

An introduction to the Concentration of Measure and Empirical Process Theory

Peadar Coyle

January 13, 2014

1 Introduction

In this thesis we will introduce some aspects from Empirical process theory and Concentration of Measure. We are interested in applications in Learning theory but do not have the time nor the scope to introduce model theory and VC-theory. Our aims are more modest, we want to introduce some concentration inequalities, show their applications in Empirical process theory and mention in passing some Statistical Learning theory. We will not be able to introduce all of these complex notions and we refer to the literature when we feel it is good to do so. Our first subsection will introduce Concentration of Measure and in passing we will see some of the celebrated theorems and inequalities. Elsewhere in thesis we will examine these theorems in more detail readers who are already familiar with the theorems - or perhaps Computer Scientists can just read the introduction and move to the end. I consider this thesis like a fine menu - take what you want from it, and don't try to digest everything at once. The style of this paper will be mostly aimed at pure mathematicians - so the standard of rigor is high. I have endeavoured to add my own opinions about what is difficult and what is easy, and elucidate some of these *unoriginal* and classical proofs. Most of what is reproduced here can be found in other textbooks or seminars, but the aim here is to condense as many of these as possible and rewrite them for an applied Statistician or Machine Learning expert audience.

Let us suppose that we have a large number of scalar random variables X_1, \dots, X_n , which have a bounded size on average (e.g. their mean and

variance would be $O(1)$). An interesting and important question is what can one say about their sum? $S = X_1 + \dots + X_n$? If each individual summand X_i varies in an interval of size $O(1)$, then their sum of course varies in an interval of size $O(n)$. However a remarkable phenomenon, known as *concentration of measure*, asserts that assuming a sufficient amount of independence between the component variables X_1, \dots, X_n , this sum sharply concentrates to a much narrower range, typically an interval of size $O(\sqrt{n})$. This phenomenon is quantified by a variety of *large deviation inequalities*¹ that give upper bounds (often exponential in nature) on the probability that such a combined random variable deviates significantly from its mean. The same phenomenon applies not only to linear expressions such as $S = X_1 + \dots + X_n$, but more generally to nonlinear combinations $F(X_1, \dots, X_n)$ of such variables, provided that the nonlinear function F is sufficiently regular - for example Lipschitz.

The basic intuition is that independent random variables find it difficult to "work together" to simultaneously pull a sum $S = X_1 + \dots + X_n$ or a more general combination $F(X_1, \dots, X_n)$ too far from its' mean. Independence is the key here; concentration of measure results typically fail if the X_i are too highly concentrated to each other. Although such results do exist see for example the work by Kontorovich on Mixing Phenomena [3] or [5] for work using time series data.

There are many applications of concentration of the concentration of measure phenomenon, such as random matrix theory, but we will mostly focus on applications in Computer Science not Physics². We wholeheartedly recommend the excellent book by Terry Tao on Random Matrix Theory [8].

Assuming that one has a sufficient amount of independence, the concentration of measure tends to be sub-gaussian in nature, this probability that one is at least λ standard deviations (s.d.) from the mean tends to drop off like $C \exp(-c\lambda^2)$ for some $C, c > 0$. In particular, one is $O(\log \frac{1}{\epsilon} n)$ standard

¹Roughly speaking, large deviations theory concerns itself with the exponential decline of probability measures of certain kinds of extreme or tail events as the number of observations grows large see [9, 12] for further details. Readers interested in a thorough introduction can read from the "*Mozart*" of large deviations theory [11]

²Ironic perhaps, since the author once studied in a Physics department

deviations from the mean with overwhelming probability. Indeed, concentration of measure is our primary tool for ensuring that various events hold with overwhelming probability.

1.1 What is in this thesis?

In section 2 we give a rapid overview of concentration of measure including some of the more advanced methods like the Herbst argument which will not be used in the thesis. In section 2.1 we will introduce some of the language from applications in statistical learning theory including in particular empirical process theory. In the following section we include some remarks about Machine Learning and Artificial Intelligence, this is just to motivate the thesis 2.2.

In 4 we introduce the major concentration inequalities which we shall use - we start with Hoeffding and then introduce a few variants and finally the celebrated McDiarmid Inequality[?] - which we will use to prove results in the applications section. In 5 we return to our applications - and give a more rigorous and complete introduction to our chosen subfield of Statistical Learning Theory - this will be of particular interest to Mathematicians who may not be familiar with the language. In ?? we apply the concentration inequalities in a applied setting. The remainder of the 5 introduces the celebrated Vapnik-Chervonenkis dimension. In 6 we introduce the meat of the applications through the technique of Rademacher averages, we end the section by elucidating on the difficulty of calculating such averages - in learning theory through the lens of empirical risk. In ?? we end the thesis with some discussion and bibliographic remarks.

2 What are other examples of Concentration of Measure?

We've introduced this loosely in the introduction. We shall introduce this in more detail. Let X be a random variable taking values in some metric space \mathbf{X} . Then we say that X (or, more correctly, the distribution of X) has the concentration property if, for any set $A \subset \mathbf{X}$ such that $\mathbb{P}(X \in A) \geq 1/2$, we have

$$\mathbb{P}(d(X, A) \leq r) \xrightarrow{r \rightarrow \infty} 1. \quad (1)$$

Here, $d(x, A)$ is the distance from the point $x \in \mathsf{X}$ to the set A :

$$d(x, A) := \inf_{y \in A} d(x, y).$$

Another way to express 1 is as follows: for any set $A \subset \mathsf{X}$ and any $r \geq 0$, define the r -blowup of A by

$$A_r := \{y \in \mathsf{X} : d(y, A) \leq r\} \equiv \{y \in \mathsf{X} : \exists x \in A \text{ such that } d(x, y) \leq r\}.$$

Then X has the concentration property if

$$\mathbb{P}(X \in A) \geq 1/2 \implies \lim_{r \rightarrow \infty} \mathbb{P}(X \in A_r) = 1.$$

In other words, X has the concentration property if any set containing X with not too small a probability can be blown up to contain X with near-certainty.

Here are two classic examples of concentration:

- **Gaussian distribution in Euclidean space.** Let $\mathsf{X} = \mathbb{R}^n$ and take $d(x, y) = \|x - y\|_2$ the usual Euclidean distance. Let X be a standard n -dimensional Gaussian random variable, i.e., $X \sim N(0, I_n)$, where I_n is the $n \times n$ identity matrix. Then for any $r \geq 0$ we have

$$\mathbb{P}(X \in A) \geq 1/2 \implies \mathbb{P}(X \in A_r) \geq 1 - e^{-r^2/2}.$$

- **Uniform distribution in Hamming space.** Let X be the Hamming cube $\{0, 1\}^n$ equipped with the normalized Hamming distance

$$d(x, y) = \frac{1}{n} \sum_{i=1}^n 1_{\{x_i \neq y_i\}}$$

that counts the fraction of bits in which $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ disagree. Let X have the uniform distribution on $\{0, 1\}^n$, i.e., $\mathbb{P}(X = x) = 2^{-n}$ for all x . Then

$$\mathbb{P}(X \in A) \geq 1/2 \implies \mathbb{P}(X \in A_r) \geq 1 - e^{-2nr^2}.$$

These two examples suggest that we should aim for hard statements in the form of sharp bounds on the concentration function

$$\alpha_X(r) := \sup_{A: \mathbb{P}(X \in A) \geq 1/2} \mathbb{P}(X \notin A_r)$$

as opposed to merely soft statements of the form $\alpha_X(r) \rightarrow 0$ as $r \rightarrow \infty$. The 64,000 pounds question is: how do we get such bounds?

There are two ways to accomplish this goal, and the main idea underlying these two ways is to replace sets with some other objects that are hopefully easier to handle. The first way is to replace sets by probability measures, the second is to replace them by functions. Here is what I mean:

Fix a set $A \subset \mathbf{X}$ with $\Pr(X \in A) > 0$. Let \mathbb{P} denote the distribution of X , and let P_A denote the conditional distribution of X given $X \in A$. That is, for any (measurable) set $B \subset \mathbf{X}$ we have

$$P_A(B) := \frac{P(A \cap B)}{P(A)}.$$

I am using the subscript notation \mathbb{P}_A instead of the more usual $\mathbb{P}(\cdot|A)$ to indicate the fact that \mathbb{P}_A is a probability measure in its own right. In this way, we can associate to each non-null set A a probability measure P_A . Now, here is a very simple observation that turns out to be very consequential:

$$D(P_A \| P) = \log \frac{1}{P(A)}. \quad (2)$$

This is very easy to prove: for any set B we have

$$P_A(B) = \frac{1}{P(A)} \int_B 1_A(x) P(dx), \quad (3)$$

so \mathbb{P}_A is absolutely continuous with respect to \mathbb{P} with the Radon-Nikodym derivative

$$dP_A/dP = 1_A/P(A)$$

. Therefore, by definition of the divergence,

$$D(P_A \| P) = \int dP_A \log \frac{dP_A}{dP} = \frac{1}{P(A)} \int_A dP \log \frac{1}{P(A)} = \log \frac{1}{P(A)}.$$

So if we are interested in bounding the probabilities of various sets A , we may hope to get some mileage out of the relationship (??).

