

# Contents

1

<b>1</b>	<b>Introduction</b>	<b>1</b>	2
1.1	Basic inequalities and theorems . . . . .	1	3
1.2	Why bother? . . . . .	2	4
1.2.1	Coin Tossings . . . . .	2	5
1.2.2	Central Limit Theorem . . . . .	2	6
1.3	Cantelli's inequality . . . . .	4	7
<b>2</b>	<b>Exponential Inequalities</b>	<b>5</b>	8
2.1	Uniform Law of Large Numbers . . . . .	8	9
<b>3</b>	<b>Application to Estimation of Data Dimension</b>	<b>13</b>	10
3.1	Chernoff-Okamoto Inequalities . . . . .	13	11
3.2	The problem . . . . .	13	12
3.3	Proofs . . . . .	16	13
<b>4</b>	<b>Applications to graph theory</b>	<b>19</b>	14
4.1	The Azuma-Hoeffding Inequality . . . . .	19	15
4.2	An heuristic algorithm for the Travelling Salesman Problem . . . . .	21	16
4.3	Three additional short examples . . . . .	23	17
<b>5</b>	<b>Applications to Vapnik–Chervonenkis theory</b>	<b>27</b>	18
5.1	Sets with Polynomial Discrimination . . . . .	27	19
5.2	Vapnik–Chervonenkis inequality . . . . .	31	20

# 1 Introduction

## 1.1 Basic inequalities and theorems

**Theorem 1.1** (Markov's inequality). For a random variable  $X$  with  $\mathbf{P}\{X < 0\} = 0$  and  $t > 0$ , we have

$$\mathbf{P}\{X \geq t\} \leq \frac{\mathbf{E} X}{t}.$$

*Proof.* In the first place, note that

$$\begin{aligned} X &= X \cdot \mathbb{1}_{\{X \geq t\}} + X \cdot \mathbb{1}_{\{X < t\}} \\ &\geq t \cdot \mathbb{1}_{\{X \geq t\}} + 0, \end{aligned}$$

and thus,

$$\mathbf{E} X \geq t \cdot \mathbf{E} \mathbb{1}_{\{X \geq t\}} = t \cdot \mathbf{P}\{X \geq t\}.$$

□

**Theorem 1.2** (Chebyshev's inequality). For  $t > 0$ , a random variable  $X$  with mean  $\mu = \mathbf{E} X$  and variance  $\sigma^2 = \mathbf{Var} X$ , we have

$$\mathbf{P}\{|X - \mu| \geq t\} \leq \frac{\sigma^2}{t^2}.$$

*Proof.* We apply Markov's inequality to the non-negative random variable  $Y = |X - \mu|^2$  in order to obtain the desired result

$$\mathbf{P}\{|X - \mu| \geq t\} = \mathbf{P}\{|X - \mu|^2 \geq t^2\} \leq \frac{\mathbf{E} [(X - \mu)^2]}{t^2} = \frac{\sigma^2}{t^2}.$$

□

## 1.2 Why bother?

The concentration inequalities are used to obtain information on how a random variable is distributed at some specific places of its domain. In the most common scenarios, these inequalities will be used to quantify how concentrated a random variable at its tails, for example,

$$\mathbf{P}\{|X - \mu| \geq t\} < f(t) \ll 1.$$

A concentration inequality is specially useful when this probability cannot be calculated at a low computational cost or estimated with high precision. The following will illustrate a case where using concentration inequalities achieves the best results.

### 1.2.1 Coin Tossings

A coin tossing game is fair if the chances of winning are equal to the chances of losing. We can verify from a sample of  $N$  games that the game is not rigged if the number of heads in the sample is not very distant from the average  $N/2$ . However, there's a chance that one may classify the coin as rigged, even when the coin is fair. By the *Law of Large Numbers*, we know that the larger the sample, the less likely it is to obtain a false positive. But let's ask ourselves how fast this probability converges to 0.

Let  $S_N \sim \text{Bi}(N, 1/2)$  denote the number of heads in a fair coin tossing game. Then,

$$\mu = \mathbf{E} S_N = \frac{N}{2}, \quad \sigma^2 = \mathbf{Var} S_N = \frac{N}{4}.$$

For a fixed  $\varepsilon > 0$ , we may classify a coin tossing game as rigged if, after  $N$  trials, the ratio of heads vs tails in the sample is greater than  $[1 + \varepsilon : 1 - \varepsilon]$ , or similarly,

$$S_N \geq \mu + \frac{\varepsilon}{2}N = \frac{1 + \varepsilon}{2}N.$$

It's clear that calculating the exact probability of the previous event for any  $N$ ,  $\varepsilon$  is a very demanding task computationally. The Chebyshev's inequality 1.2 gives us a “good-enough” result for this problem,

$$\mathbf{P}\left\{S_N \geq \mu + \frac{\varepsilon}{2}N\right\} \leq \mathbf{P}\left\{|S_N - \mu| \geq \frac{\varepsilon}{2}N\right\} \leq \sigma^2 \frac{4}{\varepsilon^2 N^2} = \frac{1}{\varepsilon^2 N}.$$

Therefore, the probability of bad events tends to 0 at least linearly with the number of games.

### 1.2.2 Central Limit Theorem

The proof of the following three theorems can be found in [Boucheron et al. \(2003\)](#)

## 1 Introduction

**Theorem 1.3.** Let  $X_i$  be a i.i.d. sample. Let  $S_N = \sum_{i=1}^N X_i$ , with mean  $\mu = \mathbf{E} S_N$  and variance  $\sigma^2 = \mathbf{Var} S_N$ . If

$$Z_N = \frac{S_N - N \cdot \mathbf{E} X_i}{\sqrt{N \cdot \mathbf{Var} X_i}} = \frac{S_N - \mu}{\sqrt{N} \sigma},$$

then,

$$Z_N \rightarrow Z \sim \mathcal{N}(0, 1), \text{ in distribution.}$$

□

**Theorem 1.4** (Tails of the Normal Distribution). Let  $Z \sim \mathcal{N}(0, 1)$ , for  $t > 0$  we have

$$\left( \frac{1}{t} - \frac{1}{t^3} \right) \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{t^2}{2} \right) \leq \mathbf{P}\{Z \geq t\} \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{t^2}{2} \right).$$

□

With that in mind, we might naively assume that better bounds can be obtained by using the previous theorem. For a large enough  $N$  we can say that for the coin tossing,

$$Z_N = \frac{S_N - N/2}{\sqrt{N/4}}$$

$$\implies \mathbf{P} \left\{ S_N \geq \frac{1+\varepsilon}{2} N \right\} = \mathbf{P} \left\{ Z_N \geq \varepsilon \sqrt{N} \right\} \sim \mathbf{P} \left\{ Z \geq \varepsilon \sqrt{N} \right\}.$$

However, this raises the question of whether we can draw the following conclusion from Theorem 1.4:

$$\mathbf{P} \left\{ S_N \geq \frac{1+\varepsilon}{2} N \right\} \leq \frac{1}{\varepsilon \sqrt{N}} \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{\varepsilon^2 \cdot N}{2} \right).$$

Unfortunately, the answer is no. The following theorem will show why.

**Theorem 1.5** (Convergence Rate for Central Limit Theorem). For  $Z_N, Z$  in Theorem 1.3, we have:

$$|\mathbf{P}\{Z_N \geq t\} - \mathbf{P}\{Z \geq t\}| = O\left(\frac{1}{\sqrt{N}}\right).$$

□

Since the approximation error of the Central Limit Theorem is of greater order than the normal bounds, the previous results cannot be taken into account.

In the context of coin tossing, this may not matter at all because the linear bound obtained using Chebyshev's inequality indicates that the probability of wrongly classifying a fair coin as a rigged coin converges at least linearly to zero. Even the Central Limit Theorem shows, in a less precise way, this convergence. However, for some specific problems in statistics, these basic tools are not precise enough to solve them. The main objective of this project is to study different ideas that improve these bounds and show examples where they can be used.

### 1.3 Cantelli's inequality

We can start with a small modification of the Chebyshev's bound for the one-sided tails

**Theorem 1.6** (Cantelli's Inequality). For  $t > 0$ , a random variable  $X$  with mean  $\mu = \mathbf{E} X$  and variance  $\sigma^2 = \mathbf{Var} X$ , we have

$$\mathbf{P}\{X - \mu \geq t\} \leq \frac{\sigma^2}{t^2 + \sigma^2}.$$

*Proof.* In the first place note that,

$$\mathbf{P}\{Y \geq s\} \leq \mathbf{P}\{Y \geq s\} + \mathbf{P}\{Y \leq s\} = \mathbf{P}\{|Y| \geq s\} = \mathbf{P}\{Y^2 \geq s^2\}. \quad (\star)$$

Let  $u \geq 0$ , define  $Y = X - \mu + u$  and  $s = t + u$  to obtain

$$\mathbf{P}\{X - \mu \geq t\} = \mathbf{P}\{X - \mu + u \geq t + u\} = \mathbf{P}\{Y \geq s\}.$$

We use  $(\star)$  and Markov's inequality (1.1) on  $Y^2$  to conclude,

$$\mathbf{P}\{Y \geq s\} \stackrel{(\star)}{\leq} \mathbf{P}\{Y^2 \geq s^2\} \stackrel{(1.1)}{\leq} \frac{\mathbf{E}[(X - \mu + u)^2]}{(t + u)^2}.$$

By linearity of expectation,

$$\mathbf{E}[(X - \mu + u)^2] = \mathbf{E}[(X - \mu)^2] + 2u \cdot \underbrace{\mathbf{E}(X - \mu)}_0 + \mathbf{E}(u^2) = \sigma^2 + u^2.$$

Finally, we choose an optimal  $u = \frac{\sigma^2}{t}$  to conclude

$$\mathbf{P}\{X - \mu \geq t\} \leq \frac{\sigma^2 + u^2}{(t + u)^2} = \frac{\sigma^2 + \sigma^4/t^2}{(t + \sigma^2/t)^2} = \frac{\sigma^2(\frac{t^2 + \sigma^2}{t^2})}{(\frac{t^2 + \sigma^2}{t})^2} = \frac{\sigma^2}{t^2 + \sigma^2}$$

□

On the other hand, the two-sided tail inequality, Cantelli's inequality is not always better than Chebyshev,

**Corollary 1.6.1** (Two-sided Cantelli inequality).

$$\mathbf{P}\{|X - \mu| \geq t\} \leq \frac{2\sigma^2}{t^2 + \sigma^2}.$$

In fact, this bound is only better than Chebyshev's  $t^2 + \sigma^2 \leq 2t^2$ , or equivalently, when  $\sigma^2 \leq t^2$ . However, in this case both formulas provide bounds greater than 1, and thus, are useless. Therefore, the conclusion is that, in general, Chebyshev's inequality is better for two-sided tails and Cantelli is for one-sided tails.

