

# Adversarial Laws of Large Numbers and Optimal Regret in Online Classification

Noga Alon\*  
Shay Moran<sup>§</sup>

Omri Ben-Eliezer<sup>†</sup>  
Moni Naor<sup>¶</sup>

Yuval Dagan<sup>‡</sup>  
Eylon Yogev<sup>||</sup>

## Abstract

Laws of large numbers guarantee that given a large enough sample from some population, the measure of any fixed sub-population is well-estimated by its frequency in the sample. We study laws of large numbers in sampling processes that can affect the environment they are acting upon and interact with it. Specifically, we consider the sequential sampling model proposed by Ben-Eliezer and Yogev (2020), and characterize the classes which admit a uniform law of large numbers in this model: these are exactly the classes that are *online learnable*. Our characterization may be interpreted as an online analogue to the equivalence between learnability and uniform convergence in statistical (PAC) learning.

The sample-complexity bounds we obtain are tight for many parameter regimes, and as an application, we determine the optimal regret bounds in online learning, stated in terms of *Littlestone's dimension*, thus resolving the main open question from Ben-David, Pál, and Shalev-Shwartz (2009), which was also posed by Rakhlin, Sridharan, and Tewari (2015).

---

\*Department of Mathematics, Princeton University, Princeton, New Jersey, USA and Schools of Mathematics and Computer Science, Tel Aviv University, Tel Aviv, Israel. Research supported in part by NSF grant DMS-1855464, BSF grant 2018267 and the Simons Foundation. Email: [nalon@math.princeton.edu](mailto:nalon@math.princeton.edu).

†Center for Mathematical Sciences and Applications, Harvard University, Massachusetts, USA. Research partially conducted while the author was at Weizmann Institute of Science, supported in part by a grant from the Israel Science Foundation (no. 950/15). Email: [omribene@cmsa.fas.harvard.edu](mailto:omribene@cmsa.fas.harvard.edu).

‡Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA. Email: [dagan@mit.edu](mailto:dagan@mit.edu).

§Department of Mathematics, Technion, Israel. Email: [smoran@technion.ac.il](mailto:smoran@technion.ac.il). Research supported in part by the Israel Science Foundation (grant No. 1225/20), by an Azrieli Faculty Fellowship, and by a grant from the United States - Israel Binational Science Foundation (BSF).

¶Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot, Israel. Supported in part by grants from the Israel Science Foundation (no. 950/15 and 2686/20) and by the Simons Foundation Collaboration on the Theory of Algorithmic Fairness. Incumbent of the Judith Kleeman Professorial Chair. Email: [moni.naor@weizmann.ac.il](mailto:moni.naor@weizmann.ac.il).

||Department of Computer Science, Boston University and Department of Computer Science, Tel Aviv University. Email: [eylony@gmail.com](mailto:eylony@gmail.com). Research supported in part by ISF grants 484/18, 1789/19, Len Blavatnik and the Blavatnik Foundation, and The Blavatnik Interdisciplinary Cyber Research Center at Tel Aviv University.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	The Adversarial Sampling Model . . . . .	1
<b>2</b>	<b>Main Results</b>	<b>2</b>
2.1	Adversarial Laws of Large Numbers . . . . .	2
2.1.1	Lower Bounds . . . . .	3
2.2	Online Learning . . . . .	4
2.3	Applications and Extensions . . . . .	4
<b>3</b>	<b>Technical Overview</b>	<b>6</b>
3.1	Upper Bounds . . . . .	6
3.1.1	Step 1: Reduction to Online Discrepancy via Double Sampling . . . . .	6
3.1.2	Step 2: From Online Discrepancy to Sequential Rademacher . . . . .	7
3.1.3	Step 3.1: Bounding Sequential Rademacher Complexity – Oblivious Case . . . . .	7
3.1.4	Step 3.2: Bounding Sequential Rademacher Complexity – Adversarial Case . . . . .	8
3.2	Lower Bounds . . . . .	10
<b>4</b>	<b>Related Work</b>	<b>11</b>
4.1	VC Theory . . . . .	11
4.2	Online Learning . . . . .	11
4.3	Streaming Algorithms . . . . .	12
<b>5</b>	<b>Preliminaries</b>	<b>12</b>
5.1	Basic Definitions: Littlestone Dimension and Sampling Schemes . . . . .	12
5.2	Notation . . . . .	13
5.3	Additional Central Definitions . . . . .	14
5.4	Sampling Without Replacement . . . . .	14
<b>6</b>	<b>Epsilon Approximations</b>	<b>15</b>
<b>7</b>	<b>Epsilon Nets</b>	<b>16</b>
<b>8</b>	<b>Double Sampling</b>	<b>18</b>
8.1	Proof of Lemma 6.2 . . . . .	20
8.2	Proof of Lemma 7.2 . . . . .	21
<b>9</b>	<b>Covering Numbers</b>	<b>21</b>
9.1	Overview . . . . .	21
9.1.1	Epsilon Approximation and Sequential Rademacher . . . . .	22
9.1.2	Epsilon Nets . . . . .	23
9.1.3	Organization . . . . .	23
9.2	Covering for Littlestone Classes . . . . .	24
9.2.1	Proof of Lemma 9.3 . . . . .	24
9.2.2	Proof of Lemma 9.5 . . . . .	27
9.3	Deriving Bounds on $\epsilon$ -Approximation via Fractional Covering Numbers . . . . .	30

9.3.1	Basic Lemmas for Deterministic Covers and Proof of Lemma 9.6 . . . . .	30
9.3.2	Basic Lemmas for Fractional Covers . . . . .	31
9.3.3	Chaining for Non-Fractional Covers . . . . .	32
9.3.4	Proof of Lemma 9.7 . . . . .	33
9.4	Bounds on $\epsilon$ -Nets via Fractional Covering Numbers . . . . .	37
<b>10</b>	<b>Reductions Between Different Sampling Schemes</b>	<b>38</b>
10.1	Intuition for the Reduction Method . . . . .	38
10.2	Abstract Reduction Method . . . . .	39
10.3	Bounds for Reservoir Sampling via Uniform Sampling . . . . .	41
10.4	Bounds for Bernoulli Sampling via Uniform Sampling . . . . .	43
10.5	Bounds for Uniform Sampling via Bernoulli Sampling . . . . .	44
10.5.1	Proof of Lemma 6.3 . . . . .	45
10.5.2	Proof of Lemma 7.3 . . . . .	49
<b>11</b>	<b>Continuous <math>\epsilon</math>-Approximation</b>	<b>50</b>
<b>12</b>	<b>Online Learning</b>	<b>51</b>
12.1	Formal Definitions . . . . .	51
12.2	Statement and Proof . . . . .	52
<b>13</b>	<b>Lower Bounds</b>	<b>52</b>
13.1	Proofs . . . . .	53
<b>A</b>	<b>Probabilistic Material</b>	<b>59</b>
A.1	Filtration and Martingales . . . . .	59
A.2	Sampling Without Replacement . . . . .	60

# 1 Introduction

When analyzing an entire population is infeasible, statisticians apply *sampling methods* by selecting a *sample* of elements from a target population as a guide to the entire population. Thus, one of the most fundamental tasks in statistics is to provide bounds on the sample size which is sufficient to soundly represent the population, and probabilistic tools are used to derive such guarantees, under a variety of assumptions. Virtually all of these guarantees are based on classical probabilistic models which assume that *the target population is fixed in advance and does not depend on the sample collected throughout the process*. Such an assumption, that the setting is *offline* (or *oblivious* or *static*), is however not always realistic. In this work we explore an abstract framework which removes this assumption, and prove that natural and efficient sampling processes produce samples which soundly represent the target population.

Situations where the sampling process explicitly or implicitly affects the target population are abundant in modern data analysis. Consider, for instance, navigation apps that optimize traffic by routing drivers to less congested routes: such apps collect statistics from drivers to estimate the traffic-load on the routes, and use these estimates to guide their users through faster routes. Thus, such apps interact with and affect the statistics they estimate. Consequently, the assumption that the measured populations do not depend on the measurements is not realistic.

Similar issues generally arise in settings involving decision-making in the face of an ever-changing (and sometimes even adversarial) environment; a few representative examples include autonomous driving [SBM<sup>+</sup>18], adaptive data analysis [DFH<sup>+</sup>15, WFRS18], security [NY15], and theoretical analysis of algorithms [CGP<sup>+</sup>18]. Consequently, there has recently been a surge of works exploring such scenarios, a partial list includes [MNS11, GHR<sup>+</sup>12, GHS<sup>+</sup>12, HW13, NY15, BJWY20, CN20, HRS20, HKM<sup>+</sup>20, WZ20]. In this work, we focus on the sequential sampling model recently proposed by Ben-Eliezer and Yogev [BEY20].

## 1.1 The Adversarial Sampling Model

We next formally describe the sampling setting and the main question we investigate. Ben-Eliezer and Yogev [BEY20] model sampling processes over a domain  $X$  as a sequential game between two players: a sampler and an adversary. The game proceeds in  $n$  rounds, where in each round  $i = 1, \dots, n$ :

- The adversary picks an item  $x_i \in X$  and provides it to the sampler. The choice of  $x_i$  might depend on  $x_1, \dots, x_{i-1}$  and on all information sent to the adversary up to this point.
- Then, the sampler decides whether to add  $x_i$  to its sample.
- Finally, the adversary is informed of whether  $x_i$  was sampled by the sampler.

The number of rounds  $n$  is known in advance to both players.<sup>1</sup> We stress that both players can be randomized, in which case their randomness is private (i.e., not known to the other player).

**Oblivious Adversaries.** In the oblivious (or static) case, the sampling process consists only of the first two bullets. Equivalently, oblivious adversaries decide on the entire stream in advance, without receiving any feedback from the sampler. Unless stated otherwise, the adversary in this paper is assumed to be adaptive (not oblivious).

---

<sup>1</sup>Though we will also consider samplers which are oblivious to the number of rounds  $n$ .

**Uniform Laws of Large Numbers.** Uniform laws of large numbers (ULLN) quantify the minimum sample size which is sufficient to *uniformly estimate multiple statistics of the data*. (Rather than just a *single* statistic, as in standard laws of large numbers.) This is relevant, for instance, in the example given above regarding the navigation app: it is desirable to accurately compute the congestion along *all* routes (paths). Otherwise, one congested route may be regarded as entirely non-congested, and it will be selected for navigation.

Given a family  $\mathcal{E}$  of subsets of  $X$ , we consider ULLNs that estimate the frequencies of each subset  $E \in \mathcal{E}$  within the adversarial stream. Formally, let  $\bar{x} = \{x_1, \dots, x_n\}$  denote the input-stream produced by the adversary, and let  $\bar{s} = \{x_{i_1}, \dots, x_{i_k}\}$  denote the sample chosen by the sampler. The sample  $\bar{s}$  is called an  $\epsilon$ -approximation of the stream  $\bar{x}$  with respect to  $\mathcal{E}$  if:

$$(\forall E \in \mathcal{E}) : \left| \frac{|\bar{s} \cap E|}{|\bar{s}|} - \frac{|\bar{x} \cap E|}{|\bar{x}|} \right| \leq \epsilon. \quad (1)$$

That is,  $\bar{s}$  is an  $\epsilon$ -approximation of  $\bar{x}$  if the *true-frequencies*  $|\bar{x} \cap E|/|\bar{x}|$  are uniformly approximated by the *empirical frequencies*  $|\bar{s} \cap E|/|\bar{s}|$ . The following question is the main focus of this work:

**Question (Main Question).** *Given a family  $\mathcal{E}$ , an error-parameter  $\epsilon > 0$ , and  $k \in \mathbb{N}$ , is there a sampler that, given any adversarially-produced input stream  $\bar{x}$ , picks a sample  $\bar{s}$  of at most  $k$  items which forms an  $\epsilon$ -approximation of  $\bar{x}$ , with high probability?*

**The Story in the Statistical Setting.** It is instructive to compare with the statistical setting in which the sample  $\bar{s}$  is drawn independently from an unknown distribution over  $X$ . Here, ULLNs are characterized by the Vapnik-Chervonenkis (VC) Theory which asserts that a family  $\mathcal{E}$  satisfies a ULLN if and only if its VC dimension,  $\text{VC}(\mathcal{E})$ , is finite [VC71].

This fundamental result became a corner-stone in statistical machine learning. In particular, *The Fundamental Theorem of PAC Learning* states that the following properties are equivalent for any family  $\mathcal{E}$ : (1)  $\mathcal{E}$  satisfies a uniform law of large numbers, (2)  $\mathcal{E}$  is PAC learnable, and (3)  $\mathcal{E}$  has a finite VC dimension. Quantitatively, the sample size required for both  $\epsilon$ -approximation and for PAC learning with excess-error  $\epsilon$  is  $\Theta((\text{VC}(\mathcal{E}) + \log(1/\delta))/\epsilon^2)$ .

**Spoiler:** Our main result (stated below) can be seen as an online/adversarial analogue of this theorem where the **Littlestone dimension** replaces the VC dimension.

## 2 Main Results

### 2.1 Adversarial Laws of Large Numbers

The main result in this paper is a characterization of adversarial uniform laws of large numbers in the spirit of VC theory and The Fundamental Theorem of PAC Learning. We begin with the following central definition.

**Definition 2.1** (Adversarial ULLN). *We say that a family  $\mathcal{E}$  satisfies an adversarial ULLN if for any  $\epsilon, \delta > 0$ , there exist  $k = k(\epsilon, \delta) \in \mathbb{N}$  and a sampler  $\mathcal{S}$  satisfying the following. For any adversarially-produced input-stream  $\bar{x}$  (of any size),  $\mathcal{S}$  chooses a sample of at most  $k$  items, which form an  $\epsilon$ -approximation of  $\bar{x}$  with probability at least  $1 - \delta$ . We denote by  $k(\mathcal{E}, \epsilon, \delta)$  the minimal such value of  $k$ .*

Note that this definition requires the sample complexity  $k = k(\epsilon, \delta)$  to be a constant independent of the stream size  $n$ . Another reasonable requirement is  $k = o(n)$ . It turns out that these two requirements are equivalent.

Which families  $\mathcal{E}$  satisfy an adversarial law of large numbers? Clearly,  $\mathcal{E}$  must have a finite VC-dimension, as otherwise, basic VC-theory implies that any sampler will fail to produce an  $\epsilon$ -approximation even against oblivious adversaries which draw the input-stream  $\bar{x}$  independently from a distribution on  $X$ . However, finite VC dimension is not enough in the fully adversarial setting: [BEY20] exhibit a family  $\mathcal{E}$  with  $\text{VC}(\mathcal{E}) = 1$  that does not satisfy an adversarial ULLN.

Our first result provides a characterization of adversarial ULLN in terms of *Online Learnability*, which is analogous to the Fundamental Theorem of PAC Learning. In this context, the role of VC dimension is played by the *Littlestone dimension*, a combinatorial parameter which captures online learnability similar to how the VC dimension captures PAC learnability. (See Section 5.1 for the formal definition.)

**Theorem 2.2** (Adversarial ULLNs – Qualitative Characterization). *Let  $\mathcal{E}$  be a family of subsets of  $X$ . Then, the following statements are equivalent:*

1.  $\mathcal{E}$  satisfies an adversarial ULLN;
2.  $\mathcal{E}$  is online learnable; and
3.  $\mathcal{E}$  has a finite Littlestone dimension.

The proof follows from Theorems 2.3 and 13.1 (and from the well-known equivalence between online learnability and finite Littlestone dimension [Lit88, BPS09]). Our quantitative upper bound for the sample-complexity  $k(\mathcal{E}, \epsilon, \delta)$ , which is the main technical contribution of this paper, is stated next.

**Theorem 2.3** (Adversarial ULLNs – Quantitative Characterization). *Let  $\mathcal{E}$  be a family with Littlestone dimension  $d$ . Then, the sample size  $k(\mathcal{E}, \epsilon, \delta)$ , which suffices to produce an  $\epsilon$ -approximation satisfies:*

$$k(\mathcal{E}, \epsilon, \delta) \leq O\left(\frac{d + \log(1/\delta)}{\epsilon^2}\right).$$

The above upper bound is realized by natural and efficient samplers; for example it is achieved by: (i) the *Bernoulli sampler*  $\text{Ber}(n, p)$  which retains each element with probability  $p = k/n$ ; (ii) the *uniform sampler*  $\text{Uni}(n, k)$  that draws a subset  $I \subseteq [n]$  uniformly at random from all the subsets of size  $k$  and selects the sample  $\{x_t : t \in I\}$ ; and (iii) the *reservoir sampler*  $\text{Res}(n, k)$  (see Section 2.3) that maintains a uniform sample continuously throughout the stream.

### 2.1.1 Lower Bounds

The upper bound in Theorem 2.3 cannot be improved in general. In particular, it is tight in all parameters for *oblivious samplers*: a sampler is called oblivious if the indices of the chosen subsample are independent of the input-stream. (The Bernoulli, Reservoir, and Uniform samplers are of this type.) A lower bound of  $\Omega((d + \log(1/\delta))/\epsilon^2)$  for oblivious samplers directly follows from VC-theory, and applies to any family  $\mathcal{E}$  for which the VC dimension and Littlestone dimension are of the same order.<sup>2</sup> For unrestricted samplers we obtain bounds of  $\Omega(d/\epsilon^2)$

---

<sup>2</sup>E.g., projective spaces, Hamming balls, lines in the plane, and others.

for  $\epsilon$ -approximation and  $\Omega(d \log(1/\epsilon)/\epsilon)$  for  $\epsilon$ -nets. We state these results and prove them in Section 13.

The above lower bound proofs hold for specific “hard” families  $\mathcal{E}$ . This is in contrast with the statistical or oblivious settings in which a lower bound of  $\Omega((\text{VC}(\mathcal{E}) + \log(1/\delta))/\epsilon^2)$  applies to any class. We do not know whether an analogous result holds in the adversarial sampling setting and leave it as an open problem. We do show, however, that the linear dependence in  $d$  is necessary for any  $\mathcal{E}$ , as part of proving Theorem 2.2.

