E \leq r_2\},$$

where $\vec{c}_1, \vec{c}_2 \in \mathfrak{R}^n$ and $r_1, r_2 \in \mathfrak{R}^+$ are fixed. Further let us express the intersection \bar{I} of surfaces of the sets \bar{B}_1, \bar{B}_2 as

$$\bar{I} \stackrel{\text{def}}{=} \left\{ \vec{y} \in \mathfrak{R}^n \mid \sum_{i=1}^n \vec{y}_i^2 - 2\vec{y}_i (\vec{c}_1)_i + (\vec{c}_1)_i^2 = r_1 \quad \text{and} \quad \sum_{i=1}^n \vec{y}_i^2 - 2\vec{y}_i (\vec{c}_2)_i + (\vec{c}_2)_i^2 = r_2 \right\}.$$

Now let us take into mind two possible situations,

1. $\bar{I} \neq \emptyset$:

In this case each vektor $\vec{y} \in \bar{I}$ satisfy

$$\sum_{i=1}^n 2\vec{y}_i ((\vec{c}_2)_i - (\vec{c}_1)_i) = r_1 - r_2 - \sum_{i=1}^n ((\vec{c}_1)_i^2 - (\vec{c}_2)_i^2) .$$

Obviously it imply that

$$\bar{I} \subset \left\{ \vec{x} \in \mathfrak{R}^n \mid \langle \vec{x} \mid \vec{c}_2 - \vec{c}_1 \rangle = \frac{r_1 - r_2 - \sum_{i=1}^n ((\vec{c}_1)_i^2 - (\vec{c}_2)_i^2)}{2} \right\} \stackrel{\text{def}}{=} \bar{H} .$$

The set \bar{H} is a hyperplane in the space \mathfrak{R}^n and this hyperplane separate sets $\bar{B}_1 - \bar{B}_2$ and $\bar{B}_2 - \bar{B}_1$, therefore it separates also the sets \bar{S}_1 and \bar{S}_2 which contradicts Radon's lemma. So we have proved that $\text{VC}_{dim}(\mathbf{BALL}_n) \leq n + 1$.

2. $\bar{I} = \emptyset$:

Obviously any hyperplane separate the sets \bar{B}_1, \bar{B}_2 and therefore especially the sets \bar{S}_1, \bar{S}_2 which contradicts Radon's lemma as in previous case.

Now let the set $\bar{S} = \{\vec{x}_1, \dots, \vec{x}_{n+1}\} \subset \mathfrak{R}^n$ be any set shattered by concept class $\mathbf{HALFSPACE}_n$. Let \bar{S}_1, \bar{S}_2 be any dichotomy of the set \bar{S} and let the hyperplane $\{\vec{x} \in \mathfrak{R}^n \mid \langle \vec{x} \mid \vec{w} \rangle = t\}$ separates the sets \bar{S}_1, \bar{S}_2 . Put $\vec{z} \stackrel{\text{def}}{=} t \cdot \frac{\vec{w}}{\|\vec{w}\|_E^2}$. Clearly \vec{z} lies on hyperplane and the sets $\bar{B}_1, \bar{B}_2 \in \mathbf{BALL}_n$ defined as

$$\bar{B}_1 = \left\{ \vec{x} \in \mathfrak{R}^n \mid \left\| \vec{x} - \left(\vec{z} - r \frac{\vec{w}}{\|\vec{w}\|_E} \right) \right\|_E \leq r \right\}$$

and

$$\bar{B}_2 = \left\{ \vec{x} \in \mathfrak{R}^n \mid \left\| \vec{x} - \left(\vec{z} + r \frac{\vec{w}}{\|\vec{w}\|_E} \right) \right\|_E \leq r \right\}$$

separates the sets \bar{S}_1, \bar{S}_2 too for sufficiently large value of radius $r > 0$ because the set \bar{S} is finite and therefore bounded.

– q. e. d. –

2.2.1 Application of Cover's lemma

Now we turn our attention to Cover's lemma. This lemma enumerate the number of dichotomies whose can be obtained by intersection of homogeneous halfspaces with a finite sets of points. To point out the main idea of its proof we start with the following two lemmas.

Definition 2.2.1 Let $\bar{A}, \bar{B} \subset \mathfrak{R}^N$ be given. Then a vector $\vec{w} \in \mathfrak{R}^N$ is **HOMOGENOUS LINEAR SEPARATOR** of sets \bar{A} and \bar{B} iff $(\forall \vec{a} \in \bar{A}) (\langle \vec{w} \mid \vec{a} \rangle > 0)$ and at the same time $(\forall \vec{b} \in \bar{B}) (\langle \vec{w} \mid \vec{b} \rangle < 0)$. In this case, the sets \bar{A} and \bar{B} are called **HOMOGENOUS LINEARLY SEPARABLE SETS**. The set of all homogenously linearly separable tuples of sets (\bar{A}, \bar{B}) will be denoted by the symbol \mathbf{HLS}_{sets} .

Lemma 2.2.4 ([Cov]) Let $\bar{A}, \bar{B} \subset \mathbb{R}^n$ be finite sets, $\vec{y} \in \mathbb{R}^n$, $\vec{y} \neq \vec{0}$. Then pair of sets $\{\bar{A}, \bar{B} \cup \{\vec{y}\}\}$ and $\{\bar{A} \cup \{\vec{y}\}, \bar{B}\}$ are simultaneously homogeneously linearly separable iff there exists hyperplane, containing zero vector and vector \vec{y} , which separate sets \bar{A} and \bar{B} .

■ *Proof:*

Let us denote $\bar{W} \stackrel{\text{def}}{=} \left\{ \vec{w} \in \mathbb{R}^n \mid (\forall \vec{a} \in \bar{A}) (\langle \vec{w} \mid \vec{a} \rangle < 0) \text{ and } (\forall \vec{b} \in \bar{B}) (\langle \vec{w} \mid \vec{b} \rangle > 0) \right\}$.

Then the sets \bar{A} and $\bar{B} \cup \{\vec{y}\}$ are homogeneously linearly separable if and only if there exists a vector $\vec{w}_a \in \bar{W}$ such that $\langle \vec{w}_a \mid \vec{y} \rangle < 0$ and the sets \bar{A} and \bar{B} are homogeneously linearly separable if and only if there exists $\vec{w}_b \in \bar{W}$ such that $\langle \vec{w}_b \mid \vec{y} \rangle > 0$. Hence if we define a vector $\vec{w}_{ab} \stackrel{\text{def}}{=} \langle \vec{w}_a \mid \vec{y} \rangle \vec{w}_b - \langle \vec{w}_b \mid \vec{y} \rangle \vec{w}_a$ then it is straightforward that the sets \bar{A} and \bar{B} are homogeneously linearly separable by hyperplane $\{\vec{x} \in \mathbb{R}^n \mid \langle \vec{x} \mid \vec{w}_{ab} \rangle = 0\}$.

Now prove the opposite implication. Let the sets \bar{A} and \bar{B} be homogeneously linearly separate by hyperplanes containing a vector \vec{y} . Then there exists $\vec{w}_{\vec{y}} \in \bar{W}$ in such a way that $\langle \vec{w}_{\vec{y}} \mid \vec{y} \rangle = 0$. By reason that \bar{W} is an open set, there exists a number $\epsilon > 0$ such that the vectors $\vec{w}_{\vec{y}} + \epsilon \vec{y}$ and $\vec{w}_{\vec{y}} - \epsilon \vec{y}$ belong to the set \bar{W} . Hence the tuple of sets $\{\bar{A}, \bar{B} \cup \{\vec{y}\}\}$ and $\{\bar{A} \cup \{\vec{y}\}, \bar{B}\}$ are both simultaneously homogenous linearly separable by hyperplanes passed by vectors $\vec{w}_{\vec{y}} + \epsilon \vec{y}$ and $\vec{w}_{\vec{y}} - \epsilon \vec{y}$.

– q. e. d. –

The following lemma will be useful in explanation of main step of Cover's lemma proof.

Lemma 2.2.5 Let $\{\vec{x}_1, \dots, \vec{x}_k\}$ is an set of linearly independent vectors of \mathbb{R}^n and let following vectors are defined as:

$$\begin{aligned} \vec{y}_i &\stackrel{\text{def}}{=} \vec{x}_i \mid \vec{x}_k^\perp, \quad \forall i \in \{1, \dots, k-1\}, \\ \vec{y}_k &\stackrel{\text{def}}{=} \vec{x}_k. \end{aligned}$$

Then vectors $\{\vec{y}_1, \dots, \vec{y}_k\}$ are linearly independent.

■ *Proof:*

We prove this lemma via contradiction. Let $\{\vec{y}_1, \dots, \vec{y}_k\}$ be linearly dependent vectors. As for all $i \in \{1, \dots, k-1\}$ holds $\vec{y}_k \perp \vec{y}_i$ vectors $\{\vec{y}_1, \dots, \vec{y}_{k-1}\}$ must be linearly dependent. Hence there exist nontrivial linear combination

$$\sum_{i=1}^{k-1} \alpha_i \vec{y}_i = \vec{0}.$$

Therefore vectors \vec{y}_i are defined as orthogonal projection along vector \vec{x}_k it holds that

$$\vec{y}_i = \vec{x}_i - \beta_i \vec{x}_k.$$

It implies that

$$\sum_{i=1}^{k-1} \alpha_i \vec{x}_i - \sum_{i=1}^{k-1} \alpha_i \beta_i \vec{x}_k = \vec{0}.$$

But $\vec{y}_i \neq \vec{0}$, $i \in \{1, \dots, k-1\}$, hence there exist at least two nonzero α_i which implies contradiction with linear independency of vectors $\{\vec{x}_1, \dots, \vec{x}_k\}$.

– q. e. d. –

Theorem 2.2.6 (Cover, [Cov]) Let $\bar{X} \stackrel{\text{def}}{=} \{\vec{x}_1, \dots, \vec{x}_d\} \subset \mathfrak{R}^N$ are linearly independent vectors. Than there exists

$$2 \sum_{k=0}^{d-1} \binom{N-1}{k}$$

mutually different disjoint splittings of the set \bar{X} into sets \bar{A} and \bar{B} whose are homogeneously linearly separable (i.e. they can be separated via hyperplane which contains zero vector).

■ *Proof:*

We prove this theorem by induction on the number of vectors d and dimension N . Let us denote the number of mutually different homogeneously separable disjoint splittings of the set \bar{X} by the symbol $\tilde{C}(d, N)$. Further let us assume that the set $\bar{Y} \stackrel{\text{def}}{=} \bar{X} \cup \{\vec{x}_{d+1}\}$ fulfills the assumption of the theorem proved and that \bar{A}, \bar{B} be arbitrary disjoint splitting of the set \bar{Y} . Because of the set \bar{X} is finite and linear homogenous separator of the set is defined via nonstrict inequality, at least one of the tuples $(\bar{A} \cup \{\vec{x}_{d+1}\}, \bar{B})$ and $(\bar{A}, \bar{B} \cup \{\vec{x}_{d+1}\})$ are homogeneously linearly separable sets. In addition, it can occur that both of these set tuples are simultaneously homogeneously linearly separable, therefore it is reasonable to define the following set:

$$\bar{S} \stackrel{\text{def}}{=} \left\{ \bar{A} \subset X \mid \bar{A} = \bar{X} - \bar{B} \wedge (\bar{A} \cup \{\vec{x}_{d+1}\}, \bar{B}) \in \mathbf{HLS}_{sets} \wedge (\bar{A}, \bar{B} \cup \{\vec{x}_{d+1}\}) \in \mathbf{HLS}_{sets} \right\}.$$

So if we recapitulate the previous explanation, we see that the number $\tilde{C}(d+1, N)$ can be expressed as the sum of $\tilde{C}(d, N)$ and the number of simultaneously linearly separable dichotomies $\{\bar{A} \cup \{\vec{x}_{d+1}\}, \bar{B}\}$ and $\{\bar{A}, \bar{B} \cup \{\vec{x}_{d+1}\}\}$ of the set $\bar{X} \cup \{\vec{x}_{d+1}\}$, e.g. $|\bar{S}|$. Applying the statement of the lemma 2.2.4 we deduce that the number of such set tuples \bar{A}, \bar{B} is equal to the number of tuples \bar{A}, \bar{B} , that are homogeneously linearly separable by hyperplane containing the vector \vec{x}_{d+1} . Apparently the sets \bar{A}, \bar{B} are separated by such a hyperplane if and only if the sets \bar{A}', \bar{B}' are linearly homogeneously separated, where \bar{A}' and \bar{B}' are orthogonal projection of the sets \bar{A} and \bar{B} into the orthogonal complement of the vector \vec{x}_{d+1} .

Owing to the lemma 2.2.5 the sets \bar{A}', \bar{B}' contains linear independent vectors (projections of vectors $\vec{x}_1, \dots, \vec{x}_d$ into orthogonal complement of the vector \vec{x}_{d+1}). So we can understand these projections as $(N-1)$ dimensional vectors (in orthogonal complement of the vector \vec{x}_{d+1}) and thus we can apply on them the induction assumption. Hence we get the following recursive formula

$$\tilde{C}(d+1, N) = \tilde{C}(d, N) + \tilde{C}(d, N-1). \quad (2.5)$$

On the base of this recursive formula we derive the number of mutually different disjoint splittings of the set $\vec{x}_1, \dots, \vec{x}_d$ in two homogeneously linearly separable parts.

Evidently $\tilde{C}(1, N) = 2$ for all $N \geq 1$, because one-point set $\{\vec{x}\}$ can be split up in two homogenously linearly separable tuples of sets $\{\emptyset, \{\vec{x}\}\}$ and $\{\{\vec{x}\}, \emptyset\}$ only.

To illustrate recursive formula for $\tilde{C}(d, N)$ let us write down these values for initial d and N in the following scheme (the each number in this scheme in a fixed position is equal to the sum of the number on the left and the number in the position which is immediately north-west):

	$d = 1$	2	3	4	5	6	7	8
$N = 1$		2						
2		2	4	8				
3		2	4	8				
4		2	4	8	16			
5		2	4	8	16	32		
6		2	4	8	16	32	64	
7		2	4	8	16	32	64	128

We can deduce from the above scheme of values $\tilde{C}(d, N)$ that the number $\tilde{C}(d, N)$ is equal to double of sum of combinatorial numbers. This observation can be easily verified using equality 2.5 (recall that $\binom{j}{i} = \binom{j-1}{i} + \binom{j-1}{i-1}$):

$$\sum_{k=0}^{N-1} \binom{d}{k} = \sum_{k=0}^{N-1} \binom{d-1}{k} + \sum_{k=0}^{N-2} \binom{d-1}{k} = \sum_{k=0}^{N-1} \binom{d-1}{k} + \sum_{k=1}^{N-1} \binom{d-1}{k-1} = \sum_{k=0}^{N-1} \binom{d}{k}.$$

- q. e. d. -

2.3 VC-dimension of composed mapping

For an arbitrary system of (elementar) functions we can define the VC-dimension of such functions as the VC-dimension of all the sets, which are inverse image of the intervals $(-\infty, \alpha)$. Next, we will show how to use the knowledge of the VC-dimension of such class concepts to determine upper estimate of VC-dimension of function sets contain a composition of those elementar functions. First of all, let us define the term composed mapping, though.

Definition 2.3.1 Let \mathbf{D} be a connected directed acyclic graph (DAG) and \bar{I} denote the set of vertices \mathbf{D} . For a vertex $j \in \bar{I}$ let the numbers d_j^+ and d_j^- denote the number of edges leading into the vertex j and leading from the vertex j , respectively. Let the following properties hold:

1. There exists exactly one vertex with $d_j^- = 0$. We call this vertex an OUTPUT VERTEX .
2. Let there exist at least one vertex with $d_j^+ = 0$. We call such vertex an INPUT VERTEX .
3. For each vertex j with $d_j^+ > 0$ let there exist a mapping $\tilde{Z}_j : \{\mathbb{R}^{d_j^+} \times W_j\} \rightarrow \mathbb{R}$, where \bar{W}_j is a given parametric space of the mapping \tilde{Z}_j .
4. Let for each noninput vertex j the value v_j of the vertex j satisfy the following:

$$v_j = \tilde{Z}_j \left(\left(v_{i_1} \cdots v_{i_{d_j^+}} \right), \bar{\mathbf{w}}_j \right),$$

where $v_{i_1} \cdots v_{i_{d_j^+}}$ are values of vertices from which an edge leads to the vertex j and $\bar{\mathbf{w}}_j \in \bar{W}_j$.

Then, $(\mathbf{D}, (\tilde{Z}_j, \bar{W}_j), j \in \bar{I})$ is called COMPOSED MAPPING . Vertex j is an INTERIOR VERTEX if j is not an input vertex or the output vertex. The number of input vertices is DIMENSION OF COMPOSED MAPPING . Let z be the number of noninput vertices. Each point in $\bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_z$, is called PARAMETRIZATION OF COMPOSED MAPPING .

Example 2.3.1 For example, let $\tilde{s} : \mathbb{R}^3 \times \bar{W}_s \times \bar{W}_f \times \bar{W}_g \times \bar{W}_h \rightarrow \mathbb{R}$, where $\bar{W}_s, \bar{W}_f, \bar{W}_g, \bar{W}_h$ are parameter spaces of mappings $\tilde{s}, \tilde{f}, \tilde{g}, \tilde{h}$, respectively, and

$$\tilde{s} \left(\tilde{f}(x, y, z, \bar{\mathbf{w}}_f), \tilde{g} \left(y, \tilde{h}(z, \bar{\mathbf{w}}_h), \bar{\mathbf{w}}_g \right), \bar{\mathbf{w}}_s \right).$$

The corresponding graph is sketched in Figure 2.1.

In the following text, we will derive basic upper estimate of the VC-dimension of composed mapping using the properties of the functions \tilde{Z}_j . For this purpose, we will have to operate with parts of composed mapping that will be defined by means of what is called proper numbering of vertices of composed mapping.

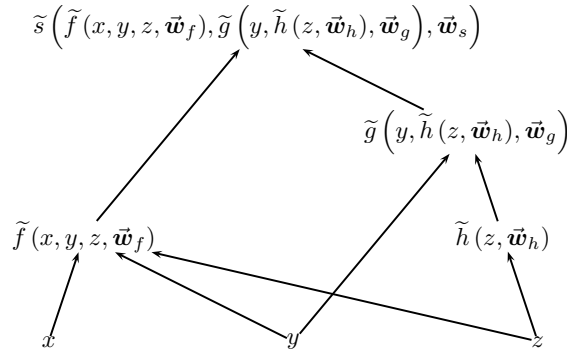


Figure 2.1: Corresponding graph for composed mapping.