## 2 Exponential Inequalities

Even if we are satisfied with the linear convergence rate provided by Chebyshev's inequality or the improvement of one sided tails given by Cantelli's inequality, there is a simple but powerful modification we can make to Markov's inequality that will greatly improve both bounds. The following result will provide the main idea from which most of the exponential inequalities are derived.

**Theorem 2.1** (MGF inequality). Let  $X_i$  be a finite sequence of independent random variables and let  $S_N := \sum_{i=1}^N a_i X_i$ . Let  $\lambda > 0$  the following inequality holds,

$$\mathbf{P}\{S_N \geq t\} \leq e^{-\lambda t} \cdot \prod_{i=1}^N \mathbf{E} e^{\lambda a_i X_i}$$

*Proof.* Let  $\lambda > 0$ , using Markov's inequality (Theorem 1.1) we assert that since  $x \mapsto e^{\lambda x}$  is a non-decreasing function,

$$\mathbf{P}\{S_N \geq t\} = \mathbf{P}\left\{e^{\lambda S_N} \geq e^{\lambda t}\right\} \leq e^{-\lambda t} \cdot \mathbf{E} \exp\left(\lambda \sum_{i=1}^N a_i X_i\right).$$

Since  $X_i$  are independent, the MGF of  $S_N$  is the product of MGFs of each  $X_i$ :

$$\begin{aligned} \mathbf{E} \exp\left(\lambda \sum_{i=1}^N a_i X_i\right) &= \prod_{i=1}^N \mathbf{E} e^{\lambda a_i X_i} \\ \implies \mathbf{P}\{S_N \geq t\} &\leq e^{-\lambda t} \cdot \prod_{i=1}^N \mathbf{E} e^{\lambda a_i X_i}. \end{aligned}$$

□

The following two theorems are examples on how we can obtain even tighter bounds than the ones we've already studied. In particular, these theorems can be obtained from the previous theorem and are be considered to be corollaries by some authors.

**Theorem 2.2** (Chernoff's inequality). Let  $X_i \sim \text{Be}(p_i)$  be independent random variables. Define  $S_N = \sum_{i=1}^N X_i$  and let  $\mu = \mathbf{E} S_N$ . Then, for  $t > \mu$ , we have

$$\mathbf{P}\{S_N \geq t\} \leq \left(\frac{\mu}{t}\right)^t e^{-\mu+t}.$$

## 2 Exponential Inequalities

*Proof.* In the first place, use Theorem 2.1 to assert that for a  $\lambda > 0$  that

$$\mathbf{P}\{S_N \geq t\} \leq e^{-\lambda t} \cdot \prod_{i=1}^N \mathbf{E} e^{\lambda X_i}$$

Now it is left to bound every  $X_i$  individually. Using the inequality  $1 + x \leq e^x$  we obtain

$$\mathbf{E} e^{\lambda X_i} = e^{\lambda p_i} + (1 - p_i) = 1 + (e^{\lambda} - 1)p_i \leq \exp(e^{\lambda} - 1)p_i.$$

Finally, we plug this inequality on the equation to conclude that

$$e^{-\lambda t} \cdot \prod_{i=1}^N \mathbf{E} e^{\lambda X_i} \leq e^{-\lambda t} \cdot \prod_{i=1}^N \exp((e^{\lambda} - 1)p_i) = e^{-\lambda t} \exp((e^{\lambda} - 1)\mu).$$

By using the substitution  $\lambda = \ln(t/\mu)$  we obtain the desired result,

$$\mathbf{P}\{S_N \geq t\} \leq \left(\frac{\mu}{t}\right)^t \exp\left(\frac{\mu t}{\mu} - \mu\right) = \left(\frac{\mu}{t}\right)^t e^{-\mu+t}.$$

□

Another exponential inequality that is derived using a similar technique is Hoeffding's inequality:

**Theorem 2.3** (Hoeffding's inequality). Let  $X_1, \dots, X_N$  be independent random variables, such that  $X_i \in [a_i, b_i]$  for every  $i = 1, \dots, N$ . Define  $S_N = \sum_{i=1}^N X_i$  and let  $\mu = \mathbf{E} S_N$ . Then, for every  $t > 0$ , we have

$$\mathbf{P}\{S_N \geq \mu + t\} \leq \exp\left(\frac{-2t^2}{\sum (a_i - b_i)^2}\right).$$

*Proof.* Since  $x \mapsto e^x$  is a convex function, it follows that, for a random variable  $X \in [a, b]$ :

$$e^{\lambda X} \leq \frac{e^{\lambda a}(b - X)}{b - a} + \frac{e^{\lambda b}(X - a)}{b - a}, \quad a \leq b.$$

Next, take expectations on both hands of the equation to obtain:

$$\mathbf{E} e^{tX} \leq \frac{(b - \mathbf{E} X) \cdot e^{\lambda a}}{b - a} - \frac{(\mathbf{E} X - a) \cdot e^{\lambda b}}{b - a}.$$

To simplify the expression, let  $\alpha = (\mathbf{E} X - a)/(b - a)$ ,  $\beta = (b - \mathbf{E} X)/(b - a)$  and  $u = \lambda(b - a)$ . Since  $a < \mathbf{E} X < b$ , it follows that  $\alpha$  and  $\beta$  are positive. Also, note that,

$$\alpha + \beta = \frac{\mathbf{E} X - a}{b - a} + \frac{b - \mathbf{E} X}{b - a} = \frac{b - a}{b - a} = 1.$$

Now,

$$\ln \mathbf{E} e^{\lambda X} \leq \ln(\beta e^{-\alpha u} + \alpha e^{\beta u}) = -\alpha u + \ln(\beta + \alpha e^u).$$

## 2 Exponential Inequalities

168 This function is differentiable with respect to  $u$ .

$$\begin{aligned}
 L(u) &= -\alpha u + \ln(\beta + \alpha e^u) \\
 L'(u) &= -\alpha + \frac{\alpha}{\alpha + \beta e^{-u}} \\
 L''(u) &= \frac{\alpha}{\alpha + \beta e^{-u}} \cdot \frac{\beta e^{-u}}{\alpha + \beta e^{-u}}.
 \end{aligned}$$

170 Note that if  $x = \frac{\alpha}{\alpha + \beta e^{-u}} \leq 1$ , then  $L''(u) = x(1 - x) \leq \frac{1}{4}$ . Remember that  $\alpha + \beta = 1$ .  
 171 Now, by expanding the Taylor series we obtain,

$$\begin{aligned}
 L(u) &= L(0) + uL'(0) + \frac{1}{2}u^2L''(u) \\
 &= \ln(\beta + \alpha) + u \left( -\alpha + \frac{\alpha}{\alpha + \beta} \right) + \frac{1}{2}u^2L''(u) \\
 &= \frac{1}{2}u^2L''(u) \\
 &\leq \frac{1}{8}\lambda^2(b - a)^2.
 \end{aligned} \tag{*}$$

173 Finally, use the inequality from Theorem 2.1 to conclude that

$$\begin{aligned}
 \mathbf{P}\{S_N - \mu \geq t\} &\leq e^{-\lambda t} \prod_{i=1}^N \mathbf{E} e^{\lambda X_i} \\
 &\stackrel{(*)}{\leq} e^{-\lambda t} \exp \left( \frac{1}{8} t^2 \sum_{i=1}^N (b_i - a_i)^2 \right)
 \end{aligned}$$

175 □

176 **Corollary 2.3.1.** Let  $X_1, \dots, X_N$  be independent random Bernoulli variables such that  
 177  $X_i \sim \text{Be}(p_i)$ , then

$$\mathbf{P} \left\{ \sum_{i=1}^N (X_i - p_i) \geq t \right\} \leq \exp \left( \frac{-2t^2}{N} \right).$$

179 □

180 Returning to the coin tossing problem, we can now make a stronger assertion of the  
 181 rate of convergence of a false negative classification using Hoeffding inequality:

$$\mathbf{P} \left\{ S_N - \frac{N}{2} \geq \frac{\varepsilon}{2} N \right\} \leq \exp(-\varepsilon^2 N).$$



## 2.1 Uniform Law of Large Numbers

For any probability measure  $P$  on the real line and  $t \in \mathbb{R}$ , define  $P_n$  as the empirical probability measure obtain from an independent sample  $X_1, \dots, X_n$  of  $P$ , that is:

$$P_n(t) = P_n(-\infty, t) = n^{-1} \cdot \sum_{i=1}^n \mathbf{1}_{\{X_i \leq t\}}.$$

From the law of large numbers we know that for a fixed  $t$ ,  $P_n(t)$  converges to  $P(t)$  with probability 1. However we can formulate a stronger statement on this convergence. The first application of concentration inequalities we are going to explore is the uniform law of large numbers, which states the following:

**Theorem 2.4** (Glivenko-Cantelli Theorem). For  $P$ ,  $P_n$  and  $t \in \mathbb{R}$ ,

$$\|P_n - P\| = \sup_{t \in \mathbb{Q}} |P_n(t) - P(t)| \xrightarrow{P} 0.$$

*Proof.* The proof, adapted from [Pollard \(2012\)](#), consists of 5 steps. At first instance, the author clarifies that we must establish the condition of  $t \in \mathbb{Q}$  to avoid problems with measurability. The author later proves that the theorem is true for any  $t \in \mathbb{R}$ , but for practical purposes, we will only prove it for rationals. Another remark the author makes is that this result from the real line can be later generalized for some classes of polynomials, and we will cover more about this in the final section.

### First Symmetrization

In the first place, define  $P'_n$  as the empirical measure obtained from an independent but identical sample  $X'_1, \dots, X'_n$  of  $P$ . Note that for any fixed  $t$ ,  $P_n(t)$  and  $P'_n(t)$  are random variables derived from their respective samples which satisfy that

$$\mathbf{E} P_n(t) = \mathbf{E} P'_n(t) = P(t), \quad \mathbf{Var} P_n(t) = \mathbf{Var} P'_n(t) = P(t).$$

We will bound the concentration of  $\|P_n - P'_n\|$  first, which will later result in a bound for  $\|P_n - P\|$  at the end of the following lemma.

