

Scombinatorics

Preliminaries	2
1 Three minimax theorems	3
1 Hall's Marriage Theorem	3
2 König's Minimax Theorem	4
3 Dilworth's Theorem	5
2 Set systems	7
1 Sperner's Theorem	7
2 The Erdős-Ko-Rado Theorem	8
3 Stable and NIP relations	10
1 Stable formulas	10
2 The Vapnik-Chervonenkis dimension	10
3 The Sauer-Shelah lemma	11
4 Law(s) of large numbers	14
1 Inequalities	14
2 Two Weak Laws of Large Numbers	18
3 A Uniform Law of Large Numbers	19
4 A Uniform Law of Large Numbers, again	22
References	25

Notation

Let \mathcal{U} and \mathcal{V} be two (large) sets. Let $\varphi(x; z)$ be a relation symbol, or a formula, whatever. We denote by $\varphi(\mathcal{U}; \mathcal{V})$ the set $\{\langle a, b \rangle \in \mathcal{U} \times \mathcal{V} : \varphi(a; b)\}$ which we call: the relation defined by $\varphi(x; z)$. Sets of the form $\varphi(\mathcal{U}; b) = \{a \in \mathcal{U} : \varphi(a; b)\}$, for some $b \in \mathcal{V}$, are called **definable** sets.

In the first chapters we always restrict the study to the **trace** of $\varphi(\mathcal{U}; \mathcal{V})$ on some finite set $A \times B$, where $A \subseteq \mathcal{U}$ and $B \subseteq \mathcal{V}$. We write $\varphi(A; B)$ for $\varphi(\mathcal{U}; \mathcal{V}) \cap A \times B$. Similarly, we write $\varphi(A; b)$ for the trace of $\varphi(\mathcal{U}; b)$ on A , that is, the set $\varphi(\mathcal{U}; b) \cap A$. We call it a definable subset of A .

We denote by $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ and $\varphi(A; b)_{b \in \mathcal{V}}$ the collection of definable sets, respectively definable subsets of A .

It would be more appropriate to call the definable sets *global types*, respectively *types over A*. But this would make the terminology (if possible) more obscure.

For $k \leq |A|$ we use following notation interchangeably

$$\binom{A}{k} = A^{(k)} = \{A' \subseteq A : |A'| = k\}$$

Chapter 1

Three minimax theorems

Though apparently unrelated, the three theorems in this chapter can be derived one each other. We prove them in an arbitrary order.

As evident from the statement, the last two theorems are minimax theorems. The first theorem less so, hence the title is only approximately correct.

1 Hall's Marriage Theorem

Let $\varphi(x; z)$ be given. Let $A \subseteq \mathcal{U}$ and $B \subseteq \mathcal{V}$ be finite sets.

We say that $A' \subseteq A$ is a **set of distinct representatives** for $\varphi(A; B)$ if

$$|\varphi(A'; b)| = |\varphi(a; B)| = 1 \quad \text{for every } a, b \in A', B,$$

or, in other words, if $\varphi(A'; B)$ is the graph of a bijection.

1.1 Hall's Marriage Theorem *For every finite $B \subseteq \mathcal{V}$, the following are equivalent*

1. $\varphi(A; B)$ has a set of distinct representatives;
2. $|B'| \leq \left| \bigcup_{b \in B'} \varphi(A; b) \right|$ for every $B' \subseteq B$.

Proof (1 \Rightarrow 2) The following holds for any set of distinct representatives A' and every set $B' \subseteq B$

$$|B'| = \left| \bigcup_{b \in B'} \varphi(A'; b) \right| \subseteq \left| \bigcup_{b \in B'} \varphi(A; b) \right|.$$

(2 \Rightarrow 1) Reason by induction on the cardinality of B . If $|B| = 1$, the claim is clear. Now assume $|B| > 1$ and consider two cases.

- a. This is the case when the inequality in 2 is strict for all nonempty $B' \subset B$. Pick any pair $a, b \in A, B$ such that $\varphi(a; b)$. Then $\varphi(A \setminus \{a\}; B \setminus \{b\})$ still satisfy 2. By induction hypothesis, it has a set of distinct representatives A' . Then $A' \cup \{a\}$ is a set of distinct representatives for $\varphi(A; B)$.
- b. Suppose instead that for some nonempty $B' \subset B$ the inequality in 2 holds with equality. Define

$$A' = \bigcup_{b \in B'} \varphi(A; b)$$

It is clear that 2 holds for $\varphi(A'; B')$. Below we prove that 2 also holds for $\varphi(A \setminus A'; B \setminus B')$. Once this claim is proved, we apply the induction hypothesis to obtain sets of distinct representatives for these two relations and note that

their union is a set of distinct representatives for $\varphi(A; B)$.

To prove the claim assume that there is a set $B'' \subseteq B \setminus B'$ that contradicts 2, then

$$\left| \bigcup_{b \in B''} \varphi(A \setminus A'; b) \right| < |B''|.$$

By the definition of A', B'

$$\begin{aligned} \bigcup_{b \in B' \cup B''} \varphi(A; b) &= \bigcup_{b \in B'} \varphi(A; b) \cup \bigcup_{b \in B''} \varphi(A; b) \\ &= A' \cup \bigcup_{b \in B''} \varphi(A; b) \\ &= A' \cup \bigcup_{b \in B''} \varphi(A \setminus A'; b) \end{aligned}$$

The two sets above are disjoint, hence

$$\left| \bigcup_{b \in B' \cup B''} \varphi(A; b) \right| > |A'| + |B''|$$

As $|A'| = |B'|$, by the choice of A', B' , we obtain that $B' \cup B''$ contradicts the inequality in 2. This prove the claim and with it the theorem. \square

2 König's Minimax Theorem

Let $\varphi(x; z)$ be given. Let $A \subseteq \mathcal{U}$ and $B \subseteq \mathcal{V}$ be finite sets.

A **matching** of $\varphi(A; B)$ is a pair of sets $A' \subseteq A$ and $B' \subseteq B$ such that $\varphi(A'; B')$ is the graph of a bijection between A' and B' in other words

$$|\varphi(A'; b)| = |\varphi(a; B')| = 1 \quad \text{for every } a, b \in A', B'.$$

Yet in other words, A' is a set of distinctive representatives for $\varphi(A; B')$

We call $|A'| = |B'|$ the cardinality of the matching. The **matching number** of $\varphi(A; B)$ is the maximal cardinality of a matching.

Note that is A' is a set of distinct representatives for $\varphi(A; B)$, then there is a $B' \subseteq B$ such that A', B' . Hence the matching number is less or equal than the cardinality of any set of distinct representatives (if it exists).

A **(vertex) cover** of $\varphi(A; B)$ is a pair of sets $A' \subseteq A$ and $B' \subseteq B$ such that $\varphi(A; B)$ is contained in $(A' \times (B \setminus B')) \cup ((A \setminus A') \times B')$. We will mainly use this property as characterized by the easy fact below.

1.2 Fact *The following are equivalent*

1. A', B' is a cover;
2. $\varphi(A; b) \subseteq A'$ for every $b \in B \setminus B'$;
3. $\varphi(a; B) \subseteq B'$ for every $a \in A \setminus A'$.

\square

We call $|A'| + |B'|$ the cardinality of the cover. the **cover number** of $\varphi(A; B)$ is the minimal cardinality of a cover.

1.3 Kőnig's Minimax Theorem *For any given $\varphi(A; B)$, matching number = cover number. That is, the maximal cardinality of a matching equals the minimal cardinality of a cover.*

Proof (\leq) We prove that $|A''| \leq |A'| + |B'|$ for every cover A', B' and every matching A'', B'' .

As $\varphi(A; b) \subseteq A'$ for every $b \in B \setminus B'$, in particular we have that $\varphi(A''; b) \subseteq A'$ for every $b \in B'' \setminus B'$. By the definition of matching, all these sets $\varphi(A''; b)$ are distinct singletons. Hence $|B''| - |B'| \leq |B \setminus B'| \leq |A'|$ is clear.