On the other hand, we may also associate to a set A with $P(A) > 0$ the function

$$f_A(x) := d(x, A) \equiv \inf_{y \in A} d(x, y).$$

This function is Lipschitz: for any $x, x' \in \mathbf{X}$,

$$f_A(x) - f_A(x') = \inf_{y \in A} d(x, y) - \inf_{y \in A} d(x', y) \leq \sup_{y \in A} [d(x, y) - d(x', y)] \leq d(x, x'),$$

where the last step is by the triangle inequality. Interchanging the roles of x and x' , we get the Lipschitz property. Moreover, let us consider the random variable $Z = f_A(X)$, where X is our original \mathbf{X} -valued random variable. Then we immediately notice two things:

For any $r \geq 0$, $\mathbb{P}(Z \leq r) = \mathbb{P}(d(X, A) \leq r) = \mathbb{P}(A_r)$. If $P(A) = \mathbb{P}(X \in A) \geq 1/2$, then 0 is a median of Z , in the sense that

$$\mathbb{P}(Z \leq 0) = P(A) \geq 1/2 \quad \text{and} \quad \mathbb{P}(Z > 0) \geq 1/2$$

(the second inequality is obviously true since Z is nonnegative with probability 1).

These two observations suggest that we may obtain concentration bounds by bounding the probability that a given Lipschitz function of X deviates from its median by more than r . In fact, it is easy to derive an alternative expression for the concentration function α_X :

$$\alpha_X(r) = \sup_{1\text{-Lipschitz } f} \mathbb{P}(f(X) > m_f + r), \quad (4)$$

where m_f denotes any median of $f(X)$. We already showed, by passing from A to $f_A = d(\cdot, A)$, that α_X is bounded from above by the quantity on the right-hand side of ((??)):

$$\alpha_X(r) = \sup_{A: P(A) \geq 1/2} \mathbb{P}(f_A(X) > \underbrace{m_{f_A}}_{=0} + r) \leq \sup_{1\text{-Lipschitz } f} \mathbb{P}(f(X) > m_f + r)$$

To prove the reverse inequality, fix any 1-Lipschitz function f and consider the set $A_f := \{x \in \mathbf{X} : f(x) \leq m_f\}$, where m_f is any median of f . Then, by definition,

$$\mathbb{P}(X \in A_f) = \mathbb{P}(f(X) \leq m_f) \geq 1/2.$$

Moreover, if we consider the r -blowup

$$[A_f]_r = \left\{ x \in \mathbf{X} : d(x, A_f) \leq r \right\},$$

then for any $x \in \mathsf{X}$ and any $y \in [A_f]_r$ we must have

$$f(x) - m_f \leq f(x) - f(y) \leq d(x, y),$$

where the last step is by the Lipschitz property of f . Consequently, by definition of the concentration function,

$$\mathbb{P}\left(f(X) > m_f + r\right) \leq \mathbb{P}\left(d(X, A_f) > r\right) = 1 - P([A_f]_r) \leq \alpha_X(r).$$

By passing to the functional viewpoint, we obtain another equivalent characterization of the concentration property: a random variable X taking values in a metric space (X, d) has the concentration property if real-valued Lipschitz functions X are nearly constant.

Lets look at the first, probabilistic viewpoint, which was born out of a 1996 breakthrough paper by Marton. Given a metric space (X, d) , let us define the L_1 Wasserstein distance (or transportation distance) between any two probability measures P and Q on it:

$$W_1(P, Q) := \inf_{X \sim P, Y \sim Q} \mathbb{E}[d(X, Y)],$$

where the infimum is over all jointly distributed random variables $X, Y \in \mathsf{X}$, such that $P_X = P$ and $P_Y = Q$. Now consider a random variable $X \in \mathsf{X}$, for which we wish to establish concentration. What Marton showed is the following: Suppose the distribution P of X satisfies the L_1 transportation inequality

$$W_1(P, Q) \leq \sqrt{2c D(Q \| P)} \tag{5}$$

for some constant $c > 0$. Then X has the concentration property, and moreover

$$P(A) \geq 1/2 \implies P(A_r) \geq 1 - \exp\left(-\frac{1}{2c} \left(r - \sqrt{2c \log 2}\right)^2\right), \forall r > \sqrt{2c \log 2}.$$

Martons proof is breathtakingly beautiful. Consider any two sets A, B with $P(A), P(B) \neq 0$. Recalling our notation for conditional distributions, we can write

$$\begin{aligned} W_1(P_A, P_B) &\leq W_1(P_A, P) + W_1(P_B, P) \\ &\leq \sqrt{2c D(P_A \| P)} + \sqrt{2c D(P_B \| P)} \\ &= \sqrt{2c \log \frac{1}{P(A)}} + \sqrt{2c \log \frac{1}{P(B)}}, \end{aligned}$$

where in the first step we have used the triangle inequality, in the second we have used the fact that P satisfies the transportation inequality (5), and in the last step we have used the formula ((?)). Now suppose that $P(A) \geq 1/2$ and let $B = A_r^c$ for some r , where c denotes set-theoretic complement. Then we can show that $W_1(P_A, P_B) \geq d(A, B) \geq r$. On the other hand,

$$\log \frac{1}{P(A)} \leq \log 2 \quad \text{and} \quad \log \frac{1}{P(B)} = \log \frac{1}{1 - P(A_r)}.$$

Combining these facts gives us the bound

$$r \leq \sqrt{2c \log 2} + \sqrt{2c \log \frac{1}{1 - P(A_r)}}$$

that holds for all r . If $r > \sqrt{2c \log 2}$, then we get

$$P(A_r) \geq 1 - \exp \left(-\frac{1}{2c} \left(r - \sqrt{2c \log 2} \right)^2 \right),$$

so we indeed have concentration and a sharp bound on $\alpha_X(r)$, at least for large enough values of r . The main message here is that, in order to study concentration, it suffices to work on the level of probability measures and to focus ones effort on showing that the distribution of X satisfies a suitable transportation inequality. Since Martons original work, there have been many refinements and extensions, which I will not go into here. One such result, due to Sergey Bobkov and Friedrich Götze, says that P satisfying a transportation inequality ((?)) is equivalent to the Gaussian concentration property

$$\alpha_X(r) \leq e^{-r^2/2c}, \quad \forall r \geq 0.$$

Now lets look at the complementary functional viewpoint. Recall that we seek tight upper bounds on deviation probabilities of the form

$$\mathbb{P} \left(f(X) \geq m_f + r \right), \quad \forall r > 0.$$

It is easier to work with means instead of medians, and indeed it can be shown that concentration around the mean is equivalent to concentration around any median. So lets focus on the mean. Let X , as before, be a random variable over some metric space (\mathbf{X}, d) , and consider a Lipschitz function $f : \mathbf{X} \rightarrow \mathbb{R}$ such that $\mathbb{E}[f(X)] = 0$. We can apply the well-known Chernoff trick: for any $r, \lambda > 0$ we have

$$\mathbb{P} \left(f(X) \geq r \right) = \mathbb{P} \left(e^{\lambda f(X)} \geq e^{\lambda r} \right) \leq e^{-\lambda r} \mathbb{E}[e^{\lambda f(X)}].$$

Now the whole affair hinges on the availability of tight upper bounds on the logarithmic moment-generating function $\Lambda(\lambda) := \log \mathbb{E}[e^{\lambda f(X)}]$. The entropy method, which was the subject of Gabor Lugosis famous work, is the name for a set of techniques for bounding $\Lambda(\lambda)$ by means of various inequalities involving the relative entropy between various tilted distributions derived from P and P itself. The entropy method has roots in the work of Michel Ledoux, who in turn distilled it from some very deep results of Michel Talagrand. The entropy method will be further examined in this thesis and proofs for the major 'named' theorems will be given.

The simplest version of the entropy method goes something like this. Let us define, for any $t \in \mathbb{R}$, the tilted distribution $P^{(t)}$ via

$$\frac{dP^{(t)}}{dP}(x) = \frac{e^{tf(x)}}{e^{\Lambda(t)}}$$

(assuming, of course, that $\Lambda(t)$ exists and is finite). Then we have

$$\begin{aligned} D(P^{(t)} \| P) &= \int dP^{(t)} [tf - \Lambda(t)] \\ &= \frac{1}{e^{\Lambda(t)}} t \mathbb{E} [f(X) e^{tf(X)}] - \Lambda(t) \\ &= t\Lambda'(t) - \Lambda(t) \\ &= t^2 \frac{d}{dt} \left(\frac{\Lambda(t)}{t} \right). \end{aligned}$$

Integrating and using the fact that $\Lambda(0) = 0$, we get

$$\Lambda(\lambda) = \lambda \int_0^\lambda \frac{D(P^{(t)} \| P)}{t^2} dt. \quad (6)$$

Now suppose that we can bound $D(P^{(t)} \| P) \leq ct^2/2$ for some $c > 0$. Then from ((?)) we have

$$\Lambda(\lambda) \leq \frac{c\lambda^2}{2}, \quad \forall t$$

which in turn gives the Gaussian tail bound

$$\mathbb{P}(f(X) \geq r) \leq \inf_{\lambda > 0} \exp \left(-\lambda r + \frac{c\lambda^2}{2} \right) = \exp \left(-\frac{r^2}{2c} \right).$$

This is the so-called Herbst argument. Of course, I have glossed over the most nontrivial part namely, showing that we can bound the relative entropy $D(P^{(t)} \| P)$ by a quadratic function of t . We shall not delve much more here, but later in the thesis we shall consider some of these bounds. I refer to [?] for further details

