

# Probability Theory

Jacob Denson

February 4, 2019

# Table Of Contents

<b>1</b>	<b>Foundations</b>	<b>2</b>
1.1	Frequentist Probability . . . . .	2
1.2	Bayesian Probability . . . . .	4
1.3	Axioms of Probability . . . . .	5
1.4	Conditional Probabilities . . . . .	11
1.5	Kolmogorov's Zero-One Law . . . . .	14
<b>2</b>	<b>Random Variables</b>	<b>15</b>
2.1	Expectation . . . . .	15
2.2	Distributions . . . . .	16
<b>3</b>	<b>Inequalities</b>	<b>17</b>
3.1	Convexity . . . . .	17
3.2	Deviation From the Mean . . . . .	18
3.3	Subgaussian Random Variables . . . . .	22
3.4	Subexponential Random Variables . . . . .	26
<b>4</b>	<b>Existence Theorems</b>	<b>28</b>
<b>5</b>	<b>Entropy</b>	<b>29</b>
<b>6</b>	<b>Appendix: Uniform Integrability</b>	<b>32</b>
<b>7</b>	<b>High Dimensional Probability</b>	<b>37</b>
7.1	Concentration of Norm . . . . .	37
7.2	Isotropic Vectors . . . . .	38
7.3	Subgaussian Vectors . . . . .	41
7.4	Concentration Without Independance . . . . .	43

7.5	The Johnson-Lindenstrauss Lemma . . . . .	46
<b>8</b>	<b>Percolation Theory</b>	<b>48</b>
8.1	Duality . . . . .	49
8.2	Boolean Functions and Sharp Thresholds . . . . .	49
8.3	Conformal Invariance . . . . .	52
<b>9</b>	<b>High Dimensional Probability</b>	<b>54</b>

# Chapter 1

## Foundations

These notes outline the basics of probability theory, the mathematical framework which allows us to interpret the statement that we are *80% more likely* to develop lung disease if you are smoker than if you are a non-smoker, or that there is a *50-50 chance* of rain on Saturday? These statements seem intuitive, and we use them naturally in everyday conversation, but a closer analysis of these statements reveals a couple difficulties with understanding such a statement. For instance, on Saturday, it will either rain, or not rain, so it is difficult to believe that there is a reasonable universal ‘chance’ of such an event happening except for absolute certainty. The mathematician has a rigorous abstraction of these statements interpreted in the language of measure theory. Nonetheless, in order to justify our intuition and to apply the theory in the sciences, we require a deeper understanding of what a probabilistic statement means. In this chapter, we will explore the two major interpretations of probability theory in real life, each of which use the same underlying mathematical theory to make judgements about the world. Regardless of which of the common interpretations you have, we will find there are certain properties your probabilistic statements possess, which can be taken as axiomatic properties of a synthetic study of probability.

### 1.1 Frequentist Probability

Classical probability theory was developed according to the intuitions of what is now known as the frequentist school of probability theory, and is

the simplest interpretation of probability to understand. It is easiest to understand from the point of view of a scientist. Suppose you are repeatedly performing some well-controlled experiment, in the sense that you do not expect the outcome of the experiment to change drastically between trials. Even under rigorously controlled conditions, the experiment will not always result in the same outcome. Slight experimental error results in slight changes in the outcome of the experiment. Nonetheless, some outcomes will occur more frequently than others.

Let us perform an experiment as often as desired, obtaining an infinite sequence of outcomes  $\omega_n$ , for  $n = 1, n = 2$ , and so on. Let  $D$  be a certain question, or *proposition* about the outcome of the experiment. For instance,  $D$  may ask whether a flipped coin lands heads up when flipping a coin repeatedly. Mathematically, we can represent the proposition as a subset of the set of all outcomes in an experiment – the outcomes for which the proposition is true. We define the *relative frequency* of  $D$  being true in  $n$  trials by the equation

$$P_n(D) := \frac{\#\{k \leq n : \omega_k \in D\}}{n}$$

The key assumption of the frequentist school of probability is that, if our experiments are suitably controlled, then regardless of the specific sequence of measured outcomes, our relative frequencies will always converge to a well defined invariant ratio, which we define to be the probability of a certain event:

$$\mathbf{P}(D) := \lim_{n \rightarrow \infty} P_n(D)$$

Let's explore some consequences of this doctrine. First,  $0 \leq P_n(D) \leq 1$  is true for any  $n$ , so for any proposition  $D$ ,  $0 \leq \mathbf{P}(D) \leq 1$ . If we let  $\Omega$  denote the set of all possible outcomes to the experiment (a proposition true for all outcomes of the experiment), then

$$P_n(\Omega) = \frac{\#\{k \leq n : \omega_k \in \Omega\}}{n} = \frac{\#\{1, 2, \dots, n\}}{n} = 1$$

Thus we conclude  $\mathbf{P}(\Omega) = 1$ . If  $A_1, A_2, \dots$  is a sequence of disjoint propositions, in the sense that no more than one outcome is true in each instance of the experiment, then

$$P_n\left(\bigcup_i A_i\right) = \frac{\#\{k \leq n : \omega_k \in \bigcup_i A_i\}}{n} = \frac{\sum_i \#\{k \leq n : \omega_k \in A_i\}}{n} = \sum_i P_n(A_i)$$

Hence  $P(\bigcup A_i) = \sum P(A_i)$ , when the events are disjoint <sup>1</sup>. There is no real generality here, because only countably many disjoint propositions can be true in the sequence of experimental outcomes, hence the probability of only countably many propositions is nonzero.

The properties we have so described turn out to be sufficient to describe all the mathematically important rules of frequentist probability. What's more, we can use these rules to *prove* that the probability of a sequence of controlled experiments eventually settles down, commonly called the strong and weak laws of probability, which justifies the thought process of the frequentist school in the first place.

## 1.2 Bayesian Probability

The frequentist school is sufficient to use probability theory to model scientific experiments, but in everyday life we make a more expansive use of probabilistic language. If you turn on the news, it's common to hear that "there is an 80% chance of downpour this evening". It is difficult to interpret this result in the frequentist definition of probability. Even if we see each night's temperament as an experimental trial, it is hard to convince yourself that these experiments are controlled enough to converge to a probabilistic result. The Bayesian school of probability redefines probability theory to be attuned to a person's individual beliefs, so that we can interpret "there is an 80% chance of downpour this evening" as an individual's belief that they think it will rain this evening rather than not rain.

You might argue that, if probability is a personal belief in an unknown event, we can choose probabilities however we want, and this would break down the logical structural required for a mathematical theory of probability. However, the probabilities that the Bayesian school studies are forced to be 'logically consistent'. Consistency can be formulated in various ways; here we discuss what is known as the Dutch book method, developed by the Italian probabilist Bruno de Finetti; if you assign to a certain unknown event  $D$  a probability  $P(D)$ , then you are willing to make a bet at  $[P(D) : 1 - P(D)]$  odds, playing the following game: If  $D$  occurs, you win  $1 - P(D)$  dollars, but if  $D$  does not occur, you have to pay up  $P(D)$  dollars. You *must* also be willing to play the game where you lose  $1 - P(D)$

---

<sup>1</sup>There are problems if we take uncountably many events  $A_i$ , rather than a countable set, but this needn't concern us now.

dollars if  $D$  occurs, and gain  $\mathbf{P}(D)$  dollars if  $D$  does not occur, so that you think the bets are ‘fair’ to both sides. For instance, you might be willing to bet a dollar against a dollar that a coin will turn up heads, which is  $[1 : 1] = [1/2 : 1/2]$  odds, so we would assign the probability that a coin will turn up heads as  $1/2$ , because then we win or lose an equal amount depending on the outcome of the bet. A person’s probability function is inconsistent if it is possible to make a series of bets that will guarantee a profit regardless of the outcome; this is known as a dutch book.

Here’s an example of how the Dutch book method can be employed to obtain general rules of probability. We claim that for any event  $D$ ,  $0 \leq \mathbf{P}(D) \leq 1$ . If a person believed that  $\mathbf{P}(D) < 0$ , then I could make a bet that person that  $D$  occurred, and I would make money regardless of the outcome. Similar results occur from betting against  $D$  if  $\mathbf{P}(D) > 1$ . It can be shown, via similar arguments, that the laws of probability hold for any logically consistent Bayesian choice of probabilities. As an aside, DeFinetti would have only allowed finitely many bets at once, which means that he would only accept  $\mathbf{P}(\bigcup A_i) = \sum \mathbf{P}(A_i)$  for finite sums, but here we allow countably many bets to be made at once. Allowing limit operations is too useful to eschew! What this means is that, regardless of whether you think that probabilities are a measure of ‘degrees of belief’ in an event happening, or the experiment frequencies of an experiment, then you still believe in the same laws of probability. Regardless of which philosophy you agree with, the fundamental principles of probability theory remain the same. We shall take the three laws we derived, and use it to make a rigorous model so that we can avoid future philosophical controversies, and this is where mathematical probability theory takes its form.

### 1.3 Axioms of Probability

Mathematically, rigorous probability theory is defined under the banner of measure theory. The framework enables us to avoid some paradoxes which can be found if we aren’t careful when analyzing experiments with infinitely many outcomes. Note, however, that the focus of probability theory is on events and random quantities (what we will soon refer to as random variables), rather than on focusing on a particular measure space under questions. Probability theorists focus on studying these *concepts*, and the framework provides the formality to understand these con-

cepts. A **probability space** is a measure space  $\Omega$  with a positive measure  $\mathbf{P}$  such that  $\mathbf{P}(\Omega) = 1$ . To the non-initiated, this means that there is a function  $\mathbf{P}$ , mapping certain subsets of  $X$  to numbers in  $[0, 1]$ , satisfying  $\mathbf{P}(\bigcup A_i) = \sum \mathbf{P}(A_i)$  for disjoint events  $A_i$ , and  $\mathbf{P}(\Omega) = 1$ .  $\Omega$  is known as the **sample space**, and  $\mathbf{P}$  is known as the **probability distribution** or **probability measure**. We interpret  $\Omega$  as the space of outcomes to some random phenomena, and  $\mathbf{P}$  measuring the likelihood of each outcome happening. Classically, probability theory was the study of certain techniques used to calculate  $\mathbf{P}(E)$  for certain outcomes  $E \subset \Omega$ , which encouraged the development of the modern fields of combinatorics and integration theory. Nowadays, probability theory tends to focus more on general principles underlying probability spaces.

**Example.** Suppose we flip a coin. There is a certain chance of flipping a heads, or flipping a tails. Since the coin is essentially symmetric, we should expect that the chance of a heads is as equally likely as a chance of tails. We can encode the set of outcomes in the sample space  $\{H, T\}$ , and then model the probability distribution as  $\mathbf{P}(H) = \mathbf{P}(T) = 1/2$ . More generally, if we have a finite sample space  $S$ , we can put a distribution on  $S$  which considers all points equally known as the **uniform distribution**, with distribution  $\mathbf{P}(s) = 1/|S|$ , for each  $s \in S$ .

**Example.** If  $\omega \in \Omega$  is fixed, the **point mass distribution**  $\delta_\omega$  at  $\omega$  is the probability distribution defined by

$$\delta_\omega(E) = \begin{cases} 1 & \omega \in E \\ 0 & \omega \notin E \end{cases}$$

The distribution represents an event where a single outcome is certain to occur, and all other situations are impossible.

We remark that there is no restriction in mathematical probability theory on *how* particular probabilistic events are obtained. All we need is that every agrees on a single value of the probability of each event, and that such values obey the laws of probability encompassed by the axioms of probability spaces. In the study of classical statistical mechanics, one can often use a discrete model consisting of placing a certain number of indistinguishable molecules in a certain number of positions. If the positions are indistinguishable, it is reasonable to assume that each molecules independently occurs in any position with equal probability, an assumption



leading to the Maxwell-Boltzmann theory of statistical mechanics. Given  $n$  points to place in  $m$  positions, combinatorics tells us the probability that  $k_1$  points occur in the first position,  $k_2$  in the second, and so on up to  $k_m$  in the  $m$ 'th position, is

$$\frac{1}{m^n} \binom{n}{k_1 \ k_2 \ \dots \ k_m} = \frac{1}{m^n} \frac{n!}{k_1! k_2! \dots k_m!}$$

Thus more evenly spread states are more likely to occur. However, in the 20th century Bose and Einstein found that in the study of certain particles, any such configuration  $(k_1, \dots, k_m)$  has an *equal* chance of occurring, leading to a completely different assignment of probabilities. To the mathematician, both theories are an equally applicable area of study.

**Example.** If  $\Omega$  is a countable set, then we can view a probability measure on  $\Omega$  as a member of the set

$$\left\{ v : \Omega \rightarrow [0, 1] : \sum_{\omega \in \Omega} v(\omega) = 1 \right\}$$

This is a convex subset of the unit ball under the  $l^\infty$  norm on  $\Omega$ , which leads to some interesting linear analysis.

The above example shows that the  $\sigma$  algebra of an at most countable probability space plays no real role in the theory. This allows us to get away with discussing most of the basic principles of probability theory without running into too many technicalities. Nonetheless, even in the study of discrete phenomena understanding probability spaces with uncountably many points becomes necessary. For instance, in the study of the limiting average of a sequence of discrete coins flips, our sample space must consist of the space of infinite sequences of coin flips  $\{0, 1\}^\omega$ , which is an infinite dimensional space. And it is often necessary in applications to select a point in an interval uniformly at random.

**Example.** Consider a five digit number  $X$  selected uniformly at random. Then the sample space consists of the numbers  $[0, 99999]$ , and the probability that a particular number is selected is one in a 100,000. There are  $10!/5! = 30240$  numbers all of whose digits are different, so the probability that a number is selected all of whose digits are different is 0.3024. If we look at the first 800 digits

of the decimal expansion of the number  $e$ , and we take each 5 digit consecutive sequence in this expansion, we end up with approximately the same frequency of unique digits, leading us to believe the digits occurring in the number  $e$  are essentially random.

**Example.** An interesting application of combinatorial probability theory is in the so called Birthday paradox. Given a number of  $n$  points to place uniformly in  $m$  boxes, the probability that no single one of them lies in the same box is  $m!/(m-n)!m^n = \prod_{k=1}^n (1 - k/m)$ . In particular, if  $m = 365$ , and  $n = 23$ , then we calculate the probability that two points lie in the same box exceeds one half. To determine this approximately, we take logarithms, using the fact that  $\log(1 - x) = -x + O(x^2)$ , so

$$\log \left( \prod_{k=1}^n (1 - k/m) \right) = \sum_{k=1}^n k/m + O(k^2/m^2) = \frac{n(n+1)}{2m} + O(n^3/m^2)$$

Thus

$$1 - \prod_{k=1}^n (1 - k/m) = 1 - \exp \left( \frac{n(n+1)}{2m} + O(n^3/m^2) \right) = \frac{n(n+1)}{2m} + O(n^4/m^2)$$

Thus we should expect it to be more likely for two points to lie in the same box then unlikely if  $n \geq \sqrt{2m}$ . In particular, if  $m = 365$ , then this estimate says we should expect that two people share the same birthday in a group of more than 27 people.