(\geq) Let A', B' be a cover of minimal cardinality. We prove that there is a matching of cardinality at least $|A'| + |B'|$.

We break $\varphi(A; B)$ into two relations, find a matching of each of these and join them together to obtain a matching of cardinality $\geq |A'| + |B'|$. Precisely, first we show that $\varphi(A \setminus A'; B')$ has a set of distinct representatives $A_1 \subseteq A \setminus A'$. Hence A_1, B' is a matching. Second, we apply the same argument shows that $\varphi(A'; B \setminus B')^*$ has a set of distinct representatives $B_2 \subseteq B \setminus B'$. Hence A', B_2 is a matching. Then $(A_1 \cup A'), (B' \cup B_2)$ is a matching of $\varphi(A; B)$. The cardinality of this matching is $|A_1| + |A'| = |B_2| + |A'| = |B'| + |A'|$.

We use Hall's Marriage Theorem to prove the first claim above. The second is proved by the symmetric argument (using 3 of the fact above in place of 2).

We need to check that $\varphi(A \setminus A'; B')$ satisfies 2 of Theorem 1.1. Suppose not. Then there is a set $B'' \subseteq B'$ such that $|A''| < |B''|$, where

$$A'' = \bigcup_{b \in B''} \varphi(A \setminus A', b)$$

Then $(A' \cup A''), (B' \setminus B'')$ would be a cover of cardinality $< |A'| + |B'|$. This contradicts the minimality of A', B' . \square

3 Dilworth's Theorem

Dilworth's Theorem is minimax theorem essentially equivalent to Kőnig's Theorem. To highlight the connection we choose to prove it using Kőnig's Theorem. Alternatively we could have proved Dilworth's Theorem directly and derived Kőnig's and Hall's Theorem from it.

Let $<$ be a strict partial order on \mathcal{U} . An **antichain** is a set $A' \subseteq \mathcal{U}$ such that $a < a'$ for every $a, a' \in A'$. A **chain** is a set $A' \subseteq \mathcal{U}$ such that $a < a' \vee a' < a$ for every distinct $a, a' \in A'$.

1.4 Dilworth's Theorem *Let $A \subseteq \mathcal{U}$ be finite. The maximal cardinality of an antichain $A' \subseteq A$ equals the minimal cardinality of a partition of A into chains.*

Proof (\leq) We prove that the cardinality of an antichain cannot exceed the cardinal-

ity of a partition of A into chains.

Let A_1, \dots, A_k be a partition of A into chains and let A' be an antichain. A chain can contain at most one element of A' , hence $|A'| \leq k$.

(\geq) Let $A' \subseteq A$ be an antichain of maximal cardinality. We prove that there is a partition A_1, \dots, A_k into chains for some $k \leq |A'|$.

Let \mathcal{V} be a disjoint copy of \mathcal{U} . Let $f : \mathcal{U} \rightarrow \mathcal{V}$ the bijection that maps each element of \mathcal{U} to its copy in \mathcal{V} . For $a, b \in \mathcal{U}$ such that $a < b$ let $\varphi(a; fb)$. Let A_1, B_1 be a cover of $\varphi(A; f[A])$. We claim that $A \setminus (A_1 \cup f^{-1}[B_1])$ is an antichain. In fact, if $a < b$ then either $a \in A_1$ or $fb \in B_1$, by the definition of cover. This proves the claim.

As A' has maximal cardinality, $|A| - |A'| \leq |A_1 \cup f^{-1}[B_1]| \leq |A_1 \cup B_1|$. If we choose a over A_1, B_1 of minimal cardinality, by König's Theorem there is a matching A'', B'' of cardinality $|A''| \geq |A_1 \cup B_1|$. Hence $|A''| \geq |A \setminus A'|$.

We construct a chain-partition of A as follow. Pick an element of $a_0 \in A''$ and construct the longest possible chain $a_0, b_0, a_1, b_1, \dots, a_m, b_m, a_{m+1}$ where $a_i \in A''$ for all $i \leq m$, and $b_i \in B''$ is the (unique) element such that $\varphi(a_i; b_i)$ and $a_{i+1} \in A$ is the copy of $b_i \in B''$. The construction halts at the first $a_{m+1} \notin A''$. Then we start a new chain from some fresh element of A'' until the chains $a_0 < a_1 < \dots < a_m < a_{m+1}$ constructed in this way cover the whole of A'' . Note that these chains are pairwise disjoint. Finally, put each element of A not covered by these chains in a chain on its own.

Notice that the elements of A'' belongs to a chain of length at least 2. Therefore the number k of chains necessary to cover A is $\leq |A| \setminus |A''| \leq |A'|$. \square

Chapter 2

Set systems

1 Sperner's Theorem

We say that $\varphi(A; b)_{b \in \mathcal{V}}$ is an **antichain** if there is no pair of distinct elements $b, b' \in \mathcal{V}$ such that $\varphi(A; b) \subset \varphi(A; b')$. Antichains are also called **Sperner systems**.

If all sets in $\varphi(A; b)_{b \in \mathcal{V}}$ are distinct and of equal cardinality, then we clearly have an antichain. If $|A| = n$, the cardinality of a collection of subsets of A , all of cardinality k , is maximal when $k = \lfloor n/2 \rfloor$ or $k = \lceil n/2 \rceil$. In this case

$$\begin{aligned} |\varphi(A; b)_{b \in \mathcal{V}}| &= \binom{n}{\lfloor n/2 \rfloor} \\ &= \binom{n}{\lceil n/2 \rceil}. \end{aligned}$$

By the following classical theorem, this bound holds for all antichain. This is one of the first results of external combinatorics (though the term has been coined a few years later).

2.1 Sperner's Theorem *Let $A \subseteq \mathcal{U}$ have cardinality n , finite. If $\varphi(A; b)_{b \in \mathcal{V}}$ is an antichain then*

$$|\varphi(A; b)_{b \in \mathcal{V}}| \leq \binom{n}{\lfloor n/2 \rfloor}.$$

Proof Clearly, $\varphi(A; b)_{b \in \mathcal{V}}$ is the disjoint union of the sets $\binom{A}{k} \cap \varphi(A; b)_{b \in \mathcal{V}}$ for k ranging over $\{0, \dots, n\}$. Then

$$|\varphi(A; b)_{b \in \mathcal{V}}| \leq \sum_{k=0}^n \left| \binom{A}{k} \cap \varphi(A; b)_{b \in \mathcal{V}} \right|.$$

As for every $k \leq n$

$$\binom{n}{k} \leq \binom{n}{\lfloor n/2 \rfloor},$$

the theorem follows immediately from the LYM inequality that we prove below. \square

The acronym LYM stands for Lubell-Yamamoto-Meshalkin.