## 2.2 Online Learning

We continue with our main application to online learning. Consider the setting of online prediction with binary labels; a learning task in this setting can be described as a guessing game between a learner and an adversary. The game proceeds in rounds  $t = 1, \dots, T$ , each consisting of the following steps:

- The adversary selects  $(x_t, y_t) \in X \times \{0, 1\}$  and reveals  $x_t$  to the learner.
- The learner provides a prediction  $\hat{y}_t \in \{0, 1\}$  of  $y_t$  and announces it to the adversary.
- The adversary announces  $y_t$  to the learner.

The goal is to minimize the number of mistakes,  $\sum_t \mathbb{1}(y_t \neq \hat{y}_t)$ . Given a class  $\mathcal{E}$ , the *regret* of the learner w.r.t.  $\mathcal{E}$  is defined as the difference between the number of mistakes made by the learner and the number of mistakes made by the best  $E \in \mathcal{E}$ :

$$\sum_t \mathbb{1}(y_t \neq \hat{y}_t) - \min_{E \in \mathcal{E}} \sum_t \mathbb{1}(y_t \neq \mathbb{1}(x_t \in E)).$$

A class  $\mathcal{E}$  is *online-learnable* if there exists an online learner whose (expected) regret w.r.t. every adversary is at most  $R(T)$ , where  $R(T) = o(T)$ . (The amortized regret  $R(T)/T$  vanishes as  $T \rightarrow \infty$ .) Ben-David, Pál, and Shalev-Shwartz [BPS09] proved that for every class  $\mathcal{E}$ , the optimal regret  $R_T(\mathcal{E})$  satisfies

$$\Omega(\sqrt{d \cdot T}) \leq R_T(\mathcal{E}) \leq O(\sqrt{d \cdot T \log T}), \quad (2)$$

where  $d$  is the Littlestone dimension of  $\mathcal{E}$ , and left closing that gap as their main open question. Subsequently, Rakhlin, Sridharan, and Tewari [RST10, RST15a, RST15b] defined the notion of *Sequential Rademacher Complexity*, proved that it captures regret bounds in online learning in a general setting, and used it to re-derive Equation (2). They also asked as an open question whether the logarithmic factor in Equation (2) can be removed and pointed on difficulties to achieve this using some known techniques [RS14, RST15b].

We show that the sequential Rademacher complexity also captures the sample-complexity of  $\epsilon$ -approximations and bound it in the proof of Theorem 2.3. This directly implies a tight bound on online learning: (See Section 12 for more details.)

**Theorem 2.4** (Tight Regret Bounds in Online Learning). *Let  $\mathcal{E}$  be a class with Littlestone dimension  $d$ . Then the optimal regret bound in online learning  $\mathcal{E}$  is  $\Theta(\sqrt{d \cdot T})$ .*

The lower bound was shown by [BPS09]. We prove the upper bound in Section 12.

## 2.3 Applications and Extensions

We next discuss applications and extensions of our results.

**Epsilon Nets.** We also provide sample complexity bounds for producing  $\epsilon$ -nets: a subsample  $\bar{s}$  of the stream  $\bar{x}$  is an  $\epsilon$ -net if whenever  $E \in \mathcal{E}$  satisfies  $|E \cap \bar{x}| \geq \epsilon n$ , then  $\bar{s} \cap E \neq \emptyset$ . I.e. the subsample  $\bar{s}$  hits every  $E \in \mathcal{E}$  which contains at least an  $\epsilon$ -fraction of the items in the stream.

Epsilon nets are a fundamental primitive in computational geometry and in learning theory. In computational geometry this notion underlies fundamental algorithmic techniques, and in learning theory it is tightly linked to the learnability in the *realizable* setting. In that sense, it is analogous to  $\epsilon$ -approximations, which correspond to learnability in the *agnostic* setting.

In Section 7 we show that, like  $\epsilon$ -approximations,  $\epsilon$ -nets are also characterized by the Littlestone dimension; and similarly, our results here provide tight sample-complexity bounds.

**Maintaining An  $\epsilon$ -Approximation Continuously.** Some natural applications require that the sampler continuously maintains an  $\epsilon$ -approximation with respect to the prefix of the stream observed thus-far. To address such scenarios we slightly modify the adversarial sampling setting by allowing the sampler to delete items from its sample. In this modified setting, we prove that the classical *Reservoir sampler* [Vit85],  $\text{Res}(n, k)$  (see Section 5 for the precise definition), enjoys similar guarantees to those of Theorem 2.3 above. Concretely, the exact same bound of Theorem 2.3 is achieved by reservoir sampling if one is only interested in  $\epsilon$ -approximation at the end of the process; for continuous  $\epsilon$ -approximation, the same bound with an added term of  $O(\log \log(n))$  in the numerator suffices (see Theorem 11.1).

Notably, allowing deletions does not add significant power to the sampler, and in particular Theorem 2.2 still applies in this setting.

**ALLNs for Real-Valued Function Classes** The adversarial sampling setting naturally extends to real-valued function classes  $\mathcal{E}$ . Moreover, much of the machinery developed in this paper readily applies in this case. In particular, the relationship with the sequential Rademacher complexity is retained. Therefore, since the sequential Rademacher complexity captures regret bounds in online learning, this allows an automatic translation of regret bounds from online learning to sample complexity bounds in adversarial ULLNs w.r.t. real-valued function classes.<sup>3</sup>

**Algorithmic Applications** Part of the reason that the Fundamental Theorem of PAC Learning became a corner-stone in machine learning theory is due to its algorithmic implications. In particular, because it justifies the *Empirical Risk Minimization Principle* (ERM), which asserts that in order to learn a VC class, it suffices to minimize the empirical loss w.r.t. a random sample. This principle reduces the learning problem (of minimizing the loss w.r.t. an unknown distribution) to an optimization problem of minimizing the loss w.r.t. the (known) input sample.

It will be interesting to explore such implications in the adversarial setting. One promising direction is to use these sampling methods to design *lazy streaming/online algorithms*. That is, algorithms that update their internal state only on a small (random) substream. Intuitively, if that substream represents the entire stream in an appropriate way, then the performance of the algorithm will be satisfactory, and the gain in efficiency can be significant. In fact, our proof of Lemma 9.5 identifies and exploits such a phenomenon in online learning: we use a *lazy online learner* that updates its predictor rarely, only in a small random subsample of examples.

---

<sup>3</sup>The reduction from bounds on  $\epsilon$ -approximations to bounds on the sequential Rademacher complexity appear in Section 6. They rely on concentration inequalities for  $\{0,1\}$  valued random variables that have analogues for  $[0,1]$  valued random variables with the same guarantees. This enables a direct extension of this reduction.



### 3 Technical Overview

We next overview the technical parts in this work. We outline the proofs of the main theorems, and try to point out which technical arguments are novel, and which are based on known techniques. A more detailed overview of particular proofs is given in the dedicated sections.

#### 3.1 Upper Bounds

We begin with the sample-complexity upper bound, Theorem 2.3 (which is the longest and most technical derivation in this work).

**Reductions Between Samplers.** Our goal is to derive an upper bound for the Bernoulli, uniform, and reservoir samplers. In order to abstract out common arguments, we develop a general framework which serves to methodically transform sample-complexity bounds between the different samplers via a type of “online reductions”. This framework allows us to bound the sample-complexity with respect to one sampler, and *automatically* deduce them for the other samplers. The reduction relies on transforming one sampling scheme into another in an online fashion, and from a technical perspective, this boils down to coupling arguments, similar to coupling techniques in Markov Chains processes [LP17]. Section 10 contains a more detailed overview followed by the formal derivations.

**Upper Bounds for The Uniform Sampler.** Thus, for the rest of this overview we focus the sampling scheme to be the uniform sampler which uniformly draws a  $k$ -index-set  $I \subseteq [n]$ , and selects the subsample  $\bar{x}_I = (x_i : i \in I)$ . Our goal is to show that with probability  $\geq 1 - \delta$ ,

$$\sup_{E \in \mathcal{E}} \left| \frac{|\bar{x}_I \cap E|}{k} - \frac{|\bar{x} \cap E|}{n} \right| \leq O\left(\sqrt{\frac{d + \log(1/\delta)}{k}}\right), \quad (3)$$

where  $d$  is the Littlestone dimension of  $\mathcal{E}$  and  $\bar{x}$  is the *adversarially* produced sequence. The proof consists of two main steps which are detailed below.

##### 3.1.1 Step 1: Reduction to Online Discrepancy via Double Sampling

The first step in the proof consists of an online variant of the celebrated *double-sampling argument* due to [VC71]. This argument serves to replace the error w.r.t. the entire population by the error w.r.t. a small *test-set* of size  $k$ , thus effectively restricting the domain to the  $2k$  items in the union of the selected sample and the test-set. In more detail, let  $J \subseteq [n]$  be a uniformly drawn *ghost* subset of size  $k$  which is disjoint from  $I$ , and is not known to the adversary. Consider the maximal deviation between the sample  $\bar{x}_I$  and the “test-set”  $\bar{x}_J$ :

$$\sup_{E \in \mathcal{E}} \left| \frac{|\bar{x}_I \cap E|}{k} - \frac{|\bar{x}_J \cap E|}{k} \right|. \quad (4)$$

The argument proceeds by showing that for a typical  $J$ , the deviation w.r.t. the entire population  $\bar{x}$  in the LHS of Equation (3) has the same order of magnitude like the deviation w.r.t. the test-set  $\bar{x}_J$  in Equation (4) above. Hence, it suffices to bound (4).

In order to bound Equation (4), consider sampling  $I, J$  according to the following process: (i) First sample the  $2k$  indices in  $I \cup J$  uniformly from  $[n]$ , and reveal these  $2k$  indices to both players (in advance). (ii) Then, the sampler draws  $I$  from these  $2k$  indices in an online fashion (i.e., the adversary does not know in advance the sample  $I$ ). Intuitively, this modified process only helps the adversary who has the additional information of a superset of size  $2k$ , which contains  $I$ . *What we gain is that the modified process is essentially equivalent to reducing the horizon from  $n$  to  $2k$ .* The case of  $n = 2k$  can be interpreted as an online variant of the well-studied *Combinatorial Discrepancy* problem, which is described next.

**Online Combinatorial Discrepancy.** The online discrepancy game w.r.t.  $\mathcal{E}$  is a sequential game played between a painter and an adversary which proceeds as follows: at each round  $t = 1, \dots, 2k$  the adversary places an item  $x_t$  on the board, and the painter colors  $x_t$  in either red or blue. The goal of the painter is that each set in  $\mathcal{E}$  will be colored in a balanced fashion; i.e., if we denote by  $I$  the set of indices of items colored red, her goal is to minimize the discrepancy

$$\text{Disc}_{2k}(\mathcal{E}, \bar{x}, I) := \max_{E \in \mathcal{E}} \left| |\bar{x}_I \cap E| - |\bar{x}_{[2k] \setminus I} \cap E| \right|.$$

One can verify that minimizing the discrepancy is equivalent to minimizing Equation (4). Moreover, each of the samplers  $\text{Ber}(2k, 1/2)$  and  $\text{Uni}(2k, k)$  corresponds to natural coloring strategies of the painter; in particular,  $\text{Uni}(2k, k)$  colors a random subset of  $k$  of the items in red (and the rest in blue.) Thus, we focus now on analyzing the performance of  $\text{Uni}(2k, k)$  in the online discrepancy problem.

### 3.1.2 Step 2: From Online Discrepancy to Sequential Rademacher

Instead of analyzing the discrepancy of  $\text{Uni}(2k, k)$ , it will be more convenient to consider the discrepancy of  $\text{Ber}(2k, 1/2)$ , which colors each item in red/blue uniformly and independently of its previous choices. Towards this end, we show that these two strategies are essentially equivalent, using the reduction framework described at the beginning of this section.

The discrepancy of  $\text{Ber}(2k, 1/2)$  connects directly to the *Sequential Rademacher Complexity* [RS15], defined as the expected discrepancy  $\text{Rad}_{2k}(\mathcal{E}) = \mathbb{E} \text{Disc}_{2k}(\mathcal{E}, \bar{x}, I)$ , where the expectation is taken according to a uniformly drawn  $I \subseteq [2k]$ . (Which is precisely the coloring strategy of  $\text{Ber}(2k, 1/2)$ .)

### 3.1.3 Step 3.1: Bounding Sequential Rademacher Complexity – Oblivious Case

In what follows, it is convenient to set  $n = 2k$ . Our goal here is to bound  $\text{Rad}_n(\mathcal{E}) \leq O(\sqrt{d \cdot n})$ . As a prelude, it is instructive to consider the oblivious setting where the items  $x_1, \dots, x_n$  are fixed in advance, before they are presented to the painter. Here, the analysis is exactly as in the standard i.i.d. setting, and the sequential Rademacher complexity becomes the standard Rademacher complexity. Consider the following three approaches, in increasing level of complexity.

**First Approach: a Union Bound.** Assume  $\mathcal{E}$  is finite. Then, for each  $E \in \mathcal{E}$  it is possible to show by concentration inequalities that with high probability, the discrepancy  $||\bar{x}_I \cap E| - |\bar{x}_{[n] \setminus I} \cap E||$  is small. By applying a union bound over all  $E \in \mathcal{E}$ , one can derive that  $\text{Rad}_n(\mathcal{E}) \leq O(\sqrt{n \log |\mathcal{E}|})$ .

**Second Approach: Sauer-Shelah-Parles Lemma.** Since  $\mathcal{E}$  can be very large or even infinite, the bound in the previous attempt may not suffice. An improved argument relies on the celebrated Sauer-Shelah-Perles (SSP) Lemma [Sau72], which asserts that the number of distinct intersection-patterns of sets in  $\mathcal{E}$  with  $\{x_1, \dots, x_n\}$  is at most  $\binom{n}{\leq \text{VC}(\mathcal{E})} \leq O(n^{\text{VC}(\mathcal{E})})$ . The proof then follows by union bounding the discrepancy over  $\{\bar{x} \cap E : E \in \mathcal{E}\}$ , resulting in a bound of

$$O\left(\sqrt{n \log(n^{\text{VC}(\mathcal{E})})}\right) \leq O\left(\sqrt{\text{VC}(\mathcal{E})n \log n}\right),$$

which is off only by a factor of  $\sqrt{\log n}$ .

**Third Approach: Using Approximate Covers and Chaining.** Shaving the extra logarithmic factor is a non-trivial task which was achieved in the seminal work by Talagrand [Tal94] using a technique called *chaining* [Dud87]. It relies on the notion of *approximate covers*:

**Definition 3.1** (Approximate Covers). *A family  $\mathcal{C}$  is an  $\epsilon$ -cover of  $\mathcal{E}$  with respect to  $x_1, \dots, x_n$  if for every  $E \in \mathcal{E}$  there exists  $C \in \mathcal{C}$  such that  $E$  and  $C$  agree on all but at most  $\epsilon \cdot n$  of the  $x_i$ 's.*

In a nutshell, the chaining approach starts by finding covers  $\mathcal{C}_0, \mathcal{C}_1, \dots$  where  $\mathcal{C}_i$  is a  $2^{-i}$ -cover for  $\mathcal{E}$  w.r.t.  $\bar{x}$ , then writing the telescopic sum

$$\text{Disc}_n(\mathcal{E}, \bar{x}, I) = \text{Disc}_n(\mathcal{C}_0, \bar{x}, I) + \sum_{i=1}^{\infty} (\text{Disc}_n(\mathcal{C}_i, \bar{x}, I) - \text{Disc}_n(\mathcal{C}_{i-1}, \bar{x}, I))$$

and bounding each summand using a union bound.

Note that the SSP Lemma provides a bound of  $|\mathcal{C}| \leq \binom{n}{\leq \text{VC}(\mathcal{E})}$  in the case of  $\epsilon = 0$ , where  $d$  is the VC-dimension of  $\mathcal{E}$ . For  $\epsilon > 0$ , a classical result by Haussler [Hau92] asserts that every family admits an  $\epsilon$ -cover of size  $(1/\epsilon)^{O(d)}$ . The latter bound allows via chaining to remove the redundant logarithmic factor and bound  $\text{Rad}_n(\mathcal{E}) \leq O(\sqrt{\text{VC}(\mathcal{E})n})$ .

### 3.1.4 Step 3.2: Bounding Sequential Rademacher Complexity – Adversarial Case

We are now ready to outline the last and most technical step in this proof. Our goal is twofold: first, we discuss how previous work [BPS09, RST10] generalized the above arguments to the adversarial (or the online learning) model, culminating in a bound of the form  $\text{Rad}_n(\mathcal{E}) = O(\sqrt{dn \log n})$ . Then, we describe the proof approach for our improved bound of  $O(\sqrt{dn})$ .

**An  $O(\sqrt{dn \log n})$  Bound via Adaptive SSP.** First, the union bound approach generalizes directly to the adversarial setting. However, the second approach, via the SSP lemma, does not. The issue is that in the adversarial setting, the stream  $\bar{x}$  can depend on the coloring that the painter chooses, and hence  $\{E \cap \{x_1, \dots, x_n\} : E \in \mathcal{E}\}$  depends on the coloring as well. In particular, it is not possible to apply a union bound over a small number of such patterns. Moreover, it is known that a non-trivial bound depending only on the VC dimension and  $n$  does not exist [RST15a]. To overcome this difficulty we use an adaptive variant of the SSP Lemma due to [BPS09], which is based on the following notion:

**Definition 3.2** (Dynamic Sets). A dynamic set  $\mathbb{B}$  is an online algorithm that operates on a sequence  $\bar{x} = (x_1, \dots, x_n)$ . At each time  $t = 1, \dots, n$ , the algorithm decides whether to retain  $x_t$  as a function of  $x_1, \dots, x_t$ . Let  $\mathbb{B}(\bar{x})$  denote the set of elements retained by  $\mathbb{B}$  on a sequence  $\bar{x}$ .<sup>4</sup>

Ben-David, Pál, and Shalev-Shwartz [BPS09] proved that any family  $\mathcal{E}$  whose Littlestone dimension is  $d$  can be covered by  $\binom{n}{\leq d}$  dynamic sets. That is, for every  $n$  there exists a family  $\mathcal{C}$  of  $\binom{n}{\leq d}$  dynamic sets such that for every sequence  $\bar{x} = (x_1, \dots, x_n)$  and for every  $E \in \mathcal{E}$  there exists a dynamic set  $\mathbb{B} \in \mathcal{C}$  which agrees with  $E$  on the sequence  $\bar{x}$ , namely,  $\mathbb{B}(\bar{x}) = E \cap \bar{x}$ .