Definition 2.3.2 Let G be an acyclic connected directed graph, and noninput vertices be numbered by natural numbers $i \in \{1, \dots, z\}$. Then, this numbering is PROPER NUMBERING iff each edge leads from a vertex with a smaller number.

Proper numbering of a connected directed acyclic graph will be attained for instance in such a way whereby, after removing all the nodes having no input edge, we will first gradually number in a graph thus originated all the nodes having no input edge (this is possible since the graph is acyclic); these nodes will be removed from the graph, and this particular procedure will be repeated on the resultant graph, while gradually numbering, using subsequent numbers.

Definition 2.3.3 Let $(\mathbf{D}, (\tilde{Z}_j, \bar{W}_j), j \in \bar{I})$ be a composed mapping with dimension n , let \mathbf{D} has proper numbering with $(z - 1)$ inner vertices, and let the output vertex have number z . Then, let us define for each $j \in \{1, \dots, z\}$ the following concept class

$$C_j^{loc} \stackrel{def}{=} \left\{ \bar{c} \subset \mathfrak{R}^{d_j^+} \mid (\exists \bar{\mathbf{w}} \in W_j) \left(\bar{c} = \left\{ \bar{\mathbf{s}} \in \mathfrak{R}^{d_j^+} \mid \tilde{Z}_j(\bar{\mathbf{w}}, \bar{\mathbf{s}}) \leq 0 \right\} \right) \right\}.$$

We will call each concept class C_j^{loc} a LOCAL CONCEPT CLASS.

Furthermore, let us denote the value of noninput vertex $j \in \{1, \dots, z\}$ in dependence of $\omega \in \bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_j$ and input values $\bar{\mathbf{x}} \in \mathfrak{R}^n$ of composed mapping as $v_{j,\omega,\bar{\mathbf{x}}}$ and let us define PARTIAL CONCEPT CLASS OF COMPOSED MAPPING as the following system of sets

$$C_j^{par} \stackrel{def}{=} \left\{ \bar{c} \subset \mathfrak{R}^n \mid (\exists \omega \in \bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_j) \left(\bar{c} = \left\{ \bar{\mathbf{x}} \in \mathfrak{R}^n \mid v_{j,\omega,\bar{\mathbf{x}}} \leq 0 \right\} \right) \right\}.$$

Then we will call the number $\text{VC}_{dim}(C_z^{par})$ as the VC-DIMENSION OF COMPOSED MAPPING.

In other words, for a given j and C_j^{loc} is a system of all subsets $\bar{c} \in \mathfrak{R}^{d_j^+}$ such that there exist parameters in \bar{W}_j in such a way that the set \bar{c} is an inverse image of the interval $(-\infty, 0)$ in mapping $\tilde{Z}_j(\bar{\mathbf{w}}, \bar{\mathbf{s}})$. Concept classes C_z^{par} are defined in a similar fashion.

Definition 2.3.4 Let $\bar{V} \stackrel{\text{def}}{=} \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_m\} \subset \mathfrak{R}^n$ and let $(\mathbf{D}, (\bar{Z}_j, \bar{W}_j), j \in \bar{I})$ be a composed mapping of the dimension n with $(n+z)$ vertices and with proper numbering. Further, let $l \in \{1, \dots, z\}$ and for each $\omega \in \bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_z$ let the matrix $\mathbf{R}^{\omega, \bar{V}, l} \in \mathfrak{R}^{l \times m}$ be defined as

$$\mathbf{R}_{i,j}^{\omega, \bar{V}, l} \stackrel{\text{def}}{=} \widetilde{\text{sgn}}(v_{j,\omega} \bar{\mathbf{x}}_i), \quad j \in \{1, \dots, l\}, i \in \{1, \dots, m\}.$$

Then, we say that two parametrizations ω_1 and ω_2 are l, \bar{V} -EQUIVALENT if and only if

$$\mathbf{R}^{\omega_1, \bar{V}, l} = \mathbf{R}^{\omega_2, \bar{V}, l}.$$

Lemma 2.3.1 The relation l, \bar{V} -equivalency defined in 2.3.4 is equivalence relation on the set $\bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_l$.

■ *Proof:*

Reflexivity and symmetry is obvious. Transitivity follows from the fact that if $\widetilde{\text{sgn}}(a) = \widetilde{\text{sgn}}(b)$ and, at the same time, $\widetilde{\text{sgn}}(b) = \widetilde{\text{sgn}}(c)$, then also $\widetilde{\text{sgn}}(a) = \widetilde{\text{sgn}}(c)$.

– q. e. d. –

Theorem 2.3.2 Let the symbol $\mathbf{S}_{l, \bar{V}}$ denote a number of equivalence classes of the l, \bar{V} -equivalence, let m be an arbitrary positive integer and $\bar{V} \stackrel{\text{def}}{=} \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_m\} \subset \mathfrak{R}^n$. Then, for all $k \in \{1, \dots, z\}$ the following estimation

$$\Pi_{\mathbf{C}_k^{\text{par}}}(m) \leq \mathbf{S}_{k, \bar{V}} \leq \Pi_{\mathbf{C}_1^{\text{loc}}}(m) \cdot \Pi_{\mathbf{C}_2^{\text{loc}}}(m) \cdots \Pi_{\mathbf{C}_k^{\text{loc}}}(m)$$

holds.

■ *Proof:*

For a fixed $k \in \{1, \dots, z\}$ let us define the following sets

$$\Gamma \stackrel{\text{def}}{=} \left\{ \mathbf{R}^{\omega, \bar{V}, l} \mid \omega \in \bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_k \right\}$$

and

$$\Gamma_j \stackrel{\text{def}}{=} \left\{ \bar{\mathbf{z}} \in \mathfrak{R}^m \mid (\exists \omega \in \bar{W}_1 \times \bar{W}_2 \cdots \times \bar{W}_z) (\forall i \in \{1, \dots, m\}) (\bar{\mathbf{z}}_i = \mathbf{R}_{i,j}^{\omega, \bar{V}, k}) \right\}, \quad j \in \{1, \dots, k\}.$$

Obviously, the number of equivalence classes is equal to $|\bar{\Gamma}|$ and $|\bar{\Gamma}| \leq |\bar{\Gamma}_1| \cdot |\bar{\Gamma}_2| \cdots |\bar{\Gamma}_k|$ ($\bar{\Gamma}_j$ contains all possible j -th columns of the matrices $\mathbf{R}^{\omega, \bar{V}, k}$). Further, it is obvious that each $\bar{\mathbf{z}} \in \Gamma_j$ define disjoint splitting \bar{V}_-, \bar{V}_+ of the set \bar{V} , where

$$\bar{V}_- = \{\bar{\mathbf{x}}_i \in \bar{V} \mid \bar{\mathbf{z}}_i \leq 0\} \quad \text{and} \quad \bar{V}_+ = \{\bar{\mathbf{x}}_i \in \bar{V} \mid \bar{\mathbf{z}}_i > 0\}.$$

But the number of such disjoint splittings is less or equal to $\left| \Pi_{\mathbf{C}_j^{\text{loc}}}(\bar{V}) \right| \leq \Pi_{\mathbf{C}_j^{\text{loc}}}(m)$, so $|\bar{\Gamma}_j| \leq \Pi_{\mathbf{C}_j^{\text{loc}}}(m)$, $j \in \{1, \dots, k\}$. Hence, $\mathbf{S}_{k, \bar{V}} \leq \Pi_{\mathbf{C}_1^{\text{loc}}}(m) \cdot \Pi_{\mathbf{C}_2^{\text{loc}}}(m) \cdots \Pi_{\mathbf{C}_k^{\text{loc}}}(m)$.

The inequality $|\Pi_{C_k^{\text{par}}}(\bar{V})| \leq S_{k,\bar{V}}$ follows directly from the definition of the number $S_{k,\bar{V}}$ (last column of the matrix $\mathbf{R}^{\omega,\bar{V},k}$ defines all possible disjoint splittings of the set \bar{V}). But this restriction holds for all \bar{V} of the size m , which follows $\Pi_{C_k^{\text{par}}}(m) \leq S_{k,\bar{V}}$.
 – q. e. d. –

Up to this point in our explication we have not used anywhere the properties of partial mappings \tilde{Z}_i of a composed mapping; so far we have used only the fact that the graph of composed mapping is solely acyclic and directed. Theorem 2.3.2 provides instruction on how to proceed in deriving upper estimate of VC-dimension of more general composed mappings since it makes it possible to reduce the process of finding such an estimate into estimates of the VC-dimension of mapping corresponding to the individual nodes of the graph only. We will illustrate this procedure on a simple example of the composition mappings. We start with the following lemma.

Lemma 2.3.3 *Let $\{\alpha_i\}_1^z$ be positive numbers and $\sum_{i=1}^z \alpha_i = 1$. Then,*

$$-\sum_{i=1}^z \alpha_i \ln(\alpha_i) \leq \ln(z). \quad (2.6)$$

■ *Proof:*

Firstly, let $\alpha_i = \frac{1}{z}$, $i \in \{1, \dots, z\}$. Then,

$$-\sum_{i=1}^z \frac{1}{z} \ln\left(\frac{1}{z}\right) = \sum_{i=1}^z \ln\left(z^{\frac{1}{z}}\right) = \ln(z),$$

e.g. for $\alpha_i = \frac{1}{z}$, $i \in \{1, \dots, z\}$, equality holds.

Further, let $\vec{s} \in \left((1, 1, \dots, 1)^T\right)^\perp$. Hence, $\sum_{i=1}^z \vec{s}_i = 0$ and $(\forall t \in \mathfrak{R}) \left(\sum_{j=1}^z \frac{1}{z} + t\vec{s}_j = 1\right)$.

Let $t_{\min}, t_{\max} \in \mathfrak{R}$ be numbers such that

$$(\forall t \in (t_{\min}, t_{\max})) (\forall i \in \{1, \dots, z\}) \left(\frac{1}{z} + t\vec{s}_i > 0\right).$$

Let us define function $\tilde{f}: (t_{\min}, t_{\max}) \rightarrow \mathfrak{R}$ as

$$\tilde{f}(t) \stackrel{\text{def}}{=} \sum_{i=1}^z \left(\frac{1}{z} + t\vec{s}_i\right) \ln\left(\frac{1}{z} + t\vec{s}_i\right).$$

Then (recall $\sum_{i=1}^z \vec{s}_i = 0$)

$$\frac{d\tilde{f}}{dt} = \sum_{i=1}^z \vec{s}_i + \sum_{i=1}^z \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) = \sum_{i=1}^z \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right).$$

Obviously for $t = 0$ the following holds:

$$\sum_{i=1}^z \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) = 0 \quad \text{e.g.} \quad \frac{d\tilde{f}}{dt}(0) = 0.$$

Because $\ln(\cdot)$ is monotonously increasing function, we can write

$$\begin{aligned} t > 0 \quad s > 0 &\Rightarrow \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) > \vec{s}_i \ln\left(\frac{1}{z}\right) \\ s < 0 &\Rightarrow \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) > \vec{s}_i \ln\left(\frac{1}{z}\right) \\ t < 0 \quad s > 0 &\Rightarrow \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) < \vec{s}_i \ln\left(\frac{1}{z}\right) \\ s < 0 &\Rightarrow \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) < \vec{s}_i \ln\left(\frac{1}{z}\right) . \end{aligned}$$

If we take into account the fact $\sum_{i=1}^z \vec{s}_i = 0$ it is straightforward that

$$t > 0 \Rightarrow \sum_{i=1}^z \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) > 0 \quad \text{and} \quad t < 0 \Rightarrow \sum_{i=1}^z \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) < 0.$$

In other words

$$(\forall t \in (t_{min}, t_{max})) \left(t \neq 0 \Rightarrow \frac{d\tilde{f}}{dt}(t) \neq 0 \right) .$$

Finally, let us show that the second derivative of the \tilde{f} in $t = 0$ is positive:

$$\frac{d}{dt} \left(\frac{d\tilde{f}}{dt} \right) = \frac{d}{dt} \left(\sum_{i=1}^z \vec{s}_i \ln\left(\frac{1}{z} + t\vec{s}_i\right) \right) = \sum_{i=1}^z \frac{\vec{s}_i^2}{\frac{1}{z} + t\vec{s}_i} .$$

$$\tilde{f}''(0) = \sum_{i=1}^z z \cdot \vec{s}_i^2 > 0 .$$

Hence, the function \tilde{f} has maximum in $t = 0$ for arbitrary vector $\vec{s} \in \left((1, 1, \dots, 1)^T \right)^\perp$.

– q. e. d. –

Theorem 2.3.4 Let $(\mathbf{D}, (\tilde{Z}_j, \bar{W}_j), j \in \bar{I})$ be a composed mapping with dimension n , let w be the number of edges, z be the number of noninput vertices and $q \stackrel{\text{def}}{=} w + z$. Further let for all noninput vertices k is $\text{VC}_{dim}(\mathbf{C}_k^{\text{loc}}) = d_k^+ + 1$. Then, for any $m > \max_{i \in \{1, \dots, z\}} \{d_i^+ + 1\}$ is

$$\Pi_{\mathbf{C}_z^{\text{par}}}(m) \leq \left(\frac{ezm}{q} \right)^q \quad (2.7)$$

and further the estimation

$$\text{VC}_{dim}(\mathbf{C}_z^{\text{par}}) < 2q \log_2(ez) \quad (2.8)$$

holds.

■ *Proof:*

■ add Proof 2.7)

Let us assume that proper numbering is defined in the graph \mathbf{D} and let d_i^+ denote the number of edges leading to the node i . According to the claim of the lemma 2.1.3 the following estimate holds for each considered i and $m \geq d_i^+ + 1$

$$\Pi_{\mathbf{C}_i}(m) \leq \left(\frac{em}{d_i^+ + 1} \right)^{d_i^+ + 1} ,$$

because by assumption is $\text{VC}_{\dim}(\mathcal{C}_k^{\text{loc}}) = d_k^+ + 1$. But then

$$\Pi_{\mathcal{C}_L}(m) \leq \Pi_{\mathcal{C}_1}(m) \cdot \Pi_{\mathcal{C}_2}(m) \cdots \Pi_{\mathcal{C}_z}(m) \leq \prod_{i=1}^z \left(\frac{em}{d_i^+ + 1} \right)^{d_i^+ + 1}. \quad (2.9)$$

Obviously it holds

$$q = \sum_{i=1}^z (d_i^+ + 1) \quad (2.10)$$

and let us set

$$\alpha_i \stackrel{\text{def}}{=} \frac{d_i^+ + 1}{q}.$$

By substituting into the estimate 2.6, we come to the expression

$$\sum_{i=1}^z \frac{d_i^+ + 1}{q} \ln \left(\frac{q}{d_i^+ + 1} \right) \leq \ln(z).$$

Hence (multiplying by q and using the equality $\ln\left(\frac{a}{b}\right) = \ln(a) - \ln(b)$)

$$\sum_{i=1}^z (d_i^+ + 1) \ln \left(\frac{1}{d_i^+ + 1} \right) \leq q \ln(z) - \left(\sum_{i=1}^z (d_i^+ + 1) \right) \ln(q) = q \ln(z) - q \ln(q).$$

Thus, finally, by delogarithmizing we obtain that

$$\prod_{i=1}^z \left(\frac{1}{d_i^+ + 1} \right)^{d_i^+ + 1} \leq \left(\frac{z}{q} \right)^q.$$

Using the equation 2.10, we can arrange the previous expression into the following form

$$\prod_{i=1}^z \left(\frac{em}{d_i^+ + 1} \right)^{d_i^+ + 1} \leq \left(\frac{z}{q} \right)^q \cdot \prod_{i=1}^z (em)^{d_i^+ + 1} = \left(\frac{z}{q} \right)^q \cdot (em)^q$$

from which the first proven estimate ensues by substituting into 2.9.

■ add Proof 2.8)

First, we will prove the inequality for $z = 1$. Hence, it obviously holds that $\text{VC}_{\dim}(\mathcal{C}_1^{\text{par}}) = w + 1$ and it holds that

$$\text{VC}_{\dim}(\mathcal{C}_1^{\text{par}}) = w + 1 < 2(1 + w) \log_2(ez).$$

Further, we will prove the inequality 2.8 or the case $z \geq 2$. For each real number $a > 4$, the following estimates ($q > 1$) hold

$$2 \log_2(a) < a \Leftrightarrow 2a \log_2(a) < a^2 \Leftrightarrow 2a \log_2(a) < 2^{2 \log_2(a)} \Leftrightarrow (2a \log_2(a))^q < 2^{2q \log_2(a)}.$$

If we substitute $a \stackrel{\text{def}}{=} ez$, $z \geq 2$, then $a > 2.71 \cdot z > 4$ and we obtain that the inequality given below always holds

$$\left(\frac{ez 2q \log_2(ez)}{q} \right)^q < 2^{2q \log_2(ez)}. \quad (2.11)$$

Let $\text{VC}_{\dim}(\mathcal{C}_z^{\text{par}}) = m$. Then, from the definition of the VC-dimension, m is a maximum number such that the equality $2^m = \Pi_{\mathcal{C}_z^{\text{par}}}(m)$ holds; furthermore, if this inequality holds for some m , it holds for all natural m which are smaller. But according to the estimate 2.7 and 2.11 it must also hold that

$$\Pi_{\mathcal{C}_l}(2q \log_2(ez)) \leq \left(\frac{ez 2q \log_2(ez)}{q} \right)^q < 2^{2q \log_2(ez)}$$

and, therefore, an set of the size $2q \log_2(ez)$ cannot be shattered by the class $\mathcal{C}_z^{\text{par}}$, hence $\text{VC}_{\dim}(\mathcal{C}_z^{\text{par}}) < 2q \log_2(ez)$.