2.2 Lemma (LYM inequality) *Let $A \subseteq \mathcal{U}$ have cardinality n , finite. If $\varphi(A; b)_{b \in \mathcal{V}}$ is an antichain then*

$$\sum_{k=0}^n \left| \binom{A}{k} \cap \varphi(A; b)_{b \in \mathcal{V}} \right| \cdot \binom{n}{k}^{-1} \leq 1.$$

Proof Let Π be uniform random variable that ranges over the set of permutations of $A = \{a_1, \dots, a_n\}$. For any $\varphi(A; b)$ of cardinality k

$$\mathbb{P}\left(\Pi\{a_1, \dots, a_k\} = \varphi(A; b)\right) = \binom{n}{k}^{-1}.$$

The events above are disjoint for distinct sets $\varphi(A; b)$, hence

$$\mathbb{P}\left(\Pi\{a_1, \dots, a_k\} \in \varphi(A; b)_{b \in \mathcal{V}}\right) = \left|\binom{A}{k} \cap \varphi(A; b)_{b \in \mathcal{V}}\right| \cdot \binom{n}{k}^{-1}.$$

As $\varphi(A; b)_{b \in \mathcal{V}}$ is an antichain, for distinct k the events above are disjoint, hence

$$\mathbb{P}\left(\bigcup_{k=0}^n \Pi\{a_1, \dots, a_k\} \in \varphi(A; b)_{b \in \mathcal{V}}\right) = \sum_{k=0}^n \left|\binom{A}{k} \cap \varphi(A; b)_{b \in \mathcal{V}}\right| \cdot \binom{n}{k}^{-1}.$$

Now, the inequality is evident. \square

Let \mathbb{P}_k be the probability measure on the subsets of A that is concentrated and uniform on $A^{(k)}$. Namely, for $A' \subseteq A$

$$\mathbb{P}_k(\{A'\}) = \begin{cases} 0 & \text{if } |A'| \neq k \\ \binom{n}{k}^{-1} & \text{if } |A'| = k \end{cases}$$

Then the the LYM inequality asserts that if $\varphi(A; b)_{b \in \mathcal{V}}$ is an antichain then

$$\sum_{k=0}^n \mathbb{P}_k(\varphi(A; b)_{b \in \mathcal{V}}) \leq 1.$$

This inequality is strict when $\varphi(A; b)_{b \in \mathcal{V}} = A^{(k)}$ for some k . In the next section we show that these are the only cases.

2 The Erdős-Ko-Rado Theorem

2.3 Lemma (Peter J. Cameron) *Let G be a 1-transitive finite graph. If G contains a clique of cardinality m , then every subgraph $H \subseteq G$ contains a clique of cardinality*

$$\geq m \frac{|H|}{|G|}.$$

Proof Let C be a clique in G of cardinality m . Let k the cardinality of the largest clique in H . Let $n = |\text{Aut}(G)|$. By 1-transitivity, the sets $\{f \in \text{Aut}(G) : fa = b\}$, for any fixed $a \in G$ and b ranging over G , have all the same cardinality. Hence, for any given pair $\langle a, b \rangle$, they have cardinality $n/|G|$.

Count the pairs $\langle a, f \rangle \in C \times \text{Aut}(G)$ such that $fa \in H$. For every $a \in C$ there are $n \cdot |H|$ automorphisms. So the number of pairs is $m \cdot n \cdot |H|/|G|$

On the other hand for each $f \in \text{Aut}(G)$ there are at most k choices of $a \in C$. So $m \cdot n \cdot |H|/|G| \leq kn$. \square

2.4 Erdős-Ko-Rado Theorem *Let $A \subseteq \mathcal{U}$ be a finite set of cardinality n . Let $k \leq n/2$. Let $\varphi(A; b)_{b \in \mathcal{V}}$ be an intersecting family of sets of cardinality k . Then*

$$\left|\varphi(A; b)_{b \in \mathcal{V}}\right| \leq \binom{n-1}{k-1}.$$

Proof Let $m = \left| \varphi(A; b)_{b \in \mathcal{V}} \right|$. Consider the graph

$$G = \binom{A}{k},$$

$$E(G) = \left\{ \{A', A''\} : A' \cap A'' \neq \emptyset \right\}.$$

Enumerate the elements of A , say $A = \{a_0, \dots, a_{n-1}\}$. Consider the following subgraph of G

$$H = \left\{ \{a_i, \dots, a_{i+k-1}\} : 0 \leq i < n \right\},$$

where the indices are intended modulo n . As $k \leq n$, the largest clique in H has cardinality k . As $\varphi(A'; b)_{b \in \mathcal{V}}$ is a clique of G , by the lemma above,

$$k \geq m \frac{|H|}{|G|} = m \cdot n \cdot \binom{n}{k}^{-1}$$

therefore

$$m \leq \binom{n-1}{k-1} \quad \square$$

Chapter 3

Stable and NIP relations

1 Stable formulas

The **ladder-dimension** of $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$, or of $\varphi(x; z)$ when \mathcal{U} and \mathcal{V} are clear, is the maximal length n of a chain of the form

$$\varphi(A; b_0) \subset \dots \subset \varphi(A; b_{n-1})$$

for some set $A \subseteq \mathcal{U}$ and some $b_0, \dots, b_{n-1} \in \mathcal{V}$. If a maximal length exists we say that $\varphi(x; z)$ is **stable** otherwise we say that $\varphi(x; z)$ is **unstable**.

2 The Vapnik-Chervonenkis dimension

If all subsets of $A \subseteq \mathcal{U}$ are definable, that is $\mathcal{P}A = \varphi(A, b)_{b \in \mathcal{V}}$ we say that A is **shattered** by $\varphi(x; z)$. The following is called the **shatter function**

$$\pi_\varphi(n) = \max \left\{ |\varphi(A, b)_{b \in \mathcal{V}}| : A \in \binom{\mathcal{U}}{n} \right\}$$

So, $\pi_\varphi(n)$ gives the maximal number of definable subsets that a set of cardinality n can have. Trivially, $\pi_\varphi(n) \leq 2^n$ for all n . Moreover, if $\pi_\varphi(n) = 2^n$ for some n , then $\pi_\varphi(k) = 2^k$ for every $k \leq n$.

The **Vapnik-Chervonenkis dimension** of $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$, abbreviated by **VC-dimension**, is the maximal cardinality of a finite set $A \subseteq \mathcal{U}$ that is shattered by $\varphi(x; z)$. Equivalently, it is the maximal k such that $\pi_\varphi(k) = 2^k$. If such a maximum does not exist, we say that the VC-dimension is infinite or that $\varphi(x; z)$ has **IP** (the independence property). Otherwise, we say that $\varphi(x; z)$ has **NIP** (not the independence property). We may also say: *is IP*, or *is NIP*.

As \mathcal{U} and \mathcal{V} are usually clear from the context, we may say VC-dimension of $\varphi(x; z)$ for the VC-dimension of $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$.

3.1 Example If $\varphi(x; z)$ is either \top or \perp , then it shatters only the empty set, therefore it has VC-dimension 0. □

3.2 Example If $\varphi(x; z)$ has ladder dimension n then it has VC-dimension at most n . Hence stable formulas are NIP. □

3.3 Example If $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ is a non trivial chain of sets, then its VC-dimension is 1. □

3.4 Example Let $\mathcal{U} = \mathbb{R}$ and $\mathcal{V} = \mathbb{R}^2$. Let $\varphi(x; z_1, z_2)$ be the formula $z_1 < x < z_2$. Then its VC-dimension 2. □

3.5 Example Let $\mathcal{U} = \mathcal{V} = \mathbb{R}^2$. Let $\varphi(x_1, x_2; z_1, z_2)$ be the formula $y < z_1 \cdot x + z_2$. Then its VC-dimension 3 (by Radon's Theorem). \square

3.6 Example If $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ is the set of all subsets of \mathcal{U} of cardinality $\leq k$. Then its VC-dimension is k and

$$\pi_\varphi(n) = \sum_{i=0}^k \binom{n}{i}. \quad \square$$

We call the VC-dimension of $\varphi(x; z)^*$ the **dual VC-dimension** of $\varphi(x; z)$.

3.7 Proposition If $\varphi(x; z)$ has VC-dimension $< k$ then its dual VC-dimension is $< 2^k$.

Proof Suppose that the VC-dimension of $\varphi(x; z)^*$ is at least 2^k . We prove that the VC-dimension of $\varphi(x; z)$ is at least k . Let $B = \{b_I : I \subseteq k\}$ be a set of cardinality 2^k shattered by $\varphi(x; z)^*$. That is, for every $\mathcal{J} \subseteq \mathcal{P}(k)$ there is $a_{\mathcal{J}}$ such that

$$\varphi(a_{\mathcal{J}}, b_I) \Leftrightarrow I \in \mathcal{J}$$

Let $a_i = a_{\{I \subseteq k : i \in I\}}$. Then from the equivalence above we obtain

$$\varphi(a_i, b_I) \Leftrightarrow i \in I$$

That is, $\varphi(x; z)$ shatters $A = \{a_i : i \in k\}$. \square

3.8 Exercise (???) Prove that there is a formula with VC-dimension k and dual dimension 2^k . \square

3 The Sauer-Shelah lemma

According to Gil Kalai in [5], Sauer-Shelah's Lemma can be described as an *eigen-theorem* because it is important in many different areas of mathematic (model theory, learning theory, probability theory, ergodic theory, Banach spaces, to name a few). No wonder it has been discovered and rediscovered may times.