Using this adaptive SSP Lemma, one can proceed to bound the discrepancy as in the oblivious case by applying a union bound over the  $(2k)^d$  dynamic sets, and bounding the discrepancy with respect to each dynamic set using Martingale concentration bounds. Implementing this reasoning yields a bound of  $\text{Rad}_n(\mathcal{E}) \leq O(\sqrt{dn \log n})$  which is off by a logarithmic factor.

**Removing the Logarithmic Factor.** To adapt the chaining argument to the adversarial setting we first need to find small  $\epsilon$ -covers. This raises the following question:

*Can every Littlestone family be  $\epsilon$ -covered by  $(1/\epsilon)^{O(d)}$  dynamic sets?*

Unfortunately, we cannot answer this question and leave it for future work. In fact, [RST15b] identified a variant of this question as a challenge towards replicating the chaining proof in the online setting. To circumvent the derivation of dynamic approximate covers, we introduce a fractional variant which we term *fractional-covers*. It turns out that any Littlestone family admits “small” approximate fractional covers and these can be used to complete the chaining argument.

**Definition 3.3** (Approximate Fractional-Covers). A probability measure  $\mu$  over dynamic sets  $\mathbb{B}$  is called an  $(\epsilon, \gamma)$ -fractional cover for  $\mathcal{E}$  if for any  $\bar{x} = (x_1, \dots, x_n)$  and any  $E \in \mathcal{E}$ ,

$$\mu(\{\mathbb{B} : E \text{ and } \mathbb{B}(\bar{x}) \text{ agree on all but at most } \epsilon n \text{ of the } x_i\}) \geq 1/\gamma.$$

The parameter  $\gamma$  should be thought of as the size of the cover. Observe that fractional-covers are relaxations of covers: indeed, if  $\mathcal{C}$  is an  $\epsilon$ -cover for  $\mathcal{E}$  then the uniform distribution over  $\mathcal{C}$  is an  $(\epsilon, \gamma)$ -fractional cover for  $\mathcal{E}$  with  $\gamma = |\mathcal{C}|$ .

**Small Approximate Fractional-Covers Exist.** We prove Lemma 9.5 which asserts that every Littlestone family  $\mathcal{E}$  admits an  $(\epsilon, \gamma)$ -fractional cover of size

$$\gamma = (O(1)/\epsilon)^d.$$

This fractional cover is essentially a mixture of non-fractional covers for subsets of the sequence  $\bar{x}$  of size  $d/\epsilon$ . In more detail, the distribution over dynamic sets is defined by the following two-step sampling process: (1) draw a uniformly random subset  $\bar{s}$  of  $\bar{x}$  of size  $d/\epsilon$ , and let  $\mathcal{C}_s$  denote the (non-fractional) cover of  $\mathcal{E}$  with respect to  $\bar{s}$ , which is promised by the dynamic variant of the SSP-Lemma. (2) Draw  $\mathbb{B}$  from the uniform distribution over  $\mathcal{C}_s$ .

We outline the proof that this is an  $(\epsilon, \gamma)$ -fractional cover with  $\gamma = O(1/\epsilon)^d$ . Fixing  $E$  and  $\bar{x}$ , our goal is to show that with probability at least  $1/\gamma$  over  $\mu$ , the drawn  $\mathbb{B}$  agrees with  $E$  on all but at most  $\epsilon \cdot n$  elements of  $\bar{x}$ . This relies on the following two arguments: (1) For every  $\bar{s}$  there

<sup>4</sup>[BPS09] refers to dynamic-sets as experts, which is compatible with the terminology of online learning.

exists  $\mathbb{B}_{\bar{s}} \in \mathcal{C}_s$  that agrees with  $E$  on  $\bar{s}$ ; and (2) it can be shown that with high probability over the selection of the subset  $\bar{s}$ ,  $\mathbb{B}_{\bar{s}}$  agrees with  $E$  on all but at most  $\epsilon n$  of the stream  $\bar{x}$ . We call such values of  $\bar{s}$  as *good*, and conclude from the two steps above that

$$\begin{aligned} \Pr_{\mathbb{B} \sim \mu} [\mathbb{B} \text{ agrees with } E \text{ on } (1 - \epsilon)n \text{ of the } x_i\text{'s}] &\geq \Pr[\bar{s} \text{ is good}] \Pr_{\mathbb{B} \sim \text{uniform}(\mathcal{C}_s)} [\mathbb{B} = \mathbb{B}_{\bar{s}}] \\ &\geq \frac{1}{2} \cdot \frac{1}{|\mathcal{C}_s|} \geq \frac{1}{2^{\binom{d}{\leq d}/\epsilon}} \geq \Omega(\epsilon)^d \geq \frac{1}{\gamma}. \end{aligned}$$

We further comment on the proof that  $\bar{s}$  is *good* with high probability: the proof relies on analyzing a *lazy* online learner that updates its internal state only once encountering elements from  $\bar{s}$ . We show that if  $\bar{s}$  is drawn uniformly, then with high probability such a learner will make  $\leq \epsilon \cdot n$  mistakes and this will imply that w.h.p.  $\mathbb{B}_{\bar{s}}$  agrees with  $E$  on  $(1 - \epsilon)n$  stream elements. We refer the reader to Section 9.2.2 for the proof.

**Chaining with Fractional Covers: Challenges and Subtleties.** Here, we discuss how approximate fractional covers are used to bound the sequential Rademacher complexity. We do so by describing how to modify the bound that uses 0-covers to use  $(0, \gamma)$ -fractional covers instead. Recall that this argument goes by two steps: (1) bounding the discrepancy for each dynamic set in the cover, and (2) arguing by a union bound that, with high probability the discrepancies of *all* dynamic sets in the cover are bounded. In comparison, with fractional covers, the second step is modified to: (2') arguing that with high probability (over the random coloring), the discrepancies of *nearly all* the dynamic sets are bounded. In particular, if more than a  $(1 - \gamma)$ -fraction of the dynamic sets have bounded discrepancies, then the discrepancies of all sets in  $\mathcal{E}$  are bounded. Indeed, this follows since every  $E \in \mathcal{E}$  is covered by at least a  $\gamma$ -fraction of the dynamic-sets, and therefore, the pigeonhole principle implies that at least one such dynamic set also has bounded discrepancy, and hence  $E$  has bounded discrepancy as well.

We note that multiple further technicalities are required to generalize the chaining technique for fractional covers and refer the reader to Section 9.3 for a short overview of this method followed by its adaptation to the adversarial setting.

## 3.2 Lower Bounds

Beyond the  $\Omega((d + \log(1/\delta))/\epsilon^2)$  lower bound for oblivious samplers, which follows immediately from the VC literature, we prove several non-trivial lower bounds in other contexts. We distinguish between two types of approaches used to derive our lower bounds, described below. As the proofs are shorter than those of the upper bounds and more self-contained, we omit the exact technical details of the proofs in this overview and refer the reader to Section 13.

**Universal Lower Bound by Adversarial Arguments.** The main lower bound in [BEY20] exhibits a separation between the static and adversarial setting by proving an adversarial lower bound for the family of one-dimensional *thresholds*. We identify that their proof implicitly constructs a tree as in the definition of the Littlestone dimension, and generalize their argument to derive an  $\Omega(d)$  lower bound for *all* families of Littlestone dimension  $d$ . For more details, see Theorem 13.1.

**Lower Bounds on the Minimum Sizes of  $\epsilon$ -Approximations/Nets.** These lower bounds actually exhibit a much stronger phenomenon, showing that small  $\epsilon$ -approximations/nets *do not exist* for some families  $\mathcal{E}$ . Thus, obviously, these cannot be captured by a sample of the same size.

It is natural to seek lower bounds of this type in the VC-literature. The main challenge is that many of the known lower bounds apply for geometric VC classes whose Littlestone dimension is unbounded. To overcome this, we present two lower bounds where  $\text{Ldim}$  can be controlled: one for  $\epsilon$ -approximation, which carefully analyzes a simple randomized construction, and another for  $\epsilon$ -nets, which combines intersection properties of lines in the projective plane with probabilistic arguments. For more details, see Theorems 13.2 and 13.3 respectively.

## 4 Related Work

### 4.1 VC Theory

As suggested by the title, the results presented by this work are inspired by uniform laws of large numbers in the statistical i.i.d. setting and in particular by VC theory. (A partial list of basic manuscripts on this subject include [VC71, VC74, Dud84, Vap98].) Moreover, the established equivalence between online learning and adversarial laws of large numbers is analogous to the the equivalence between PAC learning and uniform laws of large numbers in the i.i.d. setting. (See e.g. [VC71, BEHW89, Bou04, SSSS10, SSBD14].) From a technical perspective, our approach for deriving sample complexity upper bound is based on the chaining technique [Dud73, Dud78, Dud87], which was analogously used to establish optimal sample complexity bounds in the statistical setting [Tal94]. (The initial bounds by [VC71] are off by a  $\log(1/\epsilon)$  factor.)

From the lower bound side, our proofs are based on ideas originated from combinatorial discrepancy and  $\epsilon$ -approximations. (E.g., [MWW93]; see the book by Matoušek [Mat09] for a text-book introduction.)

### 4.2 Online Learning

The first works in online learning can be traced back to [Rob51, Bla56, Bla54, Han57]. In terms of learning binary functions, Littlestone’s dimension was first proposed in [Lit88] to characterize online learning in the realizable (noiseless) setting. The agnostic (noisy) setting was first proposed by [Hau92] in the statistical model and later extended to the online setting by [LW94] who studied function-classes of bounded cardinality and then by [BPS09] and [RST10] who provided both upper and lower bounds with only a logarithmic gap.

We note that Rakhlin, Sridharan, and Tewari [RST10, RST15a, RST15b], in the same line of work that proved the equivalence between online learning and sequential Rademacher complexity, analyzed uniform martingales laws of large numbers in the context of online learning. These laws of large numbers are conceptually different from ours: roughly, they assert uniform concentration of certain properties of martingales, where the uniformity is over a given family of martingales. In particular, in contrast with our work, there is no aspect of sub-sampling in these laws. Below, we compare their techniques to those of this paper:

- [RST10] used a symmetrization argument to reduce from Martingale quantities relating to online learning to the Rademacher complexity. This does not reduce the effective sample size, which is what we achieve using the double sampling argument.

- [RST10] developed methods for analyzing the sequential Rademacher complexity. In particular, they developed a notion of covering numbers that is generally more powerful than the *non-fractional* cover that uses dynamic sets, which was developed by [BPS09] and was the baseline for our analysis. Yet, obtaining tight bound on the sequential Rademacher of Littlestone classes remained open.
- Reductions between sampling schemes did not appear in the above work as they did not study sampling.

### 4.3 Streaming Algorithms

The streaming model of computation is useful when analyzing massive datasets [AMS99]. There is a wide variety of algorithms for solving different tasks. One common method that is useful for various approximation tasks in streaming is random sampling. To approximate a function  $f$ , each element is sampled with some small probability  $p$ , and at the end, the function  $f$  is computed on the sample. For tasks such as computing a center point of a high-dimensional dataset, where the objective is (roughly speaking) preserved under taking an  $\epsilon$ -approximation, this can result in improved space complexity and running time. Motivated by streaming applications, Ben-Eliezer and Yogev [BEY20] proposed the adversarial sampling model that we study in this paper, and proved preliminary bounds on it. Their main result, a weaker quantitative analogue of our Theorem 2.3, is an upper bound of  $O((\log(|\mathcal{E}|) + \log(1/\delta))/\epsilon^2)$  for any finite family  $\mathcal{E}$ .

Streaming algorithms in the adversarial setting is an emerging topic that is not well understood. Hardt and Woodruff [HW13] showed that linear sketches are inherently *non-robust* and cannot be used to compute the Euclidean norm of its input (where in the static setting they are used mainly for this reason). Naor and Yogev [NY15] showed that Bloom filters are susceptible to attacks by an adversarial stream of queries. On the positive side, several recent works [BJWY20, HKM<sup>+</sup>20, WZ20] present generic compilers that transform non-robust randomized streaming algorithms into efficient adversarially robust ones, for various classical problems such as distinct elements counting and  $F_p$ -sampling, among others.

## 5 Preliminaries

### 5.1 Basic Definitions: Littlestone Dimension and Sampling Schemes

**Littlestone Dimension** Let  $X$  be a domain and let  $\mathcal{E}$  be a family of subsets of  $X$ . The definition of the *Littlestone Dimension* [Lit88], denoted  $\text{Ldim}(\mathcal{E})$ , is given using mistake-trees: these are binary decision trees whose internal nodes are labelled by elements of  $X$ . Any root-to-leaf path corresponds to a sequence of pairs  $(x_1, y_1), \dots, (x_d, y_d)$ , where  $x_i$  is the label of the  $i$ 'th internal node in the path, and  $y_i = 1$  if the  $(i + 1)$ 'th node in the path is the right child of the  $i$ 'th node, and otherwise  $y_i = 0$ . We say that a tree  $T$  is shattered by  $\mathcal{E}$  if for any root-to-leaf path  $(x_1, y_1), \dots, (x_d, y_d)$  in  $T$  there is  $E \in \mathcal{E}$  such that  $x_i \in E \iff y_i = +1$ , for all  $i \leq d$ .  $\text{Ldim}(\mathcal{E})$  is the depth of the largest complete tree shattered by  $\mathcal{E}$ , with the convention that  $\text{Ldim}(\emptyset) = -1$ . See Figure 1 for an illustration.

**Sampling Algorithms** Our results are achieved by three of the simplest and most commonly used sampling procedures: Bernoulli sampling, uniform sampling, and reservoir sampling.

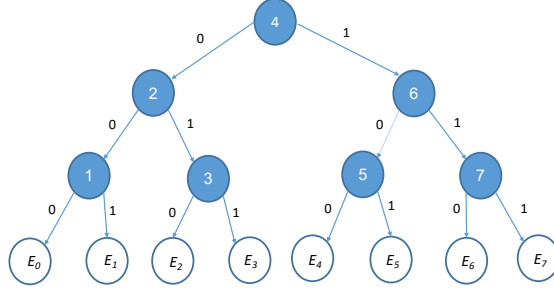


Figure 1: A tree shattered by the class  $\mathcal{E}$  containing all thresholds  $E_i \subseteq \{1, 2, \dots, 7\}$ , where  $E_i = \{1, \dots, i\}$ .

- **Bernoulli sampling:**  $\text{Ber}(n, p)$  samples the element arriving in each round  $i \in [n]$  independently with probability  $p$ .
- **uniform sampling:**  $\text{Uni}(n, k)$  randomly draws  $k$  indices  $1 \leq i_1 < \dots < i_k \leq n$  and samples the elements arriving at rounds  $i_1, \dots, i_k$ .<sup>5</sup>
- **Reservoir sampling:**  $\text{Res}(n, k)$  [Vit85] maintains a sample of size  $k$  at all times using insertions and deletions: the first  $k$  elements are always added to the sample, and for any  $i > k$ , with probability  $k/i$  the element arriving in round  $i$  is added to the sample while one of the existing elements (picked uniformly) is removed from the sample.

## 5.2 Notation

- **Random variables** are denoted in a bold font.
- **Universal constants:** let  $C, c, C', \dots$  denote universal numerical constants, that are independent of the problem parameters. Further, the values of these constants can change from the left-hand side to the right-hand side of some inequalities. We use  $C_0, C_1, \dots$  to denote universal constants whose values are fixed.
- **Sampling schemes:** We use  $I$  to denote the set of indices of elements sampled by the algorithm if they are sampled from a scheme with no deletions (e.g. Bernoulli or uniform sampling), and use  $\bar{I} = (I_1, \dots, I_n)$  to denote a sampling scheme with deletions, where  $I_j$  is the set of indices of elements retained after round  $j \in [n]$ .
- **Adversaries:** We denote the set of adversaries that generate a stream of size  $n$  by  $\text{Adv}_n$  and commonly denote adversaries by  $\mathcal{A}$ . We assume the adversary to be deterministic: since the sampler has a fixed strategy, we can make this assumption without limiting the generality of the theorem.
- **The stream and subsets of it:** Let  $\bar{x} = (x_1, \dots, x_n)$  denote the stream where  $x_i \in X$ . We set  $\bar{x}(\mathcal{A}, I)$  to denote the stream presented by the adversary when the sampler samples elements indexed by  $I$ . Notice that  $\bar{x}(\mathcal{A}, I)_t$  depends only on  $\mathcal{A}$  and  $I \cap [t-1]$ , since the  $t$ -th stream element is presented before the adversary knows if element  $t$  is added to the

<sup>5</sup>Note that the uniform sampler can be implemented efficiently in an online way: after  $i$  rounds, the probability that the next element  $x_{i+1}$  will be sampled depends only on  $i, n$ , and the number of elements sampled so far.



sample. Given a subset  $J \subseteq [n]$  we let  $\bar{x}_J = \{x_j : j \in J\}$ . By abuse of notation, we may use  $\bar{x}$  also to denote the multiset  $\{x_1, \dots, x_n\}$ , allowing operations such as set intersection.

### 5.3 Additional Central Definitions

We define the central notions used in this proof, starting with the approximation rate in  $\epsilon$ -approximations:

**Definition 5.1.** Given an arbitrary family  $\mathcal{E}$  over a domain  $X$ , an adversary  $\mathcal{A} \in \text{Adv}_n$  and a subset  $I \subseteq [n]$ , let  $\bar{x} = \bar{x}(\mathcal{A}, I)$  and define

$$\text{App}_{\mathcal{A}, I}(\mathcal{E}) = \text{App}_{\mathcal{A}, I} := \max_{E \in \mathcal{E}} \left| \frac{|E \cap \bar{x}|}{n} - \frac{|E \cap \bar{x}_I|}{k} \right|.$$

Secondly, define the notion of online discrepancy:

**Definition 5.2.** Let  $\mathcal{E}$  denote an arbitrary family over  $X$ , let  $\mathcal{A} \in \text{Adv}_n$ , let  $I \subseteq [n]$  and let  $\bar{x} = \bar{x}(\mathcal{A}, I)$ . The online discrepancy is defined by:

$$\text{Disc}_{\mathcal{A}, I}(\mathcal{E}) = \text{Disc}_{\mathcal{A}, I} := \max_{E \in \mathcal{E}} \left| |E \cap \bar{x}_I| - |E \cap \bar{x}_{[n] \setminus I}| \right|.$$

The sequential Rademacher complexity is just the expected discrepancy:

**Definition 5.3.** The sequential Rademacher complexity is defined as

$$\text{Rad}_T(\mathcal{E}) := \mathbb{E}_{I \sim \text{Ber}(n, 1/2)} [\text{Disc}_{\mathcal{A}, I}(\mathcal{E})].$$

The next definition is used in the proof for  $\epsilon$ -nets. It defines an indicator to whether there exist  $E \in \mathcal{E}$  that is well represented in the stream but not sufficiently represented in the sample.