Here is one classic example. Suppose that our function f has the bounded difference property, i.e., there exist some constants $c_1, \dots, c_n \geq 0$, such that

changing the i -th argument of f while keeping others constant will change the value of f by at most c_i :

$$\sup_{x_1, \dots, x_n} \sup_{x'_i} |f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i.$$

We can express this more succinctly as a Lipschitz property of f if we define the weighted Hamming metric

$$d(x, y) = \sum_{i=1}^n c_i 1_{\{x_i \neq y_i\}} \quad (7)$$

(we can assume without loss of generality that all the c_i s are strictly positive, because we can simply ignore those coordinates of x that do not affect the value of f). Then the bounded difference property is equivalent to saying that f is 1-Lipschitz with respect to this weighted Hamming metric. Moreover, it is possible to show that any product probability measure $P = P_1 \otimes \dots \otimes P_n$ on the product space $\mathbf{X} = \mathcal{X}^n$ satisfies the transportation inequality

$$W_1(Q, P) \leq \sqrt{2c D(Q \| P)},$$

where the Wasserstein distance is computed with respect to the weighted Hamming metric (7), and $c = \frac{1}{4} \sum_{i=1}^n c_i^2$. By the Bobkov-Götze result quoted above, this is equivalent to the concentration bound

$$\mathbb{P}(f(X) \geq \mathbb{E}[f(X)] + r) \leq \exp\left(-\frac{r^2}{2c}\right) = \exp\left(-\frac{2r^2}{\sum_{i=1}^n c_i^2}\right).$$

This is the well-known McDiarmid's inequality - which we shall use a lot in this thesis. What is important to state is that this inequality was originally derived using martingale techniques, but here we have arrived at it through a completely different route that took us back to where we started: concentration of Lipschitz functions around their means and/or medians, which (as we saw) is the same thing as the original, stochastic-geometric view of the concentration of measure phenomenon. This *equivalence* between *concentration of Lipschitz functions* and *"stochastic-geometric"* views of the concentration of measure phenomenon is a very important general phenomenon, and important for proving new theorems or finding new insights in Probability Theory or information theory. For a thorough introduction to the uses of Concentration of Measure in Information theory we recommend Raginsky [7]

2.1 What is Empirical Process theory

The simplest example of an empirical process arises when trying to estimate a probability distribution from sample data. For those interested in applications it is worth highlighting that the *data* could be from Economics - say we wanted to model when the next recession would occur based on historical time-series data [5] or it could be data from a social media website. The difference between the empirical distribution function $F_n(x)$ and the true distribution function $F(x)$ converges to zero everywhere (by the law of large numbers), and this is non-trivial the maximum difference between the empirical and true distribution functions converges to zero, too (by the Glivenko-Cantelli theorem, a uniform law of large numbers). The "empirical process" $E_n(x)$ is the re-scaled difference, $n^{1/2}[F_n(x) - F(x)]$, and it converges to a Gaussian stochastic process that only depends on the true distribution (by the functional central limit theorem).

Empirical process theory is concerned with generalizing this sort of material to other stochastic processes determined by random samples, and indexed by infinite classes (like the real line, or the class of all Borel sets on the line, or some space parameterizing a regression model). The typical objects of concern are proving uniform limit theorems, and with establishing distributional limits. (For instance, one might want to prove that the errors of all possible regression models in some class will come close to their expected errors, so that maximum-likelihood or least-squares estimation is consistent.) This endeavor is closely linked to Vapnik-Chervonenkis-style learning theory, and in fact one can see VC theory as an application of empirical process theory. So it is very difficult to untangle Vapnik-Chervonenkis(VC)-style learning theory from Concentration of Measure and Empirical Process theory. Hence throughout this thesis you will see references to VC theory, and some of the celebrated theorems proven. It is worth highlighting that most of these areas are still active research areas - and you may encounter some 'simple' questions which are not easy to prove, or are in fact open questions. But then we do live in the *Golden age* of statistics and stochastic processes! Readers interested in more about Empirical Process theory should read [10, 2, 3, 1] these all introduce the topic in more depth.

2.2 So why should we care about empirical process theory?

Well let's say that you the reader are statistically literate, and you've learned some concentration of measure, and you have an interest in model selection - that is selecting among different mathematical model which all purport to describe the same data set. The basic strategy is to find conditions under which every model in a reasonable class will, with high probability, perform about as well on sample data as they can be expected to do on new data; this involves constraining the richness or flexibility of the model class.

One of the *buzzwords* in contemporary science is Machine Learning. We shall consider statistical learning a subtype of machine learning and leave the reader to form their own impressions of what machine learning is. Let us use the Vapnik notion of statistical learning for pedagogical reasons. We want to estimate some functional which depends on an unknown distribution over a probability space X —it could be a *concept*³, regression coefficients, moments of the distribution, Shannon entropy, etc., even the distribution itself. We have a class of admissible distributions, called hypotheses, and a “loss functional,” an integral over X which tells us, for each hypothesis, how upset we should be when we guess wrong; this implicitly depends on the true distribution. Clearly we want the best hypothesis, the one which minimizes the loss functional — but to explicitly calculate that we'd need to know the true distribution. Vapnik assumes that we have access to a sequence of independent random variables, all drawn from the (stationary) true distribution. What then are we to do?

One answer which I term the *Vapnik answer* takes two parts. The first has to do with “empirical risk minimization”: approximate the true, but unknown, loss functional, which is an integral over the whole space X , with a sum over the observed data-points, and go with the hypothesis that minimizes this “empirical risk”; call this, though Vapnik doesn't, the ERM hypothesis. It's possible that the ERM hypothesis will do badly in the future, because we blundered into unrepresentative data, but we can show necessary and sufficient conditions for the loss of the ERM hypothesis to converge in probability to the loss of the best hypothesis. Moreover, we can prove that under certain very broad conditions, that if we just collect enough data-points, then the loss of the ERM hypothesis is, with high probability, within a certain additive distance (“confidence interval” — Vapnik's scare-quotes) of the loss

³In the machine learning sense

of the best hypothesis. These conditions involve the Vapnik-Chervonenkis dimension, and a related quantity called the Vapnik-Chervonenkis entropy. Very remarkably, we can even calculate how much data we need to get a given approximation, at a given level of confidence, regardless of what the true distribution is, i.e. we can calculate distribution-independent bounds. (They do, however, depend on the nature of the integrands in the loss functional.)

These results about convergence, approximation, etc. are in essence extensions of the Law of Large Numbers to spaces of functions. As such the assumption that successive data-points are independent and identically distributed is key to the whole exercise. While it is possible to consider dependant data in Statistical learning we shall not consider it in this thesis. The second part of Vapnik's procedure is an elaboration of the first: For a given amount of data, we pick the hypothesis which minimizes the sum of the empirical risk and the "confidence interval" about it. This is termed by Vapnik - *structural risk minimization* and shall not be considered in this thesis. I recommend [?] for further details.

3 Notation

We introduce some notation that is used throughout the paper. We assume that X_1, \dots, X_n are independent random variables taking values in a measurable space \mathcal{X} . Denote by X_1^n the vector of these n random variables. Let $f : \mathcal{X}^n \rightarrow \mathbb{R}$ be some measurable function.

4 Basic concentration inequalities via the martingale approach

In the following section, some basic inequalities that are widely used for proving concentration inequalities are presented, whose derivation relies on the martingale approach. Their proofs convey the main concepts of the martingale approach for proving concentration. Their presentation also motivates some further refinements that are considered in the continuation of this chapter.

4.1 The Azuma-Hoeffding inequality

The Azuma-Hoeffding inequality⁴ is a useful concentration inequality for bounded-difference martingales. It was proved in [?] for independent bounded random variables, followed by a discussion on sums of dependent random variables; this inequality was later derived in [?] for the more general setting of bounded-difference martingales. In the following, this inequality is introduced.

Theorem 4.1. [Azuma-Hoeffding inequality] *Let $\{X_k, \mathcal{F}_k\}_{k=0}^n$ be a discrete-parameter real-valued martingale sequence. Suppose that, for every $k \in \{1, \dots, n\}$, the condition $|X_k - X_{k-1}| \leq d_k$ holds a.s. for a real-valued sequence $\{d_k\}_{k=1}^n$ of non-negative numbers. Then, for every $\alpha > 0$,*

$$\mathbb{P}(|X_n - X_0| \geq \alpha) \leq 2 \exp \left(-\frac{\alpha^2}{2 \sum_{k=1}^n d_k^2} \right). \quad (8)$$

It is noted that (??) is typically interpreted as $\frac{X_n - X_0}{\sqrt{n}}$ being sub-Gaussian. The proof of the Azuma-Hoeffding inequality serves also to present the basic principles on which the martingale approach for proving concentration results is based. Therefore, we present in the following the proof of this inequality.

Proof. For an arbitrary $\alpha > 0$,

$$\mathbb{P}(|X_n - X_0| \geq \alpha) = \mathbb{P}(X_n - X_0 \geq \alpha) + \mathbb{P}(X_n - X_0 \leq -\alpha). \quad (9)$$

Let $\xi_i \triangleq X_i - X_{i-1}$ for $i = 1, \dots, n$ designate the jumps of the martingale sequence. Then, it follows by assumption that $|\xi_k| \leq d_k$ and $\mathbb{E}[\xi_k | \mathcal{F}_{k-1}] = 0$ a.s. for every $k \in \{1, \dots, n\}$.