It has been proved independently by Shelah [8], Sauer [7], and Vapnik-Chervonenkis [9] around 1970 (Shelah gives credit to Micha Perles). Saharon Shelah was working in model theory while Norbert Sauer, Vladimir Vapnik and Alexey Chervonenkis were in statistical learning theory.

3.9 Sauer-Shelah Lemma If $\varphi(x; z)$ has VC-dimension k then for every $n \geq k$

$$\pi_\varphi(n) \leq \sum_{i=0}^k \binom{n}{i}. \quad \square$$

The set system presented in Example 3.6 shows that the bound is optimal.

An alternative proof of the Sauer-Shelah Lemma derives it as corollary of a lemma by Alain Pajor [6].

3.10 Pajor's Lemma Let $A \subseteq \mathcal{U}$ be finite.

$$|\varphi(A, b)_{b \in \mathcal{V}}| \leq \left| \{C \subseteq A : C \text{ is shattered by } \varphi(x; z)\} \right|.$$

Proof If A is empty then $|\varphi(A, b)_{b \in \mathcal{V}}| = 1$ and \emptyset is the only subset of A that φ shatters, so the inequality holds trivially. Otherwise, pick an $a \in A$ and assume the lemma holds for $A' = A \setminus \{a\}$. Define

$$\psi(x; y) = \varphi(x; y) \wedge \neg \varphi(a; y) \wedge \exists y' [\varphi(a; y') \wedge \varphi(A'; y') = \varphi(A'; y)].$$

Notice that

$$|\varphi(A, b)_{b \in \mathcal{V}}| = \left| \varphi(A', b)_{b \in \mathcal{V}} \cup \left\{ \{a\} \cup \psi(A', b) : b \in \mathcal{V} \right\} \right|.$$

as the two sets in the r.h.s. are disjoint

$$|\varphi(A, b)_{b \in \mathcal{V}}| = |\varphi(A', b)_{b \in \mathcal{V}}| + |\psi(A', b)_{b \in \mathcal{V}}|.$$

By induction hypothesis,

$$|\varphi(A', b)_{b \in \mathcal{V}}| \leq \left| \{C \subseteq A' : C \text{ is shattered by } \varphi(x; z)\} \right| \quad (1)$$

and

$$\begin{aligned} |\psi(A', b)_{b \in \mathcal{V}}| &\leq \left| \{C \subseteq A' : C \text{ is shattered by } \psi(x; z)\} \right| \\ &= \left| \{C \subseteq A' : C \cup \{a\} \text{ is shattered by } \varphi(x; z)\} \right|. \end{aligned} \quad (2)$$

In fact, $C \subseteq A'$ is shattered by $\psi(x; y)$ if and only if $C \cup \{a\}$ is shattered by $\varphi(x; y)$. Clearly,

$$(1) + (2) = \left| \{C \subseteq A : C \text{ is shattered by } \varphi(x; z)\} \right|,$$

so the lemma follows. \square

Proof of the Sauer-Shelah Lemma Assume $\varphi(x; z)$ has VC-dimension k and let $n \geq k$. Then

$$\begin{aligned} \pi_\varphi(n) &= \max_{|A|=n} |\varphi(A, b)_{b \in \mathcal{V}}| \\ \pi_\varphi(n) &\leq \max_{|A|=n} \left| \{C \subseteq A : C \text{ shattered by } \varphi(x; z)\} \right| \quad \text{by Pajor's Lemma} \\ &\leq \sum_{i=0}^k \binom{n}{i} \quad \text{because } \varphi(x; z) \text{ has VC-dimension } k \end{aligned} \quad \square$$

We write $f(n) = O(g(n))$ if there is a constant C such that $|f(n)| \leq Cg(n)$ holds for all (sufficiently large) n .

The **VC-density** of $\varphi(x; z)$ is the infimum over all real number r such that $\pi_\varphi(n) \in O(n^r)$. It is infinite if no such r exist. The **dual VC-density** is defined accordingly.

By the Sauer-Shelah lemma the VC-density is at most as large as the VC-dimension. It could be smaller, however it is usually rather difficult to compute.

We conclude this section with a couple of inequalities that is useful to have at hand.

$$\sum_{i=0}^k \binom{n}{i} = \sum_{i=0}^k \frac{n!}{i! (n-i)!}$$

$$\begin{aligned}
&\leq \sum_{i=0}^k \frac{n^i}{i!} \\
&\leq \sum_{i=0}^k \frac{n^i k!}{i!(k-i)!} \\
&= \sum_{i=0}^k n^i \binom{k}{i} \\
&= (n+1)^k \quad \text{by the binomial theorem.}
\end{aligned}$$

There is a second bound, which is better when $k \geq 3$ and holds for $n > k$

$$\begin{aligned}
\sum_{i=0}^k \binom{n}{i} &\leq \left(\frac{n}{k}\right)^k \sum_{i=0}^k \left(\frac{k}{n}\right)^i \binom{n}{i} && \text{because } \frac{k}{n} < 1 \\
&\leq \left(\frac{n}{k}\right)^k \sum_{i=0}^n \left(\frac{k}{n}\right)^i \binom{n}{i} \\
&= \left(\frac{n}{k}\right)^k \left(1 + \frac{k}{n}\right)^n && \text{by the binomial theorem} \\
&\leq \left(\frac{n e}{k}\right)^k && \text{where } e \text{ is the base of the natural logarithm.}
\end{aligned}$$

Chapter 4

Law(s) of large numbers

Quoting from some unpublished notes by Carlos C. Rodríguez

What is a Law of Large Numbers? I am glad you asked! The Laws of Large Numbers, or LLNs for short, come in three basic flavors: Weak, Strong and Uniform. They all state that the observed frequencies of events tend to approach the actual probabilities as the number of observations increases. Saying it in another way, the LLNs show that under certain conditions, we can asymptotically learn the probabilities of events from their observed frequencies. To add some drama we could say that if God is not cheating and S/he doesn't change the initial standard probabilistic model too much then, in principle, we (or other machines, or even the universe as a whole) could eventually find out the Truth, the whole Truth, and nothing but the Truth.

Bull! The Devil, is in the details.

I suspect that for reasons not too different in spirit to the ones above, famous minds of the past took the slippery slope of defining probabilities as the limits of relative frequencies. They became known as “frequentist”. They wrote books and indoctrinated generations of confused students.

1 Inequalities

Throughout this and the next section we work with a given probability space \mathcal{U}, \mathbb{P} . For simplicity, the following two propositions are proved for finite \mathcal{U} , but they are easily seen to hold in general.

4.1 Definition A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is *convex* if for every tuples of real numbers p_i and x_i such that

$$\sum_{i=1}^n p_i = 1$$

we have

$$f\left(\sum_{i=1}^n p_i x_i\right) \leq \sum_{i=1}^n p_i f(x_i).$$

□

Note that, though the definition is usually given with $n = 2$, the general property above follows easily.

4.2 Jensen's Inequality Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a convex function. Then

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$$

Proof For simplicity, assume that the sample space \mathcal{U} is finite. Then the claim is obvious from the definition. \square

The following is arguably the most basic inequality in probability theory. Although it is almost trivial, it will be required several times in this chapter.