The first immediately obvious fact from the axioms is  $\mathbf{P}(E^c) = 1 - \mathbf{P}(E)$ , since  $E$  and  $E^c$  are disjoint events whose union is  $\Omega$ . A similar discussion shows that  $\mathbf{P}(E \cup F) = \mathbf{P}(E) + \mathbf{P}(F) - \mathbf{P}(E \cap F)$  because  $E \cup F$  can be written as the union of the three disjoint events  $E \cap F$ ,  $E \cap F^c$ , and  $E^c \cap F$ , and

$$\mathbf{P}(E) = \mathbf{P}(E \cap F) + \mathbf{P}(E \cap F^c) \quad \mathbf{P}(F) = \mathbf{P}(E \cap F) + \mathbf{P}(E^c \cap F)$$

This process can be generalized to unions of finitely many events. We have

$$\mathbf{P}(E \cup F \cup G) = \mathbf{P}(E) + \mathbf{P}(F) + \mathbf{P}(G) - \mathbf{P}(E \cap F) - \mathbf{P}(E \cap G) - \mathbf{P}(F \cap G) + \mathbf{P}(E \cap F \cap G)$$

which can be reasoned by looking at the number of times each element of  $E \cup F \cup G$  is ‘counted’ on the right hand side. In general, we have the inclusion-exclusion principle

$$\mathbf{P} \left( \bigcup_{k=1}^n E_k \right) = \sum_{S \subset \{1, \dots, n\}} (-1)^{|S|} \mathbf{P} \left( \bigcap_{k \in S} E_k \right)$$

This can be proven by a clumsy inductive calculation. More interestingly, but less useful, we often want to calculate the probability of an infinite union of sets  $E_k$  occurring. The inclusion-exclusion principle can be taken ‘in the limit’ to conclude that

$$\mathbf{P}\left(\bigcup_{k=1}^{\infty} E_k\right) = \lim_{n \rightarrow \infty} \mathbf{P}\left(\bigcup_{k=1}^n E_k\right) = \sum_{\substack{S \subset \mathbf{N} \\ |S| < \infty}} (-1)^{|S|} \mathbf{P}\left(\bigcap_{k \in S} E_k\right)$$

where the sum on the right is taken as the limit of the partial sums where  $S \subset \{1, \dots, n\}$  (the sum need not convergence absolutely, so it is important to take the limit in the precise ordering given).

The inclusion-exclusion formula can be tricky to calculate in real examples, so we often rely on estimates to upper bound or lower the probability of a particular event occurring. The trivial **union bound**

$$\mathbf{P}\left(\bigcup E_i\right) \leq \sum \mathbf{P}(E_i)$$

can be applied. This is a good inequality to apply if the  $E_i$  are ‘nearly disjoint’, or each have a negligible probability of occurring. On the other hand, the bound is shockingly bad if all the  $E_i$  are equal to one another.

Another useful fact to consider is that  $\mathbf{P}(E_k) \rightarrow \mathbf{P}(E)$  if the sets  $E_k$  ‘tend to’  $E$  in one form or another. If the  $E_k$  are an increasing sequence whose union is  $E$ , then we can certainly conclude  $\mathbf{P}(E_k) \rightarrow \mathbf{P}(E)$ . Similarly, if  $E_k$  is a decreasing sequence whose intersection is  $E$ , then  $\mathbf{P}(E_k) \rightarrow \mathbf{P}(E)$ . To obtain general results, we say that  $E_k \rightarrow E$  if  $\limsup E_k = \liminf E_k = E$ , where

$$\begin{aligned} \limsup_{k \rightarrow \infty} E_k &= \bigcap_{n=1}^{\infty} \bigcup_{k \geq n} E_k = \{\omega : \omega \in E_k \text{ for infinitely many } k\} \\ \liminf_{k \rightarrow \infty} E_k &= \bigcup_{n=1}^{\infty} \bigcap_{k \geq n} E_k = \{\omega : \omega \in E_k \text{ for sufficiently large } k\} \end{aligned}$$

We can then conclude that  $\mathbf{P}(E_k) \rightarrow \mathbf{P}(E)$ , since once can show

$$\limsup \mathbf{P}(E_k) \leq \mathbf{P}(\limsup E_k)$$

$$\liminf \mathbf{P}(E_k) \geq \mathbf{P}(\liminf E_k)$$

so we can apply the squeeze theorem. This already enables us to prove a very interesting theorem which can guarantee an event can ‘never occur’.

**Lemma 1.1** (Borel-Cantelli Lemma). *If  $E_1, E_2, \dots$  is a sequence of events with  $\sum \mathbf{P}(E_k) < \infty$ , then  $\mathbf{P}(\limsup E_k) = 0$ . Thus none of the events  $E_k$  can happen infinitely often.*

*Proof.* Because

$$\mathbf{P}\left(\bigcup_{k \geq n} E_k\right) \leq \sum_{k \geq n} \mathbf{P}(E_k)$$

for any  $\varepsilon > 0$  we can find an  $N$  such that for  $n \geq N$ ,  $\mathbf{P}(\bigcup_{k \geq n} E_k) < \varepsilon$ . But for any  $n$ ,  $\limsup E_k \subset \bigcup_{k \geq n} E_k$ , and so we conclude  $\mathbf{P}(\limsup E_k) < \varepsilon$ . We then let  $\varepsilon \rightarrow 0$  to conclude  $\mathbf{P}(\limsup E_k) = 0$ .  $\square$

The next example shows that the hypothesis  $\sum \mathbf{P}(E_k) < \infty$  cannot be relaxed without further analysis of the events  $E_k$  beyond their probabilities.

**Example.** *Take the Haar measure measure on  $\mathbf{T} = \mathbf{R}/\mathbf{Z}$ . Consider a sequence of positive numbers  $x_1, x_2, \dots$ , define  $S_N = \sum_{n=1}^N x_n$ , and  $E_n = [S_{n-1}, S_n]$ , considered modulo  $\mathbf{Z}$  of course. Then  $\mathbf{P}(E_n) = x_n$ , and  $\sum x_n = \infty$  happens if and only if every point in  $\mathbf{T}$  is contained in infinitely many of the  $E_n$ .*

**Theorem 1.2.** *If  $E_1, E_2, \dots$  are events with  $\inf \mathbf{P}(E_k) > 0$ , then infinitely many of the  $E_i$  occur at once with positive probability.*

*Proof.* The event that infinitely many of the  $E_1, E_2, \dots$  occur is the complement of the event that all but finitely many of the  $E_i$  do not occur, i.e.  $\liminf E_i^c$ , and it suffices to show  $\mathbf{P}(\liminf E_i^c) < 1$ . But by Fatou's lemma,

$$\mathbf{P}\left(\inf_{k \geq n} E_k^c\right) \leq \inf_{k \geq n} \mathbf{P}(E_k^c) = 1 - \sup_{k \geq n} \mathbf{P}(E_k) \leq 1 - \delta$$

and so, letting  $n \rightarrow \infty$ , we conclude  $\mathbf{P}(\liminf E_k^c) \leq 1 - \delta$ . Alternatively, if we consider the functions  $S_n = \chi_{E_1} + \dots + \chi_{E_n}$ , then

$$S_n \leq m \mathbf{I}(S_n \leq m) + n \mathbf{I}(S_n > m) = m + (n - m) \mathbf{I}(S_n > m)$$

so if  $\delta = \inf \mathbf{P}(E_i)$ , then

$$\delta n \leq \mathbf{E}(S_n) \leq m + (n - m) \mathbf{P}(S_n > m)$$

which leads to the upper bound

$$\mathbf{P}(S_n > m) \geq \frac{\delta n - m}{n - m}$$

As  $n \rightarrow \infty$ , the events on the left hand side increasing to  $\mathbf{P}(S_\infty > m)$ , where we define  $S_\infty$  as the sum of all  $\chi_{E_k}$ . Thus

$$\mathbf{P}(S_\infty > m) \geq \limsup_{n \rightarrow \infty} \frac{\delta n - m}{n - m} = \delta$$

But we can then let  $m \rightarrow \infty$  to conclude that  $\mathbf{P}(S_\infty = \infty) = \delta$ .  $\square$

## 1.4 Conditional Probabilities

In the Bayesian interpretation of probability theory, it is natural for probabilities to change over time as more information is gained about the system in question. That is, given that we know some proposition  $F$  holds over the sample space, we obtain a new probability distribution over  $\Omega$ , denoted  $\mathbf{P}(\cdot|F)$ , which represents the ratio of winnings from the bet which is only played out if  $F$  occurs. That is

- You win  $1 - \mathbf{P}(E|F)$  dollars if  $E$  occurs, and  $F$  occurs.
- You lose  $\mathbf{P}(E|F)$  dollars if  $E$  does not occur, and  $F$  occurs.
- No money exchanges hands if  $F$  does not occur.

It then follows from a dutch book argument that  $\mathbf{P}(F)\mathbf{P}(E|F) = \mathbf{P}(E \cap F)$  TODO: Fill in this argument. In the emperical interpretation,  $\mathbf{P}(E|F)$  is the ratio of times that  $E$  is true in experiments, where we only count experiments in which  $F$  also occurs. That is, we define  $\mathbf{P}(E|F)$  as the limit of the ratios

$$P_n(E|F) = \frac{\#\{k \leq n : \omega_k \in E, \omega_k \in F\}}{\#\{k \leq n : \omega_k \in F\}}$$

But it is easy to calculate, by dividing the numerator and denominator by  $n$ , that  $P_n(E|F) = P_n(E \cap F)/P_n(F)$ , so by taking limits, we find

$$\mathbf{P}(E|F) = \lim_{n \rightarrow \infty} \frac{P_n(E \cap F)}{P_n(F)} = \frac{\mathbf{P}(E \cap F)}{\mathbf{P}(F)}$$

which gives us the formula  $\mathbf{P}(F)\mathbf{P}(E|F) = \mathbf{P}(E \cap F)$ . We must of course assume that  $\mathbf{P}(F) \geq 0$ , since otherwise we are almost certain that  $F$  will never occur, and then we can almost guarantee that the limit of the values  $P_n(E|F)$  does not exist.

Thus we have motivation to define conditional probabilities by the formula  $\mathbf{P}(F)\mathbf{P}(E|F) = \mathbf{P}(E \cap F)$ , provided that  $\mathbf{P}(F) > 0$ . It enables us to model the information gained by restricting our knowledge to a particular subset of sample space. In particular, we can use the definition to identify events which contain information ‘useless’ to learning about another event. We say two events  $E$  and  $F$  are independent if  $\mathbf{P}(E \cap F) = \mathbf{P}(E)\mathbf{P}(F)$ , or, provided  $\mathbf{P}(F) > 0$ ,  $\mathbf{P}(E|F) = \mathbf{P}(E)$ ; knowledge of  $F$  gives us no foothold over knowledge of the likelihood of  $E$ .

**Example.** *The Monty Hall problem is an incredible example of how paradoxical probability theory can seem. We are on a gameshow. Suppose there are three doors in front of you. A car (brand new!) is placed uniformly randomly behind one of the doors. After we pick a door (the first door, for instance), the gameshow host then opens the second door, which you didn’t pick, revealing the car isn’t behind the door. It is important to note that he picked randomly from the remaining doors which you didn’t pick and don’t have a car behind them. What is the chance that the door you picked has the brand new car? You likely would think the two doors have a 50-50 chance of containing the car given this info, but you’d be wrong. Let  $X \in \{1, 2, 3\}$  denote the door chosen uniformly at random where the car lies, and let  $Y \in \{1, 2, 3\}$  denote the door that the host randomly chose to open. We know  $Y \neq 1$ , because the gameshow host would never open the door we picked; that would give the game away! If  $X = 1$ , then  $Y$  is picked from  $\{2, 3\}$  with uniform possibility. However, if  $X = 2$ , something interesting occurs – the gameshow is forced to open door number 3, because that’s the only door that (he thinks) won’t give any information to the player, and similarly, if  $X = 3$ , then  $Y = 2$ . Now we know that since  $X$  is chosen uniformly at random  $\mathbf{P}(X = k) = 1/3$  for each  $k$ . Similarly, we know that  $Y$  is then chosen uniformly at random from  $\{2, 3\}$ , given that  $X = 1$ , so assuming  $X$  and  $Y$  are independent, we conclude*

$$\mathbf{P}(X = 1, Y = 2) = \mathbf{P}(X = 1)\mathbf{P}(Y = 2) = 1/6$$

$$\mathbf{P}(X = 1, Y = 3) = 1/6$$

*But we also know that if  $X = 2$ , then  $Y = 3$ , so*

$$\mathbf{P}(X = 2, Y = 3) = \mathbf{P}(X = 2) = 1/3$$

$$\mathbf{P}(X = 3, Y = 2) = \mathbf{P}(X = 3) = 1/3$$

It follows that

$$\begin{aligned}
& \mathbf{P}(\text{door 1 has a car} | \text{door 2 was opened}) \\
&= \frac{\mathbf{P}(\text{door 1 has a car, door 2 was opened})}{\mathbf{P}(\text{door 2 was opened})} \\
&= \frac{\mathbf{P}(\{(X = 1, Y = 2)\})}{\mathbf{P}(\{(X = 1, Y = 2), (X = 3, Y = 2)\})} = \frac{1/6}{1/6 + 1/3} = 1/3
\end{aligned}$$

*This means we should definitely change our minds about which door we were going to pick! The argument above causes a great media uproar when it was published in 1990 in a popular magazine, because of how convincing the fallacious argument below is. The total number of possibilities is*

$$(X = 1, Y = 2), (X = 1, Y = 3), (X = 2, Y = 3), (X = 3, Y = 2)$$

*and the car seems to be in door one half of the possibilities. However, these events do not have the same probability of occurring. However, if the host changes his strategy, the conditional probabilities fall more in line with intuition – if the host always picks door number 2 to open if door number 1 was picked and had the car behind it, then the two remaining doors have an equal chance of being picked.*

We end this chapter with a final probability rule which is important in statistical analysis. If  $B$  is partitioned into a finite sequence of disjoint events  $A_1, \dots, A_n$ , then we have the formula  $\mathbf{P}(B) = \sum_i \mathbf{P}(B|A_i)\mathbf{P}(A_i)$ . This easily gives us Bayes rule

$$\mathbf{P}(A_j|B) = \frac{\mathbf{P}(B|A_j)}{\sum_i \mathbf{P}(B|A_i)\mathbf{P}(A_i)}$$

If we view  $A_j$  as a particular hypothesis from the set of all hypotheses, and  $B$  as some obtained data, then Bayes rule enables us to compute the probability that  $A_j$  is the true hypothesis from the probability that  $B$  is the data generated given the hypothesis is true. This is incredibly important if you can interpret these probabilities correctly (if you are a Bayesian), but not so useful if you are an empiricist (in which case we assume there is a ‘true’ result we are attempting to estimate from trials, so there is no probability distribution over the correctness hypothesis, other than perhaps a point mass, in which case Bayes rule gives us no information). We reiterate that Bayes rule is a theorem of probability theory, so is true in any interpretation, but can be used by Bayesians in a much more applicable way to their statistical analysis.

## 1.5 Kolmogorov's Zero-One Law

s



# Chapter 2

## Random Variables

The formality of probability theory is ironic, because even though we require the theory of measures and real analysis to place the foundations of the theory, in the probabilistic way of thinking we try to eschew as much of this foundation as possible; studying properties of random variables which aren't 'independent' of the sample space considered is avoided. As a rough approximation, if  $T : X \rightarrow Y$  is a surjective measure preserving map between probability spaces ( $X$  is an extension of the space  $Y$ , allowing more outcomes), then the random variable  $Y \circ T$  is considered the 'same' as the random variable  $Y$ , and the concepts studied in probability theory should be preserved under this extension. As we reach further and further into statistical theory, sample spaces will soon become a distant memory, brought back only for the most technical of arguments. The irony of introducing the sample space is unfortunate, because while the space is in the background, in the probabilistic way of thinking about problems we try and eschew the sample space as much as possible.