– q. e. d. –

2.4 VC-dimension of symmetric difference

Lemma 2.4.1 *Let $\bar{S}, \bar{h}_1, \bar{h}_2 \subset \bar{X}$. Then,*

$$\bar{h}_1 \cap \bar{S} = \bar{h}_2 \cap \bar{S} \quad \Leftrightarrow \quad (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} = (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}$$

■ *Proof:*

Let us assume that $\bar{h}_1 \cap \bar{S} \neq \bar{h}_2 \cap \bar{S}$. Without loss of generality, let $y \in \bar{h}_1 \cap \bar{S}$ and $y \notin \bar{h}_2 \cap \bar{S}$. It follows $y \in \bar{h}_1$, $y \in \bar{S}$ and $y \notin \bar{h}_2$. Hence, $y \notin (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S}$ and $y \in (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}$.

Therefore, $\bar{h}_1 \cap \bar{S} \neq \bar{h}_2 \cap \bar{S}$ implies $(\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} \neq (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}$, which completes the proof.

– q. e. d. –

Lemma 2.4.2 *Let $\bar{S}, \bar{c}, \bar{h}_1, \bar{h}_2 \subset \bar{X}$. Further, let*

$$(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} = (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S} \quad \text{and} \quad \bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} = \bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}. \quad (2.12)$$

Then, $\bar{h}_1 \cap \bar{S} = \bar{h}_2 \cap \bar{S}$.

■ *Proof:*

We prove the claim by contradiction. Let us assume that $\bar{h}_1 \cap \bar{S} \neq \bar{h}_2 \cap \bar{S}$. There is clearly no loss of generality in assuming $y \in \bar{h}_1 \cap \bar{S}$ and $y \notin \bar{h}_2 \cap \bar{S}$. It follows $y \in \bar{h}_1$, $y \in \bar{S}$ and $y \notin \bar{h}_2$. Hence, $y \notin (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S}$. So the left part of 2.12 implies that $y \notin (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S}$. At the same time, $y \in \bar{h}_1 \cap \bar{S}$ which follows $y \notin (\bar{X} \dot{-} \bar{c})$. So $y \in \bar{c}$.

Recalling previous argumentation, we have

$$y \in \bar{c}, \quad y \in \bar{S}, \quad y \notin \bar{X} \dot{-} \bar{h}_1, \quad y \in \bar{X} \dot{-} \bar{h}_2.$$

Put together, these equations leads to

$$y \notin \bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} \quad \text{and} \quad y \in \bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}$$

which contradicts the right part of 2.12.

– q. e. d. –

Lemma 2.4.3 Let $\bar{S}, \bar{c}, \bar{h}_1, \bar{h}_2 \subset \bar{X}$. Further, let

$$\left[(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} \right] \cup \left[\bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} \right] = \left[(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S} \right] \cup \left[\bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S} \right]. \quad (2.13)$$

Then,

$$(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} = (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S} \quad \text{and} \quad \bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} = \bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}.$$

■ *Proof:*

Let us assume that $(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} \neq (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S}$. For example, let there exist y such that $y \in (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S}$ and $y \notin (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S}$. In this case, $y \in (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S}$ implies $y \in \bar{X} \dot{-} \bar{c}$, e.g. $y \notin \bar{c}$, $y \in \bar{h}_1$ and $y \in \bar{S}$. Hence, $y \notin (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S}$ implies that $y \notin \bar{h}_2$. It means that

$$y \in \left[(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} \right] \cup \left[\bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} \right]$$

and

$$y \notin \left[(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S} \right] \cup \left[\bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S} \right]$$

which conclude the proof of the lemma.

– q. e. d. –

Theorem 2.4.4 Let $m \geq 1$, $\bar{c} \subset \bar{X}$ and $\mathbf{H} \subset 2^{\bar{X}}$, $\mathbf{R} \stackrel{\text{def}}{=} \{ \bar{h} \Delta \bar{c} \mid \bar{h} \in \mathbf{H} \}$. Then,

1. $\Pi_{\mathbf{H}}(m) = \Pi_{\mathbf{R}}(m)$ and
2. $\text{VC}_{\dim}(\mathbf{H}) = \text{VC}_{\dim}(\mathbf{R})$. (Note that \mathbf{R} depends on the concept \bar{c} .)

■ *Proof:*

■ add 1)

Let $\bar{h}_1, \bar{h}_2 \in \mathbf{H}$ and $\bar{S}, \bar{c} \subset \bar{X}$ be arbitrary but fixed. Then, by the lemma 2.4.1 is

$$\bar{h}_1 \cap \bar{S} = \bar{h}_2 \cap \bar{S} \quad \Leftrightarrow \quad (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} = (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}. \quad (2.14)$$

From 2.14 it is apparent that

$$(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} = (\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S} \quad \text{and} \quad \bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} = \bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S}. \quad (2.15)$$

Owing to the lemma 2.4.2 (see 2.12) we are going to see that simultaneous validity of equalities 2.15 implies both equalities 2.14.

• Further, it is clear and straightforward that 2.15 implies the following set equality

$$\left[(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_1 \cap \bar{S} \right] \cup \left[\bar{c} \cap (\bar{X} \dot{-} \bar{h}_1) \cap \bar{S} \right] = \left[(\bar{X} \dot{-} \bar{c}) \cap \bar{h}_2 \cap \bar{S} \right] \cup \left[\bar{c} \cap (\bar{X} \dot{-} \bar{h}_2) \cap \bar{S} \right]. \quad (2.16)$$

Similarly as in the previous text, owing to the lemma 2.4.3 (see 2.13), we can see that the opposite implication is also true, e.g. that 2.16 implies 2.15. So 2.16 is equivalent to $\bar{h}_1 \cap \bar{S} = \bar{h}_2 \cap \bar{S}$.

Now recall definition of the symmetric difference of the sets \bar{c} and \bar{h} ,

$$\bar{h} \Delta \bar{c} = (\bar{h} \dot{-} \bar{c}) \cup (\bar{c} \dot{-} \bar{h}) = (\bar{h} \cap (\bar{X} \dot{-} \bar{c})) \cup (\bar{c} \cap (\bar{X} \dot{-} \bar{h})) .$$

It follows (use $(\bar{A} \cup \bar{B}) \cap \bar{S} = (\bar{A} \cap \bar{S}) \cup (\bar{B} \cap \bar{S})$)

$$(\bar{h} \Delta \bar{c}) \cap \bar{S} = [(\bar{h} \cap (\bar{X} \dot{-} \bar{c})) \cap \bar{S}] \cup [(\bar{c} \cap (\bar{X} \dot{-} \bar{h})) \cap \bar{S}] .$$

So, finally, using the last expression, we can rewrite the equation 2.16 to a more legible form

$$(\bar{h}_1 \Delta \bar{c}) \cap \bar{S} = (\bar{h}_2 \Delta \bar{c}) \cap \bar{S} .$$

Thus, we have proved that for arbitrary $\bar{h}_1, \bar{h}_2 \in \mathbf{H}$ and $\bar{S}, \bar{c} \subset \bar{X}$ the equivalence

$$\bar{h}_1 \cap \bar{S} = \bar{h}_2 \cap \bar{S} \quad \Leftrightarrow \quad (\bar{h}_1 \Delta \bar{c}) \cap \bar{S} = (\bar{h}_2 \Delta \bar{c}) \cap \bar{S} \quad (2.17)$$

is fulfilled.

The rest of the proof goes as follows. For the fixed set \bar{S} we show that $|\Pi_{\mathbf{H}}(\bar{S})| = |\Pi_{\mathbf{R}}(\bar{S})|$. To do this let us define a mapping $\tilde{Z} : \Pi_{\mathbf{H}}(\bar{S}) \rightarrow \Pi_{\mathbf{R}}(\bar{S})$ as

$$\tilde{Z}(\bar{u}) \stackrel{\text{def}}{=} (\bar{h}_{\bar{u}} \Delta \bar{c}) \cap \bar{S}, \quad \text{where } \bar{h}_{\bar{u}} \in \{\bar{h} \in \mathbf{H} \mid \bar{h} \cap \bar{S} = \bar{u}\} .$$

First, note that the mapping \tilde{Z} is well defined. Let $\bar{u} \in \Pi_{\mathbf{R}}(\bar{S})$ and $\bar{h}'_{\bar{u}}, \bar{h}''_{\bar{u}} \in \{\bar{h} \in \mathbf{H} \mid \bar{h} \cap \bar{S} = \bar{u}\}$. Then, owing to 2.17, we get

$$\tilde{Z}(\bar{u}) = (\bar{h}'_{\bar{u}} \Delta \bar{c}) \cap \bar{S} = (\bar{h}''_{\bar{u}} \Delta \bar{c}) \cap \bar{S} .$$

Therefore, the image of the set \bar{u} is independent on the choice of the set $\bar{h}_{\bar{u}}$. Also it is straightforward that equivalence 2.17 implies injectivity of the mapping \tilde{Z} . Finally, let $\bar{q} \in \Pi_{\mathbf{R}}(\bar{S})$. Clearly, \bar{q} is rewritable in the form $\bar{q} = (\bar{h}_{\bar{q}} \Delta \bar{c}) \cap \bar{S}$ and, hence, $\tilde{Z}(\bar{h}_{\bar{q}}) = \bar{q}$. In other words, the mapping \tilde{Z} is one-to-one which follows $|\Pi_{\mathbf{H}}(\bar{S})| = |\Pi_{\mathbf{R}}(\bar{S})|$ immediately. The set \bar{S} was chosen arbitrarily, so the definition of the $\Pi_{\mathbf{H}}(\bar{S})$ follows that for any $m \geq 1$ is

$$\Pi_{\mathbf{H}}(m) = \Pi_{\mathbf{R}}(m) .$$

■ add 2)

Let us take the alternative definition $\text{VC}_{dim}(\mathbf{H}) \stackrel{\text{def}}{=} \sup \{m \in N \mid 2^m = \Pi_{\mathbf{H}}(m)\}$. Since the first part of the lemma is proved for arbitrary value of $m \in N$, it is straightforward that for the fixed set \bar{c} (note that \mathbf{R} depends on the set \bar{c}) is

$$\sup \{m \in N \mid 2^m = \Pi_{\mathbf{H}}(m)\} = \sup \{m \in N \mid 2^m = \Pi_{\mathbf{R}}(m)\} .$$

– q. e. d. –

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

Chapter 3

Sample Complexity and VC-dimension

In this chapter, we will deal with the bottom and upper estimates of the number of examples necessary for learning algorithms of the complexity $\tilde{m}(\epsilon, \delta)$ as has been introduced in [BEHW89].

One such estimate for a sufficient quantity of enquiries is given in Theorem 1.2.1. In this chapter, we will show that a dominating role in these estimates is played primarily by the concept VC-dimension studied above. Proceeding from that particular concept, we can formulate better estimates of the necessary length of the sample used for generating a hypothesis approximating the given concept.

First, we will discuss the above-mentioned estimates for a fixed chosen system of concepts and hypotheses. In the subsequent text, we will consider cases when the given concepts are distinguished according to the dimension of space whose subsets they are, and we will analyze sample complexity of learning algorithms with a view to the dimension of vector spaces containing the given concepts as their subsets. A natural requirement is to postulate a criterion guaranteeing that the length of the sample $\tilde{m}(\epsilon, \delta)$ of learning algorithm will be upper bound by polynomial in the dimension of space.

In conclusion, we will carry out a similar analysis concerning the length of words that describe the concepts from the given concept system and which, in their essence, express through their length descriptive complexity of concepts.

3.1 Estimate of the Number of Samples

In this part, we will prove a series of lemmas and theorems, which will be necessary for the proof of Theorem 3.1.12, giving the lower and upper estimate of sample complexity for (ϵ, δ) -algorithms.

Throughout this chapter and in the subsequent ones we will assume that each considered concept class, just as each considered hypothesis class, would be made up solely of Borelian sets. For the sake of comprehensiveness, let us recall the term Borelian sets ([Jar55]). Let us assume that Ω is a non-empty system of sets, which is σ -additive, i.e. for any arbitrary sequence of sets $\bar{A}_n \in \Omega$ it holds that $\bigcup_{n=1}^{\infty} \bar{A}_n \in \Omega$. Furthermore, if it holds that for any arbitrary $\bar{A}, \bar{B} \in \Omega$ there is also $\bar{A} - \bar{B} \in \Omega$, we will call the system of sets Ω

σ -additive circle. It can be illustrated that for each non-empty system of sets Ω there exists just one smallest σ -additive circle containing Ω (this is precisely the intersection of all the σ -additive circles comprising the system Ω). We will call this minimal σ -additive circle the Borelian circle over the system Ω . Let \mathcal{M} be an arbitrary metric space and let Θ be a Borelian circle over a system of all open sets in \mathcal{M} . Then, we will call the elements of this circle BORELIAN SETS of the space \mathcal{M} .

Apart from the requirement regarding the Borelian character of sets which form the individual concepts it will be necessary for other proofs to characterize in greater detail the concepts and hypotheses under consideration. This necessary characteristic can be suitably captured by the term ϵ -transversal.

Definition 3.1.1 For any arbitrary $R \subset 2^{\bar{X}}$ and for any arbitrary probability density \tilde{P} defined on \bar{X} and for an arbitrary $\epsilon > 0$ let us define $R_{\tilde{P},\epsilon} \stackrel{\text{def}}{=} \{\bar{r} \in R \mid \text{Prob}_{\tilde{P}}(\bar{r}) > \epsilon\}$. Then, we will call $\bar{T}_{\tilde{P},\epsilon} \subset \bar{X}$ ϵ -TRANSVERSAL R just when

$$\left(\forall \bar{r} \in R_{\tilde{P},\epsilon} \right) \left(\bar{r} \cap \bar{T}_{\tilde{P},\epsilon} \neq \emptyset \right)$$

Example 3.1.1 Let $\bar{a} \stackrel{\text{def}}{=} \langle 0, 1 \rangle$ and R be a system of all closed intervals on \bar{a} and let \tilde{P} be an uniform probability density on \bar{a} . Then, the set of all the points ϵk , $1 \leq k \leq \frac{1}{\epsilon}$ forms the ϵ -transversal R for any arbitrary value $\epsilon > 0$.

Example 3.1.2 If R is formed by all the open subsets of the interval \bar{a} then - for uniform probability density - there exists no finite ϵ -transversal for any ϵ .

We will examine the probability of the selection of the ϵ -transversal of the system R on the basis of arbitrary random selection of points from \bar{X} . Specifically, we will be interested in the probability of the event described in the following definition.

Definition 3.1.2 For each $m \geq 1$ and $\epsilon > 0$ let $\bar{Q}_{m,\epsilon}$ denote a set of all $\tilde{x} \in \bar{X}^m$ such that the different elements \tilde{x} do not form ϵ -transversal for R , so

$$\bar{Q}_{m,\epsilon} \stackrel{\text{def}}{=} \left\{ \tilde{x} \in \bar{X}^m \mid \left(\exists \bar{r} \in R_{\tilde{P},\epsilon} \right) \left(\tilde{x} \cap \bar{r} = \emptyset \right) \right\}.$$

Further define the set

$$\bar{J}_{\epsilon}^{2m} \stackrel{\text{def}}{=} \left\{ \left(\tilde{x}, \tilde{y} \right) \in \bar{X}^m \times \bar{X}^m \mid \left(\exists \bar{r} \in R_{\tilde{P},\epsilon} \right) \left(\tilde{x} \cap \bar{r} = \emptyset \quad \text{and} \quad \left| \tilde{y} \cap \bar{r} \right| \geq \frac{\epsilon m}{2} \right) \right\}.$$

In the proofs of the probability properties of learning algorithms we will be interested primarily in the property of the system of sets $R \stackrel{\text{def}}{=} \{\bar{h} \Delta \bar{c} \mid \bar{h} \in H\}$, where \bar{c} is some fixed concept $\bar{c} \subset \bar{X}$ and H is the hypothesis class for this particular concept. **The importance of the term ϵ -transversal is in that if the sample of the concept \bar{c} is simultaneously ϵ -transversal for R , then it contains an counter-example for each hypothesis whose error – as seen in terms of the target concept \bar{c} is greater than ϵ .** To be able to examine the properties of the system R , we will take into consideration hypothesis classes meeting the following criterion given in the definition below.

Definition 3.1.3 The hypothesis class \mathbf{H} is WELL-BEHAVED if the sets $\bar{Q}_{m,\epsilon}$ and \bar{J}_ϵ^{2m} are measurable for any arbitrary probability density \tilde{P} , any arbitrary $m \geq 1$, $\epsilon > 0$ and any arbitrary system of sets $\mathbf{R} \stackrel{\text{def}}{=} \{\bar{h} \triangle \bar{c} \mid \bar{h} \in \mathbf{H}\}$, where \bar{c} is an arbitrary Borelian sets.

An overwhelming majority of concept classes and hypothesis classes commonly used is well-behaved. One of the possibilities of verification of this particular property is to verify universal separability.

Definition 3.1.4 The hypothesis class $\mathbf{H} \subset 2^{\bar{X}}$ is UNIVERSALLY SEPARABLE, if there exists a countable subset \mathbf{T} of the class \mathbf{H} such that for all $\bar{h} \in \mathbf{H}$ there exists a sequence $\{\bar{h}_i\}_1^\infty$ of sets from \mathbf{T} such that

$$(\forall x \in \bar{X}) (\exists n \geq 1) ((\forall i \geq n) (x \in \bar{h}_i \text{ if and only if } x \in \bar{h})).$$

The implication formulated in the following theorem holds.

Theorem 3.1.1 If \mathbf{H} is universally separable, then \mathbf{H} is well-behaved.

■ *Proof:*

Let $\bar{c} \subset \bar{X}$ be a Borelian set and $\mathbf{R} \stackrel{\text{def}}{=} \{\bar{h} \triangle \bar{c} \mid \bar{h} \in \mathbf{H}\}$. We will show that the sets $\bar{Q}_{m,\epsilon}$ and \bar{J}_ϵ^{2m} are Borelian. Proof will be given only for $\bar{Q}_{m,\epsilon}$, being similar for \bar{J}_ϵ^{2m} .