4.3 Markov's Inequality Let X be a nonnegative random variable with finite mean. Then for every $\varepsilon > 0$

$$\mathbb{P}(X \geq \varepsilon) \leq \frac{\mathbb{E}[X]}{\varepsilon}$$

Proof For simplicity, assume that the sample space \mathcal{U} is finite. (The theorem holds in general, but we only need the finite case.) Define $A = \{a \in \mathcal{U} : X(a) \geq \varepsilon\}$.

$$\begin{aligned} \mathbb{E}[X] &= \sum_{a \in \mathcal{U}} \mathbb{P}(a) X(a) \\ &= \sum_{a \in A} \mathbb{P}(a) X(a) + \sum_{a \notin A} \mathbb{P}(a) X(a) \\ &\geq \sum_{a \in A} \mathbb{P}(a) X(a) \\ &\geq \varepsilon \sum_{a \in A} \mathbb{P}(a) \\ &= \varepsilon \mathbb{P}(X \geq \varepsilon) \end{aligned}$$

\square

4.4 Corollary Let X be a nonnegative random variable. If $\mathbb{E}[X^k]$ exists, then for every $\varepsilon > 0$

$$\mathbb{P}(X \geq \varepsilon) \leq \frac{\mathbb{E}[X^k]}{\varepsilon^k}$$

Proof By Markov's inequality, since $\mathbb{P}(X \geq \varepsilon) = \mathbb{P}(X^k \geq \varepsilon^k)$. \square

Chebyshev's inequality (a.k.a. Chebysheff, Chebyshev, Tchebyscheff, Tschebycheff) is a special case of the corollary above.

4.5 Chebyshev's Inequality Let X be a random variable with finite mean and variance. Then for every $\varepsilon > 0$

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq \varepsilon) \leq \frac{\text{Var}[X]}{\varepsilon^2}$$

\square

To obtain exponential bounds, we frequently apply the following trick.

4.6 Chernoff's method Let X be a random variable with finite mean. Then for every $t > 0$

$$\mathbb{P}(X \geq \varepsilon) \leq e^{-t\varepsilon} \mathbb{E}[e^{tX}]$$

Proof For every $t > 0$

$$\mathbb{P}(X \geq \varepsilon) = \mathbb{P}(e^{tX} \geq e^{t\varepsilon}) \quad \text{because } e^{tx} \text{ is increasing}$$

$$\leq e^{-t\epsilon} \mathbb{E}[e^{tX}],$$

by Markov's inequality, which we may apply since e^{tX} is always positive. \square

4.7 Hoeffding's lemma Let X be a bounded random variable, say $a \leq X \leq b$. Let $\mathbb{E}[X] = \mu$ and $d = b - a$. Then

$$\mathbb{E}[e^{t(X-\mu)}] \leq \exp\left(\frac{t^2 d^2}{8}\right).$$

Proof For clarity, assume $\mu = 0$. The general result follows easily from this special case by centralization. Recall that, by convexity, for every $x \in [a, b]$

$$e^{tx} \leq \frac{x-a}{d} e^{tb} + \frac{b-x}{d} e^{ta}$$

Then

$$e^{tX} \leq \frac{X-a}{d} e^{tb} + \frac{b-X}{d} e^{ta}$$

By the linearity of expectation,

$$\mathbb{E}[e^{tX}] \leq \frac{b e^{ta} - a e^{tb}}{d}$$

$$\log \mathbb{E}[e^{tX}] \leq \log \frac{b e^{ta} - a e^{tb}}{d}$$

taking the Taylor series expansion of the r.h.s. at $t = 0$ (the first and second derivatives vanish at 0; the second derivative is always $\leq 1/4$) we obtain

$$\log \mathbb{E}[e^{tX}] \leq \frac{t^2 d^2}{8}. \quad \square$$

4.8 Hoeffding's Inequality Let X_1, \dots, X_n be independent random variables with bounded range, say $a \leq X_i \leq b$. Define $d = b - a$.

$$M = \sum_{i=1}^n (X_i - \mathbb{E}[X_i])$$

Then for every $\epsilon > 0$

$$\mathbb{P}(M \geq \epsilon) \leq \exp\left(-\frac{2\epsilon^2}{nd^2}\right),$$

$$\mathbb{P}(M \leq -\epsilon) \leq \exp\left(-\frac{2\epsilon^2}{nd^2}\right).$$

Clearly, the two inequalities above imply the following

$$\mathbb{P}(|M| \geq \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{nd^2}\right).$$

Proof Define $\mathbb{E}[X_i] = \mu_i$. Let $t > 0$ be arbitrary.

$$\mathbb{P}(M \geq \epsilon) \leq e^{-t\epsilon} \mathbb{E}[e^{tM}] \quad \text{by Chernoff's method (4.6)}$$

$$= e^{-t\epsilon} \prod_{i=1}^n \mathbb{E}[e^{t(X_i - \mu_i)}] \quad \text{by independence.}$$

$$\leq e^{-t\epsilon} \prod_{i=1}^n \exp\left(\frac{t^2 d^2}{8}\right) \quad \text{by Hoeffding's Lemma (4.7).}$$

$$= \exp\left(\frac{nt^2d^2}{8} - t\varepsilon\right)$$

Now substitute $4\varepsilon/nd^2$ for t . □

We prove Hoeffding's lemma with a slightly weaker bound (2 for 8). The purpose is to present two clever tricks *ghost sample* and *symmetrization* which in the following section is applied in a more complex setting.

First we need the following lemma. A **random sign variable** (a.k.a. Rademacher random variable) is a random variable $\sigma \in \{-1, 1\}$ with mean 0.

4.9 Lemma *Let σ be a random sign variable. Then for every $t > 0$*

$$\mathbb{E}\left[e^{t\sigma}\right] \leq e^{t^2/2}$$

Proof Replace e^x with its Taylor expansion around $x = 0$

$$\begin{aligned} \mathbb{E}\left[e^{t\sigma}\right] &= \sum_{i=0}^{\infty} \frac{t^i \mathbb{E}[\sigma^i]}{i!} \\ &= \sum_{i=0}^{\infty} \frac{t^{2i}}{(2i)!} && \text{since } \mathbb{E}[\sigma^i] = \begin{cases} 1 & i \text{ even} \\ 0 & i \text{ odd} \end{cases} \\ &= \sum_{i=0}^{\infty} \frac{(t/2)^{2i}}{i!} \\ &= e^{t^2/2}. \end{aligned} \quad \square$$

4.10 Second proof of Hoeffding's Lemma Recall that Hoeffding's Lemma claims that, if $a \leq X \leq b$, then

$$\mathbb{E}\left[e^{t(X-\mu)}\right] \leq \exp\left(\frac{t^2d^2}{8}\right),$$

where $\mathbb{E}[X] = \mu$ and $d = b - a$.

Let X' be an independent copy of X (a.k.a. ghost sample). In particular $\mu = \mathbb{E}(X')$. Then

$$\begin{aligned} \mathbb{E}\left[e^{t(X-\mu)}\right] &= \mathbb{E}\left[e^{t(X-\mathbb{E}[X'])}\right] \\ &\leq \mathbb{E}\left[\mathbb{E}\left[e^{t(X-X')} \mid X\right]\right] && \text{by Jensen's inequality} \\ &\leq \mathbb{E}\left[e^{t(X-X')}\right] \end{aligned}$$

Let σ be a random sign variable independent of X, X' . Then $\sigma(X - X')$ has the same distribution of $X - X'$.

$$\begin{aligned} &= \mathbb{E}\left[e^{t\sigma(X-X')}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[e^{t\sigma(X-X')} \mid X, X'\right]\right] \\ &\leq \mathbb{E}\left[e^{t^2(X-X')^2/2}\right] && \text{by Lemma 4.9} \\ &\leq e^{t^2d^2/2} && \text{because } |X - X'| \leq d. \end{aligned}$$

This yields the bound above (only with 2 in place of 8). □

2 Two Weak Laws of Large Numbers

A **sample** s is a sequence s_1, \dots, s_n of elements of \mathcal{U} . Its length $|s| = n$ is also called **size** or **dimension**. We write $\text{range}(s)$ for the set $\{s_1, \dots, s_n\}$. Note that this set may have cardinality $< n$.