### 2.1 Expectation

**Theorem 2.1.** *For any  $X \geq 0$ ,*

$$\mathbf{E}[X] = \int \mathbf{P}(X \geq x) dx$$

*Proof.* Applying Fubini's theorem,

$$\begin{aligned}\int_0^\infty \mathbf{P}(X \geq x) dx &= \int_0^\infty \int_x^\infty d\mathbf{P}_*(y) dx \\ &= \int_0^\infty \int_0^y dx d\mathbf{P}_*(y) \\ &= \int_0^\infty y d\mathbf{P}_*(y) = \mathbf{E}[X]\end{aligned}$$

□

## 2.2 Distributions

Given a random variable  $X$  into  $\mathbf{R}$ . Then  $X$  induces a pushforward measure on the Borel  $\sigma$  algebra of  $\mathbf{R}$ , which we call the **law**, or **distribution** of  $X$ , denoted  $\mu_X$ . By definition, this means that

$$\mu_X(E) = \mathbf{P}(X \in E)$$

The distribution captures the size and shape of the random variable, but not it's relation to other random variables. As examples, we note that if  $X$  is uniformly distributed on  $[0, 1]$ , then  $\mu_X$  is the Lebesgue measure on  $[0, 1]$ , and if  $X$  is a *discrete random variable*, in the sense that the range of  $X$  takes on only countably many values, then  $\mu_X$  is a discrete measure with  $\mu_X(\{a\}) = \mathbf{P}(X = a)$ .

We say two random variables  $X$  and  $Y$  are **identically distributed**, sometimes written

$$X \stackrel{(d)}{=} Y$$

if  $\mu_X = \mu_Y$ . Note that  $X$  and  $Y$  need not even be defined on the same sample space, let alone be actually equal to one another. For instance, if  $\Omega = \{0, \dots, 6\}^2$ , with the uniform measure, and  $X$  and  $Y$  are random variables with  $X(a, b) = a$  and  $Y(a, b) = b$ , then  $X$  and  $Y$  are identically distributed, but they are not equal to one another.

# Chapter 3

## Inequalities

It is often to calculate explicitly the probability values of a certain random variable, but it often suffices to bound these values, especially when discussing convergence results, and doing other analytical calculations.

### 3.1 Convexity

The first classical inequality we discuss, Jensen's inequality, allows us to upper bound functions of an average with averages of a function. It depends in an essential way on the *convexity* of the function in question.

**Theorem 3.1.** *Given a convex  $f : \mathbf{R} \rightarrow \mathbf{R}$  and random  $X$ ,  $f(\mathbf{E}X) \leq \mathbf{E}(f(X))$ .*

*Proof.* Define

$$\beta = \sup_{s < \mathbf{E}X} \frac{f(\mathbf{E}X) - f(s)}{\mathbf{E}X - s}$$

Convexity shows that  $f(u) \geq f(\mathbf{E}X) - (u - \mathbf{E}X)\beta$ , and by definition, we find  $f(s) \geq f(\mathbf{E}X) - \beta(\mathbf{E}X - s)$  for every  $s < \mathbf{E}X$ . But this means the inequality holds for all  $s$ , and so, in particular,  $f(X) - f(\mathbf{E}X) - \beta(f(X) - \mathbf{E}X) \geq 0$ . Integrating both sides of this expression gives  $\mathbf{E}f(X) \geq f(\mathbf{E}X)$ , completing the proof.  $\square$

A simple consequence is obtained for the  $L^p$  norms  $\|X\|_p = (\mathbf{E}|X|^p)^{1/p}$ , where  $\|X\|_\infty$  is the smallest value such that  $X \leq \|X\|_\infty$  almost surely.

**Corollary 3.2.** *If  $p \leq q$ ,  $\|X\|_p \leq \|X\|_q$  if  $p \leq q$ .*

*Proof.* If we set  $f(t) = t^{q/p}$ , which is convex for  $p \leq q$ , then applying Jensen's inequality to  $|X|^p$ , we conclude that  $(\mathbf{E}|X|^p)^{q/p} = f(\mathbf{E}|X|^p) \leq \mathbf{E}|X|^q$ , and we can then take  $q'$ th roots. For  $q = \infty$ , we use a pointwise bound to conclude  $\mathbf{E}|X|^p \leq \|X\|_\infty^p$ , which makes the argument trivial.  $\square$

Another simple corollary is Hölder's inequality.

**Corollary 3.3.** *For  $1 \leq p, q \leq \infty$ , with  $1/p + 1/q = 1/r$ ,  $\|XY\|_r \leq \|X\|_p \|Y\|_q$ .*

*Proof.* By trading powers of coefficients in the equation above, it suffices to prove the theorem when  $r = 1$ . Assume first  $p, q < \infty$ . Furthermore, by scaling the inequality we can assume  $\mathbf{E}|X|^p = \mathbf{E}|Y|^q = 1$ , and it suffices to prove that  $\mathbf{E}|XY| \leq 1$ . By convexity, we find

$$|XY| \leq \frac{|X|^p}{p} + \frac{|Y|^q}{q}$$

Now taking expectations proves the claim. The remaining case occurs when  $p = \infty, q = 1$ , but this is trivial.  $\square$

**Corollary 3.4.** *If  $p \geq 1$ ,  $\|X + Y\|_p \leq \|X\|_p + \|Y\|_p$ .*

*Proof.* By replacing  $X$  and  $Y$  with  $\lambda X$  and  $\lambda Y$ , for an appropriate  $\lambda$ , we may assume  $(\mathbf{E}|X|^p)^{1/p} + (\mathbf{E}|Y|^p)^{1/p} = 1$ , and it suffices to show  $(\mathbf{E}|X + Y|^p)^{1/p} \leq 1$ . Now  $(X + Y)^p = (X + Y)(X + Y)^{p-1}$ , and so applying Hölder's inequality, since if  $1/p + 1/q = 1$  implies  $p = q(p - 1)$ , so

$$\begin{aligned} \mathbf{E}|X + Y|^p &\leq \mathbf{E}|X||X + Y|^{p-1} + \mathbf{E}|Y||X + Y|^{p-1} \\ &\leq ((\mathbf{E}|X|^p)^{1/p} + (\mathbf{E}|Y|^p)^{1/p})(\mathbf{E}|X + Y|^p)^{1/q} \\ &\leq (\mathbf{E}|X + Y|^p)^{1/q} \end{aligned}$$

Rearranging gives  $(\mathbf{E}|X + Y|^p)^{1/p} \leq 1$ , which completes the proof.  $\square$

## 3.2 Deviation From the Mean

The most important inequality bounds the chance that a probability will deviate from its mean.

**Theorem 3.5** (Markov's Inequality). *If  $X \geq 0$ , then  $\mathbf{P}(X \geq x) \leq \mathbf{E}(X)/x$ .*

The bound is trivial, and is therefore very rough. Nonetheless, it suffices for many purposes.

**Example.** Let  $X_1, X_2, \dots, X_N$  be a sequence of independent and identically distributed random variables, with mean  $\mu$  and finite variance  $\sigma^2$ . Then if we set  $S_N = N^{-1} \sum X_n$ , then  $S_N$  has mean  $\mu$ , and  $\mathbf{V}(S_N) = \sigma^2/N$ . It therefore follows by Hölder's inequality that

$$\mathbf{E}|S_N - \mu| = \mathbf{E} \left| \sum_{n=1}^N (1/N) \cdot (X_n - \mu) \right| \leq \sigma/\sqrt{N}$$

Thus Markov's inequality shows that  $S_N$  is highly likely to lie within a distance  $O(1/\sqrt{N})$  of  $\mu$ .

One can obtain better estimates by taking a more detailed step function bounded by  $X$ , but the payoff isn't normally that great. We obtain a somewhat sharper estimate if  $X$  has a finite variance  $\sigma$ .

**Theorem 3.6** (Chebyshev's Inequality). *If  $X$  has finite mean and variance, then  $\mathbf{P}(|X - \mathbf{E}X| \geq x) \leq (\sigma/x)^2$ . If  $Z = (X - \mu)/\sigma$ , then  $\mathbf{P}(|Z| \geq x) \leq x^{-2}$ .*

*Proof.* Applying Markov's inequality, we find

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

The latter inequality is a special case. □

We can continue this process. For instance, to obtain a bound  $\mathbf{P}(|X - \mathbf{E}X| > x) \lesssim x^{-p}$ , it suffices to bound  $\mathbf{E}|X - \mathbf{E}X|^p$ . We will refer to any of these inequalities as Chebyshev bounds. But we can do even better, in a more restricted situation. If  $\mathbf{E}e^{\lambda X}$  is finite, then we obtain a bound

$$\mathbf{P}(X \geq x) = \mathbf{P}(e^{\lambda X} \geq e^{\lambda x}) \leq \mathbf{E}(e^{\lambda X})e^{-\lambda x}$$

$$\mathbf{P}(|X - \mathbf{E}X| \geq x) \leq \mathbf{E}(e^{\lambda|X - \mathbf{E}X|^p})e^{-\lambda|x|^p}$$

These are known as Chernoff bounds, and show that a random variable  $X$  has very thin tails.

**Example.** Let  $X_1, \dots, X_N \sim \text{Ber}(p)$  by independent and identically distribution, where  $p$  is an unknown value. A good way to estimate  $p$  is via the random variable

$$\hat{p} = \frac{X_1 + \dots + X_N}{N}$$

It's mean is certainly is equal to  $p$ . We want to also quantify it's deviation from  $p$ . We find  $\hat{p}$  has mean  $p$  and variance  $p(1-p)/n$ , so we may apply Chebyshev's inequality to conclude

$$\mathbf{P}(|\hat{p} - p| \geq t) \leq \frac{p(1-p)}{Nt^2} \leq \frac{1}{4Nt^2}$$

Thus we obtain quadratic decay in our error irrespective of the value of  $p$ . Similar to the calculations in the last example, we see that  $\hat{p}$  is within a distance  $1/\sqrt{N}$  with high probability.

Hoeffding's inequality is similar to Markov's inequality, but is generally much sharper. It therefore has a more complicated formula.

**Theorem 3.7** (Hoeffding's Inequality). Let  $X_1, \dots, X_n$  be centrally distributed independent random variables, with  $a_i \leq X_i \leq b_i$ . Then

$$\mathbf{P}\left(\sum X_i \geq t\right) \leq e^{-2t^2 / \sum (b_i - a_i)^2}$$

*Proof.* For any  $\lambda > 0$ , a Chernoff bound gives that

$$\mathbf{P}\left(\sum X_i \geq t\right) \leq e^{-\lambda t} \prod \mathbf{E}(e^{\lambda X_i})$$

We can write  $X_i = \Lambda a_i + (1 - \Lambda)b_i$  for some random value  $0 \leq \Lambda \leq 1$ . Applying convexity,

$$e^{\lambda X_i} \leq \Lambda e^{\lambda a_i} + (1 - \Lambda)e^{\lambda b_i}$$

Hence

$$\mathbf{E}(e^{\lambda X_i}) \leq \mathbf{E}(\Lambda)e^{\lambda a_i} + (1 - \mathbf{E}(\Lambda))e^{\lambda b_i}$$

Now we may explicitly calculate  $\Lambda = (X_i - a_i)/(b_i - a_i)$ , so that

$$\mathbf{E}(e^{\lambda X_i}) \leq \frac{a_i}{a_i - b_i} e^{\lambda a_i} + \frac{b_i}{b_i - a_i} e^{\lambda b_i} = e^{F(\lambda(b_i - a_i))}$$

Where  $F(x) = -sx + \log(1 - s + se^x)$ , where  $s = a_i/(a_i - b_i)$ . Note that  $F(0) = F'(0) = 0$ , and  $F''(x) \leq 1/4$  for  $x > 0$ . Thus, by Taylor's theorem, there is  $y \in (0, x)$  such that  $F(x) = g''(y)x^2/2 \leq x^2/8$ . Thus  $\mathbf{E}(e^{\lambda X_i}) \leq e^{\lambda^2(b_i - a_i)^2/8}$ . We therefore conclude that

$$\mathbf{P}\left(\sum X_i \geq t\right) \leq e^{-\lambda t} \prod_{i=1}^n e^{\lambda^2(b_i - a_i)^2/8}$$

In particular, if

$$\lambda = \frac{t}{\sum (b_i - a_i)^2/4}$$

we obtain the required inequality.  $\square$

**Example.** If  $\hat{p}$  is as before, and we consider the random variables  $\Delta_n = X_n - p$ , then  $\hat{p} - p = 1/N \sum \Delta_n$ . We know  $-p \leq \Delta_n \leq 1 - p$  so Hoeffding's inequality implies that

$$\mathbf{P}(|\hat{p} - p| \geq t) = \mathbf{P}\left(\left|\sum_{n=1}^N \Delta_n\right| \geq Nt\right) \leq 2e^{-2Nt^2}$$

This is sharper than the previous bound we got for all values of  $t$ . In particular, with probability greater than  $1 - \varepsilon$ , we can guarantee  $\hat{p}$  is within  $O(\log(1/\varepsilon)^{1/2} N^{-1/2})$  of  $p$ .

**Example.** Suppose we are performing a test of some property, which succeeds in obtaining the correct answer with probability  $1/2 + \delta$ , where  $\delta$  is small. Then we obtain a stronger test by performing the test independantly  $N$  times, and taking the majority vote. If  $X_1, \dots, X_N$  denotes the  $\{-1, 1\}$  valued outcome of the particular tests, our new test is just  $\text{sgn}(S_N)$ , where  $S_N = X_1 + \dots + X_N$  (We also assume here that  $N$  is odd, so that  $S_N$  never equals zero). Without loss of generality, if the property we are testing is true, then  $\mathbf{P}(X_n = 1) = 1/2 + \delta$ , and  $\mathbf{P}(X_n = -1) = 1/2 - \delta$ . Thus the random variable has expectation  $2\delta$ . If  $\Delta_n = X_n - 2\delta$ , then  $\Delta_n$  is centrally distributed, and  $S_N = \sum \Delta_n + 2N\delta$ . By Hoeffding's inequality, we conclude that

$$\mathbf{P}(\text{Test fails}) = \mathbf{P}(S_N \leq 0) = \mathbf{P}\left(\sum \Delta_n \leq -2N\delta\right) \leq e^{-8N\delta^2}$$

In particular, if we want a bound like  $\mathbf{P}(\text{Test fails}) \leq \varepsilon$ , then we need only choose  $N = \Omega(\delta^{-2} \log(1/\varepsilon))$ . This is why when performing a statistical test, a probability of success higher than  $1/2$  is essentially comparable to a 99% probability of success.

### 3.3 Subgaussian Random Variables

Hoeffding's inequality only applies to bounded random variables. In the general case, we can't apply the inequality (which relies on the bounded intervals to use convexity), and Chebyshev's inequality often does not suffice. We should still obtain fast tail decay in most circumstances, say, for instance a Gaussian distribution with variance  $\sigma^2$ . Calculating, we find

$$\mathbf{P}(X - \mu \geq y) = \int_y^\infty \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} dx \leq \frac{1}{y\sqrt{2\pi\sigma^2}} \int_y^\infty x e^{-\frac{x^2}{2\sigma^2}} dx = \frac{\sigma e^{-\frac{y^2}{2\sigma^2}}}{y\sqrt{2\pi}}$$

This quantity is almost always better than Chebyshev's inequality, since the ratio  $1/y$ , which measures the inaccuracy of our inequality, is nullified by the exponential function. We can find similar equalities for random variables which are 'dominated' by normal distributions.