Since \mathbf{H} is universally separable, it ensues from the definition that \mathbf{R} is universally separable as well. Let $\mathbf{T} \subset \mathbf{R}$ be a set from the definition of universal separability. Let $\{\gamma_i\}_1^\infty$ be a descending sequence of positive numbers converging to zero, $\{\epsilon_i\}_1^\infty$ be a descending sequence of positive numbers converging to $\epsilon > 0$. Let us define for each $i, j \geq 1$

$$\mathbf{T}_{i,j} \stackrel{\text{def}}{=} \{\bar{t} \in \mathbf{T} \mid (\exists \bar{r} \in \mathbf{R}) (Prob_{\tilde{P}}(\bar{r}) \geq \epsilon_i \text{ and } Prob_{\tilde{P}}(\bar{t} \triangle \bar{r}) \leq \gamma_j)\}.$$

We will show that the following equality holds:

$$\bar{Q}_{m,\epsilon} = \bigcup_{i=1}^{\infty} \bigcap_{\bar{t} \in \mathbf{T}_{i,j}} \{\bar{x} \in \bar{X}^m \mid \bar{x} \cap \bar{t} = \emptyset\}.$$

■ add inclusion \supset :

for \bar{x} from the set on the right there exists such $\bar{t} \in \mathbf{T}_{i,j}$, where $\bar{x} \cap \bar{t} = \emptyset$ furthermore, $\epsilon_i - \gamma_j > \epsilon$. For each such \bar{t} there is $\bar{r} \in \mathbf{R}$ and $Prob_{\tilde{P}}(\bar{r}) > \epsilon$, therefore $\bar{x} \in \bar{Q}_{m,\epsilon}$.

■ add inclusion \subset :

For each $\bar{x} \in \bar{Q}_{m,\epsilon}$ there exists $i \geq 1$ in such a way that for some $\bar{r} \in \mathbf{R}$ there is $\bar{x} \cap \bar{r} = \emptyset$ and $Prob_{\tilde{P}}(\bar{r}) > \epsilon_i$. Since \mathbf{R} is universally separable there exists in \mathbf{T} a sequence of sets converging, in terms of points, to \bar{r} (in the sense of the definition of universal separability), hence for each $j \geq 1$ there exists $\bar{t} \in \mathbf{T}$, for which $Prob_{\tilde{P}}(\bar{t} \triangle \bar{r}) \leq \gamma_j$ and

$\bar{x} \cap \bar{t} = \bar{x} \cap \bar{r} = \emptyset$. Therefore, \bar{x} is in the set on the right.

Hence $\bar{Q}_{m,\epsilon}$ is a Borelian set and the theorem is valid.

– q. e. d. –

We will now prove the following technical lemmas providing elementary inequalities necessary for estimates of the model complexity of learning algorithms (let us establish a convention that \tilde{P}^m is probability density defined on the cartesian product \tilde{X}^m , derived from probability density \tilde{P} defined on \tilde{X}). Chebyshev's inequality will be used for the proof of the immediately following lemma ¹

Lemma 3.1.2 *Let \tilde{X} be a probability space, $m \in \mathbb{N}$, $\bar{r} \subset \tilde{X}$, such that $\text{Prob}(\bar{r}) \geq \epsilon \stackrel{\text{def}}{=} \frac{2}{m}$. Then, the probability that an m -member independent selection of points from \tilde{X} will contain at least one point from \bar{r} is greater than $\frac{1}{2}$.*

■ *Proof:*

The proof is based on the application of Chebyshev's inequality. Let \tilde{A} be a random variable equal to the number of elements belonging to \bar{r} during implementation of m independent random selections from \tilde{X} . This random variable \tilde{A} has a binomic distribution. ²

For the sake of simplicity let us assume that the set \bar{r} has probability ϵ (this can be assumed because if probability of the set \bar{r} is greater, the theorem being proved holds all the more). Then, the random variable \tilde{A} has mean value $\mu = m\epsilon$ and variance $\sigma^2 = m\epsilon(1 - \epsilon)$. Having substituted into Chebyshev's inequality, we get

$$\text{Prob}\left(\left|\tilde{A} - m\epsilon\right| \geq \lambda\right) \leq \frac{m\epsilon(1 - \epsilon)}{\lambda^2} < \frac{m\epsilon}{\lambda^2}.$$

If we substitute $m = \frac{2}{\epsilon}$ a $\lambda \stackrel{\text{def}}{=} 2$, we get

$$\text{Prob}\left(\left|\tilde{A} - 2\right| \geq 2\right) \stackrel{\text{def}}{=} \gamma < \frac{1}{2}.$$

¹Chebyshev's theorem, (see e.g. [And85])

Let \tilde{V} be an arbitrary random variable with a mean value μ and variance σ^2 . Then, it holds that

$$(\forall \lambda > 0) \left(\text{Prob}\left(\left|\tilde{V} - \mu\right| \geq \lambda\right) \leq \frac{\sigma^2}{\lambda^2} \right). \quad (3.1)$$

²Let event \tilde{A} occur with probability p , then supplementary event $\tilde{X} - \tilde{A}$ will occur with probability $1 - p$. Let $\widetilde{B}_{k,n}$ be a phenomenon whereby phenomenon \tilde{A} occurred precisely k -times out of n independent attempts. Then, the random variable $\widetilde{B}_{k,n}$ has a binomic distribution

$$\text{Prob}\left(\widetilde{B}_{k,n}\right) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

Mean value of this distribution is equal to

$$\tilde{E}\left(\widetilde{B}_{k,n}\right) = \sum_{k=0}^n k \binom{n}{k} p^k (1 - p)^{n-k} = np$$

and variance is equal to

$$\tilde{D}^2\left(\widetilde{B}_{k,n}\right) = \tilde{E}\left(\left(\widetilde{B}_{k,n} - \tilde{E}\left(\widetilde{B}_{k,n}\right)\right)^2\right) = np(1 - p).$$

Therefore, the following formula obviously holds for opposite event

$$\text{Prob} \left(\left| \tilde{A} - 2 \right| < 2 \right) = 1 - \gamma > \frac{1}{2}.$$

Since \tilde{A} assumes solely positive integer values, it ensues from the inequality $\left| \tilde{A} - 2 \right| < 2$ that \tilde{A} assumes one of the values 1, 2, 3 with a probability greater than $\frac{1}{2}$.

Before formulating the lemma, which gives an upper bound combinatorial estimate of the probability of the set \bar{J}_ϵ^{2m} , we still have to show special properties of the integral of multiple-variable functions.

Definition 3.1.5 Let π be permutation of the set $\{1, 2, \dots, n\}$ and $(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n) \in \mathfrak{R}^n$. Then define

$$\vec{x}_\pi \stackrel{\text{def}}{=} (\vec{x}_{\pi(1)}, \vec{x}_{\pi(2)}, \dots, \vec{x}_{\pi(n)}) .$$

To show special properties of positive functions with permuted variables we need the following lemma

Lemma 3.1.3 Let $n \geq 1$ be a natural number and

$$\bar{K} \stackrel{\text{def}}{=} \left\{ \vec{a} \in \mathfrak{R}^n \mid \vec{a}_i \geq 0, i \in \{1, \dots, n\}, \text{ and } \sum_{i=1}^n \vec{a}_i \leq 1 \right\} .$$

Then

$$\vec{a} \in \bar{K} \Rightarrow \sum_{i=1}^n i \cdot \vec{a}_i \leq n .$$

■ *Proof:*

Obviously, the lemma is true for all vectors $\vec{a} \in \bar{K}$ of the form $\vec{a} = (0, 0, \dots, 0, \vec{a}_n)^T$, where $0 \leq \vec{a}_n \leq 1$. Further let us assume that there exists an index $j \in \{1, \dots, n-1\}$ such that $0 < \vec{a}_j$. In such a case, $\vec{a}_n \leq 1 - \vec{a}_j < 1$. Now let us define vector $\vec{b} \in \mathfrak{R}^n$ as

$$\vec{b}_i \stackrel{\text{def}}{=} \begin{cases} \vec{a}_i & i \notin \{j, n\} \\ 0 & \text{for } i = j \\ \vec{a}_n + \vec{a}_j & i = n . \end{cases}$$

Obviously, $\vec{b} \in \bar{K}$ and

$$\sum_{i=1}^n i \cdot \vec{b}_i - \sum_{i=1}^n i \cdot \vec{a}_i = (n-j) \vec{a}_j > 0 .$$

So the maximum of the sum $\sum_{i=1}^n i \cdot \vec{a}_i$ can not be reached at any vector $\vec{a} \in \bar{K}$ except the vector $(0, 0, \dots, 0, 1)^T$. It is straightforward that this maximal value is equal to n .

Lemma 3.1.4 Let $\tilde{f} : \mathfrak{R}^n \rightarrow \{0, 1\}$ and $\tilde{p} : \mathfrak{R}^n \rightarrow \mathfrak{R}$ be a nonnegative function, $\int_{\mathfrak{R}^n} \tilde{p}(\vec{x}) dx = 1$ and $\int_{\mathfrak{R}^n} \tilde{f}(\vec{x}) dx$ exists. Let $\bar{\Omega}$ be the set of all permutations of the $\{1, 2, \dots, n\}$. Then,

$$\int_{\mathfrak{R}^n} \sum_{\tilde{\pi} \in \bar{\Omega}} \left(\tilde{f}(\vec{x}_{\tilde{\pi}}) \tilde{p}(\vec{x}_{\tilde{\pi}}) \right) d\vec{x} \leq \max_{\vec{x} \in \mathfrak{R}^n} \left\{ \sum_{\tilde{\pi} \in \bar{\Omega}} \tilde{f}(\vec{x}_{\tilde{\pi}}) \right\}. \quad (3.2)$$

■ *Proof:*

Let us define index set

$$\bar{M} \stackrel{\text{def}}{=} \left\{ 0, 1, 2, \dots, \max_{\vec{x} \in \mathfrak{R}^n} \left\{ \sum_{\tilde{\pi} \in \bar{\Omega}} \tilde{f}(\vec{x}_{\tilde{\pi}}) \right\} \right\}$$

and sets

$$\bar{D}_i \stackrel{\text{def}}{=} \left\{ \vec{x} \in \mathfrak{R}^n \mid \sum_{\tilde{\pi} \in \bar{\Omega}} \tilde{f}(\vec{x}_{\tilde{\pi}}) = i \right\} \quad \text{and} \quad \bar{\Pi}_{\vec{x}} \stackrel{\text{def}}{=} \{ \vec{y} \in \mathfrak{R}^n \mid (\exists \tilde{\pi} \in \bar{\Omega}) (\vec{y} = \vec{x}_{\tilde{\pi}}) \}.$$

for any $i \in \bar{M}$ and $\vec{x} \in \mathfrak{R}^n$.

Obviously, $\bar{D}_i \cap \bar{D}_j = \emptyset$, $i \neq j$, $i, j \in \bar{M}$, and $\cup_{i \in \bar{M}} \bar{D}_i = \mathfrak{R}^n$. Further because $\bar{\Omega}$ is the set of all permutations for any $\vec{y} \in \bar{\Pi}_{\vec{x}}$ is $\bar{\Pi}_{\vec{y}} = \bar{\Pi}_{\vec{x}}$. So if $\vec{x} \in \bar{D}_i$ then $\bar{\Pi}_{\vec{x}} \subset \bar{D}_i$.

Further define relation \mathcal{R} as

$$\vec{x} \mathcal{R} \vec{y} \Leftrightarrow \vec{y} \in \bar{\Pi}_{\vec{x}},$$

where $\vec{x}, \vec{y} \in \mathfrak{R}^n$. Let $\vec{x}, \vec{y}, \vec{z} \in \mathfrak{R}^n$. It is straightforward that $\vec{x} \mathcal{R} \vec{x}$ (identity permutation), $\vec{x} \mathcal{R} \vec{y} \Rightarrow \vec{y} \mathcal{R} \vec{x}$ (inverse permutation) and $\vec{x} \mathcal{R} \vec{y}, \vec{y} \mathcal{R} \vec{z} \Rightarrow \vec{x} \mathcal{R} \vec{z}$ (composition of corresponding permutations). In other words, \mathcal{R} is equivalence relation on \mathfrak{R}^n and $\bar{\Pi}_{\vec{x}}$ are equivalence classes of the equivalence \mathcal{R} . Further, let

$$\bar{Z}_i \stackrel{\text{def}}{=} \{ \bar{\Pi}_{\vec{x}} \mid \vec{x} \in \bar{D}_i \}, \quad i \in \bar{M}.$$

Because $\bar{\Pi}_{\vec{x}}$ are equivalence classes, all sets in \bar{Z}_i form disjoint splitting of the set \bar{D}_i , e.g.

$$\bar{D}_i = \cup_{\bar{t} \in \bar{Z}_i} \bar{t}, \quad i \in \bar{M}. \quad (3.3)$$

Now, for any equivalence class \bar{t} define

$$\bar{\Upsilon}_{\bar{t}} \stackrel{\text{def}}{=} \left\{ \vec{x} \in \bar{t} \mid \tilde{f}(\vec{x}) = 1 \right\}. \quad (3.4)$$

Clearly, if $\vec{x} \in \bar{D}_i$, then size of the set $\bar{\Upsilon}_{\bar{\Pi}_{\vec{x}}}$ is i . Hence

$$\begin{aligned} \int_{\bar{D}_i} \left[\sum_{\tilde{\pi} \in \bar{\Omega}} \tilde{f}(\vec{x}_{\tilde{\pi}}) \tilde{p}(\vec{x}_{\tilde{\pi}}) \right] d\vec{x} &\stackrel{3.3}{=} \int_{\cup_{\bar{t} \in \bar{Z}_i} \bar{t}} \left[\sum_{\tilde{\pi} \in \bar{\Omega}} \tilde{f}(\vec{x}_{\tilde{\pi}}) \tilde{p}(\vec{x}_{\tilde{\pi}}) \right] d\vec{x} \stackrel{3.4}{=} \\ &\stackrel{3.4}{=} \int_{\cup_{\bar{t} \in \bar{Z}_i} \bar{t}} \left[\sum_{\vec{x} \in \bar{\Upsilon}_{\bar{t}}} \tilde{p}(\vec{x}) \right] d\vec{x} \leq i \cdot \int_{\bar{A}_i} \tilde{p}(\vec{x}) d\vec{x} \leq i \cdot \int_{\bar{D}_i} \tilde{p}(\vec{x}) d\vec{x}, \end{aligned}$$

where the sets \bar{A}_i are defined as

$$\bar{A}_i \stackrel{\text{def}}{=} \left\{ \bar{\mathbf{x}} \in \bar{D}_i \mid \bar{\mathbf{x}} = \max_{\bar{\mathbf{y}} \in \bar{\Upsilon}_{\bar{\mathbf{n}}}\bar{\mathbf{x}}} \{ \tilde{p}(\bar{\mathbf{y}}) \} \right\} \subset \bar{D}_i, \quad i \in \bar{M}. \quad (3.5)$$

Finally, for all $i \in \bar{M} - \{0\}$ (note that $\int_{\bar{D}_0} \left(\sum_{\bar{\pi} \in \bar{\Omega}} \tilde{f}(\bar{\mathbf{x}}_{\bar{\pi}}) \tilde{p}(\bar{\mathbf{x}}_{\bar{\pi}}) d\bar{\mathbf{x}} \right)$ is zero) put

$$\alpha_i \stackrel{\text{def}}{=} \int_{\bar{D}_i} \tilde{p}(\bar{\mathbf{x}}) d\bar{\mathbf{x}}.$$

Due to the fact that \bar{D}_i form disjoint splitting of \Re^n and $\int_{\Re^n} \tilde{p}(\bar{\mathbf{x}}) d\bar{\mathbf{x}} = 1$, we have that $\alpha_i \geq 0$, $i \in \bar{M} - \{0\}$, and $\sum_{i=1}^{|\bar{M}|} \alpha_i \leq 1$. If we apply lemma 3.1.3, we get 3.2.

– q. e. d. –

Lemma 3.1.5 *Let \mathbf{R} be a non-empty class concept on \bar{X} , and \tilde{P} probability density on \bar{X} for which $\bar{Q}_{m,\epsilon}$ and \bar{J}_ϵ^{2m} are measurable for arbitrary $m \geq 1$ and $\epsilon > 0$. Then:*

1. *it holds for each $\epsilon > 0$ and $m \geq \frac{2}{\epsilon}$ that*

$$\text{Prob}_{\tilde{P}^m}(\bar{Q}_{m,\epsilon}) < 2 \text{Prob}_{\tilde{P}^m}(\bar{J}_\epsilon^{2m}),$$

2. *it holds for each $\epsilon > 0$ and $m \geq 1$ that*

$$\text{Prob}_{\tilde{P}^m}(\bar{J}_\epsilon^{2m}) \leq \Pi_{\mathbf{R}}(2m) 2^{-\frac{\epsilon m}{2}}.$$

■ *Proof:*

■ *add 1)*

We will show $\text{Prob}_{\tilde{P}^m}(\bar{J}_\epsilon^{2m}) > \frac{1}{2} \text{Prob}_{\tilde{P}^m}(\bar{Q}_{m,\epsilon})$. Let $\tilde{\chi}_{\bar{B}}$ denote a characteristic function of the set \bar{B} . Then, according to Fubini's theorem

$$\begin{aligned} \text{Prob}_{\tilde{P}^m}(\bar{J}_\epsilon^{2m}) &= \int_{\bar{X}^{2m}} \widetilde{\chi_{\bar{J}_\epsilon^{2m}}}(\tilde{x}^{2m}) d\tilde{P}^m(\tilde{x}^{2m}) = \\ &= \int_{\tilde{x} \in \bar{X}^m} \left(\int_{\tilde{y} \in \bar{X}^m} \widetilde{\chi_{\bar{J}_\epsilon^{2m}}}(\tilde{x}, \tilde{y}) d\tilde{P}^m(\tilde{y}) \right) d\tilde{P}^m(\tilde{x}) \geq \\ &\geq \int_{\tilde{x} \in \bar{Q}_{m,\epsilon}} \left(\int_{\tilde{y} \in \bar{X}^m} \widetilde{\chi_{\bar{J}_\epsilon^{2m}}}(\tilde{x}, \tilde{y}) d\tilde{P}^m(\tilde{y}) \right) d\tilde{P}^m(\tilde{x}). \end{aligned}$$