**Definition 5.4.** Fix  $n, \bar{m}, \underline{m} \in \mathbb{N}$  such that  $0 \leq \underline{m} \leq \bar{m} \leq n$ , let  $\mathcal{A} \in \text{Adv}_n$  and  $I \subseteq [n]$ . Denote  $\bar{x} = \bar{x}(\mathcal{A}, I)$ . Define

$$\text{Net}_{\mathcal{A}, I, \bar{m}, \underline{m}}^n(\mathcal{E}) = \text{Net}_{\mathcal{A}, I, \bar{m}, \underline{m}}^n = \begin{cases} 1 & \exists E \in \mathcal{E}, |\bar{x} \cap E| \geq \bar{m} \text{ and } |\bar{x}_I \cap E| \leq \underline{m} \\ 0 & \text{otherwise} \end{cases}$$

Notice that  $\text{Net}_{\mathcal{A}, I, \epsilon n, 0}^n$  is an indicator to whether  $\bar{x}_I$  fails to be an  $\epsilon$ -net for  $\bar{x}$ .

### 5.4 Sampling Without Replacement

Here we present technical probabilistic lemmas for sampling without replacement that are used in the proof. The proofs of these lemmas appear in Appendix A.2.

For a sample  $I \sim \text{Uni}(n, k)$  chosen without replacement, one would like to estimate the size of the intersection of  $I$  with any fixed set  $U \subseteq [n]$ . The next lemma bounds the variance of the intersection:

**Lemma 5.5.** Let  $n, k$  such that  $n \geq k$ , let  $U \subseteq [n]$ , and let  $I \sim \text{Uni}(n, k)$ . Then,  $\mathbb{E}[|U \cap I|] = |U|k/n$  and  $\text{Var}(|U \cap I|) \leq |U|k/n$ .

Further, exponential tail bounds can also be obtained:

**Lemma 5.6** ([Cha05],[BM15]). *Let  $n, k \in \mathbb{N}$  such that  $n \geq k$ . Let  $U \subseteq [n]$ . Then, the following holds:*

1. For any  $t \geq 0$ ,

$$\Pr_{I \sim \text{Uni}(n,k)} \left[ \left| \frac{|I \cap U|}{k} - \frac{|U|}{n} \right| \geq t \right] \leq 2 \exp(-2t^2k).$$

2. For any  $\alpha \in [0, 1]$ ,

$$\Pr_{I \sim \text{Uni}(n,k)} \left[ \left| \frac{|I \cap U|}{k} - \frac{|U|}{n} \right| \geq \alpha \frac{|U|}{n} \right] \leq 2 \exp\left(-\frac{\alpha^2 k |U|}{6n}\right).$$

## 6 Epsilon Approximations

Below, we prove Theorem 2.3. We start with a more formal statement of the theorem:

**Theorem 6.1.** *Let  $\mathcal{E}$  denote a family of Littlestone dimension  $d$ , let  $\delta, \epsilon \in (0, 1/2)$ ,  $n \in \mathbb{N}$ ,  $\mathcal{A} \in \text{Adv}_n$  and  $p \in [0, 1]$ . Define  $k = \lfloor np \rfloor$ . If  $n \geq 3k$  then, for any  $\delta > 0$ ,*

$$\Pr_I \left[ \text{App}_{\mathcal{A}, I}(\mathcal{E}) \geq C \sqrt{\frac{d + \log(1/\delta)}{k}} \right] \leq \delta,$$

where  $I$  is drawn either from  $\text{Uni}(n, k)$ ,  $\text{Ber}(n, p)$  or  $\text{Res}(n, k)$  and  $C > 0$  is a universal constant.

Note that the requirement  $n \geq 3k$  is merely technical; when  $n$  is smaller than that, one can just add all  $n$  elements to the sample and obtain a 0-approximation trivially.

Here we prove that the *uniform sample*  $\text{Uni}(n, k)$  is an  $\epsilon$ -approximation and in Section 10 we show a reduction to Bernoulli and reservoir sampling. The bound for the uniform sampler consists of three steps (see Section 3). The first step utilizes the double sampling argument, to bound the approximation error  $\text{App}_{\mathcal{A}, I}$  in terms of the discrepancy corresponding to a sampler  $\text{Uni}(2k, k)$ :

**Lemma 6.2.** *Let  $\mathcal{E}$  denote an arbitrary family of subsets from some universe. Fix  $k, n \in \mathbb{N}$  such that  $n \geq 2k$ . Then, for any  $t \geq 0$  and  $\delta \in (0, 1)$ ,*

$$\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{I \sim \text{Uni}(n,k)} [\text{App}_{\mathcal{A}, I}(\mathcal{E}) > t] \leq 2 \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Uni}(2k,k)} [\text{Disc}_{\mathcal{A}', I'}(\mathcal{E}) > tk - \sqrt{Ck}]. \quad (5)$$

The technique of double sampling is presented in Section 8 and the proof of Lemma 6.2 is in Section 8.1.

Applying Lemma 6.2, we are left with bounding  $\text{Disc}_{\mathcal{A}, I}$  where  $\mathcal{A} \in \text{Adv}_{2k}$  and  $I \sim \text{Uni}(2k, k)$ . However, it is easier to analyze a random sample  $I' \sim \text{Ber}(2k, 1/2)$  due to the independence of the coordinates. Hence, we prove the following lemma:

**Lemma 6.3.** *Let  $\mathcal{E}$  denote an arbitrary family. Then, for any  $t > 0$  and  $\delta > 0$ ,*

$$\max_{\mathcal{A} \in \text{Adv}_{2k}} \Pr_{I \sim \text{Uni}(2k,k)} [\text{Disc}_{\mathcal{A}, I}(\mathcal{E}) > t] \leq \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Ber}(2k, 1/2)} \left[ \text{Disc}_{\mathcal{A}', I'}(\mathcal{E}) > t - \sqrt{Ck \log(1/\delta)} \right] + \delta. \quad (6)$$

The technique of reducing between sampling schemes is in Section 10 and the proof of Lemma 6.3 is in Section 10.5.1.

The last step is to bound  $\text{Disc}_{\mathcal{A},I}$  for a sample  $I$  that is drawn Bernoulli  $1/2$ , for classes of bounded Littlestone dimension.

**Lemma 6.4.** *Let  $\mathcal{E}$  be of Littlestone dimension  $d$  and let  $\mathcal{A} \in \text{Adv}_{2k}$ . Then, for all  $\delta > 0$ ,*

$$\Pr_{I \sim \text{Ber}(2k, 1/2)} \left[ \text{Disc}_{\mathcal{A},I}(\mathcal{E}) > C \sqrt{k(d + \log(1/\delta))} \right] \leq \delta.$$

Lemma 6.4 is proved in Section 9.1.1. Combining Lemma 6.2, Lemma 6.3 and Lemma 6.4, We conclude that the uniform sample as an  $\epsilon$ -approximation, as summarized below:

**Theorem 6.5.** *Let  $\mathcal{E}$  denote a family of Littlestone dimension  $d$ , let  $\delta, \epsilon \in (0, 1/2)$ ,  $n \in \mathbb{N}$ ,  $\mathcal{A} \in \text{Adv}_n$  and  $k \in \mathbb{N}$  such that  $n \geq 2k$ . Then, for any  $\delta > 0$ ,*

$$\Pr_{I \sim \text{Uni}(n, k)} \left[ \text{App}_{\mathcal{A},I}(\mathcal{E}) \geq C \sqrt{\frac{d + \log(1/\delta)}{k}} \right] \leq \delta,$$

where  $C > 0$  is a universal constant.

## 7 Epsilon Nets

We prove that the three sampling schemes discussed in this paper sample  $\epsilon$ -nets with high probability, as stated below:

**Theorem 7.1.** *Let  $\mathcal{E}$  denote a family of Littlestone dimension  $d$ , let  $\epsilon \in (0, 1/2)$ ,  $n \in \mathbb{N}$ ,  $\mathcal{A} \in \text{Adv}_n$  and  $p \in [0, 1]$ . Define  $k = \lfloor np \rfloor$ . If  $n \geq 3k$  then*

$$\Pr_I [\text{Net}_{\mathcal{A}, I, \epsilon n, 0}^n(\mathcal{E}) = 1] \leq \left( \frac{Ck}{d} \right)^d \exp(-cek)$$

where  $I$  is drawn either  $\text{Uni}(n, k)$ ,  $\text{Ber}(n, p)$  or  $\text{Res}(n, k)$  and  $C > 0$  is a universal constant.

In this section, we prove Theorem 7.1 for the uniform sampler. Reductions to the other sampling schemes are given in Section 10.

Below the main lemmas are presented, starting with a reduction from a stream of size  $n$  to  $2k$ , via the technique of double sampling, presented in Section 8:

**Lemma 7.2.** *Let  $\mathcal{E}$  be some family and  $n, k \in \mathbb{N}$  be integers such that  $n \geq 2k$  and  $k \geq C/\epsilon$ . Then,*

$$\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{I \sim \text{Uni}(n, k)} [\text{Net}_{\mathcal{A}, I, \epsilon \cdot n, 0}^n(\mathcal{E}) = 1] \leq 2 \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Uni}(2k, k)} [\text{Net}_{\mathcal{A}', I', \epsilon/4 \cdot 2k, 0}^{2k}(\mathcal{E}) = 1].$$

The proof appears in Section 8.2. It would be desirable to replace the uniform sample  $\text{Uni}(2k, k)$  with a Bernoulli sample, since it selects each coordinate independently. In particular, we will show that the probability of  $\text{Uni}(2k, k)$  to fail to be an  $\epsilon$ -net is bounded in terms of the probability of  $\text{Ber}(2k, 1/8)$ . Intuitively, this follows from the fact that a sample drawn  $\text{Uni}(2k, k)$  nearly contains a sample  $\text{Ber}(2k, 1/8)$  in some sense. The formal statement is below:

**Lemma 7.3.** Let  $\mathcal{E}$  denote some family,  $\epsilon \in (0, 1)$ , and  $k \in \mathbb{N}$ . Then,

$$\begin{aligned} & \max_{\mathcal{A} \in \text{Adv}_{2k}} \Pr_{I \sim \text{Uni}(2k, k)} [\text{Net}_{\mathcal{A}, I, \epsilon \cdot 2k, 0}^{2k}(\mathcal{E}) = 1] \\ & \leq \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Ber}(2k, 1/8)} [\text{Net}_{\mathcal{A}', I', \epsilon \cdot 2k, \epsilon/16 \cdot 2k}^{2k}(\mathcal{E}) = 1] + 2 \exp(-c\epsilon k). \end{aligned}$$

The framework to reduce between sampling schemes is presented in Section 10 and the proof of Lemma 7.3 is presented in Section 10.5.2. Lastly, we bound the error probability corresponding to Bernoulli sampling, for classes of bounded Littlestone dimension, using the technique of covering numbers presented in Section 9.

**Lemma 7.4.** Let  $\mathcal{E}$  denote a family of Littlestone dimension  $d$ , let  $m, k \in \mathbb{N}$  such that  $m \leq 2k$ , and let  $p \in (0, 1)$ . Then,

$$\Pr_{I \in \text{Ber}(2k, p)} [\text{Net}_{\mathcal{A}, I, m, mp/2}^{2k}(\mathcal{E}) = 1] \leq \left(\frac{Ck}{d}\right)^d \exp(-cmp).$$

The three lemmas stated above imply the following theorem, that is a special case of Theorem 7.1 for the uniform sampler:

**Theorem 7.5.** Let  $\mathcal{E}$  denote a family of Littlestone dimension  $d$ , let  $\epsilon \in (0, 1/2)$ ,  $n \in \mathbb{N}$ ,  $\mathcal{A} \in \text{Adv}_n$  and  $k \in \mathbb{N}$  that satisfies  $n \geq 2k$ . Then

$$\Pr_{I \sim \text{Uni}(n, k)} [\text{Net}_{\mathcal{A}, I, \epsilon n, 0}^n(\mathcal{E}) = 1] \leq \left(\frac{Ck}{d}\right)^d \exp(-c\epsilon k),$$

where  $C > 0$  is a universal constant.

*Proof.* By Lemma 7.2,

$$\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{I \sim \text{Uni}(n, k)} [\text{Net}_{\mathcal{A}, I, \epsilon \cdot n, 0}^n(\mathcal{E}) = 1] \leq 2 \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Uni}(2k, k)} [\text{Net}_{\mathcal{A}', I', \epsilon/4 \cdot 2k, 0}^{2k}(\mathcal{E}) = 1].$$

Applying Lemma 7.3 while substituting  $\epsilon$  with  $\epsilon/4$ ,

$$\begin{aligned} \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Uni}(2k, k)} [\text{Net}_{\mathcal{A}', I', \epsilon/4 \cdot 2k, 0}^{2k}(\mathcal{E}) = 1] & \leq \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Ber}(2k, 1/8)} [\text{Net}_{\mathcal{A}', I', \epsilon/4 \cdot 2k, \frac{2k\epsilon}{64}}^{2k} = 1] \\ & + 2 \exp(-c\epsilon k). \end{aligned}$$

Applying Lemma 7.4 with  $p = 1/8$  and  $m = \epsilon/4 \cdot 2k$ ,

$$\max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Ber}(2k, 1/8)} [\text{Net}_{\mathcal{A}', I', \epsilon/4 \cdot 2k, \epsilon/64 \cdot 2k}^{2k} = 1] \leq \left(\frac{Ck}{d}\right)^d \exp(-c\epsilon k).$$

□

## 8 Double Sampling

Let  $n \in \mathbb{N}$  denote the stream length and assume that the sample is of size  $k \leq n/2$ . If  $n \gg k$ , it may be difficult to analyze the sample directly, since each element is selected with small probability and the universe is very large. This section presents a framework to replace the stream of size  $n$  with a stream of size  $2k$ . Then, this framework is used to prove Lemma 6.2 and Lemma 7.2 in Section 8.1 and Section 8.2, respectively.

Let  $\bar{x} \in X^n$  denote the stream and  $I \subseteq [n]$  be the index set of the sample, that has cardinality  $|I| = k$ . Let  $f: X^k \times X^n \rightarrow \{0, 1\}$  denote some function and we view  $f(\bar{x}_I, \bar{x})$  as some indicator of whether  $\bar{x}_I$  fails to approximate the complete sample  $\bar{x}$ . For example,  $f(\bar{x}_I, \bar{x})$  could indicate whether  $\bar{x}_I$  fails to be an  $\epsilon$ -approximation for  $\bar{x}$  with respect to some family  $\mathcal{E}$ . Denote by  $\bar{x} = \bar{x}(\mathcal{A}, I)$  the stream generated by the adversary  $\mathcal{A} \in \text{Adv}_n$  when the sample is indexed by  $I$  and we would like to bound  $\Pr_{I \sim \text{Uni}(n, k)}[f(\bar{x}_I, \bar{x}) = 1]$ . We would like to bound it by a different term that corresponds to only  $2k$  elements. For this purpose, let  $f': X^k \times X^k \rightarrow \{0, 1\}$  be another function, where  $f'(\bar{x}_I, \bar{x}_J)$  is an indicator of whether  $\bar{x}_I$  fails to approximate  $\bar{x}_J$ . For example,  $f'(\bar{x}_I, \bar{x}_J)$  can indicate whether  $\bar{x}_I$  fails to be an  $\epsilon$ -approximation for  $\bar{x}_J$ .

Let  $\mathcal{A}' \in \text{Adv}_{2k}$ ,  $I' \sim \text{Uni}(2k, k)$  and  $\bar{x}' = \bar{x}(\mathcal{A}', I')$  denote the stream of size  $2k$  generated by  $\mathcal{A}'$  with sample-index  $I'$ . The following lemma gives a condition under which the probability that  $f(\bar{x}_I, \bar{x}) = 1$  can be bounded in terms of the probability that  $f(\bar{x}'_{I'}, \bar{x}'_{[n] \setminus I'}) = 1$ .

**Lemma 8.1.** *Let  $I \sim \text{Uni}(n, k)$  and let  $J$  be distributed uniformly over all subsets of  $[n] \setminus I$  of size  $k$ , conditioned on  $I$ . Let  $f: X^k \times X^n \rightarrow \{0, 1\}$ . Let  $I' \sim \text{Uni}(2k, k)$  and  $f': X^k \times X^k \rightarrow \{0, 1\}$ .*

*Assume that for every  $\bar{x}$  and  $I$  that satisfy  $f(\bar{x}_I, \bar{x}) = 1$ , it further holds that*

$$\Pr_J [f'(\bar{x}_I, \bar{x}_J) = 1 \mid I = I] \geq 1/2. \quad (7)$$

Then,

$$\max_{\mathcal{A} \in \text{Adv}_n} \Pr_I [f(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)) = 1] \leq 2 \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I'} [f'(\bar{x}(\mathcal{A}', I')_{I'}, \bar{x}(\mathcal{A}', I')_{[2k] \setminus I'}) = 1].$$

To give some intuition on the condition (7), assume again that  $f(\bar{x}_I, \bar{x})$  denotes an indicator of whether  $\bar{x}_I$  fails to be an  $\epsilon$ -approximation to  $\bar{x}$  and  $f'(\bar{x}_I, \bar{x}_J)$  denotes whether  $\bar{x}_I$  fails to be an  $\epsilon'$ -approximation to  $\bar{x}_J$ , where  $\epsilon'$  is slightly larger than  $\epsilon$ . By concentration properties, if  $f(\bar{x}_I, \bar{x}) = 1$  then with high probability over  $J$ ,  $f'(\bar{x}_I, \bar{x}_J) = 1$  as well.