For now, fix a value for  $\varepsilon > 0$ , and keep in mind that  $Z = P_n - P$ ,  $Z' = P'_n - P$ ,  $\alpha = \frac{1}{2}\varepsilon$  and  $\beta = \frac{1}{2}$ . Also, for this case define  $\mathcal{A} = \{(-\infty, t) : t \in \mathbb{R}\}$

**Lemma 2.5.** Let  $\{Z(A)\}_{A \in \mathcal{A}}$  and  $\{Z'(A)\}_{A \in \mathcal{A}}$  be independent and identical functions defined under the same collection of sets  $\mathcal{A}$ . Also, assume that there exist  $\alpha, \beta > 0$  such that

$$\mathbf{P} \left\{ \sup_{A \in \mathcal{A}} |Z(A)| \leq \alpha \right\} \geq \beta.$$

It follows that, for any  $\varepsilon > 0$ ,

$$\mathbf{P} \left\{ \sup_{A \in \mathcal{A}} |Z(A)| > \varepsilon \right\} \leq \beta^{-1} \mathbf{P} \left\{ \sup_{A \in \mathcal{A}} |Z(A) - Z'(A)| > \varepsilon - \alpha \right\}.$$

## 2 Exponential Inequalities

*Proof.* Since  $Z, Z'$  are independent, it follows from the hypothesis that for any index  $B \in \mathcal{A}$ ,

$$\mathbf{P}\{|Z'(B)| \leq \alpha | Z\} = \mathbf{P}\{|Z'(B)| \leq \alpha\} \geq \mathbf{P}\left\{\sup_{A \in \mathcal{A}} |Z'(A)| \leq \alpha\right\} \geq \beta.$$

Now, fix  $B \in \mathcal{A}$  such that  $|Z(B)| > \varepsilon$  and use the previous inequality to conclude,

$$\begin{aligned} \beta \cdot \mathbf{P}\left\{\sup_{A \in \mathcal{A}} |Z(A)| > \varepsilon\right\} &\leq \mathbf{P}\{|Z'(B)| \leq \alpha\} \cdot \mathbf{P}\{|Z(B)| > \varepsilon\} \\ &= \mathbf{P}\{|Z'(B)| \leq \alpha, |Z(B)| > \varepsilon\} \\ &\leq \mathbf{P}\{|Z(B) - Z'(B)| > \varepsilon - \alpha\} \\ &\leq \mathbf{P}\left\{\sup_{A \in \mathcal{A}} |Z(A) - Z'(A)| > \varepsilon - \alpha\right\}. \end{aligned}$$

□

Using Chebyshev's inequality (1.2) we know that the hypothesis is satisfied for the values of  $\alpha$  and  $\beta$  we chose:

$$\forall t \in \mathbb{R} : \mathbf{P}\{|Z'(t)| \leq \alpha\} = \mathbf{P}\{|P_n(t) - P(t)| \leq \varepsilon\} \geq \frac{1}{2} = \beta, \quad \text{if } n \geq 8\varepsilon^{-2}.$$

Therefore, using the previous lemma, we conclude that

$$\mathbf{P}\{\|P_n - P\| > \varepsilon\} \leq 2\mathbf{P}\{\|P_n - P'_n\| > \frac{1}{2}\varepsilon\}, \quad \text{if } n \geq 8\varepsilon^{-2}. \quad (2.1.1)$$

### Second Symmetrization

The following trick will allow us to stop considering all of the  $2n$  from the previous symmetrization, and will help us to create a simpler random variable. We will initially prove the trick for unidimensional random variables, but in chapter 4, we will generalize this proof for any kind on set on  $\mathbb{R}^n$ .

**Lemma 2.6.** Let  $\sigma_1, \dots, \sigma_n$  be Rademacher random variables, that is  $\mathbf{P}\{\sigma_i = 1\} = \mathbf{P}\{\sigma_i = -1\} = 1/2$ . Let  $Y_i = \mathbb{1}_{\{X'_i \in A\}} - \mathbb{1}_{\{X_i \in A\}}$ , and note that,

$$\mathbf{P}\{Y_i = x\} = \mathbf{P}\{\sigma_i Y_i = x\}, \quad x \in \{-1, 0, 1\}$$

*Proof.* In the first place, since  $X_i$  and  $X'_i$  are two independent and identical copies of the same distribution, the following equality holds:

$$\begin{aligned} \mathbf{P}\{Y_i = 1\} &= \mathbf{P}\{X_i \in A\} \mathbf{P}\{X'_i \notin A\} \\ &= \mathbf{P}\{X'_i \in A\} \mathbf{P}\{X_i \notin A\} \\ &= \mathbf{P}\{Y_i = -1\}. \end{aligned}$$

## 2 Exponential Inequalities

On the other hand, since  $\sigma_i$  is also independent of  $Y_i$ , it follows that

$$\begin{aligned}
 \mathbf{P}\{\sigma_i Y_i = 1\} &= \mathbf{P}\{Y_i = 1, \sigma_i = 1\} + \mathbf{P}\{Y_i = -1, \sigma_i = -1\} \\
 &= \mathbf{P}\{Y_i = 1\}\mathbf{P}\{\sigma_i = 1\} + \mathbf{P}\{Y_i = -1\}\mathbf{P}\{\sigma_i = 1\} \\
 &= \frac{1}{2}\mathbf{P}\{Y_i = 1\} + \frac{1}{2}\mathbf{P}\{Y_i = 1\} \\
 &= \mathbf{P}\{Y_i = 1\} = \mathbf{P}\{Y_i = -1\} = \mathbf{P}\{\sigma_i Y_i = -1\}.
 \end{aligned}$$

Thus,

$$\mathbf{P}\{\sigma_i Y_i = \pm 1\} = \mathbf{P}\{Y_i = \pm 1\}, \quad \mathbf{P}\{\sigma_i Y_i = 0\} = \mathbf{P}\{Y_i = 0\}.$$

□

It follows that since  $P_n - P'_n = n^{-1} \sum_{i \leq n} Y_i$ ,

$$\begin{aligned}
 \mathbf{P}\{\|P_n - P'_n\| > \tfrac{1}{2}\varepsilon\} &= \mathbf{P}\left\{\sup_{t \in \mathbb{Q}} \left| n^{-1} \sum_{i=1}^n \sigma_i Y_i \right| > \tfrac{1}{2}\varepsilon \right\} \\
 &\leq \mathbf{P}\left\{\sup_{t \in \mathbb{Q}} \left| n^{-1} \sum_{i=1}^n \sigma_i \mathbf{1}_{\{X_i < t\}} \right| > \tfrac{1}{4}\varepsilon \right\} \\
 &\quad + \mathbf{P}\left\{\sup_{t \in \mathbb{Q}} \left| n^{-1} \sum_{i=1}^n \sigma_i \mathbf{1}_{\{X'_i < t\}} \right| > \tfrac{1}{4}\varepsilon \right\} \\
 &= 2\mathbf{P}\{\|P_n^\circ\| > \tfrac{1}{4}\varepsilon\}.
 \end{aligned} \tag{2.1.2}$$

where  $P_n^\circ = n^{-1} \sum_{i \leq n} \sigma_i \mathbf{1}_{\{X_i < t\}}$ . Then, from equations 2.1.1, 2.1.2 we conclude that for  $n \geq 8\varepsilon^{-2}$ ,

$$\mathbf{P}\{\|P_n - P\| > \varepsilon\} \leq 4\mathbf{P}\{\|P_n^\circ\| > \tfrac{1}{4}\varepsilon\}.$$

### Maximal Inequality

$$-\infty \xleftarrow{t_0} X_{(1)} \xrightarrow{t_1} X_{(2)} \xrightarrow{t_2} X_{(3)} \xrightarrow{t_3} \dots \xrightarrow{t_{n-1}} X_{(n)} \xrightarrow{t_n} \infty$$

For any given sample  $X = X_1, \dots, X_n$ , define  $X_{(j)}$  as the  $j$ -th observation when we order the observations, and fix  $t_j \in (X_{(j)}, X_{(j+1)}]$  for every  $j \leq n$  as the picture above shows. Note that if  $t \in (X_{(j)}, X_{(j+1)}]$ , then  $P_n^\circ(t) = P_n^\circ(t_j)$  because:

## 2 Exponential Inequalities

$$\begin{aligned}
P_n^\circ(t) &= n^{-1} \sum_{i=1}^n \sigma_i \mathbb{1}_{\{X_i < t\}}, & t \in (X_{(j)}, X_{(j+1)}] \\
&= n^{-1} \sum_{i=j+1}^n \sigma_i \mathbb{1}_{\{X_{(i)} < t\}} + n^{-1} \sum_{i=1}^j \sigma_i \mathbb{1}_{\{X_{(i)} < t\}} \\
&= n^{-1} \sum_{i=j+1}^n \sigma_i \cdot 1 & + & 0 \\
&= P_n^\circ(t_j).
\end{aligned}$$

It follows that for some  $k$ ,  $\|P_n^\circ\| = |P_n^\circ(t_k)|$ , and thus,

$$\begin{aligned}
\mathbf{P}\{\|P_n^\circ\| > \tfrac{1}{4}\varepsilon \mid X\} &\leq \sum_{j=0}^n \mathbf{P}\{|P_n^\circ(t_k)| > \tfrac{1}{4}\varepsilon \mid X\} \\
&\leq (n+1) \cdot \max_j \mathbf{P}\{|P_n^\circ(t_k)| > \tfrac{1}{4}\varepsilon \mid X\}.
\end{aligned} \tag{2.1.3}$$

### Exponential Bounds

Since for any given sample,  $\sigma \mathbb{1}_{X_i < t} \in [-1, 1]$ , we can use Hoeffding's Inequality 2.3 to obtain the following inequality

$$\mathbf{P}\{|P_n^\circ(A)| > \tfrac{1}{4}\varepsilon\} \leq 2 \exp\left(\frac{-2(n\varepsilon/4)^2}{4n}\right) = 2e^{-n\varepsilon^2/32}, \quad \forall A \in \mathcal{A}.$$

We use equation 2.1.3 to conclude

$$\mathbf{P}\{\|P_n^\circ\| > \tfrac{1}{4}\varepsilon \mid X\} \leq 2(n+1)e^{-n\varepsilon^2/32}. \tag{2.1.4}$$

### Integration

Finally, applying the formula  $P\{A\} = \mathbf{E}_X[\mathbf{P}\{A|X\}]$ , we conclude that

$$\begin{aligned}
\mathbf{P}\{\|P_n - P\| > \varepsilon\} &= \mathbf{E}[\mathbf{P}\{\|P_n - P\| > \varepsilon \mid X\}] \\
&\leq \mathbf{E}[8(n+1)e^{-n\varepsilon^2/32}] \\
&= 8(n+1)e^{-n\varepsilon^2/32}
\end{aligned} \tag{2.1.5}$$

The Borel-Cantelli states that if the probability of a sequence of events is summable, that is  $\sum_{n=1}^{\infty} \mathbf{P}\{E_n\} < \infty$ , then

$$\lim_n \mathbf{P}(E_n) \leq \mathbf{P}\left\{\bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} E_n\right\} = 0.$$

## 2 Exponential Inequalities

Since the inequality we obtain through the previous steps is exponential, the probabilities of the events  $E_n = \{\|P_n - P\| > \varepsilon\}$  are summable: 269  
270

$$\sum_{n=1}^{\infty} \mathbf{P}\{\|P_n - P\| > \varepsilon\} < \infty. \quad 271$$

Therefore, using the Borel-Cantelli lemma we conclude that 272

$$\mathbf{P}\{\|P_n - P\| > \varepsilon\} \rightarrow 0 \text{ with probability } 1. \quad 273$$

□ 274

In chapter 4 we will elaborate further on the details required to transform this powerful theorem in a more generalized version. 275  
276