Let $\tilde{x} \in \bar{Q}_{m,\epsilon}$. Then definition of $\bar{Q}_{m,\epsilon}$ follows that there exists an $\bar{r}_{\tilde{x}} \in \mathbf{R}_{\tilde{P},\epsilon}$ such that $\tilde{x} \cap \bar{r}_{\tilde{x}} = \emptyset$. Further let the set $\bar{K}_{2m,\epsilon}^{\bar{r}_{\tilde{x}}}$ be defined as

$$\bar{K}_{2m,\epsilon}^{\bar{r}_{\tilde{x}}} \stackrel{\text{def}}{=} \left\{ (\tilde{x}, \tilde{y}) \in \bar{X}^m \times \bar{X}^m \mid (\tilde{x} \cap \bar{r}_{\tilde{x}} = \emptyset \quad \text{and} \quad |\tilde{y} \cap \bar{r}_{\tilde{x}}| \geq \frac{\epsilon m}{2}) \right\}.$$

Obviously $\bar{K}_{2m,\epsilon}^{\bar{r}_x} \subset \bar{J}_\epsilon^{2m}$, and therefore

$$Prob_{\bar{P}^{2m}}(\bar{J}_\epsilon^{2m}) \geq \int_{\bar{x} \in \bar{Q}_{m,\epsilon}} \left(\int_{\bar{y} \in \bar{X}^m} \chi_{\bar{K}_{2m,\epsilon}^{\bar{r}_x}}(\bar{x}, \bar{y}) d\tilde{P}^m(\bar{y}) \right) d\tilde{P}^m(\bar{x}).$$

For each fixed $\bar{x} \in \bar{Q}_{m,\epsilon}$ the inner integral represents the probability that in m independent selections of binomially distributed variable there occurs an event with probability ϵ at least $\frac{m\epsilon}{2}$ times. According to lemma 3.1.2, for an arbitrary $m \geq \frac{2}{\epsilon}$, this probability is greater than $\frac{1}{2}$. Hence, we get

$$Prob_{\bar{P}^{2m}}(\bar{J}_\epsilon^{2m}) > \int_{\bar{x} \in \bar{Q}_{m,\epsilon}} \frac{1}{2} d\tilde{P}^m(\bar{x}) = \frac{1}{2} Prob_{\tilde{P}^m}(\bar{Q}_\epsilon^m).$$

• add 2)

Let $\bar{\Omega} \stackrel{\text{def}}{=} \{\sigma_1, \dots, \sigma_{2m}!\}$ be the set of all permutation of the set $\{1, \dots, 2m\}$. Evidently

$$Prob_{\bar{P}^{2m}}(\bar{J}_\epsilon^{2m}) = \int_{\bar{X}^{2m}} \chi_{\bar{J}_\epsilon^{2m}}(\bar{x}^{2m}) d\tilde{P}^{2m}(\bar{x}^{2m}) = \int_{\bar{X}^{2m}} \chi_{\bar{J}_\epsilon^{2m}}(\bar{x}_{\sigma_j}^{2m}) d\tilde{P}^{2m}(\bar{x}_{\sigma_j}^{2m})$$

for all the permutations $\sigma_j, j \in \{1, \dots, 2m!\}$, where the symbol $\bar{x}_{\sigma_j}^{2m}$ denotes \bar{x}^{2m} with a permuted order of elements, and $\chi(\bar{a})$ is a characteristic function of the set \bar{a} . Hence, it ensues that

$$Prob_{\bar{P}^{2m}}(\bar{J}_\epsilon^{2m}) = \int_{\bar{X}^{2m}} \frac{1}{(2m)!} \sum_{\sigma_j \in \bar{\Omega}} \chi_{\bar{J}_\epsilon^{2m}}(\bar{x}_{\sigma_j}^{2m}) d\tilde{P}^{2m}(\bar{x}_{\sigma_j}^{2m}).$$

Using the theorem contained in the formula 3.2 and the fact that $\int_{\bar{X}^{2m}} d\tilde{P}^{2m}(\bar{x}) = 1$, we get

$$\int_{\bar{X}^{2m}} \frac{1}{(2m)!} \sum_{\sigma_j \in \bar{\Omega}} \chi_{\bar{J}_\epsilon^{2m}}(\bar{x}_{\sigma_j}^{2m}) d\tilde{P}^{2m}(\bar{x}_{\sigma_j}^{2m}) \leq \frac{1}{2m!} \max_{\bar{x} \in \bar{X}^{2m}} \left\{ \sum_{\sigma_j \in \bar{\Omega}} \chi_{\bar{J}_\epsilon^{2m}}(\bar{x}_{\sigma_j}^{2m}) \right\}. \quad (3.6)$$

It is, therefore, sufficient for proof if we demonstrate the validity of the estimate

$$\frac{1}{2m!} \sum_{\sigma_j \in \bar{\Omega}} \chi_{\bar{J}_\epsilon^{2m}}(\bar{x}_{\sigma_j}^{2m}) \leq \Pi_R(2m) 2^{-\frac{\epsilon m}{2}}$$

for any arbitrary $\bar{x}^{2m} \in \bar{X}^{2m}$.

Let us take fixed $\bar{x} \in \bar{X}^{2m}$ and define the set

$$\bar{\Psi}_{\bar{x}} \stackrel{\text{def}}{=} \left\{ \bar{r} \in R_{\bar{P},\epsilon} \mid (\exists \sigma \in \bar{\Omega}) \left(\text{if } \bar{x}_\sigma^{2m} \stackrel{\text{def}}{=} (\bar{z}, \bar{y}), \text{ then } \bar{r} \cap \bar{z} = \emptyset \text{ and } |\bar{r} \cap \bar{y}| \geq \frac{m\epsilon}{2} \right) \right\}.$$

Further, let us define for each $\bar{r} \in \bar{\Psi}_{\bar{x}}$ the set

$$\bar{\Theta}_{\bar{x},\bar{r}} \stackrel{\text{def}}{=} \left\{ \sigma \in \bar{\Omega} \mid \text{if } \bar{x}_\sigma^{2m} \stackrel{\text{def}}{=} (\bar{z}, \bar{y}), \text{ then } \bar{r} \cap \bar{z} = \emptyset \text{ and } |\bar{r} \cap \bar{y}| \geq \frac{m\epsilon}{2} \right\}. \quad (3.7)$$

Let us establish cardinality of the set $\bar{\Theta}_{\bar{x}, \bar{r}}$. Let us denote $l \stackrel{\text{def}}{=} |\bar{r} \cap \bar{x}| = |\bar{r} \cap \bar{y}|$. Let us now discuss the question in how many ways it is possible to permute elements in \bar{x} in a way to meet the conditions of 3.7 (for an \bar{r} given fixed). The number of these possibilities is evidently equal to the number of ways in which l points may be placed in m positions in \bar{y} multiplied by the number of permutations $(2m - l)$ of points not belonging to \bar{r} , multiplied by the number of permutations of l points that belong to \bar{r} and which are already in selected positions \bar{y} . Hence,

$$|\bar{\Theta}_{\bar{x}, \bar{r}}| = \binom{m}{l} \cdot (2m - l)! \cdot l! . \quad (3.8)$$

Further, the number of different points in \bar{x} at most equals the number $2m$, and it follows that with the help of all sets in $\mathbf{R}_{\bar{r}, \epsilon}$ we can create at best $\Pi_{\mathbf{R}}(2m)$ of different sets in the form of $\bar{r} \cap \bar{x}$. Hence, the estimate holds (in fact it is $\left\{ \bar{v} \mid (\exists \bar{r} \in \bar{\Psi}_{\bar{x}}) (\bar{v} = \bar{r} \cap \bar{x}) \right\} \subset \Pi_{\mathbf{R}}(\bar{x})$)

$$|\bar{\Psi}_{\bar{x}}| \leq \Pi_{\mathbf{R}}(2m) .$$

We are concerned, as our target, with the number of permutations σ such for which $\bar{x}_{\sigma}^{2m} \in \bar{J}_{\epsilon}^{2m}$. However, this number equals the cardinality of the set $\bigcup_{\bar{r} \in \bar{\Psi}_{\bar{x}}} \bar{\Theta}_{\bar{x}, \bar{r}}$, so we can write

$$\sum_{\sigma_j \in \bar{\Omega}} \widetilde{\chi}_{\bar{J}_{\epsilon}^{2m}}(\bar{x}_{\sigma_j}^{2m}) \leq \left| \bigcup_{\bar{r} \in \bar{\Psi}_{\bar{x}}} \bar{\Theta}_{\bar{x}, \bar{r}} \right| \leq |\bar{\Theta}_{\bar{x}, \bar{r}}| |\bar{\Psi}_{\bar{x}}| .$$

If we recall the definition of the set \bar{J}_{ϵ}^{2m} (see 3.1.2), we can see that $l = |\bar{r} \cap \bar{y}| \geq \frac{m\epsilon}{2}$. So we get

$$\frac{1}{(2m)!} |\bar{\Theta}_{\bar{x}, \bar{r}}| |\bar{\Psi}_{\bar{x}}| \stackrel{3.8}{=} \frac{\binom{m}{l}}{\binom{2m}{l}} |\bar{\Psi}_{\bar{x}}| = \frac{m(m-1) \cdots (m-l+1)}{2m(2m-1) \cdots (2m-l+1)} |\bar{\Psi}_{\bar{x}}| \leq \frac{1}{2l} |\bar{\Psi}_{\bar{x}}| \leq \frac{1}{2} \frac{m\epsilon}{2} |\bar{\Psi}_{\bar{x}}|$$

which directly implies that

$$\frac{1}{(2m)!} \sum_{\sigma_j \in \bar{\Omega}} \widetilde{\chi}_{\bar{J}_{\epsilon}^{2m}}(\bar{x}_{\sigma_j}^{2m}) \leq \Pi_{\mathbf{R}}(2m) 2^{-\frac{m\epsilon}{2}} .$$

Because \bar{x} was chosen arbitrarily,

$$\frac{1}{2m!} \max_{\bar{x} \in X^{2m}} \left\{ \sum_{\sigma_j \in \bar{\Omega}} \widetilde{\chi}_{\bar{J}_{\epsilon}^{2m}}(\bar{x}_{\sigma_j}^{2m}) \right\} \leq \Pi_{\mathbf{R}}(2m) 2^{-\frac{m\epsilon}{2}} ,$$

which together with 3.6 completes the proof.

– q. e. d. –

The previous lemma makes it possible to obtain upper estimate of probability of an event when, based on a random sample, we do not obtain ϵ -transversal for the sets in error sets \mathbf{R} .

It is evident, on the basis of lemma 2.1.3, that if the VC-dimension of the hypothesis class H is finite, then the number $\Pi_H(m)$ is polynomially bounded in view of m . That is why the asymptotic behavior of the expression $\Pi_H(2m)2^{-\frac{\epsilon m}{2}}$ in the neighborhood of infinity is given by the exponential member, and obviously this value approaches zero for large values of m very rapidly. Now we will estimate the size m (hence length of the sample) so that the expression $2\Pi_H(2m)2^{-\frac{\epsilon m}{2}}$ is smaller than the number δ .

Lemma 3.1.6 *Let α, β, m be positive numbers, $m \geq \alpha \widetilde{\log}_2(e)^3$ and $m \geq \alpha \widetilde{\log}_2(\beta m)$. Then,*

$$(\forall t > 0) \left(m + t > \alpha \widetilde{\log}_2(\beta(m+t)) \right) . \quad (3.9)$$

■ *Proof:*

Let us define a function $\tilde{f}(m) \stackrel{\text{def}}{=} m - \alpha \widetilde{\log}_2(\beta m)$. Then,

$$\tilde{f}'(m) = 1 - \frac{\alpha \widetilde{\log}_2(e)}{m}.$$

Hence

$$\tilde{f}'(m) = 0 \quad \Rightarrow \quad m = \alpha \widetilde{\log}_2(e)$$

and obviously for all $m > \alpha \widetilde{\log}_2(e)$ the function $\tilde{f}(m)$ is monotone increasing. At the same time, $\tilde{f}(m)$ is positive.

– q. e. d. –

Lemma 3.1.7 *Let α, β, ϵ be positive numbers, $\frac{1}{\epsilon} \geq \alpha \widetilde{\log}_2(e)$ and $\frac{1}{\epsilon} \geq \alpha \widetilde{\log}_2\left(\frac{\beta}{\epsilon}\right)$. Then,*

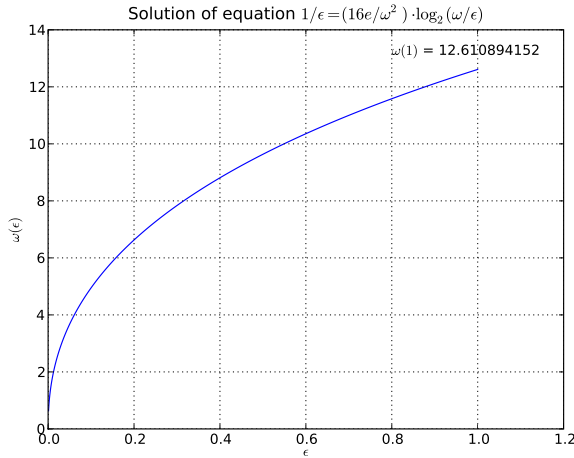
$$(\forall t \in (0, \epsilon)) \left(\frac{1}{\epsilon - t} > \alpha \widetilde{\log}_2\left(\frac{\beta}{\epsilon - t}\right) \right) . \quad (3.10)$$

■ *Proof:*

It is obvious that $\frac{1}{\epsilon - t} = \frac{1}{\epsilon} + \frac{t}{(\epsilon - t)\epsilon}$. Let $m = \frac{1}{\epsilon}$. For such m prepositions of the lemma 3.1.6 are satisfied, which follows the proof.

– q. e. d. –

${}^3\widetilde{\log}_2(e) = \frac{1}{\ln(2)} \approx 1.44269504088896$



	+ 0.0	+0.01	+0.02	+0.03	+0.04
0.00	–	1.81	2.46	2.94	3.33
0.05	3.67	3.97	4.25	4.50	4.73
0.10	4.95	5.15	5.35	5.53	5.71
0.15	5.88	6.04	6.19	6.34	6.49
0.20	6.63	6.76	6.90	7.02	7.15
0.25	7.27	7.39	7.51	7.62	7.73
0.30	7.84	7.94	8.05	8.15	8.25
0.35	8.35	8.44	8.54	8.63	8.72
0.40	8.81	8.90	8.98	9.07	9.15
0.45	9.24	9.32	9.40	9.48	9.56

Figure 3.1: Graph of the solution of equation $\frac{1}{\epsilon} = \frac{16e}{\omega^2(\epsilon)} \log_2 \left(\frac{\omega(\epsilon)}{\epsilon} \right)$. The corresponding source code in python language is listed in Appendix 4.1.1.

Lemma 3.1.8 Let m be a natural number, $0 < \epsilon < 1$ and

$$m \geq \max \left\{ \frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right\}. \quad (3.11)$$

Then, $2\Phi_{d,2m} 2^{-\frac{\epsilon m}{2}} \leq \delta$.

■ *Proof:*

The lemma 2.1.3 postulates the inequality $\Phi_{d,2m} \leq \left(\frac{2em}{d} \right)^d$, thus it suffices to show that $2 \left(\frac{2em}{d} \right)^d \leq \delta 2^{\frac{m\epsilon}{2}}$. This is equivalent to the inequality

$$\frac{m\epsilon}{2} \geq d \log_2 \left(\frac{2me}{d} \right) + \log_2 \left(\frac{2}{\delta} \right).$$

The first condition on m in 3.11 implies $\frac{m\epsilon}{4} \geq \log_2 \left(\frac{2}{\delta} \right)$, so the rest of the proof is reduced to showing

$$\frac{m\epsilon}{4} \geq d \log_2 \left(\frac{2me}{d} \right). \quad (3.12)$$

To do so, let us define (see graph on Figure 3.1)

$$\alpha_\epsilon \stackrel{\text{def}}{=} \frac{16e}{12.611^2}, \quad \beta_\epsilon \stackrel{\text{def}}{=} 12.611.$$

It is obvious that

$$(\forall \epsilon \in (0, 1)) \left(\frac{1}{\epsilon} \geq 1 > \frac{16e}{12.611^2} \log_2(e) = \alpha_\epsilon \log_2(e) \simeq 0.394539011 \right). \quad (3.13)$$

At the same time numerical calculation shows (again see graph on Figure 3.1) that for $\epsilon = 1$ is

$$\frac{1}{\epsilon} \geq \frac{16e}{12.611^2} \log_2 \left(\frac{12.611}{\epsilon} \right). \quad (3.14)$$

Inequalities 3.13 and 3.14 follows that both assumptions of the lemma 3.1.7 are true and hence the inequality 3.14 is valid for any $\epsilon \in (0, 1)$.

Further let us define

$$\alpha_m \stackrel{\text{def}}{=} \frac{4d}{\epsilon}, \quad \beta_m \stackrel{\text{def}}{=} \frac{2e}{d}, \quad m = \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right).$$

Hence we get for any $\epsilon \in (0, 1)$

$$\begin{aligned} \frac{1}{\epsilon} &\geq \frac{16e}{12.611^2} \log_2 \left(\frac{12.611}{\epsilon} \right) \Leftrightarrow 2 \log_2 \left(\frac{12.611}{\epsilon} \right) \geq \log_2 \left(\frac{16e}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right) \Leftrightarrow \\ &\Leftrightarrow \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \geq \frac{4d}{\epsilon} \log_2 \left(\frac{16e}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right) \Leftrightarrow m \geq \alpha_m \log_2 (\beta_m m). \end{aligned} \quad (3.15)$$

At the same time, for any $0 < \epsilon < 1$ it is straightforward that

$$m = \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \geq \frac{4d}{\epsilon} \log_2 (e) = \alpha_m \log_2 (e). \quad (3.16)$$

Resembly as in the previous case inequalities 3.15 and 3.16 follows that both assumptions of the lemma 3.1.6 are true and hence the last inequality 3.15 is valid for any

$$m \geq \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right).$$

But

$$m \geq \alpha_m \log_2 (\beta_m m) \Leftrightarrow \frac{m\epsilon}{4} \geq d \log_2 \left(\frac{2me}{d} \right),$$

and finally

$$\frac{m\epsilon}{2} = \frac{m\epsilon}{4} + \frac{m\epsilon}{4} \geq d \log_2 \left(\frac{2me}{d} \right) + \log_2 \left(\frac{2}{\delta} \right).$$

– q. e. d. –

Lemma 3.1.9 *Let the matrix \mathbf{Z} be a real matrix of the type $m \times n$ and let it hold that $(\forall i \in \{1, \dots, m\}) \left(\sum_{j=1}^n \frac{\mathbf{Z}_{i,j}}{n} > \alpha \right)$. Then, there exists $j_0 \in \{1, \dots, n\}$ so that*

$$\sum_{i=1}^m \frac{\mathbf{Z}_{i,j_0}}{m} > \alpha.$$

■ *Proof:*

Proof will be made by means of contradiction. It ensues from the assumption of the lemma that $\sum_{i=1}^m \sum_{j=1}^n \mathbf{Z}_{i,j} > \alpha nm$. If the lemma's theorem did not hold, then – on the contrary – $\sum_{j=1}^n \sum_{i=1}^m \mathbf{Z}_{i,j} < \alpha mn$ would have to hold, which is a contradiction.