The proof of Lemma 8.1 consists of two steps. In the first step, an index-set  $J$  is drawn uniformly at random from all the subsets of  $[n] \setminus I$  of size  $k$ , conditioned on  $I$ . The set  $J$  is called a *ghost sample*, as it is used only for the analysis and in particular, the adversary is unaware of  $J$ . The following lemma shows that under the condition (7), we can bound  $f(\bar{x}_I, \bar{x})$  in terms of  $f'(\bar{x}_I, \bar{x}_J)$ .

**Lemma 8.2.** *Let  $f: X^k \times X^n \rightarrow \{0, 1\}$ ,  $f': X^k \times X^k \rightarrow \{0, 1\}$  and  $\mathcal{A} \in \text{Adv}_n$ . Assume that for every  $\bar{x}$  and  $I$  that satisfy  $f(\bar{x}_I, \bar{x}) = 1$ , it further holds that*

$$\Pr_J [f'(\bar{x}_I, \bar{x}_J) = 1 \mid I = I] \geq 1/2. \quad (8)$$

Then,

$$\Pr_I [f(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)) = 1] \leq 2 \Pr_{I, J} [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1]. \quad (9)$$

*Proof.* By (8),

$$\begin{aligned}
\Pr_{I,J} [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1] &= \mathbb{E}_I \left[ \Pr_J [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1 \mid I] \right] \\
&\geq \mathbb{E}_I \left[ f(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) \Pr_J [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1 \mid I] \right] \\
&\geq \mathbb{E}_I [f(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) / 2] = \Pr_I [f(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1] / 2.
\end{aligned}$$

□

Notice that the right hand side of (9) corresponds to drawing two subsets of size  $k$  from a stream of size  $n$ . It is desirable to bound this with a quantity that corresponds to partitioning a sample of size  $2k$  to two subsets of size  $k$ . Essentially, this amounts to ignoring the elements out of  $I \cup J$ . Formally:

**Lemma 8.3.** *Let  $f': X^k \times X^k \rightarrow \{0, 1\}$ , let  $I$  and  $J$  be random subsets of  $[n]$  as defined above and let  $I' \sim \text{Uni}(2k, k)$ . Then,*

$$\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{I,J} [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1] \leq \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I'} [f'(\bar{x}(\mathcal{A}', I')_{I'}, \bar{x}(\mathcal{A}', I')_{[2k] \setminus I'}) = 1]. \quad (10)$$

*Proof.* Let  $\mathcal{A}$  denote the maximizer on the left hand side of (10). Let  $U \subseteq [2n]$  be a set of size  $2k$ , and we will prove that

$$\Pr_{I,J} [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) \mid I \cup J = U] \leq \max_{\mathcal{A}'} \Pr_{I'} [f'(\bar{x}(\mathcal{A}', I')_{I'}, \bar{x}(\mathcal{A}', I')_{[2k] \setminus I'})]. \quad (11)$$

The proof of Lemma 8.3 will then follow by taking an expectation over  $U$ .

The main idea to proving (11) is to match the subsets  $I \subseteq U$  with the subsets  $I' \subseteq [2k]$ . One can define an adversary  $\mathcal{A}_U \in \text{Adv}_{2k}$  that simulates the behavior of  $\mathcal{A}$  on  $U$ , hence matching the probability that  $f' = 1$ . In particular,  $\mathcal{A}_U$  simulates the selections of  $\mathcal{A}$  on the set  $U$ , while skipping all the elements not in  $U$ . Formally, denote  $U = (i_1, \dots, i_{2k})$  where  $i_1 < i_2 < \dots < i_{2k}$  and  $U_{I'} = \{i_j : j \in I'\}$ . Then  $\mathcal{A}_U$  is defined to satisfy  $\bar{x}(\mathcal{A}_U, I') := \bar{x}(\mathcal{A}, U_{I'})_U$ . This implies that

$$\bar{x}(\mathcal{A}_U, I')_{I'} = \bar{x}(\mathcal{A}, U_{I'})_{U_{I'}}; \quad \bar{x}(\mathcal{A}_U, I')_{[2k] \setminus I'} = \bar{x}(\mathcal{A}, U_{I'})_{U_{[2k] \setminus I'}}.$$

Hence,

$$f'(\bar{x}(\mathcal{A}, U_{I'})_{U_{I'}}, \bar{x}(\mathcal{A}, U_{I'})_{U_{[2k] \setminus I'}}) = f'(\bar{x}(\mathcal{A}_U, I')_{I'}, \bar{x}(\mathcal{A}_U, I')_{[2k] \setminus I'}). \quad (12)$$

Notice that the joint distribution of  $(U_{I'}, U_{[n] \setminus I'})$  taken over  $I' \sim \text{Uni}(2k, k)$ , is the same as the joint distribution of  $(I, J)$ , conditioned on  $I \cup J = U$ . In combination with (12), this implies that

$$\begin{aligned}
\Pr [f'(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)_J) = 1 \mid I \cup J = U] &= \Pr [f'(\bar{x}(\mathcal{A}, U_{I'})_{U_{I'}}, \bar{x}(\mathcal{A}, U_{I'})_{U_{[2k] \setminus I'}}) = 1] \\
&= \Pr [f'(\bar{x}(\mathcal{A}_U, I')_{I'}, \bar{x}(\mathcal{A}_U, I')_{[2k] \setminus I'}) = 1] \leq \max_{\mathcal{A}'} \Pr [f'(\bar{x}(\mathcal{A}', I')_{I'}, \bar{x}(\mathcal{A}', I')_{[2k] \setminus I'}) = 1].
\end{aligned}$$

This proves (11), and concludes the proof. □

The proof of Lemma 8.1 follow directly from Lemma 8.2 and Lemma 8.3.

## 8.1 Proof of Lemma 6.2

We start with the following auxiliary probabilistic lemma:

**Lemma 8.4.** *Let  $I \subseteq [n]$  of size  $|I| = k$ , let  $\bar{x}$  and let  $t \geq 0$  be such that*

$$\max_{E \in \mathcal{E}} \left| \frac{|E \cap \bar{x}_I|}{k} - \frac{|E \cap \bar{x}|}{n} \right| \geq t. \quad (13)$$

*Let  $J$  be distributed uniformly over all subsets of  $[n] \setminus I$  of size  $k$ . Then, with probability at least  $1/2$ ,*

$$\max_{E \in \mathcal{E}} \left| \frac{|E \cap \bar{x}_I|}{k} - \frac{|E \cap \bar{x}_J|}{k} \right| \geq t - C/\sqrt{k}.$$

*Proof.* Let  $E_0$  be a maximizer of (13). We will assume that

$$\frac{|E_0 \cap \bar{x}|}{n} \geq \frac{|E_0 \cap \bar{x}_I|}{k}. \quad (14)$$

and the proof follows similarly in the other case. Applying Lemma 5.6 (item 1) with  $U = \{i \in [n] \setminus I : x_i \in E_0\}$ ,  $n = n - k$  and  $I = J$ , we have that with probability at least  $1/2$  over  $J$ ,

$$\frac{|E_0 \cap \bar{x}_J|}{k} = \frac{|U \cap J|}{k} \geq \frac{|U|}{n - k} - C\sqrt{1/k} = \frac{|E_0 \cap \bar{x}_{[n] \setminus I}|}{n - k} - C\sqrt{1/k} \geq \frac{|E_0 \cap \bar{x}_{[n]}|}{n} - C\sqrt{1/k},$$

where the last inequality follows from (14). In particular,

$$\frac{|E_0 \cap \bar{x}_J|}{k} - \frac{|E_0 \cap \bar{x}_I|}{k} \geq t - C/\sqrt{k},$$

which concludes the proof.  $\square$

*Proof of Lemma 6.2.* First, apply Lemma 8.1 with the function  $f(\bar{x}_I, \bar{x})$  that is the indicator of the event

$$\max_{E \in \mathcal{E}} \left| \frac{|E \cap \bar{x}_I|}{k} - \frac{|E \cap \bar{x}|}{n} \right| \geq t$$

and  $f'(\bar{x}_I, \bar{x}_J)$  that is the indicator of

$$\max_{E \in \mathcal{E}} \left| \frac{|E \cap \bar{x}_I|}{k} - \frac{|E \cap \bar{x}_J|}{k} \right| \geq t - C_0/\sqrt{k},$$

where  $C_0 > 0$  corresponds to the constant  $C$  in Lemma 8.4. The condition in Lemma 8.1 is satisfied from Lemma 8.4 and one derives that

$$\begin{aligned} \max_{\mathcal{A} \in \text{Adv}_n} \Pr[\text{App}_{\mathcal{A}, I} \geq t] &= \max_{\mathcal{A} \in \text{Adv}_n} \Pr[f(\bar{x}(\mathcal{A}, I)_I, \bar{x}(\mathcal{A}, I)) = 1] \\ &\leq 2 \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr[f'(\bar{x}'(\mathcal{A}', I')_I, \bar{x}'(\mathcal{A}', I')_{[2k] \setminus I'})] \\ &= 2 \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr[\text{Disc}_{\mathcal{A}', I'} / k \geq t - C_0/\sqrt{k}]. \end{aligned}$$

$\square$

## 8.2 Proof of Lemma 7.2

We start with the following probabilistic lemma:

**Lemma 8.5.** *Let  $\epsilon \in (0, 1)$ ,  $k \geq C/\epsilon$ ,  $\bar{x} \in X^n$  and  $I \subseteq [n]$ ,  $|I| = k$ , be such that*

$$\exists E \in \mathcal{E}, |\bar{x} \cap E| \geq \epsilon n \text{ and } \bar{x}_I \cap E = \emptyset.$$

*Let  $J$  be drawn uniformly from the subsets of  $[n] \setminus I$  of size  $k$ . Then, with probability at least  $1/2$ ,*

$$\exists E \in \mathcal{E}, |\bar{x}_{I \cup J} \cap E| \geq \epsilon k/2 \text{ and } \bar{x}_I \cap E = \emptyset.$$

*Proof.* Let  $E_0 \in \mathcal{E}$  be any set such that  $|\bar{x} \cap E_0| \geq \epsilon n$  and  $\bar{x}_I \cap E_0 = \emptyset$ . It suffices to show that with probability at least  $1/2$ ,  $|\bar{x}_J \cap E_0| \geq \epsilon k/2$ . To prove this, apply Lemma 5.6 (item 2) with  $I = J$ ,  $U = \{i \in [n] \setminus I : x_i \in E_0\}$ ,  $n = n - k$  and  $\alpha = 1/2$ . Notice that  $|U| \geq \epsilon n \geq \epsilon(n - k)$  and we derive that

$$\begin{aligned} \Pr \left[ |\bar{x}_J \cap E_0| < \frac{\epsilon k}{2} \right] &= \Pr \left[ \frac{|J \cap U|}{k} < \frac{\epsilon}{2} \right] \leq \Pr \left[ \frac{|J \cap U|}{k} \leq \frac{|U|}{2(n-k)} \right] \\ &\leq \Pr \left[ \left| \frac{|J \cap U|}{k} - \frac{|U|}{n-k} \right| \geq \frac{|U|}{2(n-k)} \right] \leq 2 \exp \left( -\frac{k|U|/4}{6(n-k)} \right) \leq 2 \exp \left( -\frac{k\epsilon/4}{6} \right) \leq 1/2, \end{aligned}$$

where the last inequality follows from the assumption that  $k \geq C\epsilon$  for a sufficiently large  $C > 0$ .  $\square$

*Proof of Lemma 7.2.* We will apply Lemma 8.1 with  $f(\bar{x}_I, \bar{x})$  being the indicator of

$$\exists E \in \mathcal{E}, |\bar{x} \cap E| \geq \epsilon n \text{ and } \bar{x}_I \cap E = \emptyset$$

and  $f'(\bar{x}_I, \bar{x}_J)$  the indicator of

$$\exists E \in \mathcal{E}, |\bar{x}_{I \cup J} \cap E| \geq \epsilon k/2 \text{ and } \bar{x}_I \cap E = \emptyset.$$

The condition of Lemma 8.1 holds from Lemma 8.5 and the proof follows.  $\square$

## 9 Covering Numbers

### 9.1 Overview

While studying online algorithms, it is natural to consider the following objects:

**Definition 9.1** (Dynamic-Set.). *A dynamic set is an online algorithm  $\mathbb{B}$  which is defined on input sequences  $\bar{x} \in X^n$ . At each time-step  $t \leq n$ , the algorithm decides whether to retain  $x_t$  or to discard it. The decision whether to retain/discard  $x_t$  may depend on the elements  $x_1, \dots, x_t$  which were observed up to time  $t$ . The trace of  $\mathbb{B}$  with respect to an input sequence  $\bar{x}$  is defined by*

$$\mathbb{B}(\bar{x}) = \{x_t : x_t \text{ was retained by } \mathbb{B}, t \leq n\}.$$

*and is viewed either as a set or as an ordered set. Further, we say that  $|\mathbb{B}| \leq m$  if  $\mathbb{B}$  cannot retain more than  $m$  elements. We stress that the decision whether to retain/discard an item is not reversible: retained (discarded) items can not be discarded (retained) in the future.*



Given a family  $\mathcal{E}$ , we would like to cover it using dynamic sets, as defined below:

**Definition 9.2.** Let  $\mathcal{E}$  be some family and let  $\mathcal{N}$  denote some finite collection of dynamic sets. We say that  $\mathcal{N}$  is an  $\epsilon$ -cover for  $\mathcal{E}$  if for every input sequence  $\bar{x}$  and every  $E \in \mathcal{E}$  there exists  $\mathbb{B} \in \mathcal{N}$  such that

$$|(E \cap \{\bar{x}\}) \Delta \mathbb{B}(\bar{x})| \leq \epsilon^2 n$$

where  $\Delta$  is the symmetric difference of sets. Further, define the covering number at scale  $\epsilon$ ,  $N(\mathcal{E}, \epsilon)$ , as the smallest cardinality of an  $\epsilon$ -cover for  $\mathcal{E}$ .

We can obtain bounds on the covering numbers at scale 0 for Littlestone families, via a known argument:

**Lemma 9.3** (Covering Littlestone Families with Few Dynamic Sets). Let  $\mathcal{E}$  be a family of subsets of  $X$  with  $\text{Ldim}(\mathcal{E}) = d$  and let  $n \in \mathbb{N}$ . Then,  $N(\mathcal{E}, 0) \leq \binom{n}{\leq d}$ . Moreover, this is tight for  $n = d$ , where  $N(\mathcal{E}, 0) = 2^d = \binom{n}{\leq d}$ .

While covering numbers at scale 0 can be used to derive bounds for  $\epsilon$ -approximation and  $\epsilon$ -nets for the Bernoulli sampler  $\text{Ber}(n, p)$  with  $p$  constant, these bounds are sub-optimal for epsilon-approximations. Improved bounds can be obtained by computing covering numbers at scale  $\epsilon > 0$ . While we do not know how derive better than  $\binom{n}{\leq d}$  even for scales  $\epsilon > 0$ , it is possible to obtain improved bounds on *fractional covering numbers*, which is a notion that we define below and can replace the covering numbers:

**Definition 9.4.** Let  $\mu$  denote a probability measure over dynamic sets. We say that  $\mu$  is an  $(\epsilon, \gamma)$ -fractional cover for  $\mathcal{E}$  if for every sequence  $\bar{x}$  and every  $E \in \mathcal{E}$ ,

$$\mu(\{\mathbb{B}: |(E \cap \{\bar{x}\}) \Delta \mathbb{B}(\bar{x})| \leq \epsilon^2 n\}) \geq 1/\gamma.$$

Define the fractional covering number at scale  $\epsilon$ ,  $N'(\mathcal{E}, \epsilon)$ , as the minimal value of  $\gamma$  such that there exists an  $(\epsilon, \gamma)$  fractional cover for  $\mathcal{E}$ .

Notice that  $N'(\mathcal{E}, \epsilon) \leq N(\mathcal{E}, \epsilon)$ : if  $\mathcal{C}$  is an  $\epsilon$ -cover, then  $N'(\mathcal{E}, \epsilon) \leq |\mathcal{C}|$ , by taking  $\mu$  to be the uniform distribution over  $\mathcal{C}$ .

We can obtain the following bound on the fractional covering numbers for Littlestone classes:

**Lemma 9.5.** It holds that  $N'(\mathcal{E}, \epsilon) \leq (C/\epsilon)^{2d}$ , for some universal  $C > 0$ .

Next, we apply bounds on covering numbers for epsilon approximation and epsilon nets:

### 9.1.1 Epsilon Approximation and Sequential Rademacher

The following bound can be derived based on 0-nets:

**Lemma 9.6.** Let  $\mathcal{A} \in \text{Adv}_{2k}$ ,  $\mathbf{I} \sim \text{Ber}(2k, 1/2)$ ,  $\delta \in (0, 1/2)$  and let  $\mathcal{E}$  be any family over some universe. Then, with probability  $1 - \delta$ ,

$$\text{Disc}_{\mathcal{A}, \mathbf{I}}(\mathcal{E}) \leq C \sqrt{k(\log N(\mathcal{E}, 0) + \log 1/\delta)}.$$

The proof is via a simple union bound. In combination with the bound on the 0-cover of Littlestone classes (Lemma 9.3), this derives that with probability  $1 - \delta$ ,

$$\text{Disc}_{\mathcal{A}, I} \leq \sqrt{k(d \log k + \log(1/\delta))},$$

which implies a sample complexity of

$$O\left(\frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon^2}\right).$$

To derive sharper bounds, one can use covering numbers at scales  $\epsilon > 0$ . Since we only have fractional covering numbers for Littlestone classes, we present the following lemma that derives a bound based on them:

**Lemma 9.7.** *Let  $\mathcal{A} \in \text{Adv}_{2k}$ ,  $I \sim \text{Ber}(2k, 1/2)$ ,  $\delta \in (0, 1/2)$  and let  $\mathcal{E}$  be some family. Then, with probability  $1 - \delta$ ,*

$$\text{Disc}_{\mathcal{A}, I}(\mathcal{E}) \leq C\sqrt{k} \left( \int_0^1 \sqrt{\log N'(\mathcal{E}, \epsilon)} d\epsilon + \sqrt{\log 1/\delta} \right).$$

This has the same form as the celebrated Dudley's integral but here we extend it to fractional covering numbers.