**Theorem 3.8.** *The following are equivalent, for comparable constants  $K_i > 0$ .*

- (1)  $\mathbf{P}(|X| \geq t) \leq 2 \exp(-t^2/K_1^2)$ .
- (2)  $\|X\|_p \leq K_2 \sqrt{p}$ , for  $0 < p < \infty$ .
- (3)  $M_{X^2}(\lambda^2) \leq \exp(K_3^2 \lambda^2)$  if  $|\lambda| \leq 1/K_3$ .
- (4)  $M_{X^2}(1/K_4^2) \leq 2$ .
- (5) If, in addition,  $\mathbf{E}X = 0$ ,  $M_X(\lambda) \leq \exp(K_5^2 \lambda^2)$ .

*If any of these equations hold, we say  $X$  is subgaussian.*

*Proof.* We provide a proof that the inequalities implies one another, up to a constant change of coefficients.

- Suppose (1) holds. Then

$$\begin{aligned} \|X\|_p^p &= \int_0^\infty \mathbf{P}(|X| \geq \lambda^{1/p}) d\lambda \leq 2 \int_0^\infty \exp(-(\lambda^{1/p}/K_1)^2) d\lambda \\ &= 2K_1^p p \int_0^\infty \lambda^p \exp(-\lambda^2) \frac{d\lambda}{\lambda} = K_1^p p \cdot \Gamma(p/2) \leq K_1^p p (p/2)^{p/2} \end{aligned}$$

Thus  $\|X\|_p \leq (p^{1/p} 2^{-1/2}) \sqrt{p} K_1 \leq (e^{1/e} 2^{-1/2}) \sqrt{p} K_1$ , so we can set  $K_2 \leq 2^{-1/2} e^{1/e} K_1$ .



- Suppose (2) holds. Using Stirling's approximation, we compute

$$\begin{aligned} M_{X^2}(\lambda^2) &= \sum_{k=0}^{\infty} \frac{\mathbf{E}(X^{2k})\lambda^{2k}}{k!} = \sum_{k=0}^{\infty} \frac{\|X\|_{2k}^{2k}\lambda^{2k}}{k!} \leq \sum_{k=0}^{\infty} k^k \frac{(2K_2^2\lambda^2)^k}{k!} \\ &\leq \sum_{k=0}^{\infty} k^k \frac{(2K_2^2\lambda^2)^k}{(2\pi)^{1/2}k^{k+1/2}e^{-k}} \leq \sum_{k=0}^{\infty} (2eK_2^2\lambda^2)^k \end{aligned}$$

For  $|\lambda| \leq 1/2e^{1/2}K_2$ , since  $(1-x)^{-1} \leq \exp(2x)$  if  $0 \leq x \leq 1/2$ , we conclude

$$M_{X^2}(\lambda^2) = \frac{1}{1 - 2eK_2^2\lambda^2} \leq \exp(4eK_2^2\lambda^2)$$

Thus we can set  $K_3 \leq 2e^{1/2}K_2$ .

- The fact that (3) implies (4) is very simple, since  $M_{X^2}(\log(2)/K_3^2) \leq \exp(\log(2)) = 2$ . Thus we can set  $K_4 \leq K_3/\log(2)^{1/2}$ .
- The fact that (4) implies (1) is also pretty simple, since a Chernoff type bound gives

$$\mathbf{P}(|X| \geq t) = \mathbf{P}(e^{|X|^2/K_4^2} \geq e^{t^2/K_4^2}) \leq M_{X^2}(1/K_4^2)e^{-t^2/K_4^2} \leq 2e^{-t^2/K_4^2}$$

So we can set  $K_1 \leq K_4$ . This completes the argument that (1) through (4) are all equivalent.

- Now we show (3) implies (5), assuming  $X$  is centrally distributed. If  $|\lambda| \leq 1/K_3$ , we use the inequality  $e^{\lambda x} \leq \lambda x + e^{\lambda^2 x^2}$  to conclude

$$M_X(\lambda) \leq \lambda \mathbf{E}X + M_{X^2}(\lambda^2) = M_{X^2}(\lambda^2) \leq \exp(K_3^2\lambda^2)$$

If  $|\lambda| \geq 1/K_3$ , we use the inequality  $\lambda x \leq \lambda^2/2 + x^2/2$ , so that

$$M_X(\lambda) \leq e^{K_3^2\lambda^2/2} M_{X^2}(1/2K_3^2) \leq e^{K_3^2\lambda^2/2} e^{1/2} \leq e^{K_3^2\lambda^2}$$

Thus we may set  $K_5 \leq K_3$ .

- Suppose (5) holds. Then a Chernoff bound gives

$$\mathbf{P}(X \geq t) \leq M_X(\lambda)e^{-\lambda t} \leq e^{K_5^2\lambda^2 - \lambda t}$$

If we set  $\lambda = t/2K_5^2$ , we find  $\mathbf{P}(X \geq t) \leq e^{-t^2/4K_5^2}$ . Thus  $\mathbf{P}(|X| \geq t) \leq 2e^{-t^2/4K_5^2}$ , and so we can set  $K_1 \leq 2K_5$ .

In particular, if we check all the constants, we see that  $K_i \leq 10K_j$  for any minimal choice of  $K_i$  and  $K_j$ .  $\square$

The natural way to form a measure of being subgaussian is to come up with an Orlicz norm. Given a convex, increasing function  $\psi : [0, \infty) \rightarrow [0, \infty)$ , with  $\psi(0) = 0$ , and  $\psi(x) \rightarrow \infty$  as  $x \rightarrow \infty$ . We define the Orlicz norm  $\|X\|_\psi$  corresponding to  $\psi$  as the infimum of  $t$  such that  $\mathbf{E}(\psi(t^{-1}|X|)) \leq 1$ .

**Theorem 3.9.** *The Orlicz norm is actually a norm, and is also complete.*

*Proof.* It is obvious that  $\|\lambda X\|_\psi = |\lambda| \|X\|_\psi$ . To obtain the triangle inequality, we note that if  $\mathbf{E}(\psi(t^{-1}|X|)) \leq 1$ , and  $\mathbf{E}(\psi(s^{-1}|Y|)) \leq 1$ , then by applying convexity, we find

$$\mathbf{E} \left( \psi \left( \frac{|X+Y|}{s+t} \right) \right) \leq \mathbf{E} \left( \frac{|X|+|Y|}{s+t} \right) \leq \frac{t\mathbf{E}(\psi(t^{-1}|X|)) + s\mathbf{E}(\psi(s^{-1}|Y|))}{s+t} \leq 1$$

This gives  $\|X+Y\|_\psi \leq 1$ . If  $\|X\|_\psi = 0$ , then the increasing nature of  $\psi$  makes it easy to see that  $X = 0$  almost surely. If  $X_1, X_2, \dots$  is a Cauchy sequence with respect to the Orlicz norm, we may thin the sequence out so that  $\|X_{n+1} - X_n\|_{\psi_2} \leq 1/2^n$ . This means that  $S = \sum |X_{n+1} - X_n|$  has finite Orlicz norm, because if  $S_N = \sum_{n \leq N} |X_{n+1} - X_n|$ , then  $\|S_N\|_\psi \leq 1$ , and by the monotone convergence theorem  $\mathbf{E}(\psi(S)) = \lim \mathbf{E}(\psi(S_N)) \leq 1$ . Thus  $\|S\|_\psi \leq 1$ . This means that  $S$  is finite almost everywhere, so  $\sum X_{n+1} - X_n$  converges absolutely almost everywhere. Thus  $X_n$  converges almost everywhere to some  $X$ , and monotone convergence implies  $X_n$  converges to  $X$  in the Orlicz norm. But because the original sequence is Cauchy, this means the original sequence also converges to  $X$ .  $\square$

We can use the Orlicz norms to provide a norm measuring how subgaussian a random variable is. It is obtained by considering the convex, increasing function  $\psi_2(x) = e^{x^2}/2$ . Naturally, we let  $\|X\|_{\psi_2}$  denote this norm. Because  $\|X\|_{\psi_2} < \infty$  if and only if  $X$  is subgaussian, we conclude that the family of random variables is a vector space, and the ‘subgaussianness’ of a random variable is a normable quantity.

If  $X$  is a centered subgaussian, we can use an alternate norm. We let  $\sigma(X)$  be the smallest  $\sigma$  such that  $M_X(\lambda) \leq \exp(\sigma^2 \lambda^2 / 2)$ . A nice feature of this value is that  $\mathbf{V}(X) \leq \sigma(X)^2$ , because as  $\lambda \rightarrow 0$ ,  $M_X(\lambda) \sim 1 + \mathbf{E}(X^2) \lambda^2 / 2 = 1 + \mathbf{V}(X) \lambda^2 / 2$ , and  $\exp(\sigma^2 \lambda^2 / 2) \sim 1 + \sigma^2 \lambda^2 / 2$ . It is obvious

that this is a homogenous value. To obtain that  $\sigma(X + Y) \leq \sigma(X) + \sigma(Y)$ , we apply Hölder's inequality to conclude that if  $p^{-1} + q^{-1} = 1$ ,

$$\begin{aligned} M_{X+Y}(\lambda) &= \mathbf{E}(e^{\lambda X} e^{\lambda Y}) \leq \mathbf{E}(e^{p\lambda X})^{p^{-1}} \mathbf{E}(e^{q\lambda Y})^{q^{-1}} \\ &\leq \exp(p^{-1}(p\lambda\sigma(X))^2 + q^{-1}(q\lambda\sigma(Y))^2/2) \\ &\leq \exp((p\sigma(X)^2 + q\sigma(Y)^2)\lambda^2/2) \end{aligned}$$

Choosing  $p = 1 + \sigma(X)/\sigma(Y)$  gives  $\sigma(X + Y) \leq \sigma(X) + \sigma(Y)$ . If  $X$  and  $Y$  are independant, we actually conclude that  $X + Y$  is  $\sqrt{\sigma^2 + \tau^2}$  subgaussian, because we needn't apply Hölder's inequality, instead computing

$$M_{X+Y}(\lambda) = M_X(\lambda)M_Y(\lambda) \leq \exp(\lambda^2(\sigma^2 + \tau^2)/2)$$

Thus we get a better tail bound on the sum.

**Example.** Let  $X_1, \dots, X_N$  be independant subgaussian random variables. Let  $S = \sum X_n/N$ . Then  $\mathbf{E}(S) = \sum \mathbf{E}(X_n)/N$ , and  $\sigma(S) \leq \sqrt{\sum \sigma(X_n)^2}/N$ . We therefore conclude that

$$\mathbf{P}(|S - \mathbf{E}(S)| \geq t) \leq 2 \exp\left(\frac{-N^2 t^2}{\sum \sigma(X_n)^2}\right)$$

Thus the probability of error becomes incredibly small as  $N \rightarrow \infty$ .

If  $\sigma(X) = 0$ , then  $M_X(\lambda) \leq 1$  for all  $\lambda$ , which gives  $\mathbf{P}(|X| \geq t) \leq e^{-\lambda t}$  for all  $\lambda > 0$ , and taking  $\lambda \rightarrow \infty$  gives  $\mathbf{P}(|X| \geq t) = 0$  for all  $t > 0$ . Thus  $X = 0$  almost surely. More generally, if  $X$  is not centered, we let  $\sigma(X)$  denote  $\sigma(X - \mathbf{E}X)$ , i.e. the subgaussian measure of the noise corresponding to  $X$ . This is still a seminorm, whose kernel is the random variables constant almost surely.

**Lemma 3.10.** If  $X$  is subgaussian,  $\|X - \mathbf{E}X\|_{\psi_2} \lesssim \|X\|_{\psi_2}$ .

*Proof.* We use the triangle inequality to write

$$\|X - \mathbf{E}X\|_{\psi_2} = \|X\|_{\psi_2} + \|\mathbf{E}X\|_{\psi_2} \lesssim \|X\|_{\psi_2} + |\mathbf{E}X|$$

Thus it suffices to prove  $|\mathbf{E}X| \lesssim \|X\|_{\psi_2}$ . But by Jensen's inequality,

$$|\mathbf{E}X| \leq \mathbf{E}|X| = \|X\|_1 \leq \|X\|_{\psi_2}$$

This completes the argument. □

As a corollary, we find that  $\sigma(X) = \sigma(X - \mathbf{E}X) \lesssim \|X - \mathbf{E}X\|_{\psi_2} \lesssim \|X\|_{\psi_2}$ . The converse holds assuming that  $X$  is centrally distributed.

**Example.** If  $X$  is a symmetric Bernoulli random variable with

$$\mathbf{P}(X = 1) = \mathbf{P}(X = -1) = 1/2$$

Then  $\mathbf{E}(e^{X^2/t^2}) = e^{1/t^2}$ . This shows that  $\|X\|_{\psi_2} = (\log 2)^{-1/2}$

**Example.** If  $X$  is uniformly distributed on  $[-1, 1]$ , then

$$\mathbf{E}[X^k] = \frac{1}{2} \int_{-1}^1 x^k dx = \frac{1 - (-1)^{k+1}}{2(k+1)} = \begin{cases} \frac{1}{k+1} & k \text{ even} \\ 0 & k \text{ odd} \end{cases}$$

So

$$\mathbf{E}[e^{\lambda X}] = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k+1)(2k)!} \leq \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{2^k k!} = e^{\lambda^2/2}$$

so  $\sigma(X) \leq 1$ . Since  $\mathbf{V}(X) = 1$ , we must have  $\sigma(X) \geq 1$ , so we actually have  $\sigma(X) = 1$ . More generally, by scaling, if  $X$  is uniformly distributed on  $[-N, N]$ , then  $\sigma(X) = N$ .

**Example.** Suppose a centrally distributed random variable  $X$  satisfies  $|X| \leq M$  almost surely, then  $\sigma(X) \leq M$ . Assume without loss of generality that  $M = 1$ . Set  $Y = X + 1$ , and  $f(t) = (e^{2t} + 1)/2 - \mathbf{E}(e^{tY})$ . Since  $\mathbf{E}(Y) = 1$ ,  $f'(t) = \mathbf{E}(Y[e^{2t} - e^{tY}])$ . Since  $Y \leq 2$  almost surely,  $f'(t) \geq 0$ , and so  $f$  is increasing. In particular,  $f(0) = 1 - 1 = 0$ , so that for  $t \geq 0$ ,

$$\mathbf{E}(e^{tX}) = e^{-t} \mathbf{E}(e^{tY}) \leq \frac{e^t + e^{-t}}{2} \leq e^{t^2/2}$$

Since we can perform the same argument for  $-X$ ,  $X$  is 1 subgaussian.

### 3.4 Subexponential Random Variables

The class of subgaussian random variables is very flexible, but some distributions just don't have that thin a tail. There are obviously random variables like the Cauchy distribution, whose averages do not settle down asymptotically whatsoever, but there are still distributions which do seem rather well behaved, possessing moments of all orders, while not lying in

the category of subgaussian random variables. In this section we consider a slightly more general category of random variables for which we can get a tail bound, the subexponential random variables.

**Theorem 3.11.** *The following properties are equivalent, up to changes of constants:*

- $\mathbf{P}(|X| \geq t) \leq 2\exp(-t/K_1).$
- $\|X\|_p \leq K_2 \cdot p.$
- $\mathbf{E}(\exp(\lambda|X|)) \leq \exp(K_3 \cdot \lambda)$  for  $0 \leq \lambda \leq 1/K_3.$
- For some  $K_4$ ,  $\mathbf{E}(\exp(|X|/K_4)) \leq 2.$
- If  $\mathbf{E}X = 0$ ,  $\mathbf{E}(\exp(\lambda X)) \leq \exp(K_5^2 \lambda^2)$  for  $|\lambda| \leq 1/K_5.$

If the properties hold, we say that  $X$  is **subexponential**.

We leave the equivalence to the reader. A natural norm to place on this space is the Orlicz norm induced by  $\psi_1(x) = e^{x/t}/2.$

**Lemma 3.12.** *A variable  $X$  is subgaussian if and only if  $X^2$  is subexponential. The product of two subgaussian random variables is subexponential.*

The analogy of Hoeffding's inequality for subgaussian functions is Bernstein's inequality, which describes the tail behaviour of a sum of subexponential random variables.