– q. e. d. –

Lemma 3.1.10 Let $\bar{X} \stackrel{\text{def}}{=} \{x_1, \dots, x_d\}$ be a finite set, $k \in \{1, \dots, d\}$, and $0 \leq \omega \leq 1$. Further let

$$\mathbf{P} \stackrel{\text{def}}{=} \{\bar{a} \subset \bar{X} \mid |\bar{a}| \geq k\} \quad \text{and} \quad \mathbf{Q} \stackrel{\text{def}}{=} \{\bar{b} \subset \bar{X} \mid 1 \leq |\bar{b}| < k\}.$$

Let us randomly select a set $\bar{s} \subset \bar{X}$ and let the probability of getting $\bar{s} \in \mathbf{P}$ be equal to ω . Finally, let us denote the average size of selected sets as ρ . Then,

$$\rho \leq k \cdot (1 - \omega) + d \cdot \omega. \quad (3.17)$$

Further, let on the set \bar{X} be defined uniform probability \tilde{P} and let ϑ denotes average probability of the selected sets. Than

$$\vartheta \leq \frac{k}{d} \cdot (1 - \omega) + \omega. \quad (3.18)$$

■ *Proof:*

Denote selected sets as $\mathbf{C} \stackrel{\text{def}}{=} \{\bar{c}_1, \dots, \bar{c}_p\}$, where p is the number of trials. Obviously

$$\begin{aligned} \rho &= \frac{\sum_{i=1}^p |\bar{c}_i|}{p} = \frac{\sum_{\bar{c}_i \in \mathbf{Q}} |\bar{c}_i|}{p} + \frac{\sum_{\bar{c}_i \in \mathbf{P}} |\bar{c}_i|}{p} \leq \\ &\leq \frac{k \cdot |\{i \mid \bar{c}_i \in \mathbf{Q}\}|}{p} + \frac{d \cdot |\{i \mid \bar{c}_i \in \mathbf{P}\}|}{p} = k \cdot (1 - \omega) + d \cdot \omega. \end{aligned}$$

The proof of the second estimation is straightforward.

– q. e. d. –

Lemma 3.1.11 Let for any $\bar{z} \in \{0, 1\}^n$ the symbol $n_+(\bar{z})$ denotes the number of ones in \bar{z} . Than

$$\frac{1}{2^n} \sum_{\bar{z} \in \{0, 1\}^n} n_+(\bar{z}) = \frac{n}{2}.$$

■ *Proof:*

Let $\bar{o} \in \{0, 1\}^n$ is defined as $\bar{o} \stackrel{\text{def}}{=} (1, 1, \dots, 1)$ and let \bar{A}, \bar{B} be sets such that for all $\bar{a} \in \bar{A}$ is the vector $-(\bar{a} - \bar{o}) \in \bar{B}$ and for all $\bar{b} \in \bar{B}$ is the vector $-(\bar{b} - \bar{o}) \in \bar{A}$. Hence the sets \bar{A}, \bar{B} form disjoint splitting of the set $\{0, 1\}^n$ and in addition $|\bar{A}| = |\bar{B}|$. So we have

$$\frac{1}{2^n} \sum_{\bar{z} \in \{0, 1\}^n} n_+(\bar{z}) = \frac{1}{2^n} \sum_{\bar{a} \in \bar{A}} n_+(\bar{a}) - n_+(\bar{a} - \bar{o}) = \frac{1}{2^n} \sum_{\bar{a} \in \bar{A}} n = \frac{1}{2^n} \cdot \frac{2^n}{2} \cdot n = \frac{n}{2}.$$

– q. e. d. –

The lemmas and theorems postulated above make it possible directly to prove the fundamental theorem of this part. However, let us omit in our considerations the trivial cases described in the following definition.

Definition 3.1.6 Concept class C defined over the set \bar{X} is called **NONTRIVIAL CONCEPT CLASS** iff

$$(\exists \bar{c}_1, \bar{c}_2 \in C) \text{ such that } \bar{c}_1 \neq \bar{c}_2 \text{ and } (\bar{c}_1 \cap \bar{c}_2 \neq \emptyset \text{ or } \bar{c}_1 \cup \bar{c}_2 \neq \bar{X}).$$

Concept class is called **TRIVIAL CONCEPT CLASS** in other cases.

So what is a trivial concept class? If we negate definition condition, we get

$$(\forall \bar{c}_1, \bar{c}_2 \in C) (\bar{c}_1 = \bar{c}_2 \text{ or } (\bar{c}_1 \cap \bar{c}_2 = \emptyset \text{ and } \bar{c}_1 \cup \bar{c}_2 = \bar{X})).$$

Clearly, we can reformulate the definition 3.1.6 in the form that the concept class is trivial if and only if it consists of one concept only, or has two concepts whose form disjoint splitting of the set \bar{X} .

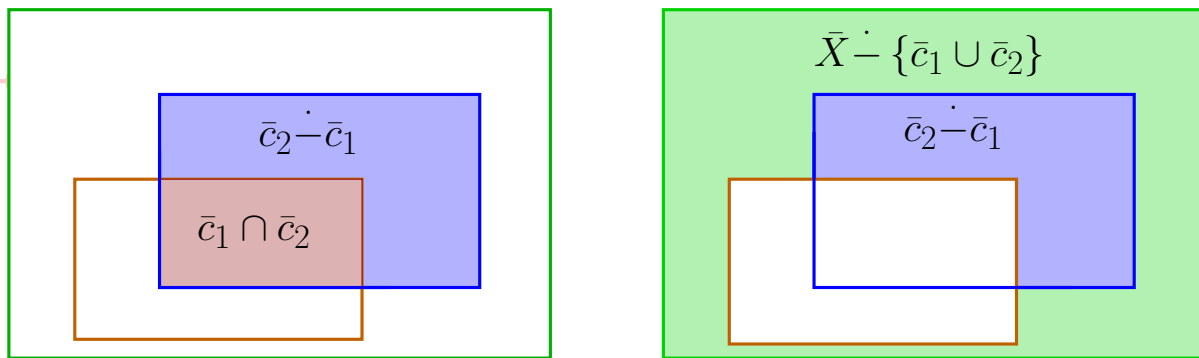


Figure 3.2: Two necessary minimal contents of nontrivial concept classes (stained areas must be nonempty sets).

It can be expected that we focus on nontrivial classes. The next theorem postulates upper bound on necessary queries that have to be provided by an environment (we mean some kind of supervision) to learning algorithm to be an (ϵ, δ) -algorithm, which is the main result of the standard PAC learning theory.

Theorem 3.1.12 Let C be a nontrivial, well-behaved concept class. Then, the following holds:

1. If $VC_{dim}(C) = d$ and $d < \infty$ then

(a) for any $0 < \epsilon < \frac{1}{2}$ there is no (ϵ, δ) -learning algorithm which exploits less than

$$\max \left(\frac{1 - \epsilon}{\epsilon} \ln \left(\frac{1}{\delta} \right), d(1 - 2(\epsilon(1 - \delta) + \delta)) \right) \quad (3.19)$$

queries.

(b) for arbitrary $0 < \epsilon < 1$, any learning algorithm using at least

$$\max \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right) \quad (3.20)$$

queries and returning a consistent hypothesis only is an (ϵ, δ) -learning algorithm.

2. \mathcal{C} is uniformly learnable if and only if $\text{VC}_{\text{dim}}(\mathcal{C}) < \infty$.

■ *Proof:*

Let \mathcal{C} be a nontrivial concept class and let $\text{VC}_{\text{dim}}(\mathcal{C}) = d$ be finite.

■ add Proof of 3.19)

Let $0 < \epsilon < \frac{1}{2}$ and $m \stackrel{\text{def}}{=} \max\left(\frac{1-\epsilon}{\epsilon} \ln\left(\frac{1}{\delta}\right), d(1 - 2(\epsilon(1 - \delta) + \delta))\right)$. We have to prove that any arbitrary learning algorithm for \mathcal{C} must exploit at least m queries. Proof will be divided into the following steps:

Proof for the first member of the estimate: Proof of the estimate given by the first member in 3.19 will be made for the following two cases:

Case 1: \mathcal{C} contains at least two different concepts \bar{c}_1 and \bar{c}_2 which have a non-empty intersection. Let A^* be a learning algorithm for \mathcal{C} and let us denote $\bar{a} \stackrel{\text{def}}{=} \bar{c}_1 \cap \bar{c}_2$ and $\bar{b} \stackrel{\text{def}}{=} \bar{c}_2 - \bar{c}_1$. Let \tilde{P} be such a probability density on \bar{X} that $\text{Prob}_{\tilde{P}}(\bar{b}) \stackrel{\text{def}}{=} \epsilon$, $\text{Prob}_{\tilde{P}}(\bar{a}) = 1 - \epsilon$ and let the probability of all the other elements from \bar{X} be zero. In view of the probability density thus defined, we can assume in further text – without loss of generality – that $\bar{X} \stackrel{\text{def}}{=} \bar{a} \cup \bar{b}$, $\mathcal{C} \stackrel{\text{def}}{=} \{\bar{a}, \bar{X}\}$ and $\mathcal{H} \stackrel{\text{def}}{=} \{\emptyset, \bar{a}, \bar{b}, \bar{X}\}$.

According to the first bound of the number of selected samples, $m \leq \frac{1-\epsilon}{\epsilon} \ln\left(\frac{1}{\delta}\right)$. Using the inequality⁴

$$\frac{1-\epsilon}{\epsilon} < \frac{-1}{\ln(1-\epsilon)},$$

we get

$$\begin{aligned} m &< \frac{-1}{\ln(1-\epsilon)} \ln\left(\frac{1}{\delta}\right), \\ m \ln(1-\epsilon) &> \ln(\delta), \\ (1-\epsilon)^m &> \delta. \end{aligned}$$

The last inequality says that the probability of selecting all the samples from the set \bar{a} is greater than δ .

All the possible learning algorithms for the concept class \mathcal{C} can obviously be divided into two groups:

⁴Power series of the function $\ln(1-\epsilon)$ is

$$\ln(1-\epsilon) = -\epsilon - \frac{\epsilon^2}{2} - \frac{\epsilon^3}{3} - \frac{\epsilon^4}{4} \dots, \quad \text{for } -1 < \epsilon < 1.$$

Taking into mind that this power series is absolute convergent (we can multiply it by any polynomial and apply arbitrary bracketing), we have

$$(1-\epsilon)\ln(1-\epsilon) = -\epsilon + \left(\epsilon^2 - \frac{\epsilon^2}{2}\right) + \left(\frac{\epsilon^3}{2} - \frac{\epsilon^3}{3}\right) + \left(\frac{\epsilon^4}{3} - \frac{\epsilon^4}{4}\right) + \dots > -\epsilon,$$

• because $0 < \epsilon < 1$, and hence all the numbers in brackets are positive. So $(1-\epsilon)\ln(1-\epsilon) > -\epsilon$, which can be expressed as

$$\frac{1-\epsilon}{\epsilon} < \frac{-1}{\ln(1-\epsilon)}.$$

1. algorithms which, in case of selecting precisely m samples from \bar{a} , generate hypothesis \bar{a} . Such algorithms produce for the target concept \bar{X} a hypothesis with an error equal to $Prob_{\tilde{P}}(\bar{b}) = \epsilon$.
2. algorithms, which – in this instance – generate hypothesis \bar{b} , \bar{X} or \emptyset . However, since – according to the assumption of the theorem $\epsilon < \frac{1}{2}$ – such algorithms produce, in case of the target concept \bar{a} , an error at least equal to ϵ .

We have, therefore, shown that there exists the concept class \mathbf{C} such that for any learning algorithm A^* there exists the target concept (\bar{x} in the group 1.) and \bar{a} in the group 2.)) for which we get the hypothesis with error greater than ϵ with probability greater than δ . Thus, there cannot exist an (ϵ, δ) -learning algorithm for \mathbf{C} , requiring less than $\frac{1-\epsilon}{\delta} \ln\left(\frac{1}{\delta}\right)$ queries.

Case 2: \mathbf{C} contains at least two different concepts \bar{c}_1 and \bar{c}_2 for which $\bar{c}_1 \cup \bar{c}_2 \neq \bar{X}$. Let us put $\bar{a} \stackrel{\text{def}}{=} \bar{X} - (\bar{c}_1 \cap \bar{c}_2)$ and $\bar{b} \stackrel{\text{def}}{=} \bar{c}_1$. Let us define probability density \tilde{P} just as in the previous case, e.g. $Prob_{\tilde{P}}(\bar{b}) \stackrel{\text{def}}{=} \epsilon$, $Prob_{\tilde{P}}(\bar{a}) = 1 - \epsilon$ and the probability of all the other elements from \bar{X} is zero. Put $\bar{X} \stackrel{\text{def}}{=} \bar{a} \cup \bar{b}$, $\mathbf{C} \stackrel{\text{def}}{=} \{\bar{b}, \emptyset\}$, $\mathbf{H} \stackrel{\text{def}}{=} \{\emptyset, \bar{a}, \bar{b}, \bar{X}\}$. Further, let us take into consideration random selection of m negative samples which are all from the set \bar{a} . The rest of the proof is in the same manner as in the preceding case.

Proof for the second member of the estimate: Now what remains to be done is to prove the validity of the second estimate in the expression 3.19. Since \mathbf{C} is a nontrivial class concept, $VC_{dim}(\mathbf{C}) = d$ is equal at least to 1. Therefore, there exists a d -element subset $\bar{\Gamma} \subset \bar{X}$, which is shattered by the class concept \mathbf{C} . Let probability density \tilde{P} be uniform on the set $\bar{\Gamma}$ and zero on its complement to the set \bar{X} (thanks to this selection the entire proof is based solely on the combinatorial properties of finite sets). In view of this probability density we can again put – without loss of to generality – $\bar{X} \stackrel{\text{def}}{=} \bar{\Gamma}$ and $\mathbf{C} \stackrel{\text{def}}{=} 2^{\bar{\Gamma}}$. Let us assume that some learning algorithm A^* has obtained a sample of the length of m elements from the set \bar{X} , and let us denote this sample (\tilde{x}, \tilde{z}) , $\tilde{x} \in \bar{X}^m$, while the number of mutually different elements x_i in the sample \tilde{x} is l . Let us define the set system of all concepts consistent with the sample (\tilde{x}, \tilde{z})

$$\mathbf{B}_{(\tilde{x}, \tilde{z})} \stackrel{\text{def}}{=} \left\{ \bar{c} \in \mathbf{C} \mid \bar{c} \text{ is consistent with } (\tilde{x}, \tilde{z}) \right\}.$$

If we subtract \tilde{x} from \bar{X} , there remains a set of the size $(d - l)$, which has 2^{d-l} subsets. Each of these subsets, unified with \tilde{x} , is a consistent hypothesis for the sample \tilde{x} . Thus, $|\mathbf{B}_{(\tilde{x}, \tilde{z})}| = 2^{d-l}$. Let us consider that the algorithm A^* always produce consistent hypotheses. Therefore

$$\left(\exists \bar{h} \in \mathbf{B}_{(\tilde{x}, \tilde{z})} \right) \left(\forall \bar{c} \in \mathbf{B}_{(\tilde{x}, \tilde{z})} \right) \left(A^* \left((\tilde{x}, \tilde{z}) \right) = \bar{h} \right),$$

in the other words, for the sample (\tilde{x}, \tilde{z}) the algorithm A^* must generate a hypothesis from $\mathbf{B}_{(\tilde{x}, \tilde{z})}$ identical one for all the target concepts from $\mathbf{B}_{(\tilde{x}, \tilde{z})}$ (the

result of the algorithm depends solely on its definition and on the sample used). The number of elements in which this hypothesis can differ from the target concept equals $d - l$. To explain this fact in a detail we can denote those elements as $\{y_1, \dots, y_{d-l}\} = \bar{X} - \bar{x}$. Obviously, all points y_i can be placed into sets $\bar{h} \cap \bar{c}$, $\bar{X} - (\bar{h} \cup \bar{c})$ and $\bar{h} \Delta \bar{c}$ arbitrarily. Therefore all subset of $\{y_1, \dots, y_{d-l}\}$ are contained inside such as outside of sets $\bar{h} \Delta \bar{c}$ for $\bar{c} \in \mathbf{B}_{(\bar{x}, \bar{z})}$. Because the probability of each element y_i is equal to the value $\frac{1}{d}$ we can estimate mean error of the generated hypothesis over all the target concepts in $\mathbf{B}_{(\bar{x}, \bar{z})}$ as (see lemma 3.1.11)

$$\frac{1}{2^{d-l}} \sum_{\bar{c} \in \mathbf{B}_{(\bar{x}, \bar{z})}} \frac{|\bar{h} \Delta \bar{c}|}{d} = \frac{1}{2^{d-l}} \sum_{\bar{z} \in \{0,1\}^{d-l}} \frac{n_+(\bar{z})}{d} \stackrel{3.1.11}{=} \frac{d-l}{2d} \geq \frac{d-m}{2d}.$$

For this moment let us define for each target concept $\bar{c} \in \mathbf{B}_{(\bar{x}, \bar{z})}$ and each hypothesis $\bar{h} \in \mathbf{B}_{(\bar{x}, \bar{z})}$ the number

$$\mathbf{Z}_{\bar{c}, \bar{h}} \stackrel{\text{def}}{=} e_{\bar{p}}(\bar{c}, \bar{h}).$$

As we explained above, obviously it holds that

$$\left(\forall \bar{h} \in \mathbf{B}_{(\bar{x}, \bar{z})} \right) \left(\sum_{\bar{c} \in \mathbf{B}_{(\bar{x}, \bar{z})}} \frac{\mathbf{Z}_{\bar{c}, \bar{h}}}{|\mathbf{B}_{(\bar{x}, \bar{z})}|} > \frac{d-m}{2d} \right).$$