Using Lemma 9.5 and Lemma 9.7, Lemma 6.4 immediately follows. Indeed,

$$\begin{aligned} \frac{\text{Disc}_{\mathcal{A}, I}}{\sqrt{k}} &\leq C \left( \int_0^1 \sqrt{\log N'(\mathcal{E}, \epsilon)} d\epsilon + \sqrt{\log 1/\delta} \right) \\ &\leq C \left( \int_0^1 \sqrt{d \log 1/\epsilon} d\epsilon + \sqrt{\log 1/\delta} \right) \leq C \left( \sqrt{d} + \sqrt{\log 1/\delta} \right). \end{aligned}$$

### 9.1.2 Epsilon Nets

We have the following statement:

**Lemma 9.8.** *Let  $\mathcal{E}$  be any family,  $\mathcal{A} \in \text{Adv}_{2k}$ ,  $I \sim \text{Ber}(2k, p)$  and let  $m \in [2k]$ . Then,*

$$\Pr_I \left[ \text{Net}_{\mathcal{A}, I, m, mp/2}^{2k} = 1 \right] \leq N(\mathcal{E}, 0) \exp(-cmp).$$

Lemma 7.4 follows directly by applying the bound on the covering numbers at scale 0 for Littlestone classes, presented in Lemma 9.3.

### 9.1.3 Organization

Section 9.2 contains the proofs of Lemma 9.3 and Lemma 9.5 on the covering numbers for Littlestone classes; Section 9.3 contains the proofs of Lemma 9.6 and Lemma 9.7 on proving  $\epsilon$ -approximations via covering numbers; and Section 9.4 contains the proof of Lemma 9.8 on proving  $\epsilon$ -nets via covering numbers.

## 9.2 Covering for Littlestone Classes

Section 9.2.1 contains the proof of Lemma 9.3 on the covering numbers at scale 0, that is based on a known arguments; and Section 9.2.2 contains the proof of Lemma 9.5, that builds on the machinery presented in Section 9.2.1.

### 9.2.1 Proof of Lemma 9.3

Before we prove Lemma 9.3, let us make a couple of comparisons to related literature.

**Remark 9.9.** *It is fair to note that the proof of this proposition exploits standard and basic ideas from the online learning literature. In particular, the constructed dynamic-sets hinge on variants of the Standard Optimal Algorithm by [Lit88], and utilize its property of being a mistake-driven algorithm<sup>6</sup>. However, for the benefit of readers who are less familiar with this literature, we provide here a self-contained proof and modify some of the terminology/notation from the language of online learning to the language of  $\epsilon$ -nets/approximations.*

**Remark 9.10.** *This comment concerns a connection with the celebrated Sauer-Shelah-Perles (SSP) Lemma [Sau72]. Note that the SSP Lemma is equivalent to a variant of Lemma 9.3 in which two quantifiers are flipped. Indeed, the SSP Lemma asserts that for every  $x_1, \dots, x_n$  there are at most  $\binom{n}{\leq d}$  sets that realize all possible intersection patterns of the sets in  $\mathcal{E}$  with  $\{x_1, \dots, x_n\}$ . That is, if one allows the sets  $\mathbb{B}_i$  to be chosen after seeing the entire input-sequence  $x_1, \dots, x_n$  then the conclusion in Lemma 9.3 extends to VC classes (which can have an unbounded Littlestone dimension, as witnessed by the class of thresholds).*

*Proof of Lemma 9.3.* We begin with the upper bound. The definition of the dynamic sets  $\mathbb{B}_i$  exploits the following property of Littlestone families. Let  $\mathcal{E}$  be a family with  $\text{Ldim}(\mathcal{E}) < \infty$ , let  $x \in X$ , and consider the two “half-families”

$$\mathcal{E}_{\not\exists x} = \{E \in \mathcal{E} : x \notin E\}, \quad \mathcal{E}_{\exists x} = \{E \in \mathcal{E} : x \in E\}.$$

The crucial observation is that if  $\mathcal{E} \neq \emptyset$  then for every  $x \in X$ :

$$\text{Ldim}(\mathcal{E}_{\not\exists x}) < \text{Ldim}(\mathcal{E}) \text{ or } \text{Ldim}(\mathcal{E}_{\exists x}) < \text{Ldim}(\mathcal{E}). \quad (15)$$

Indeed, otherwise we have  $\text{Ldim}(\mathcal{E}_{\not\exists x}) = \text{Ldim}(\mathcal{E}_{\exists x}) = \text{Ldim}(\mathcal{E}) =: d$  which implies that  $\mathcal{E}$  shatters the following tree of depth  $d + 1$ : the root is labelled with  $x$ , and the left and right subtrees of the root are trees which witness that the dimensions of  $\mathcal{E}_{\not\exists x}$  and  $\mathcal{E}_{\exists x}$  equal  $d$ . However, since  $\text{Ldim}(\mathcal{E}) = d$ , this is not possible.

**Littlestone Majority Vote.** Equation (15) allows to define a notion of majority-vote of a (possibly infinite) family  $\mathcal{E}$  with a finite Littlestone dimension. The intuition is that if  $x$  is such that  $\text{Ldim}(\mathcal{E}_{\exists x}) = d$  then by Equation (15) it must be that  $\text{Ldim}(\mathcal{E}_{\not\exists x}) < d$  and therefore  $\text{Ldim}(\mathcal{E}_{\exists x}) > \text{Ldim}(\mathcal{E}_{\not\exists x})$  which we interpret as if  $x$  is contained in a “majority” of the sets in  $\mathcal{E}$ . Similarly,  $\text{Ldim}(\mathcal{E}_{\not\exists x}) = d$  is interpreted as if most sets in  $\mathcal{E}$  do not contain  $x$ . This motivates the following definition

$$\text{Lmaj}(\mathcal{E}) = \{x : \text{Ldim}(\mathcal{E}_{\exists x}) = d\}, \quad (16)$$

---

<sup>6</sup>A mistake-driven algorithm updates its internal state only when it makes mistakes.

with the convention that  $\text{Lmaj}(\emptyset) = \emptyset$ . Observe that  $\text{Lmaj}(\mathcal{E})$  shares the following property with the standard majority-vote over finite families: let  $x \in X$  and assume  $\mathcal{E} \neq \emptyset$ . Then,

$$\left( (\forall E \in \mathcal{E}) : x \in E \right) \implies x \in \text{Lmaj}(\mathcal{E}) \quad \text{and} \quad \left( (\forall E \in \mathcal{E}) : x \notin E \right) \implies x \notin \text{Lmaj}(\mathcal{E}). \quad (17)$$

That is, if the sets in  $\mathcal{E}$  agree on  $x$  unanimously, then  $\text{Lmaj}(\mathcal{E})$  agrees with them on  $x$ . We comment that this definition is the basis of the *Standard Online Algorithm* which witnesses the online-learnability of Littlestone classes in the mistake-bound model [Lit88].

We are now ready to define the required family of  $\binom{n}{\leq d}$  dynamic sets. Each dynamic set  $\mathbb{B}_I$  is indexed by a subset  $I \subseteq [n]$  of size  $|I| \leq d$ . (Hence there are  $\binom{n}{\leq d}$  dynamic sets.) Below is the pseudo-code of  $\mathbb{B}_I$  for  $I \subseteq [n]$ .

### The Dynamic Set $\mathbb{B}_I$

Let  $\mathcal{E}$  be a family with  $\text{Ldim}(\mathcal{E}) = d$ , and let  $I \subseteq [n]$ .

Let  $x_1, \dots, x_n$  denote the (adversarially-produced) input sequence.

1. Initialize  $\mathcal{E}_0^I = \mathcal{E}$ .
2. For  $t = 1, \dots, n$ :
  - (a) If  $t \notin I$  then set  $\mathcal{E}_t^I = \mathcal{E}_{t-1}^I$ .
  - (b) Else, set

$$\mathcal{E}_t^I = \begin{cases} (\mathcal{E}_{t-1}^I)_{\not\ni x_t} & x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^I), \\ (\mathcal{E}_{t-1}^I)_{\ni x_t} & x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^I). \end{cases}$$

- (c) Retain  $x_t$  if and only if  $x_t \in \text{Lmaj}(\mathcal{E}_t^I)$ .

Observe the following useful facts regarding  $\mathbb{B}_I$ :

1. The sequence of families  $\{\mathcal{E}_t^I\}_{t=0}^n$  is a chain:  $\mathcal{E}_0^I \supseteq \mathcal{E}_1^I \supseteq \dots \supseteq \mathcal{E}_n^I$ .
2. A strict containment  $\mathcal{E}_{t-1}^I \supsetneq \mathcal{E}_t^I$  occurs only if  $t \in I$ .
3. Whenever a strict containment  $\mathcal{E}_{t-1}^I \supsetneq \mathcal{E}_t^I$  occurs then also  $\text{Ldim}(\mathcal{E}_{t-1}^I) > \text{Ldim}(\mathcal{E}_t^I)$ . (By Equations (15) and (16).)

To complete the proof it remains to show that for every  $\bar{x}$  and every  $E \in \mathcal{E}$  there exists  $I \subseteq [n]$ , with  $|I| \leq d$  such that

$$(\forall t \leq n) : x_t \in E \iff x_t \in \mathbb{B}_I(\bar{x}). \quad (18)$$

We construct the set  $I = I(E)$  in a parallel fashion to the above process:

### The Index Set $I = I(E)$

Let  $E \in \mathcal{E}$  and let  $x_1, \dots, x_n$  denote the input sequence.

1. Initialize  $\mathcal{E}_0^E = \mathcal{E}$  and  $I = \emptyset$ .
2. For  $t = 1, \dots, n$ :
  - (a) If  $E$  and  $\text{Lmaj}(\mathcal{E}_{t-1}^E)$  agree on  $x_t$  (i.e.  $x_t \in E \iff x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^E)$ ) then set  $\mathcal{E}_t^E = \mathcal{E}_{t-1}^E$ .
  - (b) Else, add  $t$  to  $I$  and set

$$\mathcal{E}_t^E = \begin{cases} (\mathcal{E}_{t-1}^E)_{\not\ni x_t} & x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^E) \wedge x_t \notin E, \\ (\mathcal{E}_{t-1}^E)_{\ni x_t} & x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^E) \wedge x_t \in E. \end{cases}$$

3. Output  $I = I(E)$ .

Note that by construction,  $\mathcal{E}_t^E = \mathcal{E}_t^I$  for every  $t \leq n$ , and  $E \in \mathcal{E}_t^E$  for all  $t$ . We need to show that the constructed set  $I$  satisfies Equation (18) and that  $|I| \leq d$ . For the first part, note that for every  $t \leq n$ :

$$\begin{aligned} x_t \in \mathbb{B}_I(\bar{x}) &\iff x_t \in \text{Lmaj}(\mathcal{E}_t^I) && \text{(by definition of } \mathbb{B}_I) \\ &\iff x_t \in \text{Lmaj}(\mathcal{E}_t^E) && \text{(since } \mathcal{E}_t^E = \mathcal{E}_t^I) \\ &\iff x_t \in E, && \text{(see below)} \end{aligned}$$

where the last step follows because all the sets  $E' \in \mathcal{E}_t^E$  agree with  $E$  on  $x_t$ . Thus, by<sup>7</sup> Equation (17) also  $\text{Lmaj}(\mathcal{E}_t^E)$  agrees with  $E$  on  $x_t$ , which amounts to the last step.

To see that  $|I| \leq d$ , consider the chain

$$\mathcal{E}_0^E \supseteq \mathcal{E}_1^E \supseteq \dots \supseteq \mathcal{E}_n^E.$$

Note that strict containments  $\mathcal{E}_{t-1}^E \supsetneq \mathcal{E}_t^E$  occurs only if  $t \in I$ , and that whenever such a strict containment occurs, we have  $\text{Ldim}(\mathcal{E}_{t-1}^E) > \text{Ldim}(\mathcal{E}_t^E)$ . Therefore, since  $\text{Ldim}(\mathcal{E}_0^E) = d$  and  $\text{Ldim}(\mathcal{E}_n^E) \geq 0$ , it follows that  $|I| \leq d$  as required.

It remains to prove the lower bound. Let  $D = \{\mathbb{B}_i : 1 \leq i < 2^d\}$  be a family of less than  $2^d$  dynamic sets. Pick a tree  $\mathcal{T}$  of depth  $d$  which is shattered by  $\mathcal{E}$  and define an adversarial sequence  $x_1, \dots, x_d$  as follows:

---

<sup>7</sup>Note that  $\mathcal{E}_t^E \neq \emptyset$  because  $E \in \mathcal{E}_t^E$ .

### The Adversarial Sequence $x_1, \dots, x_d$

1. Set  $\mathcal{T}_1 = \mathcal{T}, D_1 = D, \mathcal{E}_1 = \mathcal{E}$ , and  $i = 1$ .
2. For  $i = 1, \dots, d$ 
  - (i) Set  $x_i$  to be the item labelling the root of  $\mathcal{T}_i$ .
  - (ii) If less than half of the dynamic sets  $\mathbb{B}_j \in D_i$  retain  $x_i$  then continue the next iteration with  $\mathcal{T}_{i+1}$  being the right subtree of  $\mathcal{T}_i$  (which corresponds to the sets containing  $x_i$ ), and with  $\mathcal{E}_{i+1} = \{E \in \mathcal{E}_i : x_i \in E\}$  and  $D_{i+1} = \{\mathbb{B}_j \in D_i : x_i \in \mathbb{B}_j\}$ .
  - (iii) Else, continue to the next iteration with  $\mathcal{T}_{i+1}$  being the left subtree of  $\mathcal{T}_i$ , and with  $\mathcal{E}_{i+1} = \{E \in \mathcal{E}_i : x_i \notin E\}$  and  $D_{i+1} = \{\mathbb{B}_j \in D_i : x_i \notin \mathbb{B}_j\}$ .

Note that  $\mathcal{E}_i$  contains all the sets in  $\mathcal{E}$  that are consistent<sup>8</sup> with the path corresponding to  $x_1, \dots, x_{i-1}$ , and similarly  $D_i$  contains all dynamic sets in  $D$  which are consistent with that path. Thus, since  $|D_1| = |D| < 2^d$ , it follows by construction that  $|D_i| < 2^i$  for every  $i < d$ , and in particular that  $D_d = \emptyset$  at the end of the process. Thus, the set  $E \in \mathcal{E}$  which is consistent with the path corresponding to  $x_1, \dots, x_d$  satisfies  $E \cap \{x_1, \dots, x_d\} \neq \mathbb{B}_i(x_1, \dots, x_d)$  for every  $i < 2^d$ , as required.  $\square$

#### 9.2.2 Proof of Lemma 9.5

For convenience, let us bound  $N'(\mathcal{E}, \sqrt{\epsilon}) \leq (C/\epsilon)^d$ . We start by defining the fractional cover  $\mathcal{B}$  and then prove its validity. Let  $p = 3d/(\epsilon n)$  for a sufficiently large universal constant  $C_0$ , and  $\mathbb{B} \sim \mathcal{B}$  is sampled as follows:

1. Select a random subset  $I' \subseteq [n]$ , where each  $i \in [n]$  is selected independently with probability  $p$ .
2. Select a subset  $I \subseteq I'$  of size  $|I| \leq d$ , uniformly at random from the set of all  $\binom{I'}{\leq d}$  subsets.
3. Output  $\mathbb{B} = \mathbb{B}_I$ .

To prove that  $\mathcal{B}$  is an  $(\epsilon, (C/\epsilon)^d)$ -fractional cover, fix some  $E \in \mathcal{E}$  and sequence  $\bar{x}$ . From the proof of Lemma 9.3, for any  $I'$ , there exists  $I^* = I^*(I') \subseteq I'$  of size  $|I^*| \leq d$  such that  $(\mathbb{B}_{I^*}(\bar{x}))_{I'} = E \cap \bar{x}_{I'}$ . Denote  $I^* = I^*(I')$  where  $I'$  is distributed as above.

We will bound from below the probability that  $I^*$  satisfies

$$|\mathbb{B}_{I^*}(\bar{x}) \Delta (E \cap \bar{x})| \leq \epsilon n.$$

Further, we will bound from below the probability that  $I = I^*$ . Combining these two bounds, this will give a lower bound on the probability that

$$|\mathbb{B}_I(\bar{x}) \Delta (E \cap \bar{x})| \leq \epsilon n,$$

---

<sup>8</sup> $\mathbb{B}_j$  is consistent with the path corresponding to  $x_1, \dots, x_d$  means that  $\mathbb{B}_j(x_1, \dots, x_n)$  contains  $x_i$  if and only if  $x_{i+1}$  labels the right child of the node labelled  $x_i$ .

that suffices to complete the proof.

To begin with the first step, notice that

$$|\mathbb{B}_{I^*}(\bar{x})\Delta(E\cap\bar{x})| = |\{t \in [n] : x_t \in \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E\}|. \quad (19)$$

To analyze the right hand side of the above quantity, place each time  $t \in [n]$  in one of four categories:

1.  $x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$  and  $t \in I'$ .
2.  $x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$  and  $t \notin I'$ .
3.  $x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$  and  $t \notin I'$ .
4.  $x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$  and  $t \in I'$ .

Notice that the above properties apply:

- It holds that  $\text{Lmaj}(\mathcal{E}_t^{I^*}) \neq \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})$  if and only if  $t$  is in category (1).
  - For  $t$  in category (1) it holds that  $t \in I'$  which implies that  $x_t \notin \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$ , since  $\mathbb{B}_{I^*}$  is defined to agree with  $E$  on all  $t \in I'$ . While any  $t$  in category (1) satisfies  $x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$ , this implies that  $\text{Lmaj}(\mathcal{E}_t^{I^*}) \neq \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})$ .
  - For  $t$  in categories (2) and (3) it holds that  $t \notin I'$ . Since  $I^* \subseteq I'$ , then  $t \notin I^*$ . By definition of  $\mathbb{B}_{I^*}$  it holds that  $\text{Lmaj}(\mathcal{E}_t^{I^*}) = \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})$  whenever  $t \notin I^*$ .
  - For  $t$  in category (4), it holds that  $x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$ . Since  $\mathbb{B}_{I^*}$  agrees with  $E$  on  $I'$  and since  $t \in I'$ , it holds that  $\mathbb{B}_{I^*}$  agrees with  $E$  on  $x_t$ , namely,  $x_t \notin \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$ . This implies that  $x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta \text{Lmaj}(\mathcal{E}_t^{I^*})$ . By definition of the dynamic set  $\mathbb{B}_{I^*}$  it holds that  $\text{Lmaj}(\mathcal{E}_{t-1}^{I^*}) = \text{Lmaj}(\mathcal{E}_t^{I^*})$  if and only if  $x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta \text{Lmaj}(\mathcal{E}_t^{I^*})$ . In particular,  $\text{Lmaj}(\mathcal{E}_{t-1}^{I^*}) = \text{Lmaj}(\mathcal{E}_t^{I^*})$  as required.
- It holds that  $x_t \in \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$  if and only if  $t$  is in category (2):
  - For categories (1) and (4) it holds that  $t \in I'$ . By definition of  $I^*$ ,  $\mathbb{B}_{I^*}$  and  $E$  agree for any  $\bar{x}_t$  for  $t \in I'$ . This implies that  $x_t \notin \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$ .
  - For categories (2) and (3), it holds that  $t \notin I'$  hence  $t \notin I^*$  which implies by definition of  $\mathbb{B}_{I^*}$  that  $\text{Lmaj}(\mathcal{E}_t^{I^*}) = \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})$ . Since for category (2) we have  $x_t \in \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$ , we further have  $x_t \in \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$ . Similarly, in category (3) we have  $x_t \notin \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})\Delta E$  hence  $x_t \notin \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$ .