## 3 Application to Estimation of Data Dimension

### 3.1 Chernoff-Okamoto Inequalities

Applying Markov's Inequality to  $Y = e^{uX}$ , we can assert that

$$\mathbf{P}\{X \geq \lambda + t\} \leq e^{-u(\lambda+t)} \mathbf{E} e^{uX} = e^{-u(\lambda+t)} (1 - p + pe^u)^n.$$

The right hand equation is minimized when,

$$e^u = \frac{\lambda + t}{(n - \lambda - t)} \cdot \frac{1 - p}{p}.$$

Therefore, for  $0 \leq t \leq n - \lambda$ ,

$$\mathbf{P}\{X \geq \lambda + t\} \leq \left( \frac{\lambda}{\lambda + t} \right)^{\lambda+t} \left( \frac{n - \lambda}{n - \lambda - t} \right)^{n-\lambda-t} \quad (3.1.1)$$

**Theorem 3.1.** Let  $X$  be random variable with the binomial distribution  $\text{Bi}(n, p)$  with  $\lambda := np = \mathbf{E} X$ , then for  $t \geq 0$ ,

$$\mathbf{P}\{X \geq \lambda + t\} \leq \exp \left( -\frac{t^2}{2(\lambda + t/3)} \right) \quad (3.1.2)$$

$$\mathbf{P}\{X \leq \lambda - t\} \leq \exp \left( -\frac{t^2}{2\lambda} \right) \quad (3.1.3)$$

**Used in:** Theorem 3.3

*Proof.* (TODO I've already written the proof on paper) □

### 3.2 The problem

The article [Díaz et al. \(2019\)](#) explains how we can estimate the dimension  $d$  of a manifold  $M$  embedded on a Euclidean space of dimension  $m$ , say  $\mathbb{R}^m$ . First, we are going to introduce the method they used, and then, we will show how does the exponential inequalities can be used to prove two important results in the paper. The procedure starts with an example on a uniformly distributed sample on a  $d$ -sphere  $\mathbb{S}^{d-1} \subset \mathbb{R}^d$ , but will be later generalized for samples of any distribution on any manifold.

### 3 Application to Estimation of Data Dimension

In the first place, let  $Z_1, \dots, Z_k$  be a i.i.d. sample uniformly distributed on  $\mathbb{S}^{d-1}$ . Then, we have the following formula for the variance of the angles between  $Z_i, Z_j, i \neq j$ :

$$\beta_d := \mathbf{Var}(\arccos \langle Z_i, Z_j \rangle) = \begin{cases} \frac{\pi^2}{4} - 2 \sum_{j=1}^k (2j-1)^{-2}, & \text{if } d = 2k+1 \text{ is odd,} \\ \frac{\pi^2}{12} - 2 \sum_{j=1}^k (2j)^{-2}, & \text{if } d = 2k+2 \text{ is even.} \end{cases} \quad (3.2.1)$$

The previous formula for the angle variance is proven in [Díaz et al. \(2019\)](#). In order to give more insight on how we will be choosing an estimator  $\hat{d}$  of the dimension of the sphere, consider the following theorem.

**Theorem 3.2** (Bounds for  $\beta_d$ ). For every  $d > 1$ ,

$$\frac{1}{d} \leq \beta_d \leq \frac{1}{d-1}.$$

□

Knowing that for every  $d > 1$ ,  $\beta_d$  is in the interval  $[\frac{1}{d}, \frac{1}{d-1}]$ , one can guess the dimension of the sphere by estimating  $\beta_d$ , and then, taking  $d$  from the lower bound of the interval where our estimator is. Since  $\beta_d$  is the variance of the angles in our sphere, our best choice for an estimator is the angle's sample variance,

$$U_k = \binom{k}{2}^{-1} \sum_{i < j \leq k} \left( \arccos \langle Z_i, Z_j \rangle - \frac{\pi^2}{2} \right)^2. \quad (3.2.2)$$

In Proposition 1. of [Díaz et al. \(2019\)](#) the authors prove that it's the Minimum Variance Unbiased Estimator for  $\beta_d$  on the unit sphere.

Furthermore, the authors also prove that this result can be generalized for any manifold with samples of any distribution. Let  $X_1, \dots, X_n$  be a i.i.d. sample from a random distribution  $P$  on a manifold  $M \subset \mathbb{R}^m$ , and let  $p \in M$  denote a point on the. For  $C > 0 \in \mathbb{R}$ , let  $k = \lceil C \ln(n) \rceil$  and define  $R(n) = L_{k+1}(p)$  as the distance between  $p$  and its  $(k+1)$ -nearest neighbor. W.L.O.G. assume that  $p = 0 \in M$  and that  $X_1, \dots, X_k$  are the  $k$ -nearest neighbors of  $p$ . Additionally, for the following theorem to be true, we require that at any neighborhood of  $p$ , the probability in that neighborhood is greater than 0.

The following theorem uses a special inequality from Chernoff-Okamoto, and it's crucial in the idea behind this generalization.

### 3 Application to Estimation of Data Dimension

**Theorem 3.3** (Bound  $k$ -neighbors). For any sufficiently large  $C > 0$ , we have that, there exists  $n_0$  such that, with probability 1, for every  $n \geq n_0$ ,

$$R(n) \leq f_{p,P,C}(n) = O(\sqrt[d]{\ln(n)/n}), \quad (3.2.3)$$

where the function  $f_{p,P,C}$  is a deterministic function which depends on  $p$ ,  $P$  and  $C$ .

. □

The following theorem, although it does not require concentration inequalities, is important for connecting the idea of the previous theorem to the main frame. Let  $\pi : R^m \rightarrow T_p M$  denote the orthogonal projection on the Tangent Space of  $M$  at  $p$ . Also, define  $W_i := \pi(X_i)$  and then normalize,

$$Z_i := \frac{X_i}{\|X_i\|}, \quad \widehat{W}_i := \frac{W_i}{\|W_i\|}. \quad (3.2.4)$$

**Theorem 3.4** (Projection Distance Bounds). For any  $i < j \leq n$ ,

$$(i) \quad \|X_i - \pi(X_i)\| = O(\|\pi(X_i)\|^2) \quad (3.2.5)$$

$$(ii) \quad \|Z_i - \widehat{W}_i\| = O(\|\pi(X_i)\|) \quad (3.2.6)$$

(iii) The inner products (cosine of angles) can be bounded as it follows:

$$|\langle Z_i, Z_j \rangle - \langle \widehat{W}_i, \widehat{W}_j \rangle| \leq Kr, \quad (3.2.7)$$

for a constant  $K \in \mathbb{R}$ , whenever  $r \geq \max(\|\pi(X_i)\|, \|\pi(X_j)\|)$ .

. □

What follows is that if we know  $W_1, \dots, W_k$  are behaved similar to a uniformly distributed sample on the sphere  $\mathbb{S}^d$ , then,  $Z_1, \dots, Z_k$  (the normalized  $k$ -nearest neighbors of  $p$ ) also behave like they are uniformly distributed on  $\mathbb{S}^d$ . The following theorem is made by combining the ideas of the previous theorems.

**Theorem 3.5** (Projection's Angle Variance Statistic). For  $k = O(\ln n)$ , let

$$V_{k,n} = \binom{k}{2}^{-1} \sum_{i < j \leq k} \left( \arccos \langle \widehat{W}_i, \widehat{W}_j \rangle - \frac{\pi^2}{2} \right)^2, \quad (3.2.8)$$

and let  $U_{k,n} = U_k$  from equation 3.2.2. The following statements hold

$$(i) \quad k|U_{k,n} - V_{k,n}| \xrightarrow{n \rightarrow \infty} 0, \text{ in probability.} \quad (3.2.9)$$

$$(ii) \quad \mathbf{E} |U_{k,n} - V_{k,n}| \xrightarrow{n \rightarrow \infty} 0.$$

. □



This last theorem is as far as this document is planned to cover. However, the last result in the paper provides the main statement. It says that if we estimate  $\beta_d$  as we did with  $U_{k,n}$  from 3.5, and then, extract  $\hat{d}$  from the interval where  $U_{k,n}$  is located, it follows that,

**Theorem 3.6** (Consistency). When  $n \rightarrow \infty$ ,

$$\mathbf{P}\{\hat{d} \neq d\} \rightarrow 0.$$

### 3.3 Proofs

*Proof Theorem 3.2:* The even and the odd cases must be distinguished:

(1): When  $d = 2k + 2$  is even: In the first place, remember that,

$$\lim_{k \rightarrow \infty} \sum_{j=1}^k j^{-2} = \frac{\pi^2}{6}.$$

It follows from the equation 3.2.1 that

$$\begin{aligned} \beta_d &= \frac{\pi^2}{12} - 2 \sum_{j=1}^k (2j)^{-2} = \frac{\pi^2}{12} - \frac{1}{2} \sum_{j=1}^k j^{-2} \\ &= \frac{1}{2} \sum_{j=k+1}^{\infty} j^{-2}. \end{aligned}$$

Since  $(j^{-2})_{j \in \mathbb{N}}$  is a monotonically decreasing sequence, it follows that

$$\begin{aligned} \frac{1}{d} &= \frac{1}{2k+2} = \frac{1}{2} \int_{k+1}^{\infty} x^{-2} dx \\ &\leq \beta_d \leq \frac{1}{2} \int_{k+1/2}^{\infty} x^{-2} dx \\ &= \frac{1}{2k+1} = \frac{1}{d-1}. \end{aligned}$$

(2): When  $d = 2k + 3$  is odd: On the other hand, note that

$$\begin{aligned} \lim_{k \rightarrow \infty} \sum_{j=1}^k (2j-1)^{-2} &= \lim_{k \rightarrow \infty} \sum_{j=1}^{2k-1} j^{-2} - \sum_{j=1}^{k-1} (2j)^{-2} \\ &= \lim_{k \rightarrow \infty} \sum_{j=1}^{2k-1} j^{-2} - \frac{1}{4} \sum_{j=1}^{k-1} j^{-2} \\ &= \frac{\pi^2}{6} - \frac{\pi^2}{24} = \frac{\pi^2}{8} \end{aligned}$$

### 3 Application to Estimation of Data Dimension

Hence,

$$\begin{aligned}\beta_d &= \frac{\pi^2}{4} - 2 \sum_{j=1}^k (2j-1)^{-2} \\ &= 2 \sum_{j=k+1}^{\infty} (2j-1)^{-2}.\end{aligned}$$

Using a similar argument we conclude that

$$\begin{aligned}\frac{1}{d} &= \frac{1}{2k+1} = 2 \int_{k+1}^{\infty} (2x-1)^{-2} dx \\ &\leq \beta_d \leq 2 \int_{k+1/2}^{\infty} (2x-1)^{-2} dx \\ &= \frac{1}{2k+2} = \frac{1}{d-1}.\end{aligned}$$