⁴The Azuma-Hoeffding inequality is also known as Azuma's inequality. Since it is referred numerous times in this chapter, it will be named Azuma's inequality for the sake of brevity.

From Chernoff's inequality,

$$\begin{aligned}
& \mathbb{P}(X_n - X_0 \geq \alpha) \\
&= \mathbb{P}\left(\sum_{i=1}^n \xi_i \geq \alpha\right) \\
&\leq e^{-\alpha t} \mathbb{E}\left[\exp\left(t \sum_{i=1}^n \xi_i\right)\right], \quad \forall t \geq 0.
\end{aligned} \tag{10}$$

Furthermore,

$$\begin{aligned}
& \mathbb{E}\left[\exp\left(t \sum_{k=1}^n \xi_k\right)\right] \\
&= \mathbb{E}\left[\mathbb{E}\left[\exp\left(t \sum_{k=1}^n \xi_k\right) \mid \mathcal{F}_{n-1}\right]\right] \\
&= \mathbb{E}\left[\exp\left(t \sum_{k=1}^{n-1} \xi_k\right) \mathbb{E}[\exp(t\xi_n) \mid \mathcal{F}_{n-1}]\right]
\end{aligned} \tag{11}$$

where the last equality holds since $Y \triangleq \exp(t \sum_{k=1}^{n-1} \xi_k)$ is \mathcal{F}_{n-1} -measurable; this holds due to fact that $\xi_k \triangleq X_k - X_{k-1}$ is \mathcal{F}_k -measurable for every $k \in \mathbb{N}$, and $\mathcal{F}_k \subseteq \mathcal{F}_{n-1}$ for $0 \leq k \leq n-1$ since $\{\mathcal{F}_k\}_{k=0}^n$ is a filtration. Hence, the RV $\sum_{k=1}^{n-1} \xi_k$ and Y are both \mathcal{F}_{n-1} -measurable, and $\mathbb{E}[XY \mid \mathcal{F}_{n-1}] = Y \mathbb{E}[X \mid \mathcal{F}_{n-1}]$.

Due to the convexity of the exponential function, the straight line connecting the end points of the function over the interval $[-d_k, d_k]$ lies above this function. Since $|\xi_k| \leq d_k$ for every k (note that $\mathbb{E}[\xi_k \mid \mathcal{F}_{k-1}] = 0$), it follows that

$$\begin{aligned}
& \mathbb{E}[e^{t\xi_k} \mid \mathcal{F}_{k-1}] \\
&\leq \mathbb{E}\left[\frac{(d_k + \xi_k)e^{td_k} + (d_k - \xi_k)e^{-td_k}}{2d_k} \mid \mathcal{F}_{k-1}\right] \\
&= \frac{1}{2}(e^{td_k} + e^{-td_k}) \\
&= \cosh(td_k).
\end{aligned} \tag{12}$$

Since, for every integer $m \geq 0$,

$$(2m)! \geq (2m)(2m-2) \dots 2 = 2^m m!$$

then, due to the power series expansions of the hyperbolic cosine and exponential functions,

$$\cosh(td_k) = \sum_{m=0}^{\infty} \frac{(td_k)^{2m}}{(2m)!} \leq \sum_{m=0}^{\infty} \frac{(td_k)^{2m}}{2^m m!} = e^{\frac{t^2 d_k^2}{2}}$$

which therefore implies that

$$\mathbb{E}[e^{t\xi_k} \mid \mathcal{F}_{k-1}] \leq e^{\frac{t^2 d_k^2}{2}}.$$

Consequently, by repeatedly using the recursion in (??), it follows that

$$\mathbb{E}\left[\exp\left(t \sum_{k=1}^n \xi_k\right)\right] \leq \prod_{k=1}^n \exp\left(\frac{t^2 d_k^2}{2}\right) = \exp\left(\frac{t^2}{2} \sum_{k=1}^n d_k^2\right)$$

which then gives (see (??)) that

$$\mathbb{P}(X_n - X_0 \geq \alpha) \leq \exp\left(-\alpha t + \frac{t^2}{2} \sum_{k=1}^n d_k^2\right), \quad \forall t \geq 0.$$

An optimization over the free parameter $t \geq 0$ gives that $t = \alpha (\sum_{k=1}^n d_k^2)^{-1}$, and

$$\mathbb{P}(X_n - X_0 \geq \alpha) \leq \exp\left(-\frac{\alpha^2}{2 \sum_{k=1}^n d_k^2}\right). \quad (13)$$

Since, by assumption, $\{X_k, \mathcal{F}_k\}$ is a martingale with bounded jumps, so is $\{-X_k, \mathcal{F}_k\}$ (with the same bounds on its jumps). This implies that the same bound is also valid for the probability $\mathbb{P}(X_n - X_0 \leq -\alpha)$ and together with (??) it completes the proof of Theorem 4.1. \blacksquare

The proof of this inequality will be revisited later in this chapter for the derivation of some refined versions, whose use and advantage will be also exemplified.

Remark 4.2. In [4, Theorem 3.13], Azuma's inequality is stated as follows: Let $\{Y_k, \mathcal{F}_k\}_{k=0}^n$ be a martingale-difference sequence with $Y_0 = 0$

(i.e., Y_k is \mathcal{F}_k -measurable, $\mathbb{E}[|Y_k|] < \infty$ and $\mathbb{E}[Y_k|\mathcal{F}_{k-1}] = 0$ a.s. for every $k \in \{1, \dots, n\}$). Assume that, for every k , there exist some numbers $a_k, b_k \in \mathbb{R}$ such that a.s. $a_k \leq Y_k \leq b_k$. Then, for every $r \geq 0$,

$$\mathbb{P}\left(\left|\sum_{k=1}^n Y_k\right| \geq r\right) \leq 2 \exp\left(-\frac{2r^2}{\sum_{k=1}^n (b_k - a_k)^2}\right). \quad (14)$$

As a consequence of this inequality, consider a discrete-parameter real-valued martingale sequence $\{X_k, \mathcal{F}_k\}_{k=0}^n$ where $a_k \leq X_k - X_{k-1} \leq b_k$ a.s. for every k . Let $Y_k \triangleq X_k - X_{k-1}$ for every $k \in \{1, \dots, n\}$, so since $\{Y_k, \mathcal{F}_k\}_{k=0}^n$ is a martingale-difference sequence and $\sum_{k=1}^n Y_k = X_n - X_0$, then

$$\mathbb{P}(|X_n - X_0| \geq r) \leq 2 \exp\left(-\frac{2r^2}{\sum_{k=1}^n (b_k - a_k)^2}\right), \quad \forall r > 0. \quad (15)$$

Example 4.3. Let $\{Y_i\}_{i=0}^\infty$ be i.i.d. binary random variables which get the values $\pm d$, for some constant $d > 0$, with equal probability. Let $X_k = \sum_{i=0}^k Y_i$ for $k \in \{0, 1, \dots\}$, and define the natural filtration $\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \mathcal{F}_2 \dots$ where

$$\mathcal{F}_k = \sigma(Y_0, \dots, Y_k), \quad \forall k \in \{0, 1, \dots\}$$

is the σ -algebra that is generated by the random variables Y_0, \dots, Y_k . Note that $\{X_k, \mathcal{F}_k\}_{k=0}^\infty$ is a martingale sequence, and (a.s.) $|X_k - X_{k-1}| = |Y_k| = d$, $\forall k \in \mathbb{N}$. It therefore follows from Azuma's inequality that

$$\mathbb{P}(|X_n - X_0| \geq \alpha\sqrt{n}) \leq 2 \exp\left(-\frac{\alpha^2}{2d^2}\right). \quad (16)$$

for every $\alpha \geq 0$ and $n \in \mathbb{N}$. From the central limit theorem (CLT), since the RVs $\{Y_i\}_{i=0}^\infty$ are i.i.d. with zero mean and variance d^2 , then $\frac{1}{\sqrt{n}}(X_n - X_0) = \frac{1}{\sqrt{n}} \sum_{k=1}^n Y_k$ converges in distribution to $\mathcal{N}(0, d^2)$. Therefore, for every $\alpha \geq 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X_0| \geq \alpha\sqrt{n}) = 2Q\left(\frac{\alpha}{d}\right) \quad (17)$$

where

$$Q(x) \triangleq \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp\left(-\frac{t^2}{2}\right) dt, \quad \forall x \in \mathbb{R} \quad (18)$$

is the probability that a zero-mean and unit-variance Gaussian RV is larger than x . Since the following exponential upper and lower bounds on the Q-function hold

$$\frac{1}{\sqrt{2\pi}} \frac{x}{1+x^2} \cdot e^{-\frac{x^2}{2}} < Q(x) < \frac{1}{\sqrt{2\pi} x} \cdot e^{-\frac{x^2}{2}}, \quad \forall x > 0 \quad (19)$$

then it follows from (??) that the exponent on the right-hand side of (??) is the exact exponent in this example.