Let  $\mathbf{H}$  (hit) denote the set of all indices  $t$  that correspond to case (1) and  $\mathbf{M}$  (miss) denote the set of indices in case (2). We view  $\mathbf{H}$  and  $\mathbf{M}$  as random variables that are functions of the random variable  $I'$  (where  $\bar{x}$  and  $E$  are fixed). Since only the elements in  $t \in \mathbf{M}$  satisfy  $x_t \in \text{Lmaj}(\mathcal{E}_t^{I^*})\Delta E$ , the goal is to upper bound  $|\mathbf{M}|$ . In fact, we will upper bound its expected value and then use Markov's inequality to derive tail bounds.

Before bounding  $\mathbb{E}[|\mathbf{M}|]$ , notice that  $|\mathbf{H}| \leq d$ . This holds due to the fact that, as described above, for each  $t \in \mathbf{H}$ ,  $\text{Lmaj}(\mathcal{E}_t^{I^*}) \neq \text{Lmaj}(\mathcal{E}_{t-1}^{I^*})$ . And this can happen at most  $d$  times, since  $\text{Ldim}(\text{Lmaj}(\mathcal{E}_t^{I^*})) < \text{Ldim}(\text{Lmaj}(\mathcal{E}_{t-1}^{I^*}))$  for any such  $t$ , as described in the proof of Lemma 9.3.

We proceed with proving that  $\mathbb{E}|\mathbf{M}| \leq d/p$ . Denote  $\mathbf{H} = \{t_1, \dots, t_{|\mathbf{H}|}\}$  let  $t_0 = 0$ , and let  $y_j$  denote the number of elements of  $\mathbf{M}$  between  $t_{j-1}$  and  $t_j$ :  $y_j = |\mathbf{M} \cap \{t_{j-1} + 1, t_{j-1} + 2, \dots, t_j - 1\}|$ . Note that  $|\mathbf{M}| = \sum_{j=1}^d y_j$ .

We claim that  $\mathbb{E}[y_j] \leq 1/p$ . Define

$$\mathbf{U}_j := \{t: t > t_{j-1}, x_t \in \text{Lmaj}(\mathcal{E}_{t_{j-1}}^{I^*}) \Delta E\}; \quad \mathbf{U}_j^- = \mathbf{U}_j \cap [t - 1].$$

Notice that for any  $t$  that satisfies  $t_{j-1} < t < t_j$  it holds that  $\text{Lmaj}(\mathcal{E}_{t_{j-1}}^{I^*}) = \text{Lmaj}(\mathcal{E}_t^{I^*})$ , since, as stated above,  $\text{Lmaj}(\mathcal{E}_t^{I^*})$  only changes at iterations  $t \in \mathbf{H}$ . This will imply the following:

$$\mathbf{M} \cap \{t_{j-1} + 1, \dots, t_j - 1\} = \mathbf{U}_j^-. \quad (20)$$

For the first direction of (20), any  $t \in \mathbf{U}_j^-$  satisfies  $t \notin \mathbf{H}$  by definition, and further it satisfies  $x_t \in \text{Lmaj}(\mathcal{E}_{t_{j-1}}^{I^*}) \Delta E = \text{Lmaj}(\mathcal{E}_t^{I^*}) \Delta E$  which implies that it is in  $\mathbf{M} \cup \mathbf{H}$ . However, it cannot be in  $\mathbf{H}$  since  $\mathbf{H} = \{t_1, \dots, t_{|\mathbf{H}|}\}$ . Further, it satisfies  $t_{j-1} < t < t_j$  by definition. For the second direction, any  $t$  in the left hand side satisfies  $x_t \in \text{Lmaj}(\mathcal{E}_t^{I^*}) \Delta E = \text{Lmaj}(\mathcal{E}_{t_{j-1}}^{I^*}) \Delta E$  which implies that it is in  $\mathbf{U}_j$ . We derive (20) which implies that  $y_j = |\mathbf{U}_j^-|$ .

To estimate  $|\mathbf{U}_j^-|$ , notice that the first element of  $\mathbf{U}_j$  that is also in  $I'$  is  $t_j$ . This follows from the fact that  $\text{Lmaj}(\mathcal{E}_t^{I^*})$  changes only once an element of  $\mathbf{H}$  is observed. Further, conditioned on  $t_{j-1}$  and  $I' \cap [t_{j-1}]$ , the set  $\mathbf{U}_j$  is fixed, and conditionally, since any element of  $\mathbf{U}_j$  is in  $I'$  with probability  $p$ , the expected number of elements in  $\mathbf{U}_j$  that are encountered before the first element of  $I'$  is bounded by  $1/p$ . This quantity is exactly  $\mathbb{E}[y_j] = |\mathbf{U}_j^-| \leq 1/p$ , and we derive that  $\mathbb{E}[|\mathbf{M}|] \leq d/p$ . From Markov's inequality,  $\Pr[|\mathbf{M}| \leq 3d/p] \geq 2/3$ .

We have proved that with probability  $2/3$ ,  $|\mathbf{M}| \leq 3d/p$ . This, from (19) and from the definitions of  $\mathbf{M}$  and  $p$ , implies that with probability  $2/3$ ,  $|\mathbb{B}_{I^*}(x) \Delta (E \cap \bar{x})| \leq 3d/p \leq \epsilon n$ . Further, we want to lower bound the probability that  $\mathbf{I} = \mathbf{I}^*$ . Notice that  $\Pr[\mathbf{I} = \mathbf{I}^* | I'] = 1/\binom{|I'|}{\leq d}$ , hence it is desirable to show that  $|I'|$  is small with high probability. Indeed, since  $\mathbb{E}|I'| = np$ , by Markov's inequality,  $\Pr[|I'| \leq 3np] \geq 2/3$ . By a union bound,

$$\Pr[|I'| \leq 3np, |\mathbb{B}_{I^*}(x) \Delta (E \cap \bar{x})| \leq \epsilon n] \geq 1/3.$$

We conclude that

$$\begin{aligned} \Pr[|\mathbb{B}_I(\bar{x}) \Delta (E \cap \bar{x})| \leq \epsilon n] &\geq \Pr[|I'| \leq 3np, |\mathbb{B}_{I^*} \Delta (E \cap \bar{x})| \leq \epsilon n, \mathbf{I} = \mathbf{I}^*] \\ &= \Pr[|I'| \leq 3np, |\mathbb{B}_{I^*} \Delta (E \cap \bar{x})| \leq \epsilon n] \Pr[\mathbf{I} = \mathbf{I}^* | |I'| \leq 3np, |\mathbb{B}_{I^*} \Delta (E \cap \bar{x})|] \geq \frac{1}{3} \cdot \binom{3np}{\leq d}^{-1} \\ &= \frac{1}{3} \cdot \binom{9d/\epsilon}{\leq d}^{-1} \geq \left(\frac{C}{\epsilon}\right)^{-d}, \end{aligned}$$

using the fact that  $\mathbf{I}^*$  is a function of  $I'$ , and conditioned on any value of  $I'$ , the probability that  $\mathbf{I} = \mathbf{I}^*$  is  $\binom{|I'|}{\leq d}^{-1}$ ; and further, that  $\binom{n}{k} \leq (Cn/k)^k$  for a universal  $C > 0$ .



### 9.3 Deriving Bounds on $\epsilon$ -Approximation via Fractional Covering Numbers

In this section we prove the concentration results based on covering numbers, starting with results based on deterministic 0-covers and moving to fractional  $\epsilon$ -covers. The following definition will be useful: for any  $E \in \mathcal{E}$  define  $Y_E = |E \cap \bar{x}_I| - |E \cap \bar{x}_{[2k] \setminus I}|$ . Similarly, for any  $\mathbb{B}$ , define  $Y_{\mathbb{B}} = |\mathbb{B} \cap \bar{x}_I| - |\mathbb{B} \cap \bar{x}_{[2k] \setminus I}|$ . Notice that

$$\text{Disc}_{\mathcal{A}, I} = \max_{E \in \mathcal{E}} |Y_E|.$$

#### 9.3.1 Basic Lemmas for Deterministic Covers and Proof of Lemma 9.6

We start with concentration of a single dynamic set:

**Lemma 9.11.** *Let  $\mathbb{B}$  be a dynamic set with  $|\mathbb{B}| \leq m$ . Let  $I \sim \text{Ber}(n, 1/2)$ . Then, for any  $t \geq 0$ ,*

$$\Pr[|Y_{\mathbb{B}}| \geq t] \leq 2 \exp(-t^2/(2m)).$$

Consequently, for any  $\delta \in (0, 1/2)$ , with probability  $1 - \delta$  it holds that

$$|Y_{\mathbb{B}}| \leq C \sqrt{m \log(1/\delta)}.$$

For the proof of Lemma 9.11, we need the following Martingale lemma (notice that an overview on Martingales is given in Section A.1)

**Lemma 9.12** ([dlPn99], Theorem 6.1). *Let  $\mathbf{y}_0, \dots, \mathbf{y}_n$  be a Martingale adapted to the filtration  $F_0, \dots, F_n$ , such that  $\sum_{i=1}^n |\mathbf{y}_i - \mathbf{y}_{i-1}|^2 \leq s$  holds almost surely for some  $s > 0$ . Assume that for all  $i \in [n]$ , conditioned on  $F_{i-1}$ ,  $\mathbf{y}_i - \mathbf{y}_{i-1}$  is a symmetric random variable (namely, it has the same conditional distribution as  $\mathbf{y}_{i-1} - \mathbf{y}_i$ ). Then, for any  $t > 0$ ,*

$$\Pr[|\mathbf{y}_n - \mathbf{y}_0| \geq t] \leq 2 \exp(-t^2/(2s)).$$

*Proof of Lemma 9.11.* This follows directly from Lemma 9.12. Indeed, we apply this lemma with  $\mathbf{y}_i = |\mathbb{B} \cap I \cap [i]| - |\mathbb{B} \cap ([n] \setminus I) \cap [i]|$  and  $s = m$ .  $\square$

We are ready to prove Lemma 9.6.

*Proof of Lemma 9.6.* Notice that it suffices to prove that for any  $t \geq 0$ ,

$$\Pr[\text{Disc}_{\mathcal{A}, I} > t] \leq 2N(\mathcal{E}, 0) \exp(-t^2/(4k)).$$

Let  $\mathcal{N}$  be a minimal 0-net for  $\mathcal{E}$ . Applying Lemma 9.11 with  $m = 2k$ , for any  $\mathbb{B} \in \mathcal{N}$ ,

$$\Pr[|Y_E| > t] \leq 2 \exp(-t^2/4k)$$

Applying a union bound over  $\mathbb{B} \in \mathcal{N}$ ,

$$\Pr[\text{Disc}_{\mathcal{A}, I} > t] = \Pr[\max_{E \in \mathcal{E}} |Y_E| > t] \leq \Pr[\max_{\mathbb{B} \in \mathcal{N}} |Y_{\mathbb{B}}| > t] \leq 2N(\mathcal{E}, 0) \exp(-t^2/(4k)).$$

$\square$

### 9.3.2 Basic Lemmas for Fractional Covers

To give intuition about fractional covers, we prove a variant of Lemma 9.6. First, we start with an auxiliary lemma that replaces the union bound:

**Lemma 9.13.** *Let  $\{\mathbf{y}_j\}_{j \in J}$  denote random variables over  $\{0, 1\}$  where  $J$  is some index set, and assume that for any  $j \in J$ ,  $\Pr[\mathbf{y}_j = 1] \leq p$ , for some  $p > 0$ . Let  $\mu$  denote some probability measure over  $J$  and let  $\alpha > 0$ . Then,*

$$\Pr_{\{\mathbf{y}_j\}_{j \in J}} \left[ \mu(\{j: \mathbf{y}_j = 1\}) \geq \alpha \right] \leq p/\alpha.$$

*Proof.* Notice that by linearity of expectation,

$$\mathbb{E}[\mu(\{j: \mathbf{y}_j = 1\})] = \mathbb{E} \int \mathbf{y}_j d\mu = \int \mathbb{E}[\mathbf{y}_j] d\mu \leq \int p d\mu = p.$$

The proof follows by Markov's inequality.  $\square$

The following lemma applies Lemma 9.13 specifically for distributions over dynamic sets.

**Lemma 9.14.** *Let  $\mu$  be a probability measure over dynamic sets with  $|\mathbb{B}| \leq m$  for all  $\mathbb{B}$  in the support of  $\mu$ . Then, for any  $\delta' > 0$ , with probability at least  $1 - \delta'$  over  $I \sim \text{Ber}(n, 1/2)$  it holds that*

$$\mu \left( \left\{ \mathbb{B}: |\mathbf{Y}_{\mathbb{B}}| \leq C \sqrt{m \log(1/(\delta' \alpha))} \right\} \right) \geq 1 - \alpha.$$

*Proof.* Apply Lemma 9.13 with  $\mathbf{y}_{\mathbb{B}}$  being the indicator that  $|\mathbf{Y}_{\mathbb{B}}| \leq C \sqrt{m \log(1/(\delta' \alpha))}$  and  $p = \alpha \delta'$ . If  $C$  is a sufficiently large constant, it follows from Lemma 9.11 that  $\Pr[\mathbf{y}_{\mathbb{B}} = 1] \leq \alpha \delta'$  and Lemma 9.13 can be applied to derive the desired result.  $\square$

Using Lemma 9.14, one can derive bounds for epsilon approximation based on fractional covering numbers at scale 0:

**Lemma 9.15.** *Let  $\mathcal{A} \in \text{Adv}_n$ ,  $I \sim \text{Bin}(n, 1/2)$  and  $\delta > 0$ . Then, with probability  $1 - \delta$ ,*

$$\text{Disc}_{\mathcal{A}, I} \leq C \sqrt{n (\log N'(\mathcal{E}, 0) + \log(1/\delta))}.$$

*Proof.* Let  $\mu$  be a  $(0, N'(\mathcal{E}, 0))$ -fractional cover for  $\mathcal{E}$ . We apply Lemma 9.14 with  $m = n$ ,  $\alpha = 1/(2N'(\mathcal{E}, 0))$ ,  $\delta' = \delta$  and  $\mu = \mu$  to get that with probability  $1 - \delta$ ,

$$\mu \left( \mathbb{B}: |\mathbf{Y}_{\mathbb{B}}| \leq C \sqrt{n (\log N'(\mathcal{E}, 0) + \log(1/\delta))} \right) \geq \frac{1}{2N'(\mathcal{E}, 0)}.$$

Whenever this holds, for every  $E \in \mathcal{E}$  there exists  $\mathbb{B} \in \mathcal{E}$  that  $E \cap \bar{x} = \mathbb{B}(\bar{x})$  which implies that  $\mathbf{Y}_{\mathbb{B}} = \mathbf{Y}_E$ . Hence,

$$\text{Disc}_{\mathcal{A}, I} = \max_{E \in \mathcal{E}} |\mathbf{Y}_E| \leq \max_{\mathbb{B} \in \mathcal{N}} |\mathbf{Y}_{\mathbb{B}}| \leq C \sqrt{n (\log N'(\mathcal{E}, 0) + \log(1/\delta))}.$$

$\square$

### 9.3.3 Chaining for Non-Fractional Covers

The proof of Lemma 9.7 follows the technique of chaining. Before presenting the proof for fractional covers, we start by presenting an outline of the proof for non-fractional covers, while obtaining a similar bound with  $N(\mathcal{E}, \epsilon)$  instead of  $N'(\mathcal{E}, \epsilon)$ . The proof follows from standard techniques. Some technicalities are ignored for the sake of presentation.