□

*Proof Theorem 3.3:* The volume of a  $d$ -sphere of radius  $r$  is equal to:

$$v_d r^d = \frac{\pi^{d/2}}{\Gamma(\frac{n}{2} + 1)} r^d.$$

Where  $v_d$  is the volume of the unit  $d$ -sphere. For the assumptions we made on  $P$  and  $M$  around  $p = 0$ , we can say that for any  $r > 0$ , there's a percent (greater than 0) of the sample that is within a range  $r$  from  $p$ . This proportion is subordinated only by the volume of a  $d$ -sphere of radius  $r$  and a constant  $\alpha := \alpha(P)$  that depends on the distribution  $P$ :

$$\rho = \mathbf{P}\{X \in M : |X| < r\} \geq \alpha v_d r^d > 0.$$

We can now define a binomial process based on how many neighbors does  $p$  has within a range  $r$ . Let  $N = N_r \sim \text{Bi}(n, \rho)$  be the number of neighbors, using Theorem 3.1 with  $\lambda = n\rho$  and  $t = \frac{\lambda}{2}$  we obtain,

$$\mathbf{P}\{N \leq \lambda - t\} = \mathbf{P}\{2N \leq \lambda\} \leq \exp(-\lambda/8).$$

Since  $n(\alpha v_d r^d) \leq n\rho = \lambda$ , it follows that, by choosing  $r(n)$  such that

$$r(n) = \left( \frac{C}{\alpha v_d} \cdot \frac{\ln n}{n} \right)^{1/d} = O(\sqrt[d]{\ln(n)/n}), \quad (\star)$$

and thus,

$$C \ln n = n(\alpha v_d r(n)^d) \leq \lambda,$$

we obtain:

$$P\{2N \leq C \ln n\} \leq \mathbf{P}\{2N \leq \lambda\},$$

### 3 Application to Estimation of Data Dimension

and,

$$\exp(-\lambda/8) \leq \exp\left(\frac{-C \ln n}{8}\right) = n^{-C/8}.$$

Therefore,

$$P\{2N \leq C \ln n\} \leq n^{-C/8}.$$

Finally, with this last expression we proved that if  $k = \frac{C}{2} \ln n$ , then the  $k$ -neighbors of  $p$  are contained in the ball of radius  $r(n)$  with a probability that converges exponentially to 1.  $\square$

## 4 Applications to graph theory

### 4.1 The Azuma-Hoeffding Inequality

**Definition 4.1.** A sequence  $X_0, \dots, X_n$  of random variables is consider a **martingale** if, for every  $i \leq n$ ,

$$\mathbf{E}[X_{i+1} | X_i, \dots, X_0] = X_i$$

A random graph  $G = G(n)$  is a graph that has  $n$  labeled vertices and produces an edge between 2 of them with a probability. Let  $v_1, \dots, v_n$  denote the vertices of  $G$  and  $e_1, \dots, e_m$  all of the  $\binom{n}{2}$  potential edges that  $G$  can produce. Also, define each edge's indicator function as it follows,

$$\mathbb{1}_{e_k \in G} = \begin{cases} 1, & e_k \in G \\ 0, & \text{otherwise} \end{cases}$$

An edge exposure martingale is a sequence of random variables defined as the expected value of a function  $f(G)$  which depends on the information of the first  $j$  potential edges:

$$X_j = \mathbf{E}[f(G) | \mathbb{1}_{e_1 \in G}, \dots, \mathbb{1}_{e_j \in G}]$$

Since all of the graph information is contained in its edges, the sequence transitions from no information:  $X_0 = E(f(G))$ , to the true value of the function:  $X_m = f(G)$ . Similarly, one can define a martingale which depends on how many vertices are revealed. The vertex exposure martingale is defined as it follows,

$$X_i = \mathbf{E}[f(G) | \mathbb{1}_{\{v_k, v_j\} \in G}, k < j \leq i]$$

The following inequality is to some extent an adapted version of Hoeffding inequality 2.3 for martingale random variables. If we establish a limit for which a martingale varies from one step to another, the theorem then states that we can exponentially bound the tails of its distribution:

**Theorem 4.1** (Azuma-Hoeffding inequality). Let  $X_0, \dots, X_m$  be a martingale with  $X_0 = 0$ , and

$$|X_{i+1} - X_i| \leq 1, \quad \forall i < m.$$

Then, for  $t > 0$ ,

$$\mathbf{P}\{X_m > t\sqrt{m}\} < e^{-t^2/2}.$$

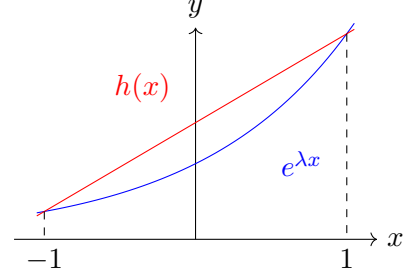
*Proof.* First, we must prove another inequality.

#### 4 Applications to graph theory

**Lemma 4.2.** Let  $Y_1, \dots, Y_m$  be random variables such that  $|Y_i| \leq 1$  and  $\mathbf{E} Y_i = 0$ , and let  $S_m = \sum_{i=1}^m Y_i$ . Then, for  $\lambda > 0$ ,

$$\mathbf{E} [e^{\lambda Y_i}] \leq e^{\lambda^2/2}. \quad 432$$

*Proof.*  $h(x) = \frac{e^\lambda + e^{-\lambda}}{2} + \frac{e^\lambda - e^{-\lambda}}{2} \cdot x,$



As the picture above shows,  $h(x)$  is the line that passes through the points  $x = -1$  and  $x = 1$  in the function  $e^{\lambda x}$ . Since  $e^{\lambda x}$  is convex ( $\lambda > 0$ ), it follows that  $h(x) \geq e^{\lambda x}$  for  $x \in [-1, 1]$ . Thus,

$$\mathbf{E} [e^{\lambda Y_i}] \leq \mathbf{E} [h(Y_i)]$$

$$\begin{aligned} (h \text{ is linear}) \quad & h(\mathbf{E} Y_i) = h(0) \\ & = \frac{e^\lambda + e^{-\lambda}}{2} = \cosh \lambda. \end{aligned} \quad 437$$

Finally,  $(2k)! \geq 2^k \cdot k!$ , for every  $k \in \mathbb{N}$ . Thus,

$$\mathbf{E} [e^{\lambda Y_i}] \leq \cosh \lambda = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} \leq \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{2^k \cdot k!} = e^{\lambda^2/2}. \quad 439$$

□ 440

Now, define  $Y_i = X_i - X_{i-1}$ . Then, by hypothesis,  $|Y_i| \leq 1$  and

$$\mathbf{E} [Y_i | X_{i-1}, \dots, X_0] = \mathbf{E} [X_i - X_{i-1} | X_{i-1}, \dots, X_0] = X_i - X_{i-1} = 0. \quad 442$$

Therefore, we can apply the previous inequality to assert,

$$\mathbf{E} [e^{\lambda Y_i} | X_{i-1}, \dots, X_0] \leq e^{\lambda^2/2}. \quad (\star) \quad 444$$

Using the formula  $E[XY] = E_X[XE[Y|X]]$  we assert that

$$\mathbf{E} e^{\lambda X_m} = \mathbf{E} \left[ \prod_{i=1}^{m-1} e^{\lambda Y_i} \cdot \mathbf{E} [e^{\lambda Y_m} | X_{m-1}, \dots, X_0] \right] \quad 446$$

447 We repeat this process  $n$  times:

$$\begin{aligned}
 & \mathbf{E} e^{\lambda X_m} = \mathbf{E} \prod_{i=1}^m e^{\lambda Y_i} \\
 & = \mathbf{E} \left[ \prod_{i=1}^{m-1} e^{\lambda Y_i} \cdot \mathbf{E} [e^{\lambda Y_m} | X_{m-1}, \dots, X_0] \right] \stackrel{(*)}{\leq} \mathbf{E} \left[ \mathbf{E} \prod_{i=1}^{m-1} e^{\lambda Y_i} \right] e^{\lambda^2/2} \\
 448 & = \mathbf{E} \left[ \prod_{i=1}^{m-2} e^{\lambda Y_i} \cdot \mathbf{E} [e^{\lambda Y_{m-1}} | X_{m-2}, \dots, X_0] \right] e^{\lambda^2/2} \stackrel{(*)}{\leq} \mathbf{E} \left[ \mathbf{E} \prod_{i=1}^{m-2} e^{\lambda Y_i} \right] e^{2\lambda^2/2} \quad (*) \\
 & = \vdots \leq \vdots \\
 & = \mathbf{E} \left[ \mathbf{E} [e^{\lambda Y_1} | X_0] \right] e^{\lambda^2/2} \leq e^{m\lambda^2/2}
 \end{aligned}$$

449 At last, by setting  $\lambda = t/\sqrt{m}$  we obtain,

$$\begin{aligned}
 & \mathbf{P}\{X_m > t\sqrt{m}\} = \mathbf{P}\{e^{\lambda X_m} > e^{\lambda t\sqrt{m}}\} \\
 & \stackrel{(\text{Markov})}{\leq} \mathbf{E} [e^{\lambda X_m}] e^{-\lambda t\sqrt{m}} \\
 450 & \stackrel{(*)}{\leq} e^{m\lambda^2/2} \cdot e^{-\lambda t\sqrt{m}} \quad (\bullet) \\
 & (\lambda = t/\sqrt{m}) = e^{t^2/2} e^{-t^2} = e^{-t^2/2}.
 \end{aligned}$$

451 □

452 **Remark.** We assumed that  $X_0 = 0$  to lighten the notation. However, we can remove  
 453 this restriction by replacing  $X_m$  with  $X_m - X_0$  in some crucial steps:

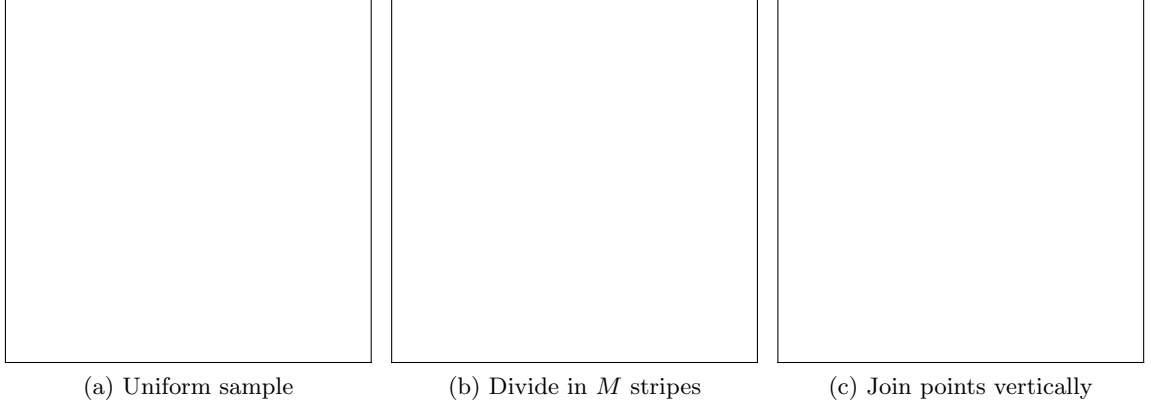
$$\begin{aligned}
 & X_m - X_0 = \sum_{i=1}^n Y_i \\
 454 & \stackrel{(*)}{\implies} \mathbf{E} e^{\lambda(X_m - X_0)} = \mathbf{E} \prod_{i=1}^m e^{\lambda Y_i} \leq e^{m\lambda^2/2} \\
 & \stackrel{(\bullet)}{\implies} \mathbf{P}\{X_m - X_0 > t\sqrt{m}\} \leq e^{-t^2/2}
 \end{aligned}$$

455 In the following section we are going to present an application of the Azuma-Hoeffding  
 456 inequality to prove the convergence to the mean of a fast (but not effective) approxima-  
 457 tion algorithm for the *Travelling Salesman Problem*.