Example 4.4. In continuation to Example 4.3, let $\gamma \in (0, 1]$, and let us generalize this example by considering the case where the i.i.d. binary RVs $\{Y_i\}_{i=0}^\infty$ have the probability law

$$\mathbb{P}(Y_i = +d) = \frac{\gamma}{1+\gamma}, \quad \mathbb{P}(Y_i = -\gamma d) = \frac{1}{1+\gamma}.$$

Hence, it follows that the i.i.d. RVs $\{Y_i\}$ have zero mean and variance $\sigma^2 = \gamma d^2$. Let $\{X_k, \mathcal{F}_k\}_{k=0}^\infty$ be defined similarly to Example 4.3, so that it forms a martingale sequence. Based on the CLT, $\frac{1}{\sqrt{n}}(X_n - X_0) = \frac{1}{\sqrt{n}} \sum_{k=1}^n Y_k$ converges weakly to $\mathcal{N}(0, \gamma d^2)$, so for every $\alpha \geq 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X_0| \geq \alpha \sqrt{n}) = 2Q\left(\frac{\alpha}{\sqrt{\gamma} d}\right). \quad (20)$$

From the exponential upper and lower bounds of the Q-function in (??), the right-hand side of (??) scales exponentially like $e^{-\frac{\alpha^2}{2\gamma d^2}}$. Hence, the exponent in this example is improved by a factor $\frac{1}{\gamma}$ as compared Azuma's inequality (that is the same as in Example 4.3 since $|X_k - X_{k-1}| \leq d$ for every $k \in \mathbb{N}$). This indicates on the possible refinement of Azuma's inequality by introducing an additional constraint on the second moment. This route was studied extensively in the probability literature, and it is the focus of Section ??.

4.2 McDiarmid's inequality

The following useful inequality is due to McDiarmid ([?, Theorem 3.1] or [4]), and its original derivation uses the martingale approach for its derivation. We

will relate, in the following, the derivation of this inequality to the derivation of the Azuma-Hoeffding inequality (see the preceding subsection).

Theorem 4.5. [McDiarmid's inequality] *Let $\{X_k\}_{k=1}^n$ be independent real-valued random variables, taking values in the set $\mathcal{X} := \prod_{k=1}^n \mathcal{X}_k \subseteq \mathbb{R}^n$. Let $g : \mathcal{X} \rightarrow \mathbb{R}$ be a measurable function such that, for some constants $\{d_k\}_{k=1}^n$,*

$$|g(\underline{x}) - g(\underline{x}')| \leq d_k, \quad \forall k \in \{1, \dots, n\} \quad (21)$$

where $\underline{x} = (x_1, \dots, x_{k-1}, x_k, x_{k+1}, \dots, x_n)$ and $\underline{x}' = (x_1, \dots, x_{k-1}, x'_k, x_{k+1}, \dots, x_n)$ are two arbitrary points in the set \mathcal{X} that may only differ in their k -th coordinate (this is equivalent to saying that the variation of the function g w.r.t. its k -th coordinate is upper bounded by d_k). Then, for every $\alpha \geq 0$,

$$\mathbb{P}(|g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)]| \geq \alpha) \leq 2 \exp\left(-\frac{2\alpha^2}{\sum_{k=1}^n d_k^2}\right). \quad (22)$$

Remark 4.6. One can use the Azuma-Hoeffding inequality for a derivation of a concentration inequality in the considered setting. However, the following proof provides in this setting an improvement by a factor of 4 in the exponent of the bound.

Proof. For $k \in \{1, \dots, n\}$, let $\mathcal{F}_k = \sigma(X_1, \dots, X_k)$ be the σ -algebra that is generated by X_1, \dots, X_k with $\mathcal{F}_0 = \{\emptyset, \Omega\}$. Define

$$\xi_k \triangleq \mathbb{E}[g(X_1, \dots, X_n) | \mathcal{F}_k] - \mathbb{E}[g(X_1, \dots, X_n) | \mathcal{F}_{k-1}], \quad \forall k \in \{1, \dots, n\}. \quad (23)$$

Note that $\mathcal{F}_0 \subseteq \mathcal{F}_1 \dots \subseteq \mathcal{F}_n$ is a filtration, and

$$\begin{aligned} \mathbb{E}[g(X_1, \dots, X_n) | \mathcal{F}_0] &= \mathbb{E}[g(X_1, \dots, X_n)] \\ \mathbb{E}[g(X_1, \dots, X_n) | \mathcal{F}_n] &= g(X_1, \dots, X_n). \end{aligned} \quad (24)$$

Hence, it follows from the last three equalities that

$$g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)] = \sum_{k=1}^n \xi_k.$$

In the following, we need a lemma:

Lemma 4.7. *For every $k \in \{1, \dots, n\}$, the following properties hold a.s.:*

1. $\mathbb{E}[\xi_k | \mathcal{F}_{k-1}] = 0$, so $\{\xi_k, \mathcal{F}_k\}$ is a martingale-difference and ξ_k is \mathcal{F}_k -measurable.
2. $|\xi_k| \leq d_k$
3. $\xi_k \in [A_k, A_k + d_k]$ where A_k is some non-positive and \mathcal{F}_{k-1} -measurable random variable.

Proof. The random variable ξ_k is \mathcal{F}_k -measurable since $\mathcal{F}_{k-1} \subseteq \mathcal{F}_k$, and ξ_k is a difference of two functions where one is \mathcal{F}_k -measurable and the other is \mathcal{F}_{k-1} -measurable. Furthermore, it is easy to verify that $\mathbb{E}[\xi_k | \mathcal{F}_{k-1}] = 0$. This proves the first item. The second item follows from the first and third items. To prove the third item, note that $\xi_k = f_k(X_1, \dots, X_k)$ holds a.s. for some function $f_k : \mathcal{X}_1 \times \dots \times \mathcal{X}_k \rightarrow \mathbb{R}$ that is \mathcal{F}_k -measurable. Let us define, for every $k \in \{1, \dots, n\}$,

$$A_k \triangleq \inf_{x \in \mathcal{X}_k} f_k(X_1, \dots, X_{k-1}, x),$$

$$B_k \triangleq \sup_{x \in \mathcal{X}_k} f_k(X_1, \dots, X_{k-1}, x)$$

which are \mathcal{F}_{k-1} -measurable, and by definition $\xi_k \in [A_k, B_k]$ holds almost surely. Furthermore, for every point $(x_1, \dots, x_{k-1}) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_{k-1}$, we

have

$$\begin{aligned}
& \sup_{x \in \mathcal{X}_k} f_k(x_1, \dots, x_{k-1}, x) - \inf_{x' \in \mathcal{X}_k} f_k(x_1, \dots, x_{k-1}, x') \\
&= \sup_{x, x' \in \mathcal{X}_k} \{f_k(x_1, \dots, x_{k-1}, x) - f_k(x_1, \dots, x_{k-1}, x')\} \\
&= \sup_{x, x' \in \mathcal{X}_k} \left\{ \mathbb{E}[g(X_1, \dots, X_n) \mid X_1 = x_1, \dots, X_{k-1} = x_{k-1}, X_k = x] \right. \\
&\quad \left. - \mathbb{E}[g(X_1, \dots, X_n) \mid X_1 = x_1, \dots, X_{k-1} = x_{k-1}, X_k = x'] \right\} \\
&= \sup_{x, x' \in \mathcal{X}_k} \left\{ \mathbb{E}[g(x_1, \dots, x_{k-1}, x, X_{k+1}, \dots, X_n)] - \mathbb{E}[g(x_1, \dots, x_{k-1}, x', X_{k+1}, \dots, X_n)] \right\} \\
&= \sup_{x, x' \in \mathcal{X}_k} \left\{ \mathbb{E}[g(x_1, \dots, x_{k-1}, x, X_{k+1}, \dots, X_n) - g(x_1, \dots, x_{k-1}, x', X_{k+1}, \dots, X_n)] \right\} \\
&\leq d_k
\end{aligned} \tag{26}$$

where (??) follows from the independence of the random variables $\{X_k\}_{k=1}^n$, and (??) follows from the condition in (??). Hence, it follows that $B_k - A_k \leq d_k$ a.s., which then implies that $\xi_k \in [A_k, A_k + d_k]$. Since $\mathbb{E}[\xi_k \mid \mathcal{F}_{k-1}] = 0$ then a.s. the \mathcal{F}_{k-1} -measurable function A_k is non-positive. It is noted that the third item of the lemma makes it different from the proof of the Azuma-Hoeffding inequality (in that case, it implies that $\xi_k \in [-d_k, d_k]$ where the length of the interval is twice larger.) \blacksquare

Applying the convexity of the exponential function (similarly to the derivation of the Azuma-Hoeffding inequality, but this time w.r.t. the interval $[A_k, A_k + d_k]$) implies that for every $k \in \{1, \dots, n\}$

$$\begin{aligned}
& \mathbb{E}[e^{t\xi_k} \mid \mathcal{F}_{k-1}] \\
&\leq \mathbb{E} \left[\frac{(\xi_k - A_k)e^{t(A_k + d_k)} + (A_k + d_k - \xi_k)e^{tA_k}}{d_k} \mid \mathcal{F}_{k-1} \right] \\
&= \frac{(A_k + d_k)e^{tA_k} - A_ke^{t(A_k + d_k)}}{d_k}.
\end{aligned}$$