According to (lemma 3.1.9) there exists target concept $\bar{b} \in \mathbf{B}_{(\bar{x}, \bar{z})}$ with average error at least $\frac{d-m}{2d}$ (in the sense that for random choice of $\bar{x} \in \bar{X}^m$ the average error of generated hypotheses is at least $\frac{d-m}{2d}$). In other words, there exists a concept in relation to which all the consistent hypotheses have mean error ρ

$$\rho \geq \frac{d-m}{2d}. \quad (3.21)$$

Finally let us define

$$\mathbf{E} \stackrel{\text{def}}{=} \left\{ \bar{h} \in \mathbf{B}_{(\bar{x}, \bar{z})} \mid e_{\bar{p}}(\bar{b}, \bar{h}) \geq \epsilon \right\}, \quad \bar{S} \stackrel{\text{def}}{=} \left\{ (\bar{x}, \bar{z}) \text{ is a sample of } \bar{b} \text{ and } \bar{x} \in \bar{X}^m \right\},$$

and

$$\mathbf{P} \stackrel{\text{def}}{=} \left\{ (\bar{x}, \bar{z}) \in \bar{S} \mid \widetilde{\mathbf{A}}^* \left((\bar{x}, \bar{z}) \right) \in \mathbf{E} \right\}, \quad \mathbf{Q} \stackrel{\text{def}}{=} \left\{ (\bar{x}, \bar{z}) \in \bar{S} \mid \widetilde{\mathbf{A}}^* \left((\bar{x}, \bar{z}) \right) \notin \mathbf{E} \right\}.$$

Let us assume, that the probability of the set \mathbf{P} be equal to $\delta \in (0, 1)$ and recall that all hypotheses in the set system \mathbf{E} are produced by the learning algorithm $\widetilde{\mathbf{A}}^*$ as a response to the samples from \mathbf{P} . So $Prob_{\bar{p}}(\mathbf{P}) = Prob_{\bar{p}}(\mathbf{E})$. Now let us assume that $Prob_{\bar{p}}(\mathbf{E}) = \delta - \gamma$, where $\gamma \in (0, \delta)$. Hence we can estimate average probability of the produced hypothesis over all samples from the set \bar{S} using the lemma 3.1.10 (note that the probability on \bar{X} is assumed to be uniform) and the 3.21 as

$$\epsilon(1-\delta) + \delta > \epsilon(1-(\delta-\gamma)) + (\delta-\gamma) \stackrel{3.18}{\geq} \rho \stackrel{3.21}{\geq} \frac{d-m}{2d}.$$

It follows that

$$m > d(1 - 2(\epsilon(1 - \delta) + \delta)) ,$$

which contradicts the second part of upperbound in the assumption 3.19. So $Prob_{\bar{P}}(\mathbf{E}) \geq \delta$.

- add Proof of the estimate 3.20)

Let $\bar{c} \in \mathbf{C}$ be fixed and let

$$m \geq \max \left\{ \frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8d}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right\}$$

and

$$\bar{Q}_{m,\epsilon} = \left\{ \tilde{x} \in \bar{X}^m \mid \left(\exists \bar{r} \in \mathbf{R}_{\bar{P},\epsilon} \right) \left(\tilde{x} \cap \bar{r} = \emptyset \right) \right\} .$$

If we recall the definition 3.1.2, the lemma 3.1.5, the theorem 2.4.4 and the lemma 3.1.8 we get

$$Prob_{\bar{P}}(\bar{Q}_{m,\epsilon}) \stackrel{3.1.5}{\leq} 2\Pi_{\mathbf{R}}(2m) 2^{-\frac{\epsilon m}{2}} \stackrel{2.4.4}{=} 2\Pi_{\mathbf{H}}(2m) 2^{-\frac{\epsilon m}{2}} \stackrel{3.1.8}{\leq} \delta ,$$

where $\mathbf{R} \stackrel{\text{def}}{=} \{ \bar{h} \triangle \bar{c} \mid \bar{h} \in \mathbf{H} \}$. Due to the fact that the algorithm A^* produces consistent hypotheses only we can write

$$\bar{h} = \widetilde{A^*} \left(\left(\tilde{x}, \bar{z} \right) \right) \quad \text{and} \quad e_{\bar{P}}(\bar{c}, \bar{h}) > \epsilon \quad \Rightarrow \quad \tilde{x} \in \bar{Q}_{m,\epsilon} .$$

But the probability of the set $\bar{Q}_{m,\epsilon}$ is less than ϵ which concludes the proof of this part.

- add Proof of the part 2)

Let us assume that $VC_{dim}(\mathbf{C}) \stackrel{\text{def}}{=} +\infty$. Obviously, in this case it is possible for every natural d to perform a sequence of steps in the proof of the second part of the estimate 3.19. However, if we select the values ϵ and δ such that the expression $(1 - 2(\epsilon(1 - \delta) + \delta))$ is positive, we get for this particular selection of ϵ, δ that there exists no (ϵ, δ) -learning algorithm for any fixed chosen length of samples.

Conversely, on the basis of the validity of the estimate 3.20 it is possible to construct an (ϵ, δ) -learning algorithm in such a way that we first well-order, in whatever fashion, the hypotheses class \mathbf{H} (each set can be well-ordered, see the Zermelo's theorem, e.g. [PBKN90], p. 44, and let us further recall that in case of well-ordering, there exist no incomparable elements in the set). Then, we assign to each sample from $\bar{S}_{\mathbf{C}}$ the first concept (in view of the actual ordering) which is consistent with it. That this sample exists ensues from the inclusion $\mathbf{C} \subset \mathbf{H}$. According to 3.20, this algorithm is an (ϵ, δ) -learning algorithm.

– q. e. d. –

Obviously, this theorem is much stronger than the theorem 1.2.1 since the value $|\mathbf{C}|$ in the upper estimate has now been replaced by the value $VC_{dim}(\mathbf{C})$ which may be substantially smaller. And in addition, this theorem is applicable in the case of infinite set \bar{X} and infinite concept class \mathbf{C} .

Example 3.1.3 Let \bar{X} , \mathbf{C} and A^* be defined in the same manner as in the example 1.2.1. Hence $VC_{dim}(\mathbf{C}) = 1$ and if

$$m \geq \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right)$$

then the learning algorithm A^* is (ϵ, δ) -learning algorithm.

Example 3.1.4 Let \bar{X} , C and A^* be defined in the same manner as in the example 1.2.2. Hence $\text{VC}_{\dim}(C) = 2n$ and if

$$m \geq \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{16n}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right)$$

then the learning algorithm A^* is (ϵ, δ) -learning algorithm.

Note that in the both cases discussed above the bellow bound on the m is independent on the number k (e.g. on the size of the space \bar{X}).

3.1.1 Delta rule learning algorithm

As an example of consistent learning algorithm, let us now study the so-called delta rule algorithm which is able to find a separation hyperplane for a given tuple of linearly separable sets of Euclidean space.

Definition 3.1.7 Let $(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)$ be a given sequence of tuples in $\mathbb{R}^n \times \{-1, +1\}$, $t \geq 1$. Further, let vector's sequence $\{\vec{w}_i\}_1^\infty$ satisfy the following recursive formulas

1. put $\vec{w}_1 \stackrel{\text{def}}{=} \vec{0}$, $k = 1$
2. let $k = k + 1$ and $\bar{J} \stackrel{\text{def}}{=} \{j \in \{1, \dots, m\} \mid \text{sgn}(\langle \vec{w}_k | \vec{x}_j \rangle) \neq y_j\}$
 - (a) if $\bar{J} = \emptyset$ put $\vec{w}_{k+1} = \vec{w}_k$ and STOP,
 - (b) else let $j_k \in \bar{J}$ be arbitrary. Then put

$$\vec{w}_{k+1} \stackrel{\text{def}}{=} \vec{w}_k + y_{j_k} \vec{x}_{j_k}$$

and REPEAT step 2).

Then we say that this sequence arose by application of DELTA RULE on $(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)$.

A basic properties of the delta rule algorithm summarizes the next theorem:

Theorem 3.1.13 Assume that sequence $\{\vec{w}_i\}_1^\infty$ arose by application of delta rule and let there exists a vector $\hat{\vec{w}}$ such that for all indexes $i \in \{1, \dots, m\}$ holds $\text{sgn}(\langle \hat{\vec{w}} | \vec{x}_i \rangle) = y_i$.

Further let

$$\alpha \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, m\}} \{\|\vec{x}_i\|^2\} \quad \text{and} \quad \beta \stackrel{\text{def}}{=} \min_{i \in \{1, \dots, m\}} \{|\langle \hat{\vec{w}} | \vec{x}_i \rangle|\} > 0.$$

Then there exists an natural number $z > 0$ satisfying $\vec{w}_{z+1} = \vec{w}_z$ and z can be estimated as

$$z \leq \frac{\alpha \|\hat{\vec{w}}\|^2}{\beta^2} + 1.$$

■ *Proof:*

By definition of delta rule it is obvious that for arbitrary $k > 1$ for which $\vec{w}_k \neq \vec{w}_{k-1}$ holds

$$\vec{w}_k = \sum_{p=1}^{k-1} \ddot{\vec{x}}_{j_p}, \quad \text{where} \quad \ddot{\vec{x}}_{j_p} \stackrel{\text{def}}{=} \begin{cases} \vec{x}_{j_p} & \text{pro } y_{j_p} = 1 \\ -\vec{x}_{j_p} & \text{pro } y_{j_p} = -1 \end{cases}.$$

Further let us assume that $k > t$. By the fact that $\widehat{sgn}(\langle \widehat{\vec{w}} | \vec{x}_{j_p} \rangle) \neq y_{j_p}$ and by definition of $\ddot{\vec{x}}_{j_p}$ the inequality $\langle \widehat{\vec{w}} | \ddot{\vec{x}}_{j_p} \rangle > 0$ is true for all vectors $\ddot{\vec{x}}_{j_p}$. Thus we can write the estimation

$$\langle \widehat{\vec{w}} | \vec{w}_k \rangle = \sum_{p=1}^{k-1} \langle \widehat{\vec{w}} | \ddot{\vec{x}}_{j_p} \rangle \geq (k-1)\beta > 0.$$

Hence, we get $\left| \langle \widehat{\vec{w}} | \vec{w}_k \rangle \right|^2 \geq (k-1)^2 \beta^2$. Further, using Schwartz inequality $\left| \langle \vec{a} | \vec{b} \rangle \right|^2 \leq \|\vec{a}\|^2 \|\vec{b}\|^2$ we finally show

$$\frac{\beta^2 (k-1)^2}{\|\widehat{\vec{w}}\|^2} \leq \|\vec{w}_k\|^2. \quad (3.22)$$

But, at the same time, for any k is

$$\vec{w}_k = \vec{w}_{k-1} + \ddot{\vec{x}}_{j_{k-1}},$$

and after square

$$\|\vec{w}_k\|^2 = \|\vec{w}_{k-1}\|^2 + 2\langle \vec{w}_{k-1} | \ddot{\vec{x}}_{j_{k-1}} \rangle + \|\ddot{\vec{x}}_{j_{k-1}}\|^2.$$

By properties of delta rule algorithm must be $\widehat{sgn}(\langle \vec{w}_{k-1} | \vec{x}_{j_{k-1}} \rangle) \neq y_{j_{k-1}}$ and therefore $\langle \vec{w}_{k-1} | \ddot{\vec{x}}_{j_{k-1}} \rangle \leq 0$. Hence, the previous inequality one could rewrite to the form

$$\|\vec{w}_k\|^2 - \|\vec{w}_{k-1}\|^2 \leq \|\ddot{\vec{x}}_{j_{k-1}}\|^2. \quad (3.23)$$

Now, let as sum up equations 3.23 over all $j \in \{1, \dots, k\}$ (remember that $\vec{w}_1 = \vec{0}$). We get

$$\|\vec{w}_k\|^2 \leq \sum_{p=1}^{k-1} \|\ddot{\vec{x}}_{j_p}\|^2 \leq (k-1)\alpha.$$

Blended together the last inequality and inequality 3.22 we can conclude

$$\frac{\beta^2 (k-1)^2}{\|\widehat{\vec{w}}\|^2} \leq (k-1)\alpha$$

and finally

$$k-1 \leq \frac{\alpha \|\widehat{\vec{w}}\|^2}{\beta^2}.$$

There is a constant independent on the value of k on the right side of the last expression. It follows that an $z > 0$ mentioned in the statement of the theorem must exist.

– q. e. d. –

Our chief interest is to illustrate by examples the PAC learning model. Reconsider now the definition of delta rule learning algorithm. Let us change step 2.(b) in definition 3.1.7 by "else let $j_k = \min \{\bar{J}\}$. Then put". To do it so, we get deterministic version of delta rule algorithm and we can apply below bound to sufficient number of examples used by (ϵ, δ) -algorithm. Obviously, each vector \vec{w} produced by delta rule algorithm is homogenous linear separator of positive and negative patterns in the sequence $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_m$. Hence corresponding concept class C is the set of all halfspaces with zero vector in its border and $C = H$. As we already shown, $VC_{dim}(C) = n$, which implies that we need at least

$$m \geq \max \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8n}{\epsilon} \log_2 \left(\frac{12.611}{\epsilon} \right) \right)$$

to keep the delta rule process to be (ϵ, δ) -algorithm.

3.1.2 Lower bound for maximal steps of delta rule algorithm.

From another point of view the statement of the Theorem 3.1.13 only says that the delta rule algorithm is finite process without any possibility to estimate the number of necessary iterations. Unfortunately, as we explain in the following text the number of delta rule iterations can not be upper bounded by polynomial in n .

Definition 3.1.8 Let $\bar{A}, \bar{B} \subset \mathcal{R}^n$. Then the tuple $(\vec{w}, t) \in \mathcal{R}^n \times \mathcal{R}$ is LINEAR SEPARATOR of sets \bar{A}, \bar{B} iff

$$(\forall \vec{a} \in \bar{A}) (\langle \vec{a} | \vec{w} \rangle < t) \quad \text{and} \quad (\forall \vec{b} \in \bar{B}) (\langle \vec{b} | \vec{w} \rangle > t) .$$

Lemma 3.1.14 Let $\bar{A} \subset \{-1, +1\}^n$, $\bar{B} \stackrel{\text{def}}{=} \{-1, +1\}^n - \bar{A}$ and (\vec{w}, t) be a linear separator of sets \bar{A} and \bar{B} . Then there exists a linear separator (\vec{w}^*, t) of sets \bar{A}, \bar{B} such that

$$(\forall \vec{x}, \vec{y} \in \{-1, +1\}^n) (\vec{x} \neq \vec{y} \Rightarrow \langle \vec{w}^* | \vec{x} \rangle \neq \langle \vec{w}^* | \vec{y} \rangle) . \quad (3.24)$$

■ *Proof:*

Let (\vec{w}, t) be an arbitrary linear separator of sets \bar{A} and \bar{B} . Further, let there exist $\vec{y}, \vec{z} \in \{-1, +1\}^n$ such that $\vec{x} \neq \vec{y}$ and $\langle \vec{w} | \vec{x} \rangle = \langle \vec{w} | \vec{y} \rangle$.

Now let us define a positive number β as

$$0 < 2\beta < \min_{\vec{r}, \vec{s} \in \{-1, +1\}^n} \{ |\langle \vec{w} | \vec{r} \rangle - \langle \vec{w} | \vec{s} \rangle| > 0 \} .$$

Because $\vec{x} \neq \vec{y}$, there exists an index k for which $\vec{x}_k \neq \vec{y}_k$. Put $\vec{w}^* = \vec{w} + \beta \vec{e}_k$, where \vec{e}_k has all coordinates zero except k -th, which is 1. Due to definition of number β the tuple (\vec{w}^*, t) is linear separator of sets \bar{A} and \bar{B} and at the same time

$$(\forall \vec{r}, \vec{s} \in \{-1, +1\}^n) (\langle \vec{w} | \vec{r} \rangle \neq \langle \vec{w} | \vec{s} \rangle \Rightarrow \langle \vec{w}^* | \vec{r} \rangle \neq \langle \vec{w}^* | \vec{s} \rangle)$$

and

$$|\langle \vec{w}^* | \vec{x} \rangle - \langle \vec{w}^* | \vec{y} \rangle| \geq 2\beta,$$

so the number of tuples \vec{y}, \vec{z} for which $\langle \vec{w} | \vec{x} \rangle = \langle \vec{w} | \vec{y} \rangle$ is decreased at least by 1. Obviously we can repeat this construction till all such tuples \vec{y}, \vec{z} are exhausted.

– q. e. d. –

Based on the previous lemma we can define recursive algorithm which allows us to postulate exponential lower bound on the number of hypercube vertices dichotomies. Thist recursive algorithm is illustrated on the Fig. 3.3.

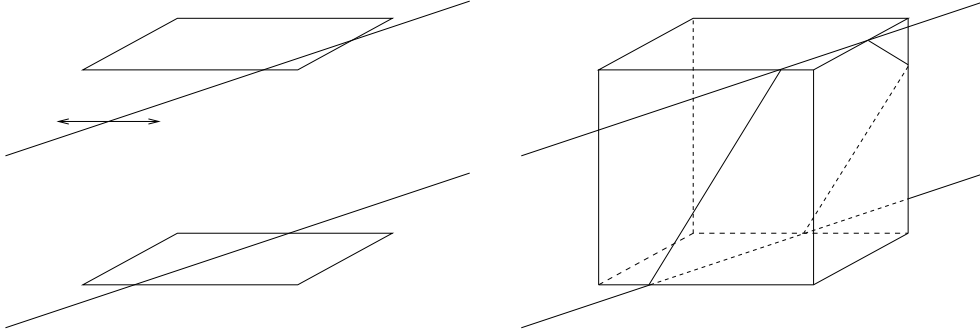


Figure 3.3: Construction of $2^k + 1$ dichotomies derived from a fixed dichotomy of $\{-1, +1\}^k$ (for $k = 2$).