**Theorem 3.13.** *If  $X_1, \dots, X_N$  are independent, centrally distributed, subexponential random variables. Then*

$$\mathbf{P}\left(\left|\frac{1}{N} \sum_{n=1}^N X_n\right| \geq t\right) \leq 2\exp\left(-c \min\left(\frac{N^2 t^2}{\sum \|X_m\|_{\psi_1}^2}, \frac{Nt}{\max \|X_m\|_{\psi_1}}\right)\right)$$

An important way to think of Bernstein's inequality is that it places the behaviour of deviation from the mean into two categories. For small deviations from the mean, the tail bound looks like that of a normal distribution. On the other hand, for large deviations, the tail becomes heavier, like an exponential distribution.

## Chapter 4

### Existence Theorems

In certain fields of probability theory, we wish to discuss collections of random variables defined over the same sample space. For instance, given a sequence  $\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_n$  of probability distributions defined over a space  $Y$ , we may want to talk about a sequence of independent random variables  $X_i : \Omega \rightarrow Y$ , such that  $\mathbf{P}(X_i \in U) = \mathbf{P}_i(U)$ . The construction here is simple; we take  $\Omega = Y^n$ , let  $X_i = \pi_i$  be the projection on the  $i$ 'th variable, and let  $\mathbf{P}$  be the probability measure induced by

$$\mathbf{P}(U_1 \times U_2 \cdots \times U_n) = \mathbf{P}_1(U_1)\mathbf{P}_2(U_2)\dots\mathbf{P}_n(U_n)$$

The construction here is simple because we have finitely many distributions, but the problem becomes much harder when we need to talk about an infinite family of distributions  $\mathbf{P}_i$ , or when we need to talk about non-independent random variables, with some specified relationships between the variables. The problem is to show there exists a sample space  $\Omega$  'big enough' for the random variables to all be defined on the space.

# Chapter 5

## Entropy

Let  $\mu$  be a probability distribution. We would like to measure the expected ‘amount of information’ contained in the distribution – in essence, the average information entropy of  $\mu$ . It was Claude Shannon who found the correct formula to measure this.

Shannon considered the problem of efficient information transfer. Suppose there was a channel of communication between two friends  $A$  and  $B$ . The friends have agreed on a standard dictionary  $X$  of possible messages, along with a probability distribution  $\mu$  over the dictionary, and we would like to encode these messages into bits, in such a way that the average length of the message is smallest. We then define this to be the information entropy of  $\mu$ . Shannon showed that if  $\mu$  is discrete with probabilities  $p_1, \dots, p_n$ , then the entropy can be calculated as

$$H(\mu) = \sum p_n \log_2 \left( \frac{1}{p_n} \right)$$

where the entropy is measured in bits, we can define the entropy in terms of the natural logarithm, in which case the entropy is said to be measured in nats. We assume that  $p_i \log 1/p_i = 0$  for  $p_i = 0$ , which makes sense by the continuity of  $x \log(1/x)$ .

The entropy of a distribution also tells us

Now suppose that we were attempting to optimize a message with respect to a discrete distribution  $\mu$ , and we instead encounter a distribution  $\nu$ . Then the policy we have used for messages will be less optimal than if we had known that  $\nu$  was the distribution in the first place. We define the relative difference in information between  $\mu$  and  $\nu$  as the difference

between the encoding of  $\nu$  with respect to  $\mu$ , and the encoding of  $\mu$  with respect to  $\mu$ . This is not a linearly ordered relation,  $\nu$  does not possess more information than  $\mu$ , just different information. If  $\mu$  takes probabilities  $p_i$  and  $\nu$  takes relative probabilities  $q_i$ , the difference in information is calculated to be

$$D(\mu, \nu) = \sum p_i \log(1/q_i) - \sum p_i \log(1/p_i) = \sum p_i \log(p_i/q_i)$$

This is known as the **Kullback Leibler distance** between  $\mu$  and  $\nu$ .

Now suppose we are viewing independent samples  $X_1, \dots, X_n$ , but we do not know where the samples are drawn from  $\mu$  or  $\nu$ . The larger  $D(\mu, \nu)$  is, the less time we should take to make an accurate decision that the distribution is  $\mu$  or  $\nu$ . Indeed, if  $p_i > 0$  and  $q_i = 0$ , then  $D(\mu, \nu) = \infty$ , and we can conclude with certainty that the distribution is  $\mu$  if we ever view the outcome corresponding to  $p_i$ .

It is necessary to define the ‘entropy’ of an arbitrary distribution, but it is then not clear how to interpret the entropy, since an encoding of uncountably many values will always have an infinite expected number of bits. However, we can define the relative entropy by performing a discretization; Let  $\mu$  and  $\nu$  be distributions on some sample space  $X$ . Consider function  $f : X \rightarrow \{1, \dots, n\}$ , and define

$$D(\mu, \nu) = \sup_f D(f_*\mu, f_*\nu)$$

where  $f_*$  pushes measures on  $X$  onto discrete measures on  $\{1, \dots, n\}$ . For a fixed  $f$ ,  $D(\mu, \nu)$  upper bounds the difference in information we expect to see over a particular discretization. One can then calculate that

$$D(\mu, \nu) = \begin{cases} \infty & : \mu \not\ll \nu \\ \int \log(d\mu/d\nu) d\mu & : \mu \ll \nu \end{cases}$$

The relative entropies of well known distributions are easy to compute. Normal distributions, for instance, have

$$D(N(\mu_1, \sigma^2), N(\mu_2, \sigma^2)) = (\mu_1 - \mu_2)^2 / 2\sigma^2$$

For Bernoulli distributions, we have

$$D(B(p), B(q)) = p \log(p/q) + (1-p) \log\left(\frac{1-p}{1-q}\right)$$



Which is true except perhaps at boundary conditions.

The Kullback Leibler distance gives us certain bounds which are essential to information theoretic lower bounds. The bound is useful, for it relates the probabilities of distributions by the difference in information contained within.

**Theorem 5.1** (The High Probability Pinsker Bound). *If  $\mu$  and  $\nu$  are probability measures on the same space  $X$ , and  $U \subset X$  is measurable, then*

$$\mu(A) + \nu(A^c) \geq \frac{1}{2} e^{-D(\mu, \nu)}$$

Suppose we have a decision procedure which attempts to distinguish between events in probability distributions. If we choose an event  $A$  upon which the decision procedure fails to make the correct decision on the measure  $\mu$ , and  $A^c$  measures the decision to fail under the measure  $\nu$ , then the bound above shows the decision procedure cannot work reliably on both  $\mu$  and  $\nu$ .

## Chapter 6

### Appendix: Uniform Integrability

Uniform integrability provides stronger conditions on controlling convergence in the  $L^1$  norm. For  $p > 1$ , inequalities often have ‘smoothing’ properties that are not apparent for the  $p = 1$  case, so uniform integrability provides particular techniques to help us. We start with a basic result in measure theory, specialized to probabilistic language.

**Lemma 6.1.** *If  $X \in L^1(\Omega)$  is a random variable, then for any  $\varepsilon > 0$ , there is  $\delta > 0$  such that for any event  $E$  with  $\mathbf{P}(E) \leq \delta$ ,*

$$\int_E |X| < \varepsilon$$

*Proof.* Suppose that there exists some  $\varepsilon$ , and events  $E_1, E_2, \dots$  with  $\mathbf{P}(E_k) \leq 1/2^k$  but with

$$\int_{E_k} |X| \geq \varepsilon$$

By taking successive unions, we may assume the  $E_i$  are a decreasing family of sets, and then

$$\int_{\bigcap_{k=1}^{\infty} E_k} |X| = \lim_{k \rightarrow \infty} \int_{E_k} |X| \geq \varepsilon$$

and  $\mathbf{P}(\bigcap E_k) = 0$ , which is impossible. □

**Corollary 6.2.** *If  $X \in L^1(\Omega)$ , and  $\varepsilon > 0$ , then there is  $K \in [0, \infty)$  with*

$$\int_{|X| > K} |X| < \varepsilon$$

A family of random variables  $\{X_\alpha\}$  is called **uniformly integrable** if given  $\varepsilon > 0$ , there is  $K \in [0, \infty)$  such that

$$\int_{|X_\alpha| > K} |X_\alpha| < \varepsilon$$

so that we can uniformly control the integral of  $X_\alpha$  over large sets. We note that

$$\mathbf{E}|X_\alpha| = \int_{|X_\alpha| > K} |X_\alpha| + \int_{|X_\alpha| \leq K} |X_\alpha| \leq \varepsilon + K$$

so a family of uniformly integrable random variables is automatically in  $L^1(\Omega)$ , and *bounded* in  $L^1(\Omega)$ . However, a family of random variables bounded in  $L^1(\Omega)$  is *not* necessarily uniformly integrable.

**Example.** Let  $\Omega$  be  $[0, 1]$  together with the Lebesgue measure. Let  $E_n = (0, 1/n)$ , and set  $X_n = n\chi_{E_n}$ . Then the  $X_n$  are bounded in  $L^1(\Omega)$ , but for  $n \geq K$ ,

$$\int_{X_n > K} X_n = 1$$

and so the random variables are not uniformly integrable.

These ‘concentrating bumps’ are essentially the only reason why we cannot always exchange expectations and limits, and require the application of the dominated convergence theorem. The condition of uniform integrability removes the ability for concentrating bumps to hide within the expectation of a family of random variables, and we find it also gives us conditions that guarantee we can exchange limits with integration. First, note that if we take a concentrated bump function  $X = n\chi_{E_n}$ , then  $\mathbf{E}|X| = 1$  is bounded uniformly over  $n$ , but  $\mathbf{E}|X|^{1+\varepsilon} = n^\varepsilon$  is unbounded, reflecting the fact that boundedness in  $L^p(\Omega)$  for  $p > 1$  removes concentrated bump functions by magnifying their effect.

**Theorem 6.3.** Suppose that  $\{X_\alpha\}$  is a class of random variables bounded in  $L^p$  for  $p > 1$ , then  $\{X_\alpha\}$  is uniformly integrable.

*Proof.* Consider some  $A \in [0, \infty)$  which gives a uniform bound  $\mathbf{E}|X_\alpha|^p < A$ . Applying Hölder’s inequality, we conclude

$$\int_{|X_\alpha| > K} |X_\alpha| \leq \int_{|X_\alpha| > K} \frac{|X_\alpha|^p}{K^{p-1}} \leq \frac{A}{K^{p-1}}$$

This is a uniform bound, and we may let  $K \rightarrow \infty$  to let the bound go to zero. Thus the family  $\{X_\alpha\}$  is uniformly integrable.  $\square$

**Corollary 6.4.** *If  $|X_\alpha| \leq Y$  is a uniform bound over a family  $\{X_\alpha\}$  of random variables, where  $Y \in L^1(\Omega)$ , then  $\{X_\alpha\}$  is uniformly integrable.*

*Proof.* We find

$$\int_{|X_\alpha| > K} |X_\alpha| \leq \int_{|X_\alpha| > K} Y \leq \int_{Y > K} Y$$

and as  $K \rightarrow \infty$ ,  $\mathbf{P}(Y > K) \rightarrow 0$ , and we can apply the continuity result to conclude that

$$\int_{Y > K} Y \rightarrow 0$$

and so we obtain a uniform bound.  $\square$

We recall that a sequence  $X_1, X_2, \dots$  of random variables **converges in probability** to a random variable  $X$  if, for every  $\varepsilon$ ,  $\mathbf{P}(|X_n - X| > \varepsilon) \rightarrow 0$ . If  $X_i \rightarrow X$  almost surely, then  $X_i \rightarrow X$  in probability, because we can apply the reverse Fatou lemma to conclude

$$\begin{aligned} 0 &= \mathbf{P}\left(\liminf_{n \rightarrow \infty} |X_n - X| > \varepsilon\right) \\ &= \mathbf{P}(\limsup\{|X_n - X| > \varepsilon\}) \geq \limsup \mathbf{P}(|X_n - X| > \varepsilon) \end{aligned}$$

Hence  $\mathbf{P}(|X_n - X| > \varepsilon) \rightarrow 0$ . The bounded convergence theorem links  $L^1$  convergence to convergence in probability using uniform integrability.

**Theorem 6.5.** *If  $X_n$  is a sequence of bounded random variables which tend to a random variable  $X$  in probability, then  $X_n \rightarrow X$  in the  $L^1$  norm.*

*Proof.* Let us begin by proving that if  $|X_n| \leq K$ , then  $|X| \leq K$  almost surely. This follows because for any  $k$ ,

$$\mathbf{P}(|X| > K + 1/k) \leq \mathbf{P}(|X - X_n| > 1/k) \rightarrow 0$$

so  $\mathbf{P}(|X| > K + 1/k) = 0$ , and letting  $k \rightarrow \infty$  gives  $\mathbf{P}(|X| > K) = 0$ . Let  $\varepsilon > 0$  be given. Then if we choose  $n$  large enough that

$$\mathbf{P}(|X_n - X| > \varepsilon) \leq \varepsilon$$

then

$$\begin{aligned}\mathbf{E}|X_n - X| &= \int_{|X_n - X| > \varepsilon} |X_n - X| + \int_{|X_n - X| \leq \varepsilon} |X_n - X| \\ &\leq 2K\varepsilon + \varepsilon\end{aligned}$$

we can then let  $\varepsilon \rightarrow 0$  to obtain  $L^1$  convergence.  $\square$

All this discussion concludes with a sufficient condition for  $L^1$  convergence, showing that uniform integrability is really the right condition which removes the pathologies which prevent us from exchanging expectation with pointwise limits.

**Theorem 6.6.** *Let  $X_n$  be a sequence of integrable random variables, and  $X$  is another integrable random variable. Then  $X_n \rightarrow X$  in the  $L^1$  norm if and only if  $X_n \rightarrow X$  in probability, and  $\{X_n\}$  is uniformly integrable.*

*Proof.* Fix  $K > 0$ , and consider

$$f_K(x) = \begin{cases} K & : x > K \\ x & : |x| \leq K \\ -K & : x < -K \end{cases}$$

Then for every  $\varepsilon > 0$ , we can choose  $K$  such that  $\|f_K(X_n) - X_n\|_1 \leq \varepsilon$ ,  $\|f_K(X) - X\|_1 \leq \varepsilon$  uniformly across  $n$  (adding a single variable to a uniformly integrable random variable keeps it uniformly integrable). But it is easy to see that  $f_K(X_n) \rightarrow f_K(X)$  in probability also, so by the bounded dominated convergence theorem, we conclude that  $\|f_K(X_n) - f_K(X)\| \rightarrow 0$ . A triangle inequality result gives the general result because the behaviour of  $X$  for large values is bounded by the uniform integrability.