Let $P_k \triangleq -\frac{A_k}{d_k} \in [0, 1]$, then

$$\begin{aligned}
& \mathbb{E}[e^{t\xi_k} \mid \mathcal{F}_{k-1}] \\
& \leq P_k e^{t(A_k+d_k)} + (1 - P_k) e^{tA_k} \\
& = e^{tA_k} (1 - P_k + P_k e^{td_k}) \\
& = e^{H_k(t)}
\end{aligned} \tag{27}$$

where

$$H_k(t) \triangleq tA_k + \ln(1 - P_k + P_k e^{td_k}), \quad \forall t \in \mathbb{R}. \tag{28}$$

Since $H_k(0) = H'_k(0) = 0$ and the geometric mean is less than or equal to the arithmetic mean then, for every t ,

$$H''_k(t) = \frac{d_k^2 P_k (1 - P_k) e^{td_k}}{(1 - P_k + P_k e^{td_k})^2} \leq \frac{d_k^2}{4}$$

which implies by Taylor's theorem that

$$H_k(t) \leq \frac{t^2 d_k^2}{8} \tag{29}$$

so, from (??),

$$\mathbb{E}[e^{t\xi_k} \mid \mathcal{F}_{k-1}] \leq e^{\frac{t^2 d_k^2}{8}}.$$

Similarly to the proof of the Azuma-Hoeffding inequality, by repeatedly using the recursion in (??), the last inequality implies that

$$\mathbb{E} \left[\exp \left(t \sum_{k=1}^n \xi_k \right) \right] \leq \exp \left(\frac{t^2}{8} \sum_{k=1}^n d_k^2 \right) \tag{30}$$

which then gives from (??) that, for every $t \geq 0$,

$$\begin{aligned}
& \mathbb{P}(g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)] \geq \alpha) \\
& = \mathbb{P} \left(\sum_{k=1}^n \xi_k \geq \alpha \right) \\
& \leq \exp \left(-\alpha t + \frac{t^2}{8} \sum_{k=1}^n d_k^2 \right).
\end{aligned} \tag{31}$$

An optimization over the free parameter $t \geq 0$ gives that $t = 4\alpha (\sum_{k=1}^n d_k^2)^{-1}$, so

$$\mathbb{P}(g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)] \geq \alpha) \leq \exp\left(-\frac{2\alpha^2}{\sum_{k=1}^n d_k^2}\right). \quad (32)$$

By replacing g with $-g$, it follows that this bound is also valid for the probability

$$\mathbb{P}(g(X_1, \dots, X_n) - \mathbb{E}[g(X_1, \dots, X_n)] \leq -\alpha)$$

which therefore gives the bound in (??). This completes the proof of Theorem 4.5. ■

4.3 Hoeffding's inequality, and its improved version (the Kearns-Saul inequality)

In the following, we derive a concentration inequality for sums of independent and bounded random variables as a consequence of McDiarmid's inequality. This inequality is due to Hoeffding (see [?, Theorem 2]). An improved version of Hoeffding's inequality, due to Kearns and Saul [?], is also introduced in the following.

Theorem 4.8 (Hoeffding). *Let $\{U_k\}_{k=1}^n$ be a sequence of independent and bounded random variables such that, for every $k \in \{1, \dots, n\}$, $U_k \in [a_k, b_k]$ holds a.s. for some constants $a_k, b_k \in \mathbb{R}$. Let $\mu_n \triangleq \sum_{k=1}^n \mathbb{E}[U_k]$. Then,*

$$\mathbb{P}\left(\left|\sum_{k=1}^n U_k - \mu_n\right| \geq \alpha\sqrt{n}\right) \leq 2 \exp\left(-\frac{2\alpha^2 n}{\sum_{k=1}^n (b_k - a_k)^2}\right), \quad \forall \alpha \geq 0. \quad (33)$$

Proof. Apply Theorem 4.5 to $g(\underline{u}) \triangleq \sum_{k=1}^n u_k$ for every $\underline{u} \in \prod_{k=1}^n [a_k, b_k]$. ■

5 Statistical Learning Theory

We want to introduce some notions from Statistical Learning Theory (an application in Computer Science/ Statistics) to describe how risk is controlled in predictive models, i.e., their expected inaccuracy on new data from the same source as that used to fit the model. Or to say this in another way, the goal of statistical learning theory is to study, in a statistical framework, the properties of learning algorithms. In this chapter, I summarize the basic forms of these results in the literature, sacrificing some rigour for brevity.

5.1 The Traditional Setup

Consider predictors $X \in \mathcal{X}$ and responses $Y \in \mathcal{Y}$. Let \mathcal{F} be a class of functions $f : \mathcal{X} \rightarrow \mathcal{Y}$ which take predictors as inputs. Define a loss function $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ which measures the cost of making poor predictions. Throughout this chapter I make the following assumption on the loss function.

Assumption 1. $\forall f \in \mathcal{F}$

$$0 \leq \downarrow(y, y') \leq M < \infty$$

Then, I can define the risk of any predictor $f \in \mathcal{F}$.

Definition 5.1 (Risk or generalization error).

$$R(f) := \int \downarrow(f(X), Y) d\mathbb{P} = \mathbb{E}_{\mathbb{P}}[\downarrow(f(X), Y)], \quad (34)$$

where $(X, Y) \sim \mathbb{P}$

The risk or generalization error measures the expected cost of using f to predict Y from X given a new observation. Just to emphasize, the expectation is taken with respect to the distribution \mathbb{P} of the test point (X, Y) which is independent of f ; the risk is a deterministic function of f with all the randomness in the data averaged away.

Since the true distribution \mathbb{P} is unknown, so is $R(f)$, but one can attempt to estimate it based on only the observed data. Suppose that I observe a random sample $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ so that $(X_i, Y_i) \stackrel{i.i.d}{\sim} \mathbb{P}$, i.e. $D_n \sim \mathbb{P}$. Define the *training error or empirical risk* of f as follows.

Definition 5.1 (Training error or empirical risk).

$$\widehat{R}_n(f) := \frac{1}{n} \sum_{i=1}^n \downarrow(f(X_i), Y_i). \quad (35)$$

In other words, the in-sample training error, $\widehat{R}_n(f)$, is the average loss over the actual training points. It is easy to see that, because the training data D_n and the test point (X, Y) are IID, then given some fixed function f (chosen independently of the sample D_n),

$$\widehat{R}_n(f) = R(f) + \gamma_n(f) \quad (36)$$

where $\gamma_n(f)$ is a mean-zero noise variable that reflects how far the training sample departs from being perfectly representative of the data-generating distribution. Here I should emphasize that $\widehat{R}_n(f)$ is random enough through the training sample D_n . By the law of large numbers, for such fixed f , $\gamma_n(f) \rightarrow 0$ as $n \rightarrow \infty$, so, with enough data, one has a good idea of how well any given function will generalize to new data.

However, one is rarely interested in the performance of a single function f without adjustable parameters fixed for them in advance by theory. Rather, researchers are interested in a class of plausible functions \mathcal{F} , possibly indexed by some possibly infinite parameter $\theta \in \Theta$, which I refer to as a model. One function (one particular parameter point) is chosen from the model class minimizing some criterion function. Maximum likelihood, Bayesian maximization *a posteriori*, least squares, regularized methods, and empirical risk minimization (ERM) all have this flavor as do many other estimation methods. In these cases, one can define the empirical risk minimizer for an appropriate loss function \downarrow .

Definition 5.2. [*Empirical Risk Minimizer*]

$$\widehat{f} := \operatorname{argmin}_{f \in \mathcal{F}} \widehat{R}_n(f) = \operatorname{argmin}_{f \in \mathcal{F}} (R(f) + \gamma_n(f)). \quad (37)$$

It is important to note that \widehat{f} is random and measurable with respect to the empirical risk process $\widehat{R}_n(f)$ for $f \in \mathcal{F}$. Choosing a predictor \widehat{f} by empirical risk minimization (tuning the adjustable parameters so that \widehat{f} fits the training data well) conflates predicting future data well (low $R(\widehat{f})$, the true risk) with exploiting the accidents and noise of the training data (large negative $\gamma_n(\widehat{f})$, finite-sample noise). The true risk of \widehat{f} will generally

be bigger than its in-sample risk precisely because I picked it to match the data well. In doing so, \hat{f} ends up reproducing some of the noise in the data and therefore will not generalize well. The difference between the true and apparent risk depends on the magnitude of the sampling fluctuations:

$$R(\hat{f}) - \hat{R}_n(\hat{f}) \leq \sup_{f \in \mathcal{F}} \|\gamma_n(f)\| = \Gamma_n(\mathcal{F}) \quad (38)$$

In ((??)), $R(\hat{f})$ is random and measurable with respect to \hat{f} .

The main goal of statistical learning theory is to control $\Gamma_n(\mathcal{F})$ while making minimal assumptions about the data generating process - i.e. to provide bounds on over-fitting. Using more flexible models (allowing more general distributions, adding parameters, etc.) has two contrasting effects. On the one hand, it improves the best possible accuracy, lowering the minimum of the true risk. On the other hand, it increases the ability to, as it were, memorize noise for any fixed sample n . This qualitative observation - a generalization of the bias-variance trade-off from estimation theory - can be made use-fully precise by quantifying the complexity of model classes. A typical result is a confidence bound on Γ_n (and hence on over-fitting), which says that with probability at least $1 - \nu$,

$$\Gamma_n(\mathcal{F}) \leq \Theta(\Delta(\mathcal{F}), n, \nu), \quad (39)$$

where $\Delta(\cdot)$ is some suitable measure of the complexity of the model \mathcal{F} . To give specific forms of $\Theta(\cdot)$, I need to show that for a particular f , $R(f)$ and $\hat{R}_n(f)$ will be close to each other for any fixed n without knowledge of the distribution of the data. Furthermore, I need the complexity $\Delta(\mathcal{F})$, to claim that $R(f)$ and $\hat{R}_n(f)$ will be close, not only for a particular f , but uniformly over all $f \in \mathcal{F}$. Together these two results will allow me to show, despite little knowledge of the data generating process, how bad the \hat{f} which I choose will be at forecasting future observations.