To a sample s of size n we associate a finite probability measure on the subsets of \mathcal{U} namely, for any event $A \subseteq \mathcal{U}$, we define the empirical frequency of A given s

$$\text{Fr}(s, A) = \frac{1}{n} \cdot |\{i : s_i \in A\}|.$$

It is convenient to rewrite it using indicator functions

$$\text{Fr}(s, A) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{s_i \in A}.$$

We are interested in the *existence* of samples that approximate the probability within ε . Suppose that, for a given event A , we can prove that

$$\mathbb{P}(s \in \mathcal{U}^n : |\text{Fr}(s, A) - \mathbb{P}(A)| \geq \varepsilon) \leq \text{some_bound}(\varepsilon, n)$$

and that, for n large enough, $\text{some_bound}(\varepsilon, n) < 1$. Then a sample of size $\leq n$ that approximate the probability within ε is guaranteed to exist.

Random variables are convenient formalism to discuss these probabilities. We say **random element** of \mathcal{U} for a random variables S such that $\mathbb{P}(S \in A) = \mathbb{P}(A)$ for every $A \subseteq \mathcal{U}$. A **random sample** from \mathcal{U} is a tuple $S = S_1, \dots, S_n$ be of independent random elements of \mathcal{U} . Then $\mathbb{I}_{S_i \in A}$ as a Bernoulli random variable with probability of success $\mathbb{P}(A)$ and

$$\text{Fr}(S, A) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{S_i \in A}$$

is (up to the factor $1/n$) a binomial random variable.

4.11 Weak Law of Large Numbers For every event $A \subseteq \mathcal{U}$ and every $n > 0$

$$\mathbb{P}(|\text{Fr}(s, A) - \mathbb{P}(A)| \geq \varepsilon) \leq \frac{1}{n\varepsilon^2}.$$

Proof Let S_1, \dots, S_n be a random sample from \mathcal{U} . Up to the factor $1/n$, the distribution of $\text{Fr}(S, A)$ is binomial with parameters n and $\mathbb{P}(A)$. Therefore it has expected value $\mathbb{P}(A)$ and variance $\leq 1/n$. By Chebyshev's inequality we obtain

$$\mathbb{P}(|\text{Fr}(S, A) - \mathbb{P}(A)| \geq \varepsilon) \leq \frac{1}{n\varepsilon^2}$$

which proves the theorem. □

Sometime we are interested in the minimal size of a sample that approximates the probability up to a given ε .

4.12 Corollary Assume \mathcal{U} is finite (of arbitrary cardinality, though). For every $A \subseteq \mathcal{U}$ and every $\varepsilon > 0$ there is a sample s of size

$$|s| = \left\lfloor \frac{1}{\varepsilon^2} + 1 \right\rfloor$$

such that

$$\left| \text{Fr}(s, A) - \mathbb{P}(A) \right| < \varepsilon.$$

Proof By the Weak Law of Large Numbers above, a sample of size n exists if

$$\frac{1}{n\varepsilon^2} < 1 \quad \square$$

In the following section we need a better bound for the Weak Law of Large Numbers. This is obtained with a similar proof.

4.13 Weak Law of Large Numbers (with exponential bound) For every event $A \subseteq \mathcal{U}$ and every $n > 0$

$$2e^{-2n\varepsilon^2} \geq \mathbb{P}\left(s \in \mathcal{U}^n : \left| \text{Fr}(s, A) - \mathbb{P}(A) \right| \geq \varepsilon\right).$$

Proof Let S_1, \dots, S_n be a random sample from \mathcal{U} . Define

$$M = \sum_{i=1}^n \left(\mathbb{I}_{S_i \in A} - \mathbb{E}[\mathbb{I}_{S_i \in A}] \right)$$

As $\mathbb{E}[\mathbb{I}_{S_i \in A}] = \mathbb{P}(A)$, the inequality we have to prove can be rewritten as

$$2e^{-2n\varepsilon^2} \geq \mathbb{P}\left(|M| \geq n\varepsilon\right)$$

and this follows immediately from Hoeffding inequality. \square

Using the exponential bounds above, we can improve (by a constant factor) the size of the minimal sample size that approximates the probability obtained in Corollary 4.12.

4.14 Corollary For every $A \subseteq \mathcal{U}$ and every $\varepsilon > 0$ there is a sample s of size n where

$$n = \left\lfloor \frac{\log 2}{2\varepsilon^2} + 1 \right\rfloor$$

such that

$$\left| \text{Fr}(s, A) - \mathbb{P}(A) \right| < \varepsilon. \quad \square$$

3 A Uniform Law of Large Numbers

Throughout this section we work with a fixed family of definable subsets $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ that are events of the sample space \mathcal{U}, \mathbb{P} . It is convenient to introduce some abbreviations

$$\mathbb{P}(b) = \mathbb{P}\left(\varphi(\mathcal{U}; b)\right)$$

$$\text{Fr}(s, b) = \text{Fr}\left(s, \varphi(\mathcal{U}; b)\right)$$

$$\mathbb{I}_{s, b} = \mathbb{I}_{\varphi(s; b)}$$

An **ε -approximation** is a sample s such that

$$\left| \text{Fr}(s, b) - \mathbb{P}(b) \right| < \varepsilon \quad \text{for every } b \in \mathcal{V}.$$

We are interested in estimating the minimal size of an ε -approximation.

The main theorem of this section is this famous result of Vapnik-Chervonenkis [9].

4.15 Vapnik-Chervonenkis Inequality *Let $\pi_\varphi(n)$ be the shatter function of $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$. Let $S = S_1, \dots, S_n$ be a random sample from \mathcal{U} . Then, for every $b \in \mathcal{V}$*

$$\mathbb{P} \left(\left| \text{Fr}(S, b) - \mathbb{P}(b) \right| \geq \varepsilon \right) \leq 6 \pi_\varphi(n) \exp \left(- \frac{n\varepsilon^2}{32} \right).$$

N.B. Some technical hypothesis of measurability are necessary when $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ is uncountable. This have been omitted in the statement above, and will be discussed below.

Proof Let $S' = S'_1, \dots, S'_n$ be an independent copy of S . Then, by the triangular inequality

$$\mathbb{P} \left(\left| \text{Fr}(S, b) - \mathbb{P}(b) \right| \geq \varepsilon \right) \leq \mathbb{P} \left(\left| \text{Fr}(S, b) - \text{Fr}(S', b) \right| \geq \frac{\varepsilon}{2} \right) + (*)$$

where, by the Weak Law of Large Numbers 4.13,

$$\begin{aligned} (*) &= \mathbb{P} \left(\left| \text{Fr}(S', b) - \mathbb{P}(b) \right| \geq \frac{\varepsilon}{2} \right) \\ &\leq 2 e^{-n\varepsilon^2/2} \end{aligned}$$

Let $\sigma = \sigma_1, \dots, \sigma_n$ be a tuple of independent sign random variables. Then

$$\mathbb{P} \left(\left| \text{Fr}(S, b) - \text{Fr}(S', b) \right| \geq \frac{\varepsilon}{2} \right) = \mathbb{P} \left(\sum_{i=1}^n |\mathbb{I}_{S_i, b} - \mathbb{I}_{S'_i, b}| \geq \frac{n\varepsilon}{2} \right)$$

Then, again by the triangular inequality

$$\leq 2 \mathbb{P} \left(\left| \sum_{i=1}^n \sigma_i \mathbb{I}_{S_i, b} \right| \geq \frac{n\varepsilon}{4} \right)$$