To verify the reverse condition, note that if  $\mathbf{E}|X_n - X| \rightarrow 0$ , then Markov's inequality gives

$$\mathbf{P}(|X_n - X| \geq K) \leq \frac{\mathbf{E}|X_n - X|}{K} \rightarrow 0$$

to obtain uniform integrability, note that for each  $n$ ,  $\{X_1, \dots, X_n, X\}$  is uniformly integrable, then for each  $\varepsilon > 0$ , there is  $\delta$  such that if  $\mathbf{P}(E < \delta)$ ,

$$\int_E |X_n| < \varepsilon \quad \int_E |X| < \varepsilon$$

Since the entire set of  $X_n$  are bounded in  $L^1(\Omega)$ , we can choose  $K$  such that  $\sup \mathbf{E}|X_k| < \delta K$ , and then for  $m > n$ ,  $\mathbf{P}(|X_m - X| > K) < \delta$ , and so

$$\int_{|X_m| > K} |X_m| \leq \int_{|X_m| > K} |X| + \mathbf{E}|X - X_m| \leq 2\varepsilon$$

where we assume we have chosen  $n$  large enough that  $\mathbf{E}|X - X_m| \leq \varepsilon$ . The fact that for  $m \leq n$ ,

$$\int_{|X_m| > K} |X_m| \leq \varepsilon$$

follows from uniform integrability of the family  $\{X, X_1, \dots, X_n\}$ , so we have shown the entire infinite sequence is uniformly integrable.  $\square$

# Chapter 7

## High Dimensional Probability

In this chapter, we study the problems and phenomena that arise when studying random phenomena in high dimensional spaces. For instance, one often studies the behaviours of random vectors  $X \in \mathbf{R}^n$ , when  $n$  is a very large number. The exponential increase in room leads to concentration of the vector in unlikely places.

### 7.1 Concentration of Norm

We now see that the norm of a random vector  $X$  in a high dimensional space is essentially guaranteed to be close to it's mean. If  $X$  is a random Gaussian vector in  $\mathbf{R}^n$  with a standard normal distribution, then

$$\mathbf{E}|X|^2 = \sum \mathbf{E}X_i^2 = n$$

Thus we should expect  $|X|$  to be close to  $\sqrt{n}$ .

**Theorem 7.1.** *Let  $X$  be a random vector with independant, subgaussian coordinates and  $\mathbf{E}(X_m^2) = 1$ . Then  $\||X| - \sqrt{n}\|_{\psi_2} \lesssim K^2$ , where  $K = \max \|X_n\|_{\psi_2}$ .*

*Proof.* Assume  $K \geq 1$  (TODO: Prove this is sufficient). The  $X_m^2$  are subexponential, with  $\|X_m^2\|_{\psi_1} = \|X_m\|_{\psi_2}^2$ . Furthermore, by centering, we know

$\|X_m^2 - 1\|_{\psi_1} \lesssim \|X_m^2\|_{\psi_1}$ . Applying Bernstein's inequality shows that

$$\begin{aligned} \mathbf{P}(|X|^2 - n \geq t) &\leq 2 \exp \left( -c \min \left( \frac{t^2}{\sum \|X_n\|_{\psi_2}^4}, \frac{t}{\max \|X_n\|_{\psi_2}^2} \right) \right) \\ &\leq 2 \exp \left( -c \cdot \min \left( t^2/K^4, t/K^2 \right) \right) \\ &\leq 2 \exp(-c/K^4 \cdot \min(t^2, t)) \end{aligned}$$

Where we have used  $K \geq 1$  so that  $-1/K^2 \leq -1/K^4$ . This is a good concentration bound for  $|X|^2$ , and we now need to turn it into a concentration bound for  $|X|$ . To do this, we note that if  $x \geq 0$ , and  $|x - \sqrt{n}| \geq t$ , then

$$|x^2 - n| = |x - \sqrt{n}|(x + \sqrt{n}) \geq t(x + \sqrt{n}) \geq \max(t\sqrt{n}, t^2)$$

Since  $\min(\max(t\sqrt{n}, t^2)^2, \max(t\sqrt{n}, t^2)) \geq t^2$ ,

$$\mathbf{P}(|X| - \sqrt{n} \geq t) \leq \mathbf{P}(|X|^2 - n \geq \max(t\sqrt{n}, t^2)) \leq 2 \exp(-c/K^4 \cdot t^2)$$

This is equivalent to  $\|X| - \sqrt{n}\|_{\psi_2} \lesssim K^2$ . □

This theorem says that with probability independent of  $n$ ,  $X$  lies within a distance  $O(K)$  from the sphere of radius  $\sqrt{n}$ .

$$\mathbf{P}(|X| - \sqrt{n} \geq t) \leq \exp(-ct^2/K^2)$$

An intuitive explanation of this is that  $|X|^2$  has mean  $n$ , and standard deviation  $O(\sqrt{n})$ . If  $|X|^2 = n + O(\sqrt{n})$ , then

$$|X| = \sqrt{n + O(\sqrt{n})} = \sqrt{n} + O(1)$$

So we should expect  $|X|$  to deviate from  $\sqrt{n}$  by a constant distance, independent of  $n$ . This is precisely what the theorem above says.

## 7.2 Isotropic Vectors

There is one nice way of generalizing this behaviour to random vectors whose coordinates need not be independent. Given a centrally distributed random vector  $X$  in  $\mathbf{R}^n$ , we define the  $n \times n$  covariance matrix  $\Sigma(X)$  by the



formula  $\Sigma(X)_{ij} = \mathbf{E}(X_i X_j)$ . This is a symmetric, positive semi-definite matrix, since  $v^T \Sigma(X) v = \mathbf{E}((v \cdot x)^2) \geq 0$ . The spectral decomposition theorem implies that there is a basis of normalized eigenvectors  $u_1, \dots, u_n$  for  $\Sigma(X)$ , with eigenvalues  $\lambda_1, \dots, \lambda_n$ . We assume that they have been arranged so that the eigenvalues are placed in decreasing order. If  $Y_i = u_i \cdot X$ , then this means  $\mathbf{E}(Y_i^2) = \lambda_i$ , and  $\mathbf{E}(Y_i Y_j) = 0$  if  $i \neq j$ . In the case where  $X$  is normally distributed, then this implies that the  $Y_i$  are independent Gaussian vectors, with variance  $\lambda_i$ . If we throw away the vectors  $Y_i$  with  $\mathbf{E}(Y_i^2) = 0$  (so  $Y_i = 0$  almost everywhere), then normalize, we end up with what is essentially an independent set of vectors.

We say a random vector  $X$  is isotropic if  $\Sigma(X)$  is the identity matrix. This is equivalent to saying that for each vector  $x$ ,  $\mathbf{E}(X \cdot x)^2 = |x|^2$ . This is because  $\mathbf{E}(X \cdot x)^2 = x^T \Sigma(X) x$ , and this is  $|x|^2$  for all  $x$  if and only if  $\Sigma(X) = I$ . Thus the one dimensional projections of an isotropic vector  $X$  onto each coordinate axis have unit variance, so the vector is extended evenly in all directions.

**Example.** Let  $K$  be a convex set with nonempty interior in  $\mathbf{R}^n$ . If  $X$  is a uniformly chosen point in  $K$ , which by translation we may assume to have mean zero, and has covariance matrix  $\Sigma$ , then  $\Sigma$  is positive definite, because if  $\Sigma$  has a zero eigenvalue, there would be a vector  $a$  such that  $X$  is almost surely orthogonal to  $a$ , which is impossible since  $K$  has non-empty interior. Thus  $\Sigma^{-1/2} X$  is isotropic, and thus  $\Sigma^{-1/2} K$  can be seen as a convex set ‘uniformly in each direction’.

**Lemma 7.2.** If  $X$  is isotropic, then  $\mathbf{E}|X|^2 = n$ . More generally, if  $X$  and  $Y$  are independent and isotropic, then  $\mathbf{E}(X \cdot Y)^2 = n$ .

*Proof.* First we show  $\mathbf{E}|X|^2 = n$ . To do this, we write

$$\begin{aligned} \mathbf{E}|X|^2 &= \mathbf{E}X^T X = \mathbf{E}\text{tr}(X^T X) \\ &= \mathbf{E}\text{tr}(XX^T) = \text{tr}(\mathbf{E}(XX^T)) = \text{tr}(I) = n \end{aligned}$$

Next, given  $Y$ , we find

$$\mathbf{E}((X \cdot Y)^2 | Y) = \sum Y_i Y_j \mathbf{E}(X_i X_j | Y) = \sum Y_i Y_j \mathbf{E}(X_i X_j) = \sum Y_i^2 = |Y|^2$$

But this means that

$$\mathbf{E}((X \cdot Y)^2) = \mathbf{E}(\mathbf{E}((X \cdot Y)^2 | Y)) = \mathbf{E}|Y|^2 = n$$

Thus  $|X|$  and  $|Y|$  must essentially be equal to  $\sqrt{n}$ . □

Thus  $|X \cdot Y| = \sqrt{n}$ . If we normalize, then we find that  $(X \cdot Y)/|X||Y|$  must be approximately  $1/\sqrt{n}$ , so that in high dimensional spaces, two independent isotropic vectors are with overwhelming probability almost orthogonal to one another. This is very different from in two dimensions, where two independent unit vectors chosen at random are on average at an angle  $\pi/4$  from one another.

**Example.** A fundamental example of an isotropic random vector is a vector chosen uniformly randomly on the sphere of radius  $\sqrt{n}$ . For  $i \neq j$ ,  $(X_i, X_j)$  is identically distributed to  $(X_i, -X_j)$ , so  $\mathbf{E}(X_i X_j) = -\mathbf{E}(X_i X_j)$ , implying  $\mathbf{E}(X_i X_j) = 0$ . Since  $\mathbf{E}|X|^2 = n$ , and the  $\mathbf{E}|X_i|^2$  are independent of  $i$ , this implies  $\mathbf{E}|X_i|^2 = 1$ .

It is good to remember that the coordinates of an isotropic vector need not be independent. A uniformly random point  $X$  on the radius  $\sqrt{n}$  sphere need not be independent, because the points must satisfy the relationship  $X_1^2 + \dots + X_n^2 = n$ .

**Example.** A discrete isotropic distribution is given by the Bernoulli distribution, which takes the values  $-1$  and  $1$  with equal likelihood. In this case the coordinates are independent from one another, and they are centrally distributed with variance one.

**Example.** For the most extreme example of a discrete isotropic distribution, we can take a vector uniformly randomly from  $\{\sqrt{n}e_1, \dots, \sqrt{n}e_n\}$ . Then  $X_i^2$  is  $\{0, 1\}$  valued, with  $\mathbf{P}(X_i^2 = 1) = 1/n$ . This gives  $\mathbf{E}(X_i^2) = 1/n$ . On the other hand,  $X_i X_j = 0$  for  $i \neq j$ , so  $\mathbf{E}(X_i X_j) = 0$ . Note that the mean of  $\mathbf{E}(X_i)$  is non-zero, it is actually equal to  $1/\sqrt{n}$ .

We obtain a family of discrete isotropic random vectors by considering uniformly distributions over discrete families of vectors used most notably in signal processing, known as **frames**. A **frame** is a set  $\{v_1, \dots, v_n\}$  which obeys the approximate  $A|x| \leq \sum (v_i \cdot x)^2 \leq B|x|^2$  for all vectors  $x$ . if  $A = B$ , the frame is called **tight**. A frame is tight if and only if  $\sum v_i^T v_i = AI_n$ .

**Example.** An example of a tight frame which isn't an orthonormal basis is the 'Mercedes Benz' frame, three equidistance points on a circle in the plane. If

$$v_1 = (0, 1) \quad v_2 = \left( -\frac{1 + \sqrt{3}}{2\sqrt{2}}, -\frac{\sqrt{3} - 1}{2\sqrt{2}} \right) \quad v_3 = \left( \frac{1 + \sqrt{3}}{2\sqrt{2}}, -\frac{\sqrt{3} - 1}{2\sqrt{2}} \right)$$

then

$$v_1^T v_1 = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \quad v_2^T v_2 = \begin{pmatrix} \frac{2+\sqrt{3}}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{2-\sqrt{3}}{4} \end{pmatrix} \quad v_3^T v_3 = \begin{pmatrix} \frac{2+\sqrt{3}}{4} & \frac{-1}{4} \\ \frac{-1}{4} & \frac{2-\sqrt{3}}{4} \end{pmatrix}$$

$$v_1 = (0, 1) \quad v_2 = (-\sqrt{3/4}, -\sqrt{1/4}) \quad v_3 = (\sqrt{3/4}, -\sqrt{1/4})$$

and this means that

$$(x, v_1)^2 + (x, v_2)^2 + (x, v_3)^2 = x_2^2 + (1/4)(\sqrt{3}x_1 + x_2)^2 + (1/4)(\sqrt{3}x_1 - x_2)^2$$

$$= 3/2|x|^2$$

So the frame is tight with  $A = B = 3/2$ .

**Theorem 7.3.** If  $\{u_1, \dots, u_N\}$  is tight, then a uniformly random choice of a frame element, scaled by  $\sqrt{n/A}$  is isotropic. Conversely, if  $X$  is isotropic, and takes only finitely many values  $\{u_1, \dots, u_N\}$  with probability  $p_i$ , then  $\sqrt{p_i}u_i$  is a tight frame with  $A = 1$ .

*Proof.* If  $X = \sqrt{n/A} \cdot \text{Uni}(u_1, \dots, u_n)$ , then

$$\mathbf{E}(X \cdot x)^2 = \frac{(n/A) \sum (u_i \cdot x)^2}{n} = (1/A)A|x|^2 = |x|^2$$

which means precisely that the frame is isotropic. To prove the converse, we find that

$$\sum (p_i^{1/2} u_i \cdot x)^2 = \sum p_i (u_i \cdot x)^2 = \mathbf{E}((X \cdot x)^2) = |x|^2$$

which is precisely the definition of a tight frame.  $\square$

## 7.3 Subgaussian Vectors

We say a random vector  $X$  is subgaussian if the  $X \cdot x$  is subgaussian for all  $x$ . Then there is a smallest value  $\|X\|_{\psi_2} < \infty$  if  $\|X \cdot x\|_{\psi_2} \leq \|X\|_{\psi_2} |x|$ . It is easy to verify that if  $X$  has independent, subgaussian coordinates  $X_1, \dots, X_n$ , then  $X$  is subgaussian and  $\|X\|_{\psi_2} \lesssim \max \|X_k\|_{\psi_2}$ . Even if the coordinates are not independent,  $X$  is subgaussian, but we need not have  $\|X\|_{\psi_2} \lesssim \max \|X_k\|_{\psi_2}$ . Examples include Gaussian vectors, and independent  $\{-1, 1\}$  Bernoulli distributions.

**Example.** The random vector  $X$  chosen uniformly randomly from  $\{\sqrt{n}e_n\}$  is technically subgaussian, but not ‘subgaussian enough’ for most purposes. Since  $\mathbf{P}(X_k \geq \sqrt{n}) = 1/n$ , the subgaussian norm for  $X_k$  is such that  $1/n \leq \exp(-cn/\|X\|_{\psi_2}^2)$ , or  $\|X\|_{\psi_2} \gtrsim (n/\log n)^{1/2}$ . We will also show that  $\|X\|_{\psi_2} \lesssim (n/\log n)^{1/2}$ . To do this, we must show there is a small constant  $c$  such that for any  $a_i$  with  $\sum a_i^2 = 1$ ,  $\mathbf{P}(\sum a_i X_i \geq \sqrt{n}t) \leq 2e^{-c \log n t^2}$ . If we assume  $a_1 \geq a_2 \geq \dots \geq a_n$ , then the discreteness of the random variable we are working with makes it sufficient to prove that

$$\mathbf{P}\left(\sum a_i X_i \geq \sqrt{n} \cdot a_k\right) = k/n \leq 2n^{-ca_k^2} = 2e^{-ca_k^2 \log n}$$

This is equivalent to showing that  $k \leq 2n^{1-ca_k^2}$ . Since  $1 \geq a_1^2 + \dots + a_k^2 \geq ka_k^2$ ,  $a_k^2 \geq 1/k$ , and it suffices to prove that  $k \leq 2n^{1-c/k}$ , or  $\log k/(1-c/k) \leq \log 2 + \log n$ . But as  $k \rightarrow \infty$ ,

$$\frac{\log k}{1-c/k} = \log k(1 + O(c/k)) \leq \log n + c \cdot o(1)$$

Thus there is  $N$ , independent of  $n$  and  $c$ , such that for  $k \geq N$ ,  $\log k(1 - c/k)^{-1} \leq \log n + c \log 2$ . Choosing  $c \leq 1$  completes the argument in this case. For  $k = \{1, \dots, N\}$ , we can just choose  $c$  small enough that the inequality holds here, and because  $N$  is independent of  $n$  and  $c$ , the choice of  $c$  made here is independent of  $n$ . Thus the general inequality is established.