## 458 4.2 An heuristic algorithm for the Travelling Salesman Problem

459 Let  $X_1, \dots, X_N$  be a sample of  $N$  uniformly distributed points in a compact square  
 460  $[0, L] \times [0, L]$ . The algorithm divides this square in  $M$  stripes of width  $L/M$  each. Then,

#### 4 Applications to graph theory



it connects each of the points in each of the stripes vertically and connects the top-most of one stripe with the top-most of the next one (or viceversa as the image below shows).

In the paper [Gzyl et al. \(1990\)](#) the authors found that the optimal number of stripes is  $M^* = \lfloor 0.58N^{1/2} \rfloor$ . If  $t_N$  is the TSP solution distance for our sample and  $d_N$  is the algorithm's answer with the optimal  $M^*$ , then the error is asymptotically:

$$\frac{d_N - t_N}{t_N} \approx 0.23.$$

The result that we are going to show is that  $d_n$  is very concentrated around its mean. In order to prove this, some modifications must be made to the algorithm's trajectory. Let  $e_N$  be the distance of a new trajectory that satisfies the following conditions:

- For any empty stripe in the plane we sum the length of its diagonal  $\sqrt{L^2 + L^2/M^2}$  and then it skips the empty stripe.
- When there are no empty stripes,  $e_N = d_N$

Since the probability that any given stripe is empty converges exponentially to 0,

$$\begin{aligned} (1 - 1/M)^N &= (1 - 0.58^{-1}N^{-1/2})^N \\ &= \left( (1 - 1/M)^M \right)^{0.58^{-1}N^{1/2}} \\ &\sim \exp(-0.58^{-1}N^{1/2}). \end{aligned}$$

Let  $\mathcal{A}_i := \sigma\{X_1, \dots, X_i\}$  denote the sigma algebra corresponding to revealing the first  $i$  points,  $\mathcal{A}_0 = \{\emptyset, [0, L]^2\}$ . The expected value of the trajectory  $e_N$  given that we only know the positions of the first  $i$  points in the sample is  $\mathbf{E}(e_N|\mathcal{A}_i)$ . Define

$$Z_i = \mathbf{E}(e_N|\mathcal{A}_i) - \mathbf{E}(e_N|\mathcal{A}_{i-1}),$$

As the difference of this expectations when we reveal 1 more point. Note that since

$$\mathbf{E}(Z_i|\mathcal{A}_i) = \mathbf{E}(e_N|\mathcal{A}_i, \mathcal{A}_i) - \mathbf{E}(e_N|\mathcal{A}_{i-1}, \mathcal{A}_i) = \mathbf{E}(e_N|\mathcal{A}_i) - \mathbf{E}(e_N|\mathcal{A}_i) = 0,$$

## 4 Applications to graph theory

481  $Z_1, \dots, Z_N$  is the difference sequence of a vertex exposure martingale.

482 Define  $e_N^{[i]}$  as the distance of the trajectory when we remove the  $i$ -th point from the  
 483 sample. Intuitively from the triangle inequality, we can obtain the following inequalities:

$$484 \quad e_N^{[i]} \leq e_N \leq e_N^{[i]} + 2L/M,$$

485 meaning that revealing one point cannot increase more than 2 widths the distance of  
 486 the trajectory. Thus,

$$487 \quad \|Z_i\|_\infty = \sup_{X_1, \dots, X_N} \|\mathbf{E}(e_N | \mathcal{A}_i) - \mathbf{E}(e_N | \mathcal{A}_{i-1})\| \leq 2L/M.. \quad (\star)$$

488 On the other hand,

$$489 \quad e_N - \mathbf{E} e_N = \mathbf{E}(e_N | \mathcal{A}_N) - \mathbf{E}(e_N | \mathcal{A}_0) = \sum_{i=1}^N Z_i.$$

490 Therefore, by the Azuma-Hoeffding inequality,

$$491 \quad \mathbf{P}\{|e_N - \mathbf{E} e_N| > t\} \leq 2 \exp\left(\frac{-t^2}{2} \sum_{i=1}^N \|Z_i\|_\infty^2\right).$$

492 Finally,

$$493 \quad \sum_{i=1}^N \|Z_i\|_\infty^2 \leq \frac{4NL^2}{M^2},$$

494 which implies that

$$495 \quad \mathbf{P}\{|e_N - \mathbf{E} e_N| > t\} \leq 2 \exp\left(\frac{-t^2}{2} \sum_{i=1}^N \frac{4NL^2}{M^2}\right) \sim e^{-t^2 KN},$$

496 for some  $K \in \mathbb{R}^+$ .

### 497 4.3 Three additional short examples

498 Three examples from [Alon and Spencer \(2016\)](#) will be exposed to illustrate some ideas  
 499 that can be associated with the main inequality of this chapter. Furthermore, the use-  
 500 fulness of the Azuma-Hoeffding inequality in the study of graphs and metric spaces can  
 501 be used in a more general frame. Let  $\Omega = A^B$  be the set of all functions  $g : B \rightarrow A$  for  
 502 which a probability measure is assigned

$$503 \quad \mathbf{P}\{g(b) = a\} = p(a, b), \quad \sum_{a \in A} p(a, b) = 1.$$

504 All the values  $g(b)$  are mutually independent. Now, fix a chain of sets

$$505 \quad \emptyset = B_0 \subset B_1 \subset \dots \subset B_m = B, \quad \mathcal{B} = \{B_i\}_{i=0}^m$$



## 4 Applications to graph theory

and let  $L : A^B \rightarrow \mathbb{R}$  be a functional. The martingale sequence  $X_0, \dots, X_m$  associated with  $L$  and  $\mathcal{B}$  is defined as it follows: For a fixed  $h \in A^B$ :

$$X_i(h) = \mathbf{E} [L(g) \mid g(b) = h(b), \forall b \in B_i].$$

What this means is that, given that we know the values in  $B_i$  of a function  $h$ , the martingale at the  $i$ -th step predicts the outcome of  $L(h)$  based only on this information. The following definition and theorem have the purpose to make our lives easier when talking about the ‘boundness’ of a martingale.

**Definition 4.2.** A functional  $L$  is said to satisfy the Lipschitz condition if for every  $i < m$ : Whenever two functions  $g(b) \neq g'(b)$  only on  $B_{i+1} - B_i$ ,

$$|L(g) - L(g')| \leq 1.$$

When we say that the outcome of  $L$  won’t change by more than 1 unit from one revelation to another, it means that it has the Lipschitz condition. The following theorem will connect this idea to Azuma’s inequality:

**Theorem 4.3.** The martingale associated with a functional  $L$  with the Lipschitz condition satisfies:

$$|X_{i+1}(g) - X_i(g)| \leq 1, \quad \forall g \in A^B, \forall i < m.$$

*Proof.* The proof is adapted from [Alon and Spencer \(2016\)](#) chapter 7. In the original proof, the author skips many steps that I believe are not trivial. Thus, I decided to restructure the proof in three parts:

**Rewriting**  $X_{i+1}$

Fix  $h \in A^B$ ,  $i \in \mathbb{N}$  and define  $H \subset A^B$  to be the set of functions  $h'$  in which  $h(b) = h'(b)$  for every  $b \in B_{i+1}$ . Let

$$p_{h'} = \mathbf{P}\{g = h' \mid g(b) = h(b), \forall b \in B_{i+1}\}.$$

Note that if  $h' \notin H$  and we are given that  $g(b) = h(b)$  for  $b \in B_{i+1}$ , then it would be imposible for  $g$  to be equal to  $h'$  because there would exist  $b^* \in B_{i+1}$  such that  $h'(b^*) \neq h(b^*) = g(b^*)$ . Thus,  $p_{h'} = 0$  if  $h' \notin H$  and,

$$\begin{aligned} X_{i+1}(h) &= \mathbf{E} [L(g) \mid g(b) = h(b), \forall b \in B_{i+1}]. \\ &= \sum_{h' \in A^B} L(h') \cdot \mathbf{P}\{g = h' \mid g(b) = h(b), \forall b \in B_{i+1}\} \\ &= \sum_{h' \in H} L(h') \cdot p_{h'} \end{aligned}$$

## 4 Applications to graph theory

533 **Rewriting**  $X_i$

534 From the definition,

$$535 \quad X_i(h) = \mathbf{E} [L(g) \mid g(b) = h(b), \forall b \in B_i]$$

536 Define  $H[h']$  to be the collection of functions  $h^*$  in which it is guaranteed that  $h^*$  agrees  
537 with  $h'$  everywhere but  $B_{i+1} - B_i$  (it can possibly agree there too). Let

$$538 \quad q_{h^*} = \mathbf{P}\{g(b) = h^*(b), \forall b \in B_{i+1} \mid g(b) = h(b) \forall b \in B_i\}.$$

539 For  $h' \in H$ ,  $h'(b) = h(b)$  for every  $b \in B_{i+1}$ , thus

$$540 \quad \mathbf{E} [L(g) \mid g(b) = h'(b), \forall b \in B_{i+1}] = \mathbf{E} [L(g) \mid g(b) = h(b), \forall b \in B_{i+1}] \mathbf{P}\{\}$$

541 □

542 Let  $g \in [n]^n$  be a random vector (uniformly chosen) with  $n$  entries, in which every  
543 entry is in  $[n] = \{1, \dots, n\}$ . Define  $L(g)$  to be the amount of number that are not included  
544 in the vector,

$$545 \quad L(g) = \#\{k : g_i \neq k, \forall i \in [n]\} = \sum_{k=1}^n \mathbb{1}_{k \notin g}$$

546 For example,

$$547 \quad L(\underset{g_1}{1}, \underset{g_2}{3}, \underset{g_3}{1}, \underset{g_4}{6}, \underset{g_5}{4}, \underset{g_6}{3}) = 2. \quad (2 \text{ and } 5 \text{ are missing})$$

548 We can understand the process of choosing  $g$  as independently assigning a random  
549 number in each of its coordinates. Thus, for a number  $k \in [n]$ , the probability of that  
550 number to not be in any of the entries of the vector is

$$551 \quad \mathbf{E} \mathbb{1}_{k \notin g} = \mathbf{P}\{g_i \neq k, \forall i\} = \prod_{i=1}^n P\{g_i \neq k\} = \left(1 - \frac{1}{n}\right)^n.$$