Putting together the inequalities above we obtain

$$(1) \quad \mathbb{P} \left(\sup_{b \in \mathcal{V}} \left| \text{Fr}(S, b) - \mathbb{P}(b) \right| \geq \varepsilon \right) \leq 2 \mathbb{P} \left(\sup_{b \in \mathcal{V}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{S_i, b} \right| \geq \frac{n\varepsilon}{4} \right) + 2 e^{-n\varepsilon^2/2}$$

Let $s = s_1, \dots, s_n$ be a possible realization of S .

$$\mathbb{P} \left(\sup_{b \in \mathcal{V}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{S_i, b} \right| \geq \frac{n\varepsilon}{4} \right) = \mathbb{P} \left(\sup_{b \in \mathcal{V}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{S_i, b} \right| \geq \frac{n\varepsilon}{4} \mid S = s \right)$$

Finally, note that the r.h.s. of (2) only depends on $\varphi(\text{range}(s), b)$. Hence the $\sup_{b \in \mathcal{V}}(\cdot)$ on the r.h.s. is actually a maximum among $m = \pi_\varphi(n)$ events. Say, we can choose $b_1, \dots, b_m \in \mathcal{V}$ such that

$$(2) \quad = \mathbb{P} \left(\sup_{j \leq m} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{S_i, b_j} \right| \geq \frac{n\varepsilon}{4} \right)$$

Note that, in general, for any real random variables X_1, \dots, X_m we have

$$\mathbb{P} \left(\sup_{i \leq m} X_i \geq \varepsilon \right) = \mathbb{P} \left(\bigcup_{i=0}^m X_i \geq \varepsilon \right)$$

$$\leq \sum_{i=1}^m \mathbb{P}(X_i \geq \varepsilon)$$

Hence, continuing from (2) we obtain

$$\begin{aligned} &\leq \sum_{j=1}^m \mathbb{P}\left(\left|\sum_{i=1}^n \sigma_i \mathbb{I}_{s_i, b_j}\right| \geq \frac{n\varepsilon}{4}\right) \\ &\leq 2\pi_\varphi(n) \exp\left(-\frac{n\varepsilon^2}{32}\right), \end{aligned}$$

where the last inequality is obtained from Hoeffding's Inequality 4.8. In fact, Hoeffding's Inequality, applied to $X_i = \sigma_i \mathbb{I}_{s_i, b}$ with $n\varepsilon/4$ for ε , yields

$$\mathbb{P}\left(\left|\sum_{i=1}^n \sigma_i \mathbb{I}_{s_i, b}\right| \geq \frac{n\varepsilon}{4}\right) \leq 2 \exp\left(-\frac{n\varepsilon^2}{32}\right).$$

Unsurprisingly, the bound does not depend on b . Finally, proceeding from (1) we obtain

$$\begin{aligned} &\leq 4\pi_\varphi(n) \exp\left(-\frac{n\varepsilon^2}{32}\right) + 2 \exp\left(-\frac{n\varepsilon^2}{2}\right) \\ &\leq 6\pi_\varphi(n) \exp\left(-\frac{n\varepsilon^2}{32}\right), \end{aligned}$$

which finally proves the theorem. \square

4.16 Corollary *Let $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ have finite VC-dimension. For every $\varepsilon > 0$ there is a finite sample s such that*

$$\left| \text{Fr}(s, b) - \mathbb{P}(b) \right| < \varepsilon \quad \text{for every } b \in \mathcal{V}.$$

Proof It suffices to require that $n = |s| = |\mathcal{S}|$ is large enough to guarantee

$$\mathbb{P}\left(\left| \text{Fr}(S, b) - \mathbb{P}(b) \right| \geq \varepsilon\right) < 1 \quad \text{for every } b \in \mathcal{V}.$$

By the Vapnik-Chervonenkis inequality 4.15, it suffices that

$$(3) \quad 6\pi_\varphi(n) \exp\left(-\frac{n\varepsilon^2}{32}\right) < 1$$

By the Sauer-Shelah Lemma 3.9, $\pi_\varphi(n)$ grows polynomially. Hence the inequality holds for n large enough. \square

The corollary above is sufficient for our intended applications. For completeness, the following proposition gives an explicit bound.

4.17 Proposition *There is a sample s as in the corollary above of size*

$$c \frac{k}{\varepsilon^2} \log \frac{k}{\varepsilon^2}$$

where c is an absolute constant and k is the VC-dimension of $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$.

Proofsketch By (3) in the proof above and the inequality proved after the Sauer-Shelah Lemma 3.9 it suffices that $n = |s|$ satisfies

$$\log 6 + k \log(n+1) < \frac{n\varepsilon^2}{32},$$

which is the case if n satisfies the following inequality

$$c' \frac{k}{\varepsilon^2} < \frac{n}{\log n},$$

for some absolute constant c' . Finally, the latter inequality is satisfied if

$$c \frac{k}{\varepsilon^2} \log \frac{k}{\varepsilon^2} < n$$

for some absolute constant c . □

4 A Uniform Law of Large Numbers, again

We prove a second version of the Vapnik-Chervonenkis Inequality. Which, I conjecture, is due to Devroye and Lugosi [3].

4.18 Vapnik-Chervonenkis Inequality (2) *Let $\pi_\varphi(n)$ be the shatter function of $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$. Let $S = S_1, \dots, S_n$ be a random sample from \mathcal{U} . Then, for every $b \in \mathcal{V}$*

$$\mathbb{E} \left| \text{Fr}(S, b) - \mathbb{P}(b) \right| \leq 2 \sqrt{\frac{\log(2 \pi_\varphi(n))}{n}}. \quad \square$$

The same caveat on measurability apply as for Inequality 4.15.

We note that the bound is not optimal, using a clever techniques called *chaining*, Dudley could prove that

$$\mathbb{E} \left| \text{Fr}(S, b) - \mathbb{P}(b) \right| < c \sqrt{\frac{k}{n}},$$

where k is the VC-dimension and c is absolute constant.

Before embarking in the proof of the theorem above, we prove the following (easy, although mysterious) lemma, which also has independent interest.

4.19 Lemma *Let X_1, \dots, X_m be real valued random variables. Let c be such that*

$$\mathbb{E}[e^{tX_i}] \leq e^{c^2 t^2 / 2} \quad \text{for every } i \leq m \text{ and every } t > 0.$$

Then

$$\mathbb{E} \left[\max_{i \leq m} X_i \right] \leq c \sqrt{2 \log m}.$$

If in addition

$$\mathbb{E}[e^{-tX_i}] \leq e^{c^2 t^2 / 2} \quad \text{for every } i \leq m \text{ and every } t > 0,$$

then

$$\mathbb{E} \left[\max_{i \leq m} |X_i| \right] \leq c \sqrt{2 \log(2m)}.$$

Proof By Jensen's inequality,

$$\begin{aligned} \exp \left(t \cdot \mathbb{E} \left[\max_{i \leq m} X_i \right] \right) &\leq \mathbb{E} \left[\exp \left(\max_{i \leq m} t X_i \right) \right] \\ &= \mathbb{E} \left[\max_{i \leq m} e^{t X_i} \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left[\sum_{i \leq m} e^{tX_i} \right] \\
&= \sum_{i \leq m} \mathbb{E} [e^{tX_i}] \\
&= m e^{c^2 t^2 / 2}
\end{aligned}$$

Taking the logarithm of both sides and replacing t with $\frac{\sqrt{2 \log m}}{c}$, we obtain the first inequality of the lemma.