**Example.** Suppose  $X$  is an isotropic random vector supported on a finite set with  $N$  elements. If  $\|X\|_{\psi_2} = O(1)$ , then we find  $N \geq e^{cn}$  for some  $c$ . In particular, the only frames which have good subgaussian properties must have exponentially many vectors in them, which makes them fairly useless in practice for generating good subgaussian results.

Let  $\mathbf{P}(X = v^i) = p_i$ . Then there is a constant  $c$  such that for any choice of  $a$  with  $|a| = 1$ ,  $\mathbf{P}(a \cdot X \geq t) \leq 2\exp(-ct^2)$ . We have  $\mathbf{E}((a \cdot X)^2) = 1$  for any  $a$  with  $|a| = 1$ . TODO

The uniform point on the sphere is a well behaved subgaussian random variable for which the coordinates are not independent of one another.

**Example.** Hoeffding’s inequality says that a random  $X$  chosen from  $\{-1, 1\}$  is subgaussian, with  $\|X\|_{\psi_2}$  bounded independent of  $n$ . This makes the subgaussian bound very useful to use.

**Theorem 7.4.** *If  $X$  is a uniformly random vector on the sphere with radius  $\sqrt{n}$ , then  $\|X\|_{\psi_2}$  is bounded independantly of  $n$ .*

*Proof.* Let  $Y$  be a Gaussian vector. Then  $\sqrt{n}Y/|Y|$  is uniformly distributed on the sphere. By rotation invariance, showing that  $\|x \cdot X\|_{\psi_2}$  is bounded, it suffices to show that  $\|X_1\|_{\psi_2}$  is bounded. It suffices to obtain tail bounds for  $X_1$  only if  $t \leq \sqrt{n}$ , for they are trivial for  $t \geq \sqrt{n}$ . We know  $\mathbf{P}(|Y| - \sqrt{n} \geq t/2) \leq 2\exp(-Ct^2)$ , so

$$\begin{aligned} \mathbf{P}(Y_1/|Y| \geq t/\sqrt{n}) &\leq 2\exp(-Ct^2) + \mathbf{P}(Y_1 \geq t(1 - t/2\sqrt{n})) \\ &= 2\exp(-Ct^2) + \mathbf{P}(Y_1 \geq t/2) \\ &= 2\exp(-Ct^2) + 2\exp(-Ct^2) = 4\exp(-Ct^2) \end{aligned}$$

This is enough to obtain the required result.  $\square$

**Corollary 7.5.** *A uniform vector  $X$  in the unit ball of radius  $\sqrt{n}$  has  $\|X\|_{\psi_2}$  bounded independantly of  $n$ .*

*Proof.* If  $Y = X/|X|$ , then  $Y$  is uniformly distributed on the unit sphere, and  $Y \geq X$ , so  $\|X\|_{\psi_2} \leq \|Y\|_{\psi_2}$ , which is bounded.  $\square$

One way to view this theorem is that the marginals  $X \cdot x$  of the uniform distribution on the sphere become very close to normal in high dimensional space. i.e. if a unit vector  $x$  is fixed, then  $X \cdot x$  converges to the normal distribution as the dimension increases. The theorem above provides a concentration bound for how close  $X \cdot x$  lies to a normal distribution.

## 7.4 Concentration Without Independance

Let  $X$  be a subgaussian vector, and  $f$  a real valued function. A natural question to ask is when  $f(X)$  concentrates about it's mean  $\mathbf{E}f(X)$ . For linear functions  $f$ , this question is easy. We cannot expect this if  $f$  is an arbitrary function. Nonetheless, if  $f$  does not oscillate too much, the theorem is true.

The main theorem of this section gives the result for Lipschitz functions on the sphere of radius  $\sqrt{n}$ , where  $X$  is chosen uniformly at random on the sphere. We will prove that if  $\|f(X) - \mathbf{E}f(X)\|_{\psi_2} \leq \|f\|_{\text{Lip}}$ . To do this, we compare the nonlinear  $f$  with a linearized version by looking at spherical caps on the sphere, and using a remarkable geometric principle called the isoperimetric inequality.

**Theorem 7.6.** *Let  $\varepsilon > 0$ . Then among all sets  $E$  with surface area  $A$ , the minimizers of the surface area of  $E_\varepsilon$  are obtained by setting  $E$  to be a spherical cap.*

This can be compared to the classical isoperimetric inequality on  $\mathbf{R}^n$ , which says that the sets  $E$  which minimize the volume of  $E_\varepsilon$  are obtained by setting  $E$  to be a ball. This principle implies a ‘blow up’ phenomenon in high dimension. We let  $\sigma$  denote the normalized spherical measure on the sphere.

**Theorem 7.7.** *Let  $E$  be a subset of the radius  $\sqrt{n}$  sphere. If  $\sigma(E) \gg 2\exp(-cs^2)$ , then for every  $t \geq s$ ,  $\sigma(E_{2t}) \geq 1 - 2\exp(-ct^2)$ .*

*Proof.* Assume first that  $\sigma(E) \geq 1/2$ . Let  $H$  denote the set of all points  $x$  on the sphere with  $x_1 \leq 0$ . Since  $\sigma(E) \geq 1/2 = \sigma(H)$ , the isoperimetric inequality on the sphere implies that  $\sigma(E_t) \geq \sigma(H_t)$ . And if  $X$  is uniformly chosen on the sphere, then

$$\sigma(H_t) = \mathbf{P}(X \leq t/\sqrt{2}) \geq 1 - 2\exp(-ct^2)$$

This gives the required inequality for  $E_t$ .

Now suppose that  $\sigma(E) > 2\exp(-cs^2)$ . Then  $\sigma(E_s) > 1/2$ . To see why, assume that  $\sigma(E_s) < 1/2$ . Then  $F = E_s^c$  has  $\sigma(F) \geq 1/2$ , so  $\sigma(F_s) \geq 1 - 2\exp(-cs^2)$ ,  $F_s$  is disjoint from  $E$ , and  $\sigma(F_s) + \sigma(E_s) > 2\exp(-cs^2) + 1 - 2\exp(-cs^2) = 1$ , which gives a contradiction.  $\square$

**Theorem 7.8.** *Let  $X$  be a random vector on the sphere of radius  $\sqrt{n}$ . if  $f$  is Lipschitz, then  $\|f(X) - \mathbf{E}f(X)\|_{\psi_2} \lesssim \|f\|_{\text{Lip}}$ .*

*Proof.* We may assume by scaling that  $\|f\|_{\text{Lip}} = 1$ . Let  $M$  be a median for  $f$ , satisfying  $\mathbf{P}(f(X) \leq M) \geq 1/2$  and  $\mathbf{P}(f(X) \geq M) \geq 1/2$ . Let  $E$  denote the level set  $\{x : f(x) \leq M\}$ . But then  $\sigma(E_t) \geq 1 - 2 \cdot \exp(-ct^2)$ , and because  $f$  is Lipschitz,  $E_t$  is also a subset of the level set  $\{x : f(x) \leq M + t\}$ . Thus  $\mathbf{P}(f(X) \leq M + t) \geq 1 - 2 \cdot \exp(-ct^2)$ . But by symmetry, we also know  $\mathbf{P}(f(X) \geq M - t) \geq 1 - 2 \cdot \exp(-ct^2)$ . This gives the concentration bound

$$\mathbf{P}(|f(X) - M| \geq t) \leq 4 \cdot \exp(-ct^2)$$

so  $\|f(X) - M\|_{\psi_2} \lesssim 1$ . To replace  $M$  with  $\mathbf{E}(f(X))$ , we just need to recenter. The expectation of  $f(X) - M$  is  $\mathbf{E}f(X) - M$ , so the recentering bound gives

$$\|f(X) - \mathbf{E}f(X)\|_{\psi_2} = \|(f(X) - M) - (\mathbf{E}f(X) - M)\|_{\psi_2} \lesssim \|f(X) - M\|_{\psi_2} \lesssim 1$$

This completes the argument.  $\square$

If we instead work on the unit sphere, then we find  $\|f(X) - \mathbf{E}(f(X))\|_{\psi_2} \lesssim \|f\|_{\text{Lip}} n^{-1/2}$ . Equivalently, this means that

$$\mathbf{P}(|f(X) - \mathbf{E}(f(X))| \geq t) \leq 2 \exp(-cnt^2 \|f\|_{\text{Lip}}^{-2})$$

which gives a strong concentration bound in high dimensions.

**Example.** If  $X$  lies on the sphere of radius 1, then we know  $\|X\|_{\psi_2} \leq 1/\sqrt{n}$ . Thus for any  $x$  lying on the sphere,  $\|X \cdot x\|_{\psi_2} \leq 1/\sqrt{n}$ , and so

$$\mathbf{P}(|X \cdot x| \geq \varepsilon) \leq 2 \exp(-cn\varepsilon^2)$$

For a geometric application, we show that while there are at most  $n$  orthogonal vectors in  $\mathbf{R}^n$ , we can have exponentially many almost orthogonal vectors. Two unit vectors  $x$  and  $y$  are almost orthogonal if  $|x \cdot y| \leq \varepsilon$ . We construct such a set inductively. Consider unit vectors  $e_1, \dots, e_N$ , which are almost orthogonal. For each  $k$  we can consider  $E_k = \{x \in S^{n-1} : |(x \cdot e_k)| \leq \varepsilon\}$ . Then  $\sigma(E_k) \geq 1 - 2 \exp(-cn\varepsilon^2)$ , and so  $\sigma(E_1 \cap \dots \cap E_N) \geq 1 - 2N \exp(-cn\varepsilon^2)$ .  $N < \exp(cn\varepsilon^2)/2$ , this is positive, so there certainly exists a unit vector simultaneously orthogonal to all other vectors. Adding this to the list and continuing, we can work up to the point where  $N \geq \exp(cn\varepsilon^2)/2$ .

There is nothing really special to the sphere here. Given any other metric measure space with an isoperimetric inequality, with a blow up in mass in high dimensions, we can obtain the same result. Here we also consider the examples of concentration of Gaussian measures, and concentration of mass on the Hamming cube.

**Example.** Consider  $\mathbf{R}^n$  equipped with the Gaussian measure, which has the Gaussian distribution as a density function. It is non-obvious, but the minimizers of measure expansion are achieved by the half plane. Thus we can deduce that if  $X$  is a Gaussian vector, and  $f$  is Lipschitz, then  $\|f(X) - \mathbf{E}f(X)\|_{\psi_2} \lesssim \|f\|_{\text{Lip}}$ . We should expect the Gaussian result to look essentially the same as the result based on the uniform distribution on spheres, because in high dimensions, the two results are essentially the same.

**Example.** Next, we look at the uniform distribution on  $\{-1, 1\}^n$ . The isoperimetric inequality here is minimized by Hamming balls, which are neighbourhoods of points with respect to the Hamming distance  $d(x, y)$ , which gives the number of  $i$  with  $x_i \neq y_i$ . Thus we conclude that  $\|f(X) - \mathbf{E}f(X)\|_{\psi_2} \lesssim \|f\|_{\text{Lip}} n^{-1/2}$ . Similar techniques work for the Hamming distance on the symmetric group, and give the same equation.

**Example.** If  $M$  is a Riemannian manifold, we can consider the arclength distance, as well as normalized volume of  $M$  inducing a probability distribution  $X$  chosen uniformly at random on  $M$ . If  $c(M)$  denotes the infimum of the Ricci curvature tensor over all points, and  $c(M) > 0$ , then we have a concentration result  $\|f(X) - \mathbf{E}f(X)\|_{\psi_2} \lesssim \|f\|_{\text{Lip}}/\sqrt{c(M)}$ . Important applications of this example occur on the special orthogonal groups  $SO_n$ , the Grassmanians  $G_{n,m}$ .

## 7.5 The Johnson-Lindenstrauss Lemma

Suppose we have  $N$  data points in  $\mathbf{R}^n$ , where  $n$  is very large. We would like to reduce the dimension of the data, while still preserving the geometric properties of the data points. The simplest data reduction is to project the data points onto a lower dimensional subspace. A natural question the smallest dimension we can project the points, while still approximately preserving the distance between points. The Johnson-Lindenstrauss lemma says the distances will be approximately preserved when projecting into a space with dimension  $\log N$ .

**Lemma 7.9.** Let  $\pi$  be a randomly chosen projection onto an  $m$  dimensional subspace of  $G_{n,m}$ . If  $z \in \mathbf{R}^n$  is fixed, and  $\varepsilon > 0$ , then  $\mathbf{E}|\pi z|^2 \leq (m/n)|z|^2$ , and with probability greater than  $1 - 2\exp(-c\varepsilon^2 m)$ ,

$$(m/n)^{1/2}(1 - \varepsilon)|z| \leq |\pi z| \leq (1 + \varepsilon)(m/n)^{1/2}|z|$$

*Proof.* Without loss of generality, assume that  $|z| = 1$ . Then, instead of considering a random projection  $\pi$ , we can consider a fixed projection acting on a random unit vector  $z$ , since the distribution of  $\pi(z)$  is the same. Using rotation invariance, we may assume that  $\pi$  is the projection onto the first  $m$  coordinates. Thus

$$\mathbf{E}|\pi(z)|^2 = \sum_{i=1}^m \mathbf{E}z_i^2 = m/n$$

Thus the first part of the lemma is proven. Next, we apply the concentration result for Lipschitz functions on a sphere. if  $f(x) = |\pi(x)|$ , then  $\|f\|_{\text{Lip}} = 1$ . Thus

$$\mathbf{P}\left(\left|Px| - (m/n)^{1/2}\right| \geq t\right) \leq 2\exp(-cnt^2)$$

Picking  $t = \varepsilon\sqrt{m/n}$  completes the proof.  $\square$



**Theorem 7.10.** *Let  $E$  denote a uniformly randomly chosen  $m$  dimensional subspace of  $\mathbf{R}^n$ . If  $X$  is a set of  $N$  points in  $\mathbf{R}^n$ ,  $\varepsilon > 0$ , and  $m \gtrsim \log N / \varepsilon^2$ , then with probability  $1 - 2\exp(-c\varepsilon^2 m)$ , the projection  $\pi$  of  $X$  onto  $E$  satisfies*

$$(1 - \varepsilon)|x - y| \leq (n/m)^{1/2}|\pi(x) - \pi(y)| \leq (1 + \varepsilon)|x - y|$$

*for all  $x$  and  $y$ , so the geometry of  $X$  is scaled and approximately preserved in lower dimensions.*

*Proof.* The idea is just to bound how the random projection  $\pi$  behaves on a pair  $x, y$ , and then take a union bound. For each  $x, y$ , setting  $z = x - y$  and applying the last lemma gives the result with overwhelming probability. And since we have  $N^2$  choices of  $x, y$ ... TODO.  $\square$

# Chapter 8

## Percolation Theory

Let us consider the two dimensional theory of percolation. The two examples we have in mind are the lattice  $\mathbf{Z}^2$ , and the triangular lattice  $\mathbf{T}$ . For any  $p \in [0, 1]$ , we define a graph structure on  $\mathbf{Z}^2$ , adding an edge between two adjacent elements of the lattice with independent probability  $p$ . On  $\mathbf{T}$ , we instead consider *site percolation*, where we keep a hexagon with probability  $p$ .