552 Hence,

$$553 \quad \mathbf{E} L(g) = \sum_{k=1}^n \mathbf{P}\{g_i \neq k, \forall i\} = n \left(1 - \frac{1}{n}\right)^n \sim \frac{n}{e}.$$

554 Now, define

$$\begin{aligned} X_0(g) &= \mathbf{E} L(g) \sim \frac{n}{e} \\ X_1(g) &= \mathbf{E} [L(g) | g_1] \\ &\vdots \\ 555 \quad X_j(g) &= \mathbf{E} [L(g) | g_1, \dots, g_j] \\ &\vdots \\ X_n(g) &= \mathbf{E} [L(g) | g_1, \dots, g_n] = L(g) \end{aligned}$$

#### 4 Applications to graph theory

$X_k(g)$  is the martingale that exposes one coordinate of  $g$  at a time. The value of  $L(g)$  556  
can vary at most by 1 for each coordinate we reveal, so  $L(g)$  has the Lipschitz condition. 557  
Then, we use theorem 4.3 and Azuma-Hoeffding inequality to conclude that 558

$$\mathbf{P}\{|L(g) - \frac{n}{e}| > t\sqrt{n}\} < 2e^{-t^2/2}. \quad 559$$

## 5 Applications to Vapnik–Chervonenkis theory

### 5.1 Sets with Polynomial Discrimination

The version of the Glivenko–Cantelli inequality we showed on chapter 2 can be generalized in multiple ways. First, we have to make some modifications in the proof of this theorem to make it work not just on intervals of the real line. The idea is to extend this property to a specific class of sets for which the final inequality will still be satisfied:

$$\mathbf{P}\{\|P_n - P\| > \varepsilon\} \leq p(n) \cdot e^{-n\varepsilon^2/32}, \text{ for a polynomial } p(n). \quad (5.1.1)$$

Remember from chapter 2 that:

- $X_i$  is a i.i.d. sample from a probability measure  $P$ .
- $P_n(A) = n^{-1} \sum \mathbb{1}_{X_i \in A}$  is the empirical measure given by  $n$  sample points.
- $\sigma_i$  is a Rademacher random variable.

In chapter 2 we assumed that  $P$  is only defined on real intervals  $(-\infty, t)$ . Then, in the section maximal inequality, we strategically defined  $(n + 1)$  different disjoint intervals when ordering the sample

$$A_0 = (-\infty, X_{(1)}], A_1 = (X_{(1)}, X_{(2)}], \dots, A_{n-1} = (X_{(n-1)}, X_{(n)}], A_n = (X_{(n)}, \infty].$$

In each one of these intervals, we fixed a representative  $t_j \in A_j$  so the function

$$P_n^\circ(B) = n^{-1} \sum_{i=1}^n \sigma_i \mathbb{1}_{X_i \in B},$$

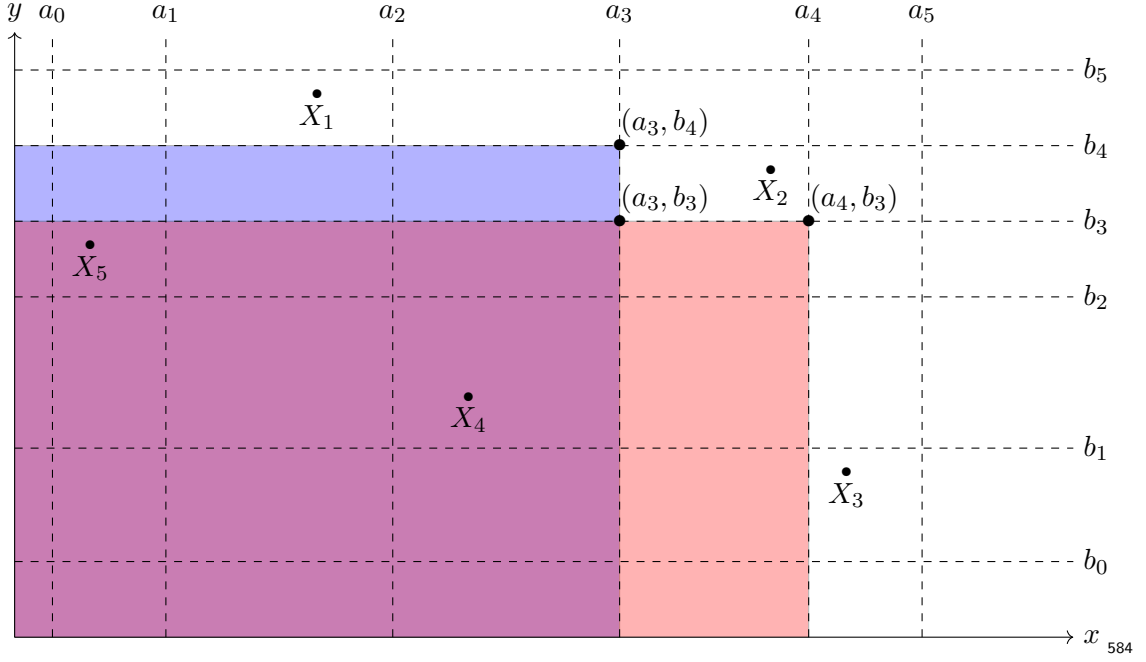
reaches its supremum in one of the sets  $B_k = (-\infty, t_k)$ :

$$\implies \exists k \leq n : \|P_n^\circ\| = |P_n^\circ(B_k)|.$$

Therefore, the  $(n + 1)$  term appears in the equation 2.1.3.

### Quadrants in $\mathbb{R}^2$

Now, imagine that instead of  $(n + 1)$  intervals we take  $(n + 1)^2$  quadrants in the form  $(-\infty, a_i) \times (-\infty, b_j) \subseteq \mathbb{R}^2$ :



Let  $A_{i,j} = (-\infty, a_i) \times (-\infty, b_j)$  be the quadrants described previously. In this example, we choose  $a_i$  and  $b_i$  in such way that the  $a_i$ 's separate the sample horizontally and vertically (similar to how we did with the  $t_j$ 's in the 1-D case). Now, let  $\mathcal{A}_n = \{A_{i,j}\}_{i,j \leq n}$ , and let  $\mathcal{A}$  be the collection of all quadrants in  $\mathbb{R}^2$ . We will see that even though  $\mathcal{A}_n \subset \mathcal{A}$  is finite, it contains all of the information of  $P_n^\circ$ .

Let  $X_j^i$  be the  $i$ -th coordinate of the point  $X_j$ , the formula for  $P_n^\circ$  at a point  $(x, y) \in \mathbb{R}^2$  is:

$$P_n^\circ(x, y) = P_n^\circ((-\infty, x) \times (-\infty, y)) = n^{-1} \sum_{k=1}^n \sigma_k \mathbb{1}_{X_k^1 < x} \cdot \mathbb{1}_{X_k^2 < y}$$

Then, because of the way we chose  $a_i$  and  $b_j$ , there exists  $i, j$  such that  $x \in (a_{i-1}, a_i)$  and  $y \in (b_{j-1}, b_j)$ . Thus,

$$\forall k \leq n : \begin{cases} \mathbb{1}_{X_k^1 < x} = \mathbb{1}_{X_k^1 < a_i} \\ \mathbb{1}_{X_k^2 < y} = \mathbb{1}_{X_k^2 < b_j} \end{cases}.$$

It follows that all the relevant information of  $\mathcal{A}$  is contained in  $\mathcal{A}_n$  since  $P_n^\circ(x, y) = P_n^\circ(a_i, b_j) = P_n(A_{i,j})$  for some  $i, j \in \mathbb{N}$ . Thus, there exist  $k_1, k_2 \in \mathbb{N}$  such that

$$\|P_n^\circ\|_{\mathcal{A}} = \max_{A \in \mathcal{A}_n} |P_n^\circ(A)| = |P_n(A_{k_1, k_2})|.$$

Hence,

$$\begin{aligned} \mathbf{P}\{\|P_n^\circ\|_{\mathcal{A}} > \tfrac{1}{4}\varepsilon \mid X\} &\leq \sum_{i,j \leq n} \mathbf{P}\{|P_n^\circ(A_{i,j})| > \tfrac{1}{4}\varepsilon \mid X\} \\ &\leq (n+1)^2 \cdot \mathbf{P}\{|P_n^\circ(A_{k_1,k_2})| > \tfrac{1}{4}\varepsilon \mid X\}. \end{aligned} \quad (5.1.2)$$

The rest of the steps in the proof of the Glivenko–Cantelli theorem (2.4) never depended on the fact that we used intervals (we will elaborate further in the next section). Therefore, the formula 5.1.1, should be changed to:

$$\mathbf{P}\{\|P_n - P\|_{\mathcal{A}} > \varepsilon\} \leq (n+1)^2 \cdot e^{-n\varepsilon^2/32} \quad (5.1.3)$$

$$\implies \mathbf{P}\{\|P_n - P\|_{\mathcal{A}} > \varepsilon\} \xrightarrow{p} 0.$$

Note that the reason why the uniform convergence worked in the previous example, was because the geometry of the collection  $\mathcal{A}$  allowed us to find a suitable sub-collection whose cardinality grows as polynomial of  $n$ . Otherwise, if we take, for instance,  $\mathcal{A} = \mathcal{R}^2$  as the collection of all the open sets in  $\mathbb{R}^2$ , then, there are at least  $2^n$  different sets in  $\mathcal{A}$  because, since  $\mathcal{R}^2$  is a metric space, we can always separate  $k$  of the sample points from the rest of the sample. Thus, the Glivenko–Cantelli inequality won’t hold anymore:

$$\mathbf{P}\{\|P_n - P\|_{\mathbb{R}^2} > \varepsilon\} \leq 2^n \cdot e^{-n\varepsilon^2/32} = e^{n(\log 2 - \varepsilon^2/32)}, \quad (5.1.4)$$

which diverges to  $\infty$  when  $\varepsilon \leq \sqrt{\log 2^{32}}$ . This will introduce us to the definition we’re looking for.

**Definition 5.1.** A collection of sets  $\mathcal{A}$  of some space  $S$  is said to have a polynomial discrimination of degree  $v$  if there exists a polynomial  $p(\cdot)$  such that:

- For any given  $n$  points  $X_1, \dots, X_n \in S$ , there exists a sub-collection  $\mathcal{A}_n$  such that for any set  $A \in \mathcal{A}$ , there exists  $B \in \mathcal{A}_n$  that satisfies  $\mathbb{1}_{X_i \in A} = \mathbb{1}_{X_i \in B}$  for every  $i \leq n$ .
- The size of  $\mathcal{A}_n$  is at most  $p(n)$ :  $\#\mathcal{A}_n \leq p(n) = O(n^v)$ .