To prove the second inequality, apply the first one to $X_1, \dots, X_m, -X_1, \dots, -X_m$. (N.B. note that independence is not assumed.) \square

Proof of the Vapnik-Chervonenkis inequality Let $S' = S'_1, \dots, S'_n$ be an independent copy of S . We claim that

$$(1) \quad \mathbb{E} \left[\sup_{b \in \mathcal{V}} |\text{Fr}(S, b) - \mathbb{P}(b)| \right] \leq \mathbb{E} \left[\sup_{b \in \mathcal{V}} |\text{Fr}(S, b) - \text{Fr}(S', b)| \right]$$

In fact,

$$\begin{aligned}
\text{Fr}(S, b) - \mathbb{P}(b) &= \text{Fr}(S, b) - \mathbb{E} [\text{Fr}(S', b)] \\
&= \mathbb{E} [\text{Fr}(S, b) - \text{Fr}(S', b) \mid S].
\end{aligned}$$

Now, apply Jensen's inequality to the absolute value function, then use that

$$(2) \quad \sup_{b \in \mathcal{V}} \mathbb{E} [\dots] \leq \mathbb{E} [\sup_{b \in \mathcal{V}} (\dots)].$$

Write \mathbb{I}_b for the indicator function of $\varphi(\mathcal{U}; b)$. Then

$$|\text{Fr}(S, b) - \text{Fr}(S', b)| = \frac{1}{n} \left| \sum_{i=1}^n (\mathbb{I}_{S_i, b} - \mathbb{I}_{S'_i, b}) \right|$$

Let $\sigma = \sigma_1, \dots, \sigma_n$ be a tuple of independent sign random variable. The random variable $\mathbb{I}_{S_i, b} - \mathbb{I}_{S'_i, b}$ has the same distribution of $\sigma_i (\mathbb{I}_{S_i, b} - \mathbb{I}_{S'_i, b})$ hence

$$= \frac{1}{n} \mathbb{E} \left[\sum_{i=1}^n \sigma_i (\mathbb{I}_{S_i, b} - \mathbb{I}_{S'_i, b}) \mid S, S' \right]$$

Inserting this into (1) we obtain

$$\mathbb{E} \left[\sup_{b \in \mathcal{V}} |\text{Fr}(S, b) - \mathbb{P}(b)| \right] \leq \frac{1}{n} \mathbb{E} \left[\sup_{b \in \mathcal{V}} \mathbb{E} \left[\sum_{i=1}^n \sigma_i (\mathbb{I}_{S_i, b} - \mathbb{I}_{S'_i, b}) \mid S, S' \right] \right]$$

Let s, s' be a generic realization of S, S'

$$\leq \frac{1}{n} \sup_{s, s'} \sup_{b \in \mathcal{V}} \mathbb{E} \left[\sum_{i=1}^n \sigma_i (\mathbb{I}_{s_i, b} - \mathbb{I}_{s'_i, b}) \right]$$

and, by what remarked in (2)

$$\leq \frac{1}{n} \sup_{s, s'} \mathbb{E} \left[\sup_{b \in \mathcal{V}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{s_i, b} - \mathbb{I}_{s'_i, b}) \right| \right]$$

Observe that once s, s' is fixed, $\sup_{b \in \mathcal{V}}$ is actually a maximum among $\pi_\varphi(2n)$ sets, in fact, $\pi_\varphi(2n)$ is the number of definable subsets of $A = \{s_1, \dots, s_n, s'_1, \dots, s'_n\}$. Then, by Lemma 4.19 (the second inequality, with $m = \pi_\varphi(2n)$ and i ranging over the definable subsets of A), for an appropriate constant c ,

$$\leq \frac{1}{n} \sup_{s,s'} c \sqrt{2 \log (2 \pi_{\varphi}(2n))}.$$

As the r.h.s. does not depend on s, s' ,

$$\leq \frac{c}{n} \sqrt{2 \log (2 \pi_{\varphi}(2n))}$$

Finally, as $\pi_{\varphi}(2n) \leq \pi_{\varphi}(n)^2$,

$$\leq \frac{2c}{n} \sqrt{\log (2 \pi_{\varphi}(n))}$$

The Vapnik-Chervonenkis inequality is proved if we can show the assumption of Lemma 4.19 holds with $c = \sqrt{n}$.

$$\mathbb{E} \left[\exp \left(t \sum_{i=1}^n \sigma_i (\mathbb{I}_{s_i, b} - \mathbb{I}_{s'_i, b}) \right) \right] = \prod_{i=1}^n \mathbb{E} \left[\exp \left(t \sigma_i (\mathbb{I}_{s_i, b} - \mathbb{I}_{s'_i, b}) \right) \right]$$

As $\sigma_i (\mathbb{I}_{s_i, b} - \mathbb{I}_{s'_i, b})$ takes values in $\{-1, 1\}$ with mean 0, by Lemma 4.9

$$\leq e^{nt^2/2}$$

and the same holds for $-\sigma_i (\mathbb{I}_{s_i, b} - \mathbb{I}_{s'_i, b})$. □

As an application we prove the Glivenko-Cantelli Theorem, an important theorem of mathematical statistics. The theorem says that the empirical cumulative distribution function converges uniformly to the true one. We prove an informative variant which gives the rate of convergence (though, this is not optimal).

4.20 Glivenko-Cantelli Theorem *Let $X = X_1, \dots, X_n$ be i.i.d. random variables. Let $F(z) = \mathbb{P}(X_i \leq z)$ be their common cumulative distribution function. Let $F_e(x)$ be the empirical cumulative distribution function, that is,*

$$F_e(X, z) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{X_i \leq z}.$$

Then, for every $z \in \mathbb{R}$,

$$\mathbb{E} \left| F_e(X, z) - F(z) \right| \leq 2 \sqrt{\frac{\log (2n+2)}{n}}. \quad \square$$

Proof Take $\mathcal{U} = \mathcal{V} = \mathbb{R}$ and $\varphi(x; z) = x < z$. Assign to $\varphi(\mathcal{U}; b) = (-\infty, b]$ the probability measure $\mathbb{P}(X_i \leq b)$. Then, if $S = S_1, \dots, S_n$ is a random sample from \mathcal{U}

$$\mathbb{E} \left| F_e(X, z) - F(z) \right| = \mathbb{E} \left| \text{Fr}(S, b) - \mathbb{P}(b) \right|$$

hence, by the Vapnik-Chervonenkis inequality 4.15

$$\leq \sqrt{2 \frac{\log (2 \pi_{\varphi}(n))}{n}}.$$

It is clear that $\varphi(\mathcal{U}; b)_{b \in \mathcal{V}}$ has VC-dimension 1. Hence, by the Sauer-Shelah Lemma 3.9 and the inequalities proved thereafter, $\pi_{\varphi}(n) \leq n + 1$. The theorem follows. □

References

- [1] R. P. Anstee, Lajos Rónyai, and Attila Sali, *Shattering news*, Graphs Combin. **18** (2002), no. 1, 59–73.
- [2] Luc Devroye, László Györfi, and Gábor Lugosi, *A probabilistic theory of pattern recognition*, Springer-Verlag, 1996.
- [3] Luc Devroye and Gábor Lugosi, *Combinatorial methods in density estimation*, Springer Series in Statistics, Springer-Verlag, 2001.
- [4] Timothy Gowers, *Dimension arguments in combinatorics*, Gowers’s Weblog (2008).
- [5] Gil Kalai, *Extremal Combinatorics III: Some Basic Theorems*, Combinatorics and more (2008).
- [6] Alain Pajor, *Sous-espaces l_1^n des espaces de Banach*, Travaux en Cours [Works in Progress], vol. 16, 1985.
- [7] N. Sauer, *On the density of families of sets*, J. Combinatorial Theory Ser. A **13** (1972), 145–147.
- [8] Saharon Shelah, *A combinatorial problem; stability and order for models and theories in infinitary languages*, Pacific J. Math. **41** (1972), 247–261.
- [9] V. N. Vapnik and A. Ya. Chervonenkis, *On the uniform convergence of relative frequencies of events to their probabilities*, Measures of complexity, Springer, Cham, 2015, pp. 11–30. Reprint of Theor. Probability Appl. **16** (1971), 264–280.