**Theorem 8.1** (Russo-Seymour-Welsh). *If  $p = 1/2$ , then for any  $a, b > 0$ , there exists  $c$  such that if  $A_n$  denotes the event that we can travel from the left edge to the right edge of the lattice  $[0, a \cdot n] \times [0, b \cdot n] \cap \mathbf{Z}^2$ , then  $c < \mathbf{P}(A_n) < 1 - c$ .*

One of the main problems in percolation theory is to determine how likely it is to find an infinite connected set of vertices, or cluster, in the randomly selected graph. As the probability of each edge becomes more likely, the graph becomes more and more connected. We find that for  $p > 1/2$ , there is almost surely an infinite cluster, and for  $p < 1/2$ , there is almost surely *not* a cluster. The value  $p = 1/2$  is therefore called the *phase transition* point. A very related value to the phase transition problem is the percolation density function  $\theta$ , which for each  $p$ , gives the probability  $\theta(p)$  of the origin being in an infinite cluster of the graph. As an example, it is known that on  $\mathbf{T}$ ,  $\theta(p) = (p - 1/2)^{5/36 + o(1)}$ , as  $p \downarrow 1/2$ . Determining phase transition points is the main focus of this chapter's notes.

## 8.1 Duality

Note an important duality in these geometric scenarios. Given any graph on  $\mathbf{Z}^2$ , we can obtain another graph, the dual graph, by taking the vertices as unit squares with corners on  $\mathbf{Z}^2$ , and with an edge between adjacent squares if there is no edge separating the two squares. Then the probability that there is an edge between two squares is the same as the probability that we do *not* select the corresponding separating edge, i.e. with probability  $1 - p$ . This will be useful.

The book says the dual graph of  $T$  is  $T$ , but I don't quite understand why?

## 8.2 Boolean Functions and Sharp Thresholds

If  $G$  is a graph, then the family of all graph structures on these vertices can be identified with  $\{0,1\}^E = \{0,1\}^{O(V^2)}$ . Thus a function  $f$  on the set of graphs can be identified with a boolean function, and we can apply boolean function techniques. In our case, the natural graphs will be the subgraphs of  $\mathbf{Z}^2$  of the form  $[0, a \cdot n] \times [0, b \cdot n]$ . An example of a Boolean function on graphs is obtained by setting

$$f_n(G) = \mathbf{I}(\text{There is a path from left to right in } G)$$

Boolean analysis tells us the main features of  $G$ . Here we introduce some basics, which will help us get the job done.

If  $f$  is a boolean function, we say it is monotone if  $x_i \leq y_i$  for each  $i$  implies  $f(x) \leq f(y)$ . An index  $i$  is pivotal for an input  $x$  if  $f(x) \neq f(x^i)$ , where  $x^i$  is  $x$  with the bit flipped at the  $i$ 'th position. The influence  $\mathbf{I}_i(f)$  of  $f$  in the variable  $i$  is then the probability that for a randomly chosen  $x$ ,  $i$  is pivotal. If we instead choose an input to be equal to one with probability  $p$ , and zero with probability  $1 - p$ , then the probability that  $i$  is pivotal is denoted  $\mathbf{I}_i^p(f)$ . The sum of the influence over all influences  $i$  is known as the *total influence*. Now if  $E$  is a monotone event, then it is obvious that as  $p$  increases,  $\mathbf{P}(E)$  should increase. The degree to which it increases is quantified by the Margulis-Russo formula.

**Theorem 8.2** (Margulis-Russo). *Let  $E$  be monotone. Then  $d\mathbf{P}(E)/dp = \mathbf{I}^p(E)$ .*

*Proof.* Temporarily, let  $\mathbf{I}_i^{p_1, \dots, p_n}(E)$  denote the chance that  $X^i \neq X$ , where the  $X_j \in \{0, 1\}$  are chosen uniformly at random with  $\mathbf{P}(X_j = 1) = p_j$ . Define  $\mathbf{I}^{p_1, \dots, p_n}(E) = \sum \mathbf{I}_i^{p_1, \dots, p_n}(E)$ . It suffices to show  $\partial \mathbf{P}(E)/\partial p_i = \mathbf{I}_i^{p_1, \dots, p_n}(E)$ , from which we can use the chain rule. We can write  $E$  as the union of two disjoint events  $E_0$  and  $E_1$ , where  $E_0 = E \cap \{\chi_E(X^i) \neq \chi_E(X)\}$ , and  $E_1 = E \cap \{\chi_E(X^i) = \chi_E(X)\}$ . Now  $E_1$  does not depend on the value of  $X_1$  at all, so  $\mathbf{P}(E_1)$  is independent of  $p_i$ , and so

$$\frac{\partial \mathbf{P}(E_1)}{\partial p_i} = 0$$

On the other hand, by monotonicity,  $E_0$  then equals the probability that  $\chi_E(X^i) \neq \chi_E(X)$ , intersected with the event  $X_i = 1$ . These two events are independent, so  $\mathbf{P}(E_0) = p_i \mathbf{P}(\chi_E(X^i) \neq \chi_E(X))$ . The latter probability does not depend on the index  $i$ , so

$$\frac{\partial \mathbf{P}(E_0)}{\partial p_i} = \mathbf{P}(\chi_E(X^i) \neq \chi_E(X)) = \mathbf{I}_i^{p_1, \dots, p_n}(E)$$

This completes the proof.  $\square$

To analyze the critical exponent, we rely on two results we will prove later on using Fourier analysis, which allow us to upper bound the influence by the variance of a function.

**Theorem 8.3** (Bourgain, Kahn, Kalai). *For any  $f$  and  $p$ , there exists  $i$  such that  $\mathbf{I}_i^p(f) \gtrsim \mathbf{V}_p(f)[\log n/n]$ , and  $\mathbf{I}^p(f) \gtrsim \mathbf{V}_p(f) \log(1/\max \mathbf{I}_i^p(f))$ .*

We also rely on a fact that, for exponentially many boolean inputs, the probability that a monotone event happens jumps from being unlikely to being likely in a very small range of  $p$  values, i.e. of length approximately  $1/\log n$ . Think like the majority function. Once  $p > 1/2$ , the chance of a vote passing grows rapidly in  $p$ .

**Theorem 8.4** (Friedgut, Kalai). *If  $A$  is a monotone event, whose influences are the same for each index, and for  $p = p_0$ , if  $\mathbf{P}(A) > \varepsilon$ , then for  $p \geq p_0 + c \log(1/\varepsilon)/\log n$ ,  $\mathbf{P}(A) > 1 - \varepsilon$ .*

*Proof.* If all the influences are the same, we find the total variance is

$$\mathbf{I}^p(E) \gtrsim \mathbf{V}_p(\chi_E) \log n = \min(\mathbf{P}(E), 1 - \mathbf{P}(E)) \log n$$

Now the Margulis-Russo formula yields that if  $\mathbf{P}(E) \leq 1/2$ ,

$$\frac{d(\log \mathbf{P}(E))}{dp} = \frac{\mathbf{I}^p(E)}{\mathbf{P}(E)} \gtrsim \log n$$

Thus if  $p \geq p_0 + c/\log n$ ,  $\mathbf{P}(E) \geq 1/2$ . Now if  $\mathbf{P}(E) \geq 1/2$ , then

$$\frac{d(\log(1 - \mathbf{P}(E)))}{dp} = -\frac{\mathbf{I}^p(E)}{1 - \mathbf{P}(E)} \lesssim -\log n$$

In order to make  $\mathbf{P}(E) \geq 1 - \varepsilon$ , we need  $\log(1 - \mathbf{P}(E)) \leq \log \varepsilon$ . To move from  $\log(1/2)$  to  $\log \varepsilon$ , we need  $p \geq p_0 + c/\log n + c \log(1/\varepsilon)/\log n = p_0 + c \log(1/\varepsilon)/\log n$ .  $\square$

**Theorem 8.5** (FKG). *Any two increasing events are positively correlated.*

We now try and prove the critical exponent for the square lattice is  $1/2$ . First, we show  $\theta(1/2) = 0$ . Consider the ‘square annulus’ between  $4^n$  and  $3 \cdot 4^n$ , and let  $E_n$  be the event that there is a ring in this annulus. Because the edge sets involved in each event are disjoint, the events are independent. Furthermore, the probability that there is a cross on each side of the annulus is bounded below by some constant  $c > 0$ . Adding more edges doesn’t hurt the probability that an event happens, so they are monotonic, and the FKG inequality says that  $\mathbf{P}(E_n) = \mathbf{P}(A \cap B \cap C \cap D) \geq \mathbf{P}(A)\mathbf{P}(B)\mathbf{P}(C)\mathbf{P}(D) \geq c^4$ . Thus infinitely many occur almost surely, so  $\theta(1/2) = 0$ , and in fact, there is almost surely no infinite cluster anywhere.

To form a contradiction, we assume the critical point is  $1/2 + \delta$  instead of  $1/2$ . Given  $n$ , we form a  $2n \times n$  box, and we let  $J_n$  denote the event that there is a crossing in the box, i.e a path from left to right and an edge from bottom to top.

**Lemma 8.6.** *As  $n \rightarrow \infty$ ,  $\max \mathbf{I}_e^p(J_n) \rightarrow 0$ , where the maximum is taken over  $1/2 \leq p \leq 1/2 + \delta/2$  and all edges  $e$ .*

*Proof.* If  $e$  is pivotal on a given input for  $J_n$ , that means there is a path from a vertex adjacent to  $e$  to a vertex a distance  $n/2$  away. Thus by translation invariance,  $\mathbf{I}_e^p(J_n)$  is bounded by half the probability that there is a path of length  $n/2$  from the origin, which is uniformly bounded for  $p \leq 1/2 + \delta/2$ . As  $\theta(1/2 + \delta/2) = 0$ , these values must tend to zero as  $n \rightarrow \infty$ .  $\square$

**Lemma 8.7.** *For some  $n$ , and with  $p = 1/2 + \delta/2$ ,  $\mathbf{P}(J_n) \geq 0.98$ .*

*Proof.* We have already seen that  $\inf \mathbf{P}(J_n) > 0$  when  $p = 1/2$ . If  $\mathbf{P}(J_n) < 0.98$  for all  $n$ , then BKS shows that  $d\mathbf{P}(J_n)/dp$  must tend to  $\infty$  uniformly for all  $1/2 \leq p \leq 1/2 + \delta/2$  as  $n \rightarrow \infty$ , by setting  $\varepsilon < 0.02$ . And this means that the probability of  $J_n$  increases massively over the range of  $p$  from  $1/2$  to  $1/2 + \delta/2$ .  $\square$

Now we show this implies that we almost surely get an infinite cluster. If we can cross a  $2n \times n$  box with probability  $1 - \varepsilon$ , crossing lines gives a probability of  $1 - 5\varepsilon$  of cross a  $4n \times n$  box. Thus the probability that we cross a  $4n$  by  $2n$  box is greater than the probability that we cross the top and bottom of the graph. Thus

$$\begin{aligned} \mathbf{P}(\text{Top} \cup \text{Bottom}) &= \mathbf{P}(\text{Top}) + \mathbf{P}(\text{Bottom}) - \mathbf{P}(\text{Top})\mathbf{P}(\text{Bottom}) \\ &\geq 2(1 - 5\varepsilon) - (1 - 5\varepsilon)^2 \geq 1 - 5\varepsilon^2 \end{aligned}$$

If  $\varepsilon$  is small, this error is REALLY small. Thus almost surely all but finitely many of the  $J_n$  occur. Thus we just have to put these lines together in a way which guarantees

### 8.3 Conformal Invariance

Brownian motion is the limit of a random walk, and in the plane, is conformally invariant, in the sense that the image of a path of Brownian motion under an analytic map looks like the path of Brownian motion, up to a change in the measurement of time. Thus a random walk is conformally invariant ‘in the limit’. In some sense, percolation should also look ‘asymptotically’ conformally invariant in the limit.

Let us describe what this principle should look like when discretized. We consider a conformal map  $\phi$  from the unit disk to some other simply connected domain  $D$  fixing the origin, and with  $\phi'(0) > 0$ . Now for any  $\delta$ , we can consider the lattice  $\delta\mathbf{Z}^2$ , restricted to the interior of the unit disk. If  $C_\delta$  is the cluster around the origin in the interior. Similarly, we define  $C'_\delta$  to be the cluster around the origin in  $\delta\mathbf{Z}^2$  restricted to  $D$ . Now  $\phi(C_\delta)$  and  $C'_\delta$  don’t even lie on the same lattice, but as  $\delta \rightarrow 0$ , they should still asymptotically describe the same law on space.

The simplest precise statement of conformal invariance was proved by Smirnov in 2001. Scale the hexagonal percolation problem by  $\delta$  at the critical percolation value  $p = 1/2$ . If four points  $A, B, C$ , and  $D$  are chosen on the boundary of  $D$ , then the probability that there is a path from points on the boundary of  $D$  between  $A$  and  $B$  to points on the boundary of  $D$  between  $C$  and  $D$  converges as  $\delta \rightarrow 0$ . Furthermore, this convergent value is invariant under conformal mappings. In the case where  $D$  is a sidelength one equilateral triangle,  $A, B$ , and  $C$  are the three corner points, and  $D$  is on the line between  $A$  and  $C$  with distance  $x$  from  $C$ , the probability converges to  $x$ . By conformal invariance, this gives a general way to calculate the limiting probability. On  $\mathbf{Z}^2$ , even this statement is still open.

# Chapter 9

## High Dimensional Probability

**Theorem 9.1.** *Let  $f$  be  $L$  Lipschitz. Then*

$$\gamma(x \in \mathbf{R}^n : |f(x) - M| > t) \leq 2 \exp(-t^2/2L^2)$$

*Proof.* Let  $\mathbf{H}$  be a half space with Gaussian measure  $1/2$ , i.e. the standard upper half plane  $\mathbf{H} = \{x : x_1 \leq 0\}$ . Then  $\gamma(H_\varepsilon^c) = \mathbf{P}(N(0,1) \geq \varepsilon) \leq e^{-\varepsilon^2/2}$ . Thus  $\gamma(H_\varepsilon) \geq 1 - e^{-\varepsilon^2/2}$ .

Next, if  $A = \{x : f(x) \geq M\}$ , we obtain

$$\gamma(x : f(x) \geq M - L\varepsilon) \geq \gamma(A_\varepsilon) \geq \gamma(\mathbf{H}_\varepsilon) \geq 1 - e^{-\varepsilon^2/2}$$

□

**Theorem 9.2** (Isoperimetric Inequality in Gaussian Space). *Among all measurable sets  $A$  in  $\mathbf{R}^n$  with the same Gaussian measure, half spaces minimize the measure  $\gamma_n(A_\varepsilon)$ .*

*Remark.* We can replace  $M$  by  $\mathbf{E}f$  using sub Gaussian centering.

**Theorem 9.3.** *Let  $X$  be sub Gaussian. Then the centered version satisfies  $\|X - \mathbf{E}X\|_{\psi_2} \leq 2\|X\|_{\psi_2}$ .*

*Proof.* We find

$$\|X - \mathbf{E}\|_p \leq \|X\|_p + \|\mathbf{E}X\|_p$$

Now  $\mathbf{E}\|X\| \leq \|X\|_p$ .

□

The normal distribution concentrates on an annulus of radius  $\sqrt{d}$  and width  $O(1)$ . Also, it concentrates on half spaces  $\{X_1 \geq c\}$ .



# Bibliography

- [1] Larry Wasserman, *All of Statistics*
- [2] Walter Rudin, *Real and Complex Analysis*