An equivalent way to express this definition is to say that for any subspace  $S_n = \{X_1, \dots, X_n\} \subset S$ , there are at most  $p(n)$  different sets with the form  $A \cap S_n$  for  $A \in \mathcal{A}$ :

$$\max_{X_1, \dots, X_n \in S} \#\{A \cap \{X_1, \dots, X_n\} \mid A \in \mathcal{A}\} \leq p(n) \leq 2^n$$

**Remark.** For any collection  $\mathcal{A}$  and a sample  $X_1, \dots, X_n$  there exists a sub-collection  $\mathcal{A}_n$  such that

$$\#\mathcal{A}_n = \#\{A \cap \{X_1, \dots, X_n\} \leq 2^n.$$

Define the equivalence relationship  $\simeq$  as it follows,

$$A \simeq B \iff \forall i \leq n : \mathbb{1}_{X_i \in A} = \mathbb{1}_{X_i \in B},$$

which is in turn equivalent to

$$A \simeq B \iff \forall i \leq n : A \cap \{X_1, \dots, X_n\} = B \cap \{X_1, \dots, X_n\}.$$

This equivalence proves that both of the definitions are the same. Then, in order to construct  $\mathcal{A}_n$  take one representative in each of the  $\#\{A \cap \{X_1, \dots, X_n\}\}$  different equivalence classes  $[A]_{\simeq}$ ,  $A \in \mathcal{A}$ .

Another important fact from the previous remark is that, for any collection  $\mathcal{A}$ , and any given sample  $X_1, \dots, X_n$ , since for every set  $A \in \mathcal{A}$  there exists a set  $B \in \mathcal{A}_n$  such that  $\mathbb{1}_{X_i \in A} = \mathbb{1}_{X_i \in B}$ ,  $\forall i \leq n$  and  $\#\mathcal{A}_n \leq 2^n$ , it follows that  $\|P_n^\circ\|_{\mathcal{A}}$  exists and,

$$\exists A^* \in \mathcal{A}_n : \sup_{A \in \mathcal{A}} \|P_n^\circ(A)\| = \max_{B \in \mathcal{A}_n} |P_n^\circ(B)| = |P_n^\circ(A^*)|$$

Similar to the quadrants example in the equations 5.1.2 and 5.1.3, we conclude that if  $\mathcal{A}$  has a polynomial discrimination, then

$$\begin{aligned} \mathbf{P}\{\|P_n^\circ\| > \tfrac{1}{4}\varepsilon \mid X\} &\leq \sum_{A \in \mathcal{A}_n} \mathbf{P}\{|P_n^\circ(A^*)| > \tfrac{1}{4}\varepsilon \mid X\} \\ &= \#\mathcal{A}_n \cdot \mathbf{P}\{|P_n^\circ(A^*)| > \tfrac{1}{4}\varepsilon \mid X\}. \\ &\leq p(n) \cdot \mathbf{P}\{|P_n^\circ(A^*)| > \tfrac{1}{4}\varepsilon \mid X\}. \end{aligned} \tag{5.1.5}$$

$$\begin{aligned} \implies \mathbf{P}\{\|P_n - P\|_{\mathcal{A}} > \varepsilon\} &\leq p(n) \cdot e^{-n\varepsilon^2/32} \\ \implies \mathbf{P}\{\|P_n - P\|_{\mathcal{A}} > \varepsilon\} &\xrightarrow{p} 0. \end{aligned} \tag{5.1.6}$$

It's clear that  $\mathcal{R}^2$  doesn't have polynomial discrimination. Another example of a class of sets without discrimination degree is the collection of closed convex sets on  $\mathbb{S}^1 \subset \mathbb{R}^2$ . For every of the  $2^n$  subsets of any  $n$  points on the sphere, we can find a convex polygon that captures  $k$  of the points and excludes the rest. We are going to show how this works for  $n = 5$ :

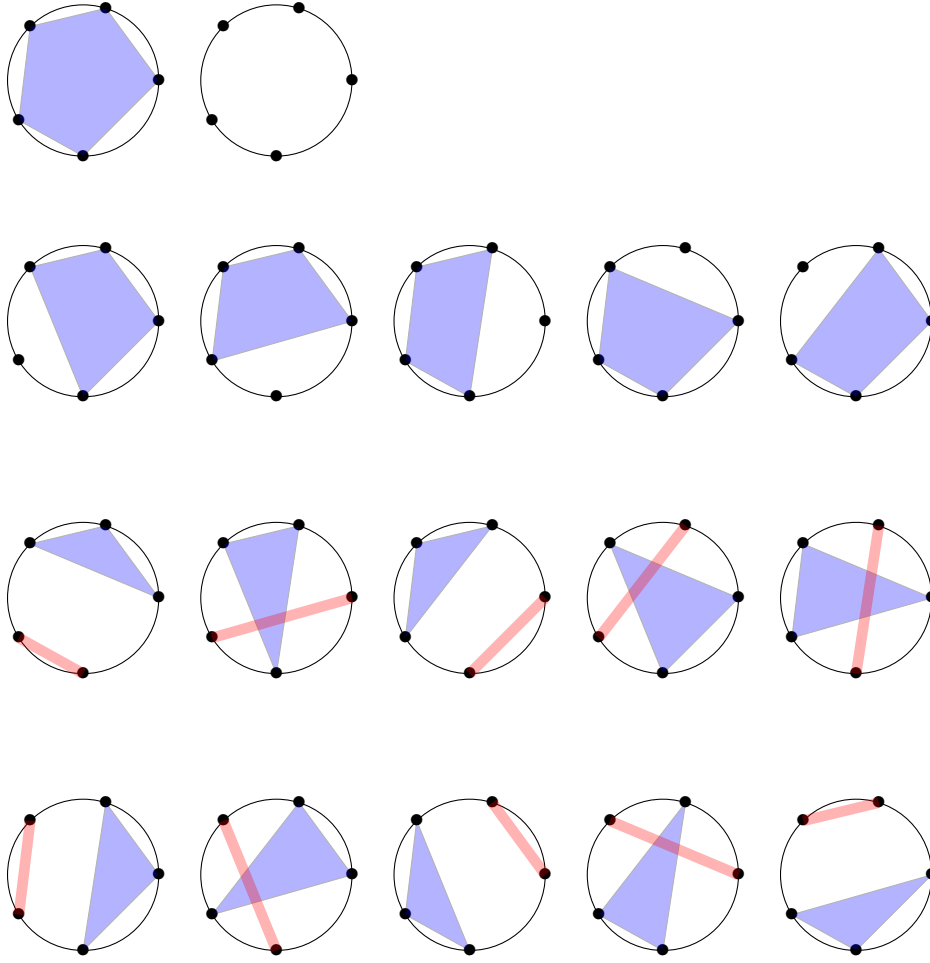


Figure 5.1: All 32 unique subsets of 5 points on  $\mathbb{S}^1$

## 5.2 Vapnik–Chervonenkis inequality

In the previous section we conclude that the uniform law of large numbers is satisfied for collections of sets with polynomial discrimination.

**Definition 5.2.** Let  $N_{\mathcal{A}}(X_1, \dots, X_n)$  be the number of different sets with the form  $\{X_1, \dots, X_n\} \cap A$  for  $A \in \mathcal{A}$

$$N_{\mathcal{A}} = \#\{\{X_1, \dots, X_n\} \cap A ; A \in \mathcal{A}\}.$$

The  $n$ -th shatter coefficient of the collection  $\mathcal{A}$  is the maximum of  $N_{\mathcal{A}}$  over all possible points in  $S$ :

$$s(\mathcal{A}, n) = \max_{X_1, \dots, X_n \in S} N_{\mathcal{A}}(X_1, \dots, X_n) \leq 2^n.$$



Finally, the Vapnik–Chervonenkis dimension is defined as the largest integer  $k$  for which  $s(\mathcal{A}, n) = 2^k$ ,

$$V_A = \operatorname{argmax}_{k \in \mathbb{N}} \{s(\mathcal{A}, k) = 2^k\} = \operatorname{argmin}_{k \in \mathbb{N}} \{s(\mathcal{A}, k) < 2^k\} - 1.$$

If  $s(\mathcal{A}, n) = 2^n$  for every  $n \in \mathbb{N}$  or equivalently if  $\mathcal{A}$  doesn't have polynomial discrimination, we say that  $V_A = \infty$ .

**Theorem 5.1** (Vapnik–Chervonenkis inequality).

$$\mathbf{P}\{\|P_n - P\|_{\mathcal{A}} > \varepsilon\} \leq s(\mathcal{A}, n) \cdot e^{-n\varepsilon^2/32}$$

*Proof.* It might be anti-climatic to tell the reader that there's no work left in this proof. But let's recapitulate everything we've done so far:

- **First Symmetrization:** Using lemma 2.5 and Chebyshev's inequality we concluded that for an identical independent copy of the empirical measure  $P'_n$  we have

$$\mathbf{P}\{\|P_n - P\|_{\mathcal{A}} > \varepsilon\} \leq 2 \mathbf{P}\{\|P_n - P'_n\|_{\mathcal{A}} > \tfrac{1}{2}\varepsilon\}, \quad \text{for } n \geq \frac{8}{\varepsilon^2}.$$

- **Second Symmetrization:** We build another distribution  $P_n^\circ(A) = n^{-1} \sum \sigma_i \mathbf{1}_{X_i \in A}$  and concluded from lemma 2.6 equation 2.1.2 that

$$\mathbf{P}\{\|P_n - P'_n\|_{\mathcal{A}} > \tfrac{1}{2}\varepsilon\} \leq 2 \mathbf{P}\{\|P_n^\circ\|_{\mathcal{A}} > \tfrac{1}{4}\varepsilon\}$$

- **Maximal Inequality:** This was the step in which we had to be most careful. In the rest of the steps it never really mattered if we worked with intervals or any other class of sets on any space. In this step the task is, for any given a sample  $X_1, \dots, X_n$ , to find a sub-collection  $\mathcal{A}_n \subset \mathcal{A}$  such that

$$\#\mathcal{A}_n = \#\{\{X_1, \dots, X_n\} \cap A; A \in \mathcal{A}\} = N_{\mathcal{A}}(X_1, \dots, X_n).$$

We proved the existence of this set in the previous theorem.

□

## 679 Bibliography

- 680 Noga Alon and Joel H Spencer. *The probabilistic method*. John Wiley & Sons, 2016.
- 681 Stéphane Boucheron, Gábor Lugosi, and Olivier Bousquet. Concentration inequalities.  
682 In *Summer school on machine learning*, pages 208–240. Springer, 2003.
- 683 Mateo Díaz, Adolfo J Quiroz, and Mauricio Velasco. Local angles and dimension esti-  
684 mation from data on manifolds. *Journal of Multivariate Analysis*, 173:229–247, 2019.
- 685 H Gzyl, R Jiménez, and AJ Quiroz. The physicist’s approach to the travelling salesman  
686 problem—ii. *Mathematical and Computer Modelling*, 13(7):45–48, 1990.
- 687 David Pollard. *Convergence of stochastic processes*. Springer Science & Business Media,  
688 2012.