Theorem 3.1.15 *The number of linearly separable dichotomies of the cube $\{-1, +1\}^n$ is larger than $2^{\frac{n(n-1)}{2}}$.*

■ *Proof:*

We proceed the proof by induction on n . The case $n = 1$ is clear because there exist 4 linearly separable dichotomies of $\{-1, +1\}$ exactly. So we suppose that $n \geq 1$ and that statement holds for this n . Let (\bar{A}, \bar{B}) be a fixed linearly separable dichotomy of the cube $\{-1, +1\}^n$. As follows from the lemma 3.1.14 we can assume that corresponding linear separator (\vec{w}^*, t) fullfils the condition 3.24. So there exists 2^n mutually different real numbers $\omega_i, i \in \{1, \dots, 2^n\}$, such that $\omega_i < \omega_{i+1}, i \in \{1, \dots, 2^n - 1\}$, and

$$\{\omega_1, \omega_2, \dots, \omega_{2^n}\} = \{\omega | \omega = \langle \vec{w}^* | \vec{x} \rangle, \vec{x} \in \{-1, +1\}^n\}$$

Let us define

$$t_0 \stackrel{\text{def}}{=} \omega_1 - 1, \quad t_i \stackrel{\text{def}}{=} \frac{\omega_i + \omega_{i+1}}{2}, \quad i \in \{1, \dots, 2^n - 1\}, \quad t_{2^n} \stackrel{\text{def}}{=} \omega_{2^n} + 1,$$

$$\bar{A}^- \stackrel{\text{def}}{=} \{(\vec{a}, -1) \in \{-1, +1\}^{n+1} | \vec{a} \in \bar{A}\}, \quad \bar{B}^- \stackrel{\text{def}}{=} \{(\vec{b}, -1) \in \{-1, +1\}^{n+1} | \vec{b} \in \bar{B}\},$$

and let for all $i \in \{0, \dots, 2^n\}$ is

$$\bar{A}_i^+ \stackrel{\text{def}}{=} \{(\vec{x}, +1) \in \{-1, +1\}^{n+1} | \langle \vec{w}^* | \vec{x} \rangle < t_i\}, \quad \bar{B}_i^+ \stackrel{\text{def}}{=} \{(\vec{x}, +1) \in \{-1, +1\}^{n+1} | \langle \vec{w}^* | \vec{x} \rangle > t_i\}$$

and

$$\bar{C}_i^A = \bar{A}^- \cup \bar{A}_i^+ \quad , \quad \bar{C}_i^B = \bar{B}^- \cup \bar{B}_i^+.$$

Obviously, a tuple $(\bar{C}_i^A, \bar{C}_i^B)$ is a dichotomy of the cube $\{-1, +1\}^{n+1}$. Further we show that

$$\left\langle \left(\bar{\mathbf{w}}^*, \frac{t-t_i}{2} \right), \frac{t+t_i}{2} \right\rangle. \quad (3.25)$$

is a linear separator of $(\bar{C}_i^A, \bar{C}_i^B)$.

Let $\bar{\mathbf{x}} \in \bar{A}^-$. Then $(\exists \bar{\mathbf{a}} \in \bar{A}) (\bar{\mathbf{x}} = (\bar{\mathbf{a}}, -1))$. Hence

$$\left\langle \left(\bar{\mathbf{w}}^*, \frac{t-t_i}{2} \right) | (\bar{\mathbf{a}}, -1) \right\rangle = \langle \bar{\mathbf{w}}^* | \bar{\mathbf{a}} \rangle - \frac{t-t_i}{2} < t - \frac{t-t_i}{2} = \frac{t+t_i}{2}.$$

Let $\bar{\mathbf{x}} \in \bar{A}_i^+$. Then $(\exists \bar{\mathbf{y}} \in \{-1, +1\}^n) (\bar{\mathbf{x}} = (\bar{\mathbf{y}}, 1) \text{ and } \langle \bar{\mathbf{w}}^* | \bar{\mathbf{y}} \rangle < t_i)$. Hence

$$\left\langle \left(\bar{\mathbf{w}}^*, \frac{t-t_i}{2} \right) | (\bar{\mathbf{y}}, 1) \right\rangle = \langle \bar{\mathbf{w}}^* | \bar{\mathbf{y}} \rangle + \frac{t-t_i}{2} < t_i + \frac{t-t_i}{2} = \frac{t+t_i}{2}.$$

Let $\bar{\mathbf{x}} \in \bar{B}^-$. Then $(\exists \bar{\mathbf{b}} \in \bar{B}) (\bar{\mathbf{x}} = (\bar{\mathbf{b}}, -1))$. Hence

$$\left\langle \left(\bar{\mathbf{w}}^*, \frac{t-t_i}{2} \right) | (\bar{\mathbf{b}}, -1) \right\rangle = \langle \bar{\mathbf{w}}^* | \bar{\mathbf{b}} \rangle - \frac{t-t_i}{2} > t - \frac{t-t_i}{2} = \frac{t+t_i}{2}.$$

Let $\bar{\mathbf{x}} \in \bar{B}_i^+$. Then $(\exists \bar{\mathbf{y}} \in \{-1, +1\}^n) (\bar{\mathbf{x}} = (\bar{\mathbf{y}}, 1) \text{ and } \langle \bar{\mathbf{w}}^* | \bar{\mathbf{y}} \rangle > t_i)$. Hence

$$\left\langle \left(\bar{\mathbf{w}}^*, \frac{t-t_i}{2} \right) | (\bar{\mathbf{y}}, 1) \right\rangle = \langle \bar{\mathbf{w}}^* | \bar{\mathbf{y}} \rangle + \frac{t-t_i}{2} > t_i + \frac{t-t_i}{2} = \frac{t+t_i}{2}.$$

The previous four inequalities follows, that the tuple 3.25 is a linear separator of the $(\bar{C}_i^A, \bar{C}_i^B)$. Thus for a given dichotomy (\bar{A}, \bar{B}) of $\{-1, +1\}^n$ we obtained $2^n + 1$ mutually different dichotomies of the $\{-1, +1\}^{n+1}$. In addition, it is obvious that for different dichotomies (\bar{A}', \bar{B}') of $\{-1, +1\}^n$ derived dichotomies of $\{-1, +1\}^{n+1}$ differs in corresponding sets \bar{A}^- and \bar{B}^- .

Finally, let K_n denotes the number of all linearly separable dichotomies of the set $\{-1, +1\}^n$. Hence we proved recursive formula $K_{n+1} \geq (2^n + 1)K_n$. As we mentioned $K_1 = 4$, so

$$K_n \geq \prod_{i=0}^{n-1} (2^i + 1) > \prod_{i=0}^{n-1} 2^i = 2^{\frac{n(n-1)}{2}}.$$

– q. e. d. –

The previous theorem follows the fact that there exists linearly separable splitting of the cube $\{-1, +1\}^n$ with integer linear separator whose indices can not be upper bounded by any polynomial in n .

Theorem 3.1.16 *There exists a linearly separable dichotomy of the $\{-1, +1\}^n$ such that any integer linear separator $(\bar{\mathbf{w}}, t)$ of this dichotomy satisfies estimation*

$$2^{\frac{n-2}{2}} \leq \sum_{k=1}^n |\bar{\mathbf{w}}_k| + |t|.$$

■ *Proof:*

Let us prove bellow bound. Denote

$$\bar{P}_n \stackrel{\text{def}}{=} \{(\bar{A}, \bar{B}) \mid (\bar{A}, \bar{B}) \text{ is linearly separable splitting of } \{-1, +1\}^n\}.$$

Further let for all $(\bar{A}, \bar{B}) \in \bar{P}_n$ is

$$\bar{W}_{(\bar{A}, \bar{B})} \stackrel{\text{def}}{=} \{(\vec{w}, t) \in \mathfrak{R}^n \times \mathfrak{R} \mid (\vec{w}, t) \text{ is integer linear separator of } (\bar{A}, \bar{B})\}.$$

Finally define integer function $\tilde{\zeta}(n)$ as

$$\tilde{\zeta}(n) \stackrel{\text{def}}{=} \max_{(\bar{A}, \bar{B}) \in \bar{P}_n} \left\{ \min_{(\vec{w}, t) \in \bar{W}_{(\bar{A}, \bar{B})}} \left\{ \left\lceil \log_2 \left(\sum_{k=1}^n |\vec{w}_k| + |t| \right) \right\rceil \right\} \right\}.$$

Definition of the function $\tilde{\zeta}(n)$ follows that for any dichotomy $(\bar{A}, \bar{B}) \in \bar{P}_n$ there exists an integer linear separator $(\vec{w}, t) \in \bar{W}_{(\bar{A}, \bar{B})}$ which has all indices bounded to the interval $(-2^{\tilde{\zeta}(n)}, 2^{\tilde{\zeta}(n)})$.

Using binary notation we need at most $\tilde{\zeta}(n) + 1$ bits (including sign) to store each indices of the linear separator. Hence, to store whole linear separator, e.g. $n + 1$ indices, $(n + 1)(\tilde{\zeta}(n) + 1)$ bits is sufficient.

This amount of bits allow to store at most $2^{(n+1)(\tilde{\zeta}(n)+1)}$ different linear separators. At the same time, by the theorem 3.1.15, there exists more than $2^{\frac{n(n-1)}{2}}$ different linearly separable dichotomies. Hence we get inequality

$$2^{(n+1)(\tilde{\zeta}(n)+1)} \geq 2^{\frac{n(n-1)}{2}}.$$

Therefore $(n + 1)(\tilde{\zeta}(n) + 1) \geq \frac{n(n-1)}{2}$ and

$$\tilde{\zeta}(n) + 1 \geq \frac{n(n-1)}{2(n+1)} = \frac{n-2}{2} + \frac{1}{n+1} > \frac{n-2}{2}.$$

– q. e. d. –

Example 3.1.5 Let $\bar{X} \stackrel{\text{def}}{=} \{-1, +1\}^n$ and

$$\mathcal{C} = \left\{ \bar{A} \subset \{-1, +1\}^n \mid \bar{A} \text{ and } (\{-1, +1\}^n \setminus \bar{A}) \text{ are linearly separable} \right\}.$$

Let $\bar{c} \in \mathcal{C}$ be such a set, that any integer linear separator (\vec{w}, t) of the sets \bar{c} and $\{-1, +1\}^n \setminus \bar{c}$ satisfies $2^{\frac{n-2}{2}} \leq \sum_{k=1}^n |\vec{w}_k| + |t|$. Let us use delta rule as a learning algorithm to find (\vec{w}, t) . The algorithm stops when return a consistent hypothesis and this happens after finite number of steps. But in each step of delta rule algorithm absolute values of the weight vector \vec{w} can increase by one only, so we need at least $\frac{2^{\frac{n-2}{2}}}{n+1}$ steps of delta rule algorithm.

3.1.3 Linear separation and linear programming

Definition 3.1.9 (Mangasarian LP) Let $\bar{A} \stackrel{\text{def}}{=} \{\vec{a}_1, \dots, \vec{a}_i\}$ and $\bar{B} \stackrel{\text{def}}{=} \{\vec{b}_1, \dots, \vec{b}_j\}$ be a finite subsets of the \mathfrak{R}^n . Then MANGASARIAN LINEAR PROBLEM is defined as the problem find vectors $\vec{y} \in \mathfrak{R}^i$, $\vec{z} \in \mathfrak{R}^j$, $\vec{w} \in \mathfrak{R}^n$ and $t \in \mathfrak{R}$ that minimizes

$$\sum_{\alpha=1}^i \vec{y}_\alpha + \sum_{\beta=1}^j \vec{z}_\beta$$

subject to

$$\begin{aligned} \vec{y}_\alpha + \langle \vec{w} | \vec{a}_\alpha \rangle - t &\geq 1 \quad \text{for } \alpha \in \{1, \dots, i\} \\ \vec{z}_\beta - \langle \vec{w} | \vec{b}_\beta \rangle + t &\geq 1 \quad \text{for } \beta \in \{1, \dots, j\} \\ \vec{y}_\alpha &\geq 0 \quad \text{for } \alpha \in \{1, \dots, i\} \\ \vec{z}_\beta &\geq 0 \quad \text{for } \beta \in \{1, \dots, j\}. \end{aligned}$$

Theorem 3.1.17 Let $\bar{A} \stackrel{\text{def}}{=} \{\vec{a}_1, \dots, \vec{a}_i\}$ and $\bar{B} \stackrel{\text{def}}{=} \{\vec{b}_1, \dots, \vec{b}_j\}$ be a finite subsets of the \mathfrak{R}^n . Then

1. There exists a linear separator of the sets \bar{A} and \bar{B} if and only if the optimal value of the corresponding Mangasarian LP is zero.
2. If the optimal value of the corresponding Mangasarian LP is zero and $(\vec{y}^*, \vec{z}^*, \vec{w}^*, t^*)$ is optimal solution, than (\vec{w}^*, t^*) is linear separator of the sets \bar{A} and \bar{B} .

■ Proof:

■ add 1 \Rightarrow)

Let (\vec{w}, t) be linear separator of the sets \bar{A} and \bar{B} . It means that $\langle \vec{w} | \vec{a}_\alpha \rangle - t > 0$, $\alpha \in \{1, \dots, i\}$, and $\langle \vec{w} | \vec{b}_\beta \rangle - t < 0$, $\beta \in \{1, \dots, j\}$. Let

$$\omega \stackrel{\text{def}}{=} \min \left\{ \min \{ \langle \vec{w} | \vec{a}_\alpha \rangle - t \mid \alpha \in \{1, \dots, i\} \}, \min \{ -\langle \vec{w} | \vec{b}_\beta \rangle + t \mid \beta \in \{1, \dots, j\} \} \right\}.$$

Clearly, $\omega > 0$ and $\langle \vec{w} | \vec{a}_\alpha \rangle - t \geq \omega$, $\alpha \in \{1, \dots, i\}$, and $\langle \vec{w} | \vec{b}_\beta \rangle - t \leq \omega$. So $\left\langle \frac{\vec{w}}{\omega} | \vec{a}_\alpha \right\rangle - \frac{t}{\omega} \geq 1$, $\alpha \in \{1, \dots, i\}$, $\left\langle \frac{\vec{w}}{\omega} | \vec{b}_\beta \right\rangle - \frac{t}{\omega} \leq 1$. So $(\vec{0}, \vec{0}, \frac{\vec{w}}{\omega}, \frac{t}{\omega})$ is the solution of the corresponding Mangasarian LP with optimal value 0.

■ add 1 \Leftarrow and 2)

Let $(\vec{y}^*, \vec{z}^*, \vec{w}^*, t^*)$ is an optimal solution with zero optimal value. It follows that $\vec{y} = \vec{0}$ and $\vec{z} = \vec{0}$. Hence a tuple (\vec{w}, t) is a linear separator of the sets \bar{A} and \bar{B} .

– q. e. d. –

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

© F. Hák, ICS CAS, Prague, Tech. Rep. №1227, Dec 2015

Chapter 4

Appendices

4.1 Source codes

4.1.1 Python code for $\omega(\epsilon)$ solver

The python code listed bellow were used to compute table and graph on the Fig. 3.1.

```
#!/usr/bin/python
import string, os, os.path, time, re, random, shutil, sys
import matplotlib.pyplot as plt
import numpy as np
import math

omega = []
epsilon = []

eps = 0.001
step = 0.001
e = 2.71828182845905
log2e = 1.442695040888963387

result = 1.0

while eps <= 1.0 + step :
    epsilon.append(eps)
    iter = result
    residum = 1000000000
    while residum > 0.00000000001 : # Newton iterations method  $x_{n+1}=x_n-f(x)/f'(x)$ 
        value = 1/eps - log2e*math.log(iter/eps)*16*e/iter/iter
        derivative = (2*math.log(iter/eps) - 1.0)*log2e*16*e/iter/iter/iter
        iter = iter - value/derivative

        residum = abs(1/eps - log2e*math.log(iter/eps)*16*e/iter/iter)

    result = iter
    omega.append(result)
    eps += step
```

```

##### COMPUTE TABLE #####
def get_index(r,c) : # r row , c column
    value = c*0.01 + r*0.05
    index = int(value*1000)
    return index-1

text_row = u'0.00 '
col = 0
while col < 5 :
    text_row = text_row + " & %8s"%(col*0.01)
    col += 1
print "%s \\\ \hline"%(text_row)

text_row = u'0.00 '
row = 0
col = 0
while row < 10 :
    while col < 5 :
        index = get_index(row,col)
        res = unicode("%.2f"%(omega[index]))
        if col == 0 and row == 0 : res = u'--'
        text_row = text_row + " & %8s"%(res)
        col += 1
    col = 0
    print "%s \\\ \"%(text_row)
    text_row = "%.2f  "%((float(row)+1)/20)
    row += 1
print "\\hline"

##### MAKE PLOT #####
plt.title(u'Solution of equation  $1/\epsilon=(16e/\omega^2)\cdot\log_2(\omega/\epsilon)$ ')

plt.xlabel(u'$\epsilon$')
plt.ylabel(u'$\omega(\epsilon)$')
plt.text(0.8,13.0,u'$\omega(1)$ = 12.610894152')
plt.grid(True)
plt.plot(epsilon,omega)

plt.savefig(u'omega-as-function-epsilon.pdf',
    dpi=75, facecolor='w', edgecolor='w', orientation='portrait',
    papertype='a4', format=None, transparent=False)

plt.show()

```

Bibliography

- [AB92] Martin Anthony and Norman Biggs. *Computational Learning Theory*. Press Syndicate of the University of Cambridge, 1992.
- [And85] Jiří Anděl. *Matematická statistika*. SNTL - ALFA Bratislava, 1985.
- [BEHW89] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Learnability and the Vapnik-Chervonenkis Dimension. *Journal of the Association for Computing Machinery*, 36:929–965, oct 1989.
- [Cov] Thomas M. Cover. Geometrical and Statistical Properties of Systems of Linear Inequalities with Applications in Pattern Recognition. *IEEE Transaction on Electronic Computers*, pages 326–334.
- [Jar55] V. Jarník. *Integrální počet II*. NČAV Praha, 1955.
- [Lej85] K. Lejchtvějs. *Vypuklité množstava*. Moskva Nauka, 1985.
- [PBKN90] Ladislav Procházka, Ladislav Bican, Tomáš Kepka, and Petr Němec. *Algebra*. ACADEMIA Praha, 1990.
- [Sau72] N. Sauer. On the Density of Families of Sets. *Journal of the Association for Computing Machinery*, 13:145–147, feb 1972.