5.2 Concentration

The first step to controlling the difference between the empirical and expected risk is to develop concentration results for fixed functions. We have already introduced various inequalities. McDiarmid Inequality and Hoeffding's Inequality. These results are extremely important in Statistical learning theory, but it is beyond the scope of this thesis to go into much more detail. In the remainder of this section, I will show how to obtain concentration for the training error around the risk for two different choices of the random

variables Z_i . This will lead to two different ways of controlling Γ_n and hence the generalization error of prediction functions.

5.3 Contol by Counting

Let us assume that we let Z_i be the loss of the i^{th} training point for some fixed function f . Then by Hoeffding's inequality we get the following remarkable result

$$\mathbb{P}^n \left(\|R(f) - \widehat{R}_n(f)\| \geq \epsilon \right) \leq 2 \exp \left\{ -\frac{2n\epsilon^2}{M^2} \right\} \quad (40)$$

This result is quite powerful, it says that the probability of observing data which will result in a training error much different from the expected risk goes to zero exponentially with the size of the training set. The only assumption necessary was $0 \leq \downarrow(y, y') \leq M < \infty$.

5.4 Capacity

For “small” models, we can just count the number of functions in the class and take the union bound. Suppose that $\mathcal{F} = \{f_1, \dots, f_N\}$. Then we have

$$\mathcal{P} \left(\sup_{1 \leq i \leq N} |R(f_i) - \widehat{R}_n(f_i)| > \epsilon \right) \leq \sum_{i=1}^N \mathcal{P} \left(|R(f_i) - \widehat{R}_n(f_i)| > \epsilon \right) \quad (41)$$

$$\leq N \exp \left\{ -\frac{2n\epsilon^2}{K} \right\}, \quad (42)$$

by Theorem ???. Most interesting models are not small in this sense, but similar results hold when model size is measured appropriately.

There are a number of measures for the size or capacity of a model. Algorithmic stability [?, ?, ?] quantifies the sensitivity of the chosen function to small perturbations to the data. Similarly, maximal discrepancy [?] asks how different the predictions could be if two functions are chosen using two separate data sets. A more direct, functional-analytic approach partitions \mathcal{F} into equivalence classes under some metric, leading to covering numbers [?, ?]. Rademacher complexity [?] directly describes a model's ability to fit random noise. We focus on a measure which is both intuitive and powerful: Vapnik-Chervonenkis (VC) dimension [?, ?].

VC dimension starts as an idea about collections of sets.

Definition 5.3. Let \mathbb{U} be some (infinite) set and S a finite subset of \mathbb{U} . Let \mathcal{C} be a family of subsets of \mathbb{U} . We say that \mathcal{C} shatters S if for every $S' \subseteq S$, $\exists C \in \mathcal{C}$ such that $S' = S \cap C$.

Essentially, \mathcal{C} can shatter a set S if it can pick out every subset of points in S . This says that the collection \mathcal{C} is very complicated or flexible. The cardinality of the largest set S that can be shattered by \mathcal{C} is the latter's VC dimension.

Definition 5.4 (VC dimension). The Vapnik-Chervonenkis (VC) dimension of a collection \mathcal{C} of subsets of \mathbb{U} is

$$\text{VCD}(\mathcal{C}) := \sup\{|S| : S \subseteq \mathbb{U} \text{ and } S \text{ is shattered by } \mathcal{C}\}. \quad (43)$$

To see why this is a “dimension”, we need one more notion.

Definition 5.5 (Growth function). The growth function $G(\mathcal{C}, n)$ of a collection \mathcal{C} of subsets of \mathbb{U} is the maximum number of subsets which can be formed by intersecting a set $S \subset \mathbb{U}$ of cardinality n with \mathcal{C} ,

$$G(n, \mathcal{C}) := \sup_{S \subset \mathbb{U} : |S|=n} |S \wedge \mathcal{C}| \quad (44)$$

The growth function counts how many *effectively* distinct sets the collection contains, when we can only observe what is going on at n points, not all of \mathbb{U} . If $n \leq \text{VCD}(\mathcal{C})$, then from the definitions $G(n, \mathcal{C}) = 2^n$. If the VC dimension is finite, however, and $n > \text{VCD}(\mathcal{C})$, then $G(n, \mathcal{C}) < 2^n$, and in fact it can be shown [?] that

$$G(n, \mathcal{C}) \leq (n + 1)^{\text{VCD}(\mathcal{C})}. \quad (45)$$

This polynomial growth of capacity with n is why VCD is a “dimension”.

Using VC dimension to measure the capacity of function classes is straightforward. Define the indicator function $\mathbf{1}_A(x)$ to take the value 1 if $x \in A$ and 0 otherwise. Suppose that $f \in \mathcal{F}$, $f : \mathbb{U} \rightarrow \mathbb{R}$. Each f corresponds to the set

$$C_f = \{(u, a) : \mathbf{1}_{(0, \infty)}(f(u) - b) = 1, \quad u \in \mathbb{U}, \quad b \in \mathbb{R}\}, \quad (46)$$

so \mathcal{F} corresponds to the class $\mathcal{C}_{\mathcal{F}} := \{C_f : f \in \mathcal{F}\}$. Essentially, the growth function $G(n, \text{VCD}(\mathcal{F}))$ counts the effective number of functions in \mathcal{F} , i.e.,

how many can be told apart using only n observations. When $\text{VCD}(\mathcal{F}) < \infty$, this number grows only polynomially with n . This observation lets us control the risk over the entire model, providing one of the pillars of statistical learning theory.

Theorem 5.6 ([?]). *Suppose that $\text{VCD}(\mathcal{F}) < \infty$ and $0 \leq \ell(y, y') \leq K < \infty$. Then,*

$$\mathcal{P} \left(\sup_{f \in \mathcal{F}} |R(f) - \widehat{R}_n(f)| > \epsilon \right) \leq 4(2n+1)^{\text{VCD}(\mathcal{F})} \exp \left\{ -\frac{n\epsilon^2}{K_1^2} \right\}, \quad (47)$$

where K_1 depends only on K and not n or \mathcal{F} .

The proof of this theorem has a similar flavor to the union bound argument given in (42).

This theorem has as an immediate corollary a bound for the out-of-sample risk. Since $\sup_{f \in \mathcal{F}}$ is inside the probability statement in (47), it applies to both pre-specified and to data-dependent functions, including any \widehat{f} chosen by fitting a model or minimizing empirical risk.

Corollary 5.7. *When Theorem 5.6 applies, for any $\eta > 0$ and any $f \in \mathcal{F}$, with probability at least $1 - \eta$,*

$$R(f) \leq \widehat{R}_n(f) + K_1 \sqrt{\frac{\text{VCD}(\mathcal{F}) \log(2n+1) + \log 4/\eta}{n}}. \quad (48)$$

The factor K_1 can be calculated explicitly but is unilluminating and we will not need it. Conceptually, the right-hand side of this inequality resembles standard model selection criteria, like AIC or BIC, with in-sample fit plus a penalty term which goes to zero as $n \rightarrow \infty$. Here however, the bound holds with high probability despite lack of knowledge of \mathcal{P} and it has nothing to do with asymptotic convergence: it holds for each n . It does however hold *only* with high \mathcal{P} probability, not always.

VC dimension is well understood for some function classes. For instance, if $\mathcal{F} = \{\mathbf{x} \mapsto \boldsymbol{\gamma} \cdot \mathbf{x} : \boldsymbol{\gamma} \in \mathbb{R}^p\}$ then $\text{VCD}(\mathcal{F}) = p+1$, i.e. it is the number of free parameters in a linear regression plus 1. VC dimension does not always have such a nice relation to the number of free parameters however; the classic example is the model $\mathcal{F} = \{x \mapsto \sin(\omega x) : \omega \in \mathbb{R}\}$, which has only one free

parameter, but $\text{VCD}(\mathcal{F}) = \infty$.⁵ At the same time, there are model classes (support vector machines) which may have infinitely many parameters but finite VC dimension [?]. This illustrates a further difference between the statistical learning approach and the usual information criteria, which are based on parameter-counting.

The concentration results in Theorem 5.6 and Theorem 5.7 work well for independent data. The first shows how quickly averages concentrate around their expectations: exponentially fast in the size of the data. The second result generalizes the first from a single function to entire function classes. Both results, as stated, depend critically on the independence of the random variables. For time series, we must be able to handle dependent data. In particular, because time-series data are dependent, the length n of a sample path Y_1, \dots, Y_n exaggerates how much information it contains. Knowing the past allows forecasters to predict future data (at least to some degree), so actually observing those future data points gives less information about the underlying process than in the IID case. Thus, while in ?? the probability of large discrepancies between empirical means and their expectations decreases exponentially in n , in the dependent case, the effective sample size may be much less than n resulting in looser bounds.

6 Applications

The word *applications* means different things to different people. If we wear our Mathematician hats we are always looking for a new tool to prove some theorem. Or reprove a celebrated theorem in a new way. Computer Scientists are always concerned with algorithms and performance analysis and statisticians or *data scientists* are after getting a handle on the asymptotics of estimators and samplers. In short, each specialist wants to know what a given tool will contribute to his field.

The remarkable thing about concentration of measure is that it uses span a wide range - from something practical as decoding neural signals to esoteric topics such as analyzing convex bodies in Banach spaces. It is too much to go beyond one or two applications. So I'll cite some here

⁵This result follows if we can show that for any positive integer J and any binary sequence (r_1, \dots, r_J) , there exists a vector (x_1, \dots, x_J) such that $\mathbf{1}_{[0,1]}(\sin(\omega x_i)) = r_i$. If we choose $x_i = 2\pi 10^{-i}$, then one can show that taking $\omega = \frac{1}{2} \left(\sum_{i=1}^J (1 - r_i) 10^i + 1 \right)$ solves the system of equations.

- Milman's proof of Dvoretzky's theorem on sections of convex bodies
- a widely cited lemma of Johnson and Lindenstrauss concerning low-distortion dimensionality reduction in \mathbb{R}^n by random projections.
- statistics and empirical processes and machine learning[?]

In the interest of keeping this document short we shall only look at the last example.

6.1 Rademacher Processes

A key technique in the theory of empirical processes is *Rademacher symmetrization*. This was first introduced into empirical processes in a classical paper by Gine [?] so we'll show how this applies in the context of Talagrand's inequality.

Let $\epsilon_i, i = 1, \dots, n$, be i.i.d Rademacher random signs (taking values -1,1 with probability 1/2), independent of the X_i 's, defined on a large product probability space with product probability \Pr , denote the joint expectation by E , and the E_ϵ and E_X the corresponding expectations w.r.t the ϵ_i 's X_i 's, respectively. The following symmetrization inequality holds for random variables in arbitrary normed spaces, but we state it for the supremum norm relevant in empirical process theory: For \mathcal{F} a class of functions on (S, \mathcal{A}) , define $\|H\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |H(f)|$.

Lemma 6.1. *Let \mathcal{F} be a uniformly bounded P -centered class of functions defined on a measurable space (S, \mathcal{A}) . Let ϵ_i be i.i.d. Rademachers as above, and let $a_i, i = 1, \dots, n$ be any sequence of real numbers. Then*

$$\frac{1}{2} E \left\| \sum_{i=1}^n \epsilon_i f(X_i) \right\|_{\mathcal{F}} \leq E \left\| \sum_{i=1}^n f(X_i) \right\|_{\mathcal{F}} \leq \left\| \sum_{i=1}^n \epsilon_i (f(X_i) + a_i) \right\|_{\mathcal{F}} \quad (49)$$

Proof. Let us assume for simplicity that \mathcal{F} is countable (so that we can neglect measurability problems). Since $E_X f(X_i) = 0$ for every f, i , the first inequality follows from

$$E \left\| \sum_{i=1}^n \epsilon_i (f(X_i)) \right\|_{\mathcal{F}} = E_\epsilon E_X \leq E_\epsilon E_X \left\| \sum_{i:\epsilon_i=-1} f(X_i) + E_X \sum_{i:\epsilon_i=1} f(X_i) \right\|_{\mathcal{F}} + E_\epsilon E_X \left\| \sum_{i:\epsilon_i=1} f(X_i) + E_X \sum_{i:\epsilon_i=-1} f(X_i) \right\|_{\mathcal{F}}$$

where in the last inequality we have used Jensen's inequality and convexity of the norm. To prove the second inequality, let $X_{n+i}, i = 1, \dots, n$ be an independent copy of X_1, \dots, X_n . Then proceeding as above, $E\|\sum_{i=1}^n f(X_i)\|_{\mathcal{F}} = E\|\sum_{i=1}^n (f(X_i) - Ef(X_{n+i}))\|_{\mathcal{F}} \leq E\|\sum_{i=1}^n (f(X_i + a_i) - \sum_{i=1}^n (f(X_{n+i} + a_i))\|_{\mathcal{F}}$ which clearly equals

$$E_{\epsilon} E_X \left\| \sum_{i:\epsilon_i=1} \epsilon_i (f(X_i) + a_i - f(X_{n+i}) - a_i) - \sum_{i:\epsilon_i=-1} \epsilon_i (f(X_i) + a_i - f(X_{n+i}) - a_i) \right\|_{\mathcal{F}}$$

Now \Pr being a product probability measure with identical coordinates, it is invariant by permutations of the coordinates, so that we may exchange $f(X_i)$ and $f(X_{n+i})$ for the i 's where $\epsilon_i = -1$ in the last expectation. This gives that the quantity in the last equation equals

$$E_{\epsilon} E_X \left\| \sum_{i=1}^n \epsilon_i (f(X_i) + a_i - f(X_{n+i}) - a_i) \right\|_{\mathcal{F}} \leq 2E \left\| \sum_{i=1}^n \epsilon_i (f(X_i) + a_i) \right\|_{\mathcal{F}}$$

which completes the proof. ■

This simple but very useful result says that we can always compare the size of the expectation of the supremum of an empirical process to a symmetrized process. The idea usual is that the symmetrized *Rademacher* process has conditional on the X_i 's a very simple structure. One can then derive results of the Rademacher process and integrate the results over the distribution of the X_i 's

Rademacher chaos and Rademacher averages are quantities that play an important role in empirical process theory[2] and in the theory of Banach spaces[6].

6.2 Rademacher averages

Let B denote a separable Banach space and let X_1, \dots, X_n be independent and identically distributed bounded B -valued random variables. Without loss of generality we assume that $\|X_1\| \leq 1$ almost surely. The quantity of interest is the conditional Rademacher average $Z = \mathbb{E} [\|\sum_{i=1}^n \epsilon_i X_i\| | X_1^n]$ where the ϵ_i are independent centered 1, -1-valued random variables. We offer the following concentration inequalities for Z :

Theorem 6.2. *For any $t > 0$,*

$$\mathbb{P}[Z \geq \mathbb{E}Z + t] \leq \exp\left[-\frac{t^2}{2\mathbb{E}Z + 2t/3}\right] \quad (50)$$

and

$$\mathbb{P}[Z \leq \mathbb{E}Z - t] \leq \exp\left[-\frac{t^2}{2\mathbb{E}Z}\right] \quad (51)$$

We are particularly interested in deriving (upper and lower) tail inequalities for conditional Rademacher averages. This is significant in statistical applications.

6.3 Rademacher Chaos

In this section \mathcal{F} denotes a collection of $n \times n$ symmetric matrices M , and $\epsilon_1, \dots, \epsilon_n$ are i.i.d. Rademacher variables. We assume that if $M \in \mathcal{F}$, then $-M \in \mathcal{F}$. To avoid problems with measurability we assume that \mathcal{F} is a finite set. For convenience assume that the matrices M have zero diagonal, that is, $M(i,i) = 0$ for all $M \in \mathcal{F}$ and $i = 1, \dots, n$. We investigate concentration of the random variable

$$Z = \sup_{M \in \mathcal{F}} \sum_{i,j \leq n} \epsilon_i \epsilon_j M(i,j)$$

Suppose the supremum of the L_2 operator norm of matrices $(M)_{M \in \mathcal{F}}$ is finite, and w.l.o.g we assume that this supremum equals one, that is,

$$\sup_{M \in \mathcal{F}} \sup_{\alpha: \sum_{i=1}^n \alpha_i^2 \leq 1} \alpha^\dagger M \alpha = 1 \quad (52)$$

where α^\dagger denotes the transpose of the vector $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$. We introduce an important theorem which follows from the previous result.

7 References

References

- [1] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. OUP Oxford, 2013.

- [2] R.M. Dudley. *Uniform Central Limit Theorems*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 1999.
- [3] L. Kontorovich. Metric and Mixing Sufficient Conditions for Concentration of Measure. *ArXiv Mathematics e-prints*, October 2006.
- [4] C. McDiarmid. *Concentration*,. Probabilistic Methods for Algorithmic Discrete Mathematics. Springer, 1998.
- [5] D. J. McDonald, C. Rohilla Shalizi, and M. Schervish. Time series forecasting: model evaluation and selection using nonparametric risk bounds. *ArXiv e-prints*, December 2012.
- [6] M.Ledoux. *The Concentration of Measure Phenomenon*, volume 89 of *Mathematical Surveys and Monographs*. American Mathematical Society, 2001.
- [7] Maxim Raginsky and Igal Sason. Concentration of measure inequalities in information theory, communications and coding. *CoRR*, abs/1212.4663, 2012.
- [8] T. Tao. *Topics in Random Matrix Theory*. Graduate studies in mathematics. American Mathematical Society, 2012.
- [9] Hugo Touchette. The large deviation approach to statistical mechanics. *Physics Reports*, 478(13):1 – 69, 2009.
- [10] V. Vapnik. *The Nature of Statistical Learning Theory*. Information Science and Statistics. Springer, 2000.
- [11] S. R. S. Varadhan. Large deviations. *The Annals of Probability*, 36(2):397–419, 03 2008.
- [12] Peter Whittle. Large-deviation theory. In *Probability via Expectation*, Springer Texts in Statistics, pages 306–316. Springer New York, 2000.