Let  $1 = \epsilon_0 > \epsilon_1 > \dots$  be a non-increasing sequence of values with  $\lim_{i \rightarrow \infty} \epsilon_i = 0$ . We take nets  $\mathcal{N}_0, \mathcal{N}_1, \dots$ , where each  $\mathcal{N}_i$  is an  $\epsilon_i$  net for  $\mathcal{E}$ . Each  $E \in \mathcal{E}$  we approximate using elements from the different nets: for any  $E \in \mathcal{E}$ ,  $i \geq 0$  and  $\bar{x} \in X^n$ , let  $\mathbb{B}_{E,i,\bar{x}}$  denote an arbitrarily chosen dynamic set  $\mathbb{B} \in \mathcal{N}_i$  that satisfies  $|\mathbb{B}_{E,i,\bar{x}} \Delta (E \cap \bar{x})| \leq \epsilon_i^2 n$ . Further, define the random variable over dynamic sets  $\mathbf{B}_{E,i} = \mathbb{B}_{E,i,\bar{x}}$  and notice that  $|\mathbf{B}_{E,i} \Delta (E \cap \bar{x})| \leq \epsilon_i^2 n$ . Since  $\epsilon_i \rightarrow 0$ , we have  $\mathbf{Y}_{\mathbf{B}_{E,i}} \rightarrow \mathbf{Y}_E$ , hence,

$$|\mathbf{Y}_E| = \left| \mathbf{Y}_{\mathbf{B}_{E,0}} + \sum_{i=1}^{\infty} (\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}) \right| \leq |\mathbf{Y}_{\mathbf{B}_{E,0}}| + \sum_{i=1}^{\infty} |\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}|. \quad (21)$$

The hope is for the sum in the right hand side of Equation (21) to converge. Indeed, notice that

$$|\mathbf{B}_{E,i} \Delta \mathbf{B}_{E,i-1}| \leq |\mathbf{B}_{E,i} \Delta E| + |E \Delta \mathbf{B}_{E,i-1}| \leq (\epsilon_i^2 + \epsilon_{i-1}^2) n \leq 2\epsilon_{i-1}^2 n, \quad (22)$$

hence, as  $i$  increases, the differences  $|\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}|$  tend to decrease. Taking a maximum over  $E \in \mathcal{E}$  in (21), we have

$$\max_{E \in \mathcal{E}} |\mathbf{Y}_E| \leq \max_{E \in \mathcal{E}} |\mathbf{Y}_{\mathbf{B}_{E,0}}| + \sum_{i=1}^{\infty} \max_{E \in \mathcal{E}} |\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}|. \quad (23)$$

We show how to bound the summand corresponding to any  $i \geq 1$  while term  $\max_{E \in \mathcal{E}} |\mathbf{Y}_{\mathbf{B}_{E,0}}|$  can be similarly bounded. Since  $\mathbf{Y}_{\mathbf{B}_{E,i}} \in \mathcal{N}_i$  and  $\mathbf{Y}_{\mathbf{B}_{E,i-1}} \in \mathcal{N}_{i-1}$ , there can be at most  $|\mathcal{N}_{i-1}| |\mathcal{N}_i| = N(\mathcal{E}, \epsilon_{i-1}) N(\mathcal{E}, \epsilon_i)$  distinct differences  $\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}$ . Intuitively, as  $i$  increases, the maximum is taken over more elements, however, the individual differences are smaller, and the hope is that the sum in Equation (23) would converge.

Using Equation (22), we can apply Lemma 9.11 with  $m = 2\epsilon_{i-1}^2 n$  to obtain that each distance  $\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}$  is bounded by  $C\epsilon_i \sqrt{n \log(1/\delta')}$  with probability  $1 - \delta'$ . Taking  $\delta'$  smaller than  $1/N(\mathcal{E}, \epsilon_{i-1})N(\mathcal{E}, \epsilon_i)$  and applying a union bound over at most  $N(\mathcal{E}, \epsilon_{i-1})N(\mathcal{E}, \epsilon_i)$  elements, we derive that with high probability,

$$\max_{E \in \mathcal{E}} |\mathbf{Y}_{\mathbf{B}_{E,i}} - \mathbf{Y}_{\mathbf{B}_{E,i-1}}| \leq O\left(\epsilon_{i-1} \sqrt{n \log N(\mathcal{E}, \epsilon_i)}\right). \quad (24)$$

We further take a union bound over  $i \geq 0$  and derive by (23) and (24) that w.h.p,

$$\max_{E \in \mathcal{E}} |\mathbf{Y}_E| \leq O\left(\sqrt{n} \sum_{i=1}^{\infty} \epsilon_{i-1} \sqrt{\log N(\mathcal{E}, \epsilon_i)}\right).$$

A standard choice for  $\epsilon_i$  is  $\epsilon_i = 2^{-i}$ , and this yields

$$\max_{E \in \mathcal{E}} |\mathbf{Y}_E| \leq O\left(\sqrt{n} \sum_{i=1}^{\infty} 2^{-i} \sqrt{\log N(\mathcal{E}, 2^{-i})}\right) \leq O\left(\sqrt{n} \int_{\epsilon=0}^1 \sqrt{\log N(\mathcal{E}, \epsilon)}\right), \quad (25)$$

where the last inequality is by approximating the sum with an integral. This, in fact, is the celebrated *Dudley's integral*.

**Remark 9.16.** The choice of  $\epsilon_i = 2^{-i}$  can generally only yield bounds on  $\mathbb{E} \max_{E \in \mathcal{E}} |\mathbf{Y}_E|$ , rather than high probability bounds. It is common in literature to obtain high probability bounds by first bounding the expectation  $\mathbb{E} \max_{E \in \mathcal{E}} |\mathbf{Y}_E|$  and then showing that this maximum concentrates around its expectation, via McDiarmid-like inequalities. However, such concentration inequalities cannot be applied in the adversarial setting. Hence we use instead a different well-studied choice of  $\epsilon_i$  that directly gives high probability bounds.

### 9.3.4 Proof of Lemma 9.7

The proof is by the standard technique of chaining, with adaptations to handle fractional covering numbers. Our goal is to bound  $\max_{E \in \mathcal{E}} \mathbf{Y}_E$ . The main idea is to create finer and finer approximations for  $\mathcal{E}$  using fractional epsilon nets. Formally, let  $\epsilon_j$  denote the minimal value of  $\epsilon$  such that  $N'(\mathcal{E}, \epsilon) \leq 2^{2^j}$ . We will derive the following bound: with probability  $1 - \delta$ ,

$$\max_{E \in \mathcal{E}} |\mathbf{Y}_E| \leq C\sqrt{n} \left( \sqrt{\log(1/\delta)} + \sum_{j=0}^{\infty} \epsilon_j \sqrt{\log N'(\mathcal{E}, \epsilon_j)} \right) \leq C\sqrt{n} \left( \sqrt{\log(1/\delta)} + \sum_{j=0}^{\infty} \epsilon_j 2^{j/2} \right). \quad (26)$$

This series in the right hand side is known to be equal, up to constant factors, to the following integral, which is known as Dudley's integral. We include the proof for completeness.

**Lemma 9.17.**

$$\sum_{j=0}^{\infty} \epsilon_j 2^{j/2} \leq \frac{\sqrt{2}}{\sqrt{2}-1} \int_0^1 \sqrt{\log N'(\mathcal{E}, \epsilon)} d\epsilon.$$

*Proof.* Notice that

$$\begin{aligned} \sum_{j=0}^{\infty} \epsilon_j 2^{j/2} &\leq \sum_{j=0}^{\infty} \sum_{i=j}^{\infty} (\epsilon_i - \epsilon_{i+1}) 2^{j/2} = \sum_{i=0}^{\infty} \sum_{j \leq i} (\epsilon_i - \epsilon_{i+1}) 2^{j/2} \\ &= \sum_{i=0}^{\infty} (\epsilon_i - \epsilon_{i+1}) \frac{\sqrt{2}^{i+1} - 1}{\sqrt{2} - 1} \leq \frac{\sqrt{2}}{\sqrt{2} - 1} \sum_{i=0}^{\infty} (\epsilon_i - \epsilon_{i+1}) 2^{i/2}. \end{aligned}$$

Further, by definition of  $\epsilon_i$  we have that for any  $\epsilon < \epsilon_i$ ,  $N'(\mathcal{E}, \epsilon) > 2^{2^i}$ , hence

$$\sqrt{\log N'(\mathcal{E}, \epsilon)} > 2^{i/2}.$$

This implies that

$$\sum_{i=0}^{\infty} (\epsilon_i - \epsilon_{i+1}) 2^{i/2} \leq \sum_{i=0}^{\infty} \int_{\epsilon_{i+1}}^{\epsilon_i} \sqrt{\log N'(\mathcal{E}, \epsilon)} \leq \int_0^1 \sqrt{\log N'(\mathcal{E}, \epsilon)}.$$

□

We start with some definitions; for any dynamic sets  $\mathbb{B}$  and  $\mathbb{B}'$  and  $E \in \mathcal{E}$ :

- Let  $\mathbb{B} \setminus \mathbb{B}'$  be defined by  $(\mathbb{B} \setminus \mathbb{B}')(x_1, \dots, x_n) = \mathbb{B}(x_1, \dots, x_n) \setminus \mathbb{B}'(x_1, \dots, x_n)$ . Notice that  $\mathbb{B} \setminus \mathbb{B}'$ , as defined, is a dynamic set.

- Define  $\mathbb{B}_{\leq m}$  as the dynamic sets that simulates  $\mathbb{B}$  up to the point where  $\mathbb{B}$  has  $m$  elements, and then it stops adding elements.
- Let  $\bar{x} = \bar{x}(\mathcal{A}, I)$  denote the stream that is output by the adversary  $\mathcal{A}$  in interaction with the sampler that samples the coordinates from  $I$ .

Let  $\mu_j$  be a probability measure over dynamic sets that is a fractional  $(\epsilon_j, N'(\mathcal{E}, \epsilon_j))$ -cover for  $\mathcal{E}$ . We will approximate each  $E \in \mathcal{E}$  using dynamic sets  $\pi_{E,0}, \pi_{E,1}, \pi_{E,2}, \dots$ , where  $\pi_{E,j} \in \text{support}(\mu_j)$  is an  $\epsilon$ -approximation for  $E$ , namely,

$$|(E \cap \bar{x}) \Delta \pi_{E,j}(\bar{x})| \leq \epsilon^2 n,$$

Notice that by definition of (fractional) covers,  $\pi_{E,j}$  may depend on the stream  $\bar{x}$ , which is a random variable, hence  $\pi_{E,j}$  is also a random variable. The following lemma shows that for a sufficiently large  $j$ ,  $\pi_{E,j}(\bar{x})$  equals  $E \cap \bar{x}$ .

**Lemma 9.18.** *Assume that  $N'(\mathcal{E}, \epsilon) < \infty$  for all  $\epsilon > 0$ . Then there exists  $j_1 > 0$  such that for all  $E \in \mathcal{E}$  and all  $j \geq j_1$ ,  $E \cap \bar{x} = \pi_{E,j}(\bar{x})$  holds with probability 1 over  $\bar{x}$ . Consequently,  $\mathbf{Y}_E = \mathbf{Y}_{\pi_{E,j}}$  for all  $j \geq j_1$ .*

*Proof.* The assumption of the lemma implies that for some  $j, \epsilon_j < n^{-1/2}$  and let  $j_1$  be the minimal such value. By definition of  $\pi_{E,j}$  we have that for all  $j \geq j_1$ ,  $|(E \cap \bar{x}) \Delta \pi_{E,j}(\bar{x})| < 1$ , hence  $E \cap \bar{x} = \pi_{E,j}(\bar{x})$  as required.  $\square$

We can assume that  $N'(\mathcal{E}, \epsilon) < \infty$  for all  $\epsilon > 0$  otherwise Dudley's integral (appearing in Lemma 9.7) would diverge. Hence, by Lemma 9.18, for any  $j_0 \geq 0$ ,

$$|\mathbf{Y}_E| := \left| \mathbf{Y}_{\pi_{E,j_0}} + \sum_{j=j_0}^{\infty} (\mathbf{Y}_{\pi_{E,j+1}} - \mathbf{Y}_{\pi_{E,j}}) \right| \leq \left| \mathbf{Y}_{\pi_{E,j_0}} \right| + \sum_{j=j_0}^{\infty} \left| \mathbf{Y}_{\pi_{E,j+1}} - \mathbf{Y}_{\pi_{E,j}} \right| \quad (27)$$

$$= \left| \mathbf{Y}_{\pi_{E,j_0}} \right| + \sum_{j=j_0}^{\infty} \left| \mathbf{Y}_{\pi_{E,j+1} \setminus \pi_{E,j}} - \mathbf{Y}_{\pi_{E,j} \setminus \pi_{E,j+1}} \right| \leq \left| \mathbf{Y}_{\pi_{E,j_0}} \right| + \sum_{j=j_0}^{\infty} \left| \mathbf{Y}_{\pi_{E,j+1} \setminus \pi_{E,j}} \right| + \sum_{j=j_0}^{\infty} \left| \mathbf{Y}_{\pi_{E,j} \setminus \pi_{E,j+1}} \right|. \quad (28)$$

We bound the supremum over  $E \in \mathcal{E}$  by taking the supremum over each term separately:

$$\sup_{E \in \mathcal{E}} |\mathbf{Y}_E| \leq \sup_{E \in \mathcal{E}} \left| \mathbf{Y}_{\pi_{E,j_0}} \right| + \sum_{j=j_0}^{\infty} \sup_{E \in \mathcal{E}} \left| \mathbf{Y}_{\pi_{E,j+1} \setminus \pi_{E,j}} \right| + \sum_{j=j_0}^{\infty} \sup_{E \in \mathcal{E}} \left| \mathbf{Y}_{\pi_{E,j} \setminus \pi_{E,j+1}} \right|. \quad (29)$$

Each supremum will be bounded using the generalized union bound Lemma 9.13.

Next, we define measures over differences of dynamic sets, that will be used to bound the right hand side of (29). For any  $j \geq 1$  we let  $\mu_{j,0}$  denote a probability measure over dynamic sets such that  $\mathbb{B} \sim \mu_{j,0}$  is drawn by first drawing  $\mathbb{B}_j \sim \mu_j$  and  $\mathbb{B}_{j+1} \sim \mu_{j+1}$  and then outputting  $\mathbb{B} = (\mathbb{B}_j \setminus \mathbb{B}_{j+1})_{\leq 2\epsilon_j^2}$ . Similarly, let  $\mu_{j,1}$  denote the measure that outputs  $(\mathbb{B}_{j+1} \setminus \mathbb{B}_j)_{\leq 2\epsilon_j^2 n}$ . By the generalized union bound, we have the following:

**Lemma 9.19.** *Let  $j_0 = \lceil \log_2 \log_2(1/\delta) \rceil$ . With probability at least  $1 - \delta$ , the following holds:*

- It holds that

$$\mu_{j_0} \left( \left\{ \mathbb{B} \in \text{support}(\mu_{j_0}) : |\mathbf{Y}_{\mathbb{B}}| > C_0 \sqrt{\log(1/\delta)} \sqrt{n} \right\} \right) \leq \frac{1}{3 \cdot 2^{2j_0}}.$$

- For all  $j \geq j_0$  and all  $b \in \{0, 1\}$ ,

$$\mu_{j,b} \left( \left\{ \mathbb{B} \in \text{support}(\mu_{j,b}) : |Y_{\mathbb{B}}| > C_0 \epsilon_j 2^{j/2} \sqrt{n} \right\} \right) \leq \frac{1}{18 \cdot 2^{2j} \cdot 2^{2^{j+1}}}.$$

*Proof.* The proof follows directly from Lemma 9.14. First, we show that the first item holds with probability  $1 - \delta/2$ : it follows by substituting in Lemma 9.14 the values  $m = n$ ,  $\delta' = \delta/2$  and  $\alpha = \frac{1}{3 \cdot 2^{2^{j_0}}}$ , and notice that

$$\log_2(1/\alpha) = \log_2 3 + 2^{j_0} \leq C \cdot \log(1/\delta),$$

by definition of  $j_0$ .

For the second item, we show that the term corresponding to a specific  $j \geq j_0$  and  $b \in \{0, 1\}$  holds with probability  $1 - 2^{-j-3}\delta$ . Indeed, we can substitute  $m = 2\epsilon_j^2 n$ ,  $\delta' = 2^{-j-3}\delta$  and  $\alpha = \frac{1}{18 \cdot 2^{2j} \cdot 2^{2^{j+1}}}$ , and notice that

$$\log_2 \frac{1}{\delta'} = \log_2 \frac{1}{\delta} + j + 3 \leq C \cdot 2^j,$$

since  $j \geq j_0$  and by definition of  $j_0$ . Further,

$$\log_2 \frac{1}{\alpha} = \log_2 18 + 2^{j+2} \leq C \cdot 2^j.$$

By a union bound, the failure probability is bounded by

$$\delta/2 + 2 \cdot \sum_{j=j_0}^{\infty} 2^{-j-3}\delta \leq \delta/2 + 2^{-j_0-1}\delta \leq \delta/2 + \delta/2 = \delta.$$

□

For the remainder of the proof, we fix some stream  $\bar{x}$  such that the condition in Lemma 9.19 holds when  $\bar{x} = \bar{x}$ . This fixes values  $\{Y_E\}_{E \in \mathcal{E}}$  and  $\{Y_{\mathbb{B}}\}_{\mathbb{B} \text{ dynamic set}}$  such that  $Y_E = Y_E$  and  $Y_{\mathbb{B}} = Y_{\mathbb{B}}$  for all  $E$  and  $\mathbb{B}$ . We will show that for any  $E \in \mathcal{E}$

$$|Y_E| \leq C \sqrt{n \log(1/\delta)} + C \sum_{j=0}^{\infty} 2^{j/2} \epsilon_j \sqrt{n},$$

which suffices to complete the proof from Lemma 9.17. Fix  $E \in \mathcal{E}$ ; we show how to define  $\pi_{E,j}$ . First, for any  $j \geq j_0$ , let  $A_j$  denote the set of elements  $\mathbb{B} \in \text{support}(\mu_j)$  such that  $|(E \cap \bar{x}) \Delta \mathbb{B}(\bar{x})| \leq \epsilon_j^2 n$ . By the property of the fractional cover  $\mu_j$ , we know that

$$\mu_j(A_j) \geq 1/2^{2^j}. \quad (30)$$

The set  $A_j$  contains the candidates for  $\pi_{E,j}$ . Notice that in order to bound (29), we would like to bound  $Y_{\pi_{E,j} \setminus \pi_{E,j-1}}$  and  $Y_{\pi_{E,j-1} \setminus \pi_{E,j}}$ . For this purpose, we now define for any  $j \geq j_0$  the function  $R: \text{support}(\mu_j) \times \text{support}(\mu_{j+1}) \rightarrow \{0, 1\}$ , that indicates which pairs of elements  $\mathbb{B}_j, \mathbb{B}_{j+1}$  are not suitable to be defined as  $\pi_{E,j}$  and  $\pi_{E,j+1}$ :

$$R(\mathbb{B}_j, \mathbb{B}_{j+1}) = \begin{cases} 1 & \max \left( |Y_{\mathbb{B}_j \setminus \mathbb{B}_{j+1}}|, |Y_{\mathbb{B}_{j+1} \setminus \mathbb{B}_j}| \right) > C_0 \epsilon_j 2^{j/2} \sqrt{n}, \mathbb{B}_j \in A_j, \mathbb{B}_{j+1} \in A_{j+1}, \\ 0 & \text{otherwise} \end{cases},$$

where  $C_0$  is the constant from Lemma 9.19. Next, we further restrict the set of candidates by creating a set  $A'_j \subseteq A_j$ , that contains only dynamic sets  $\mathbb{B}_j \in A_j$  such that for many elements  $\mathbb{B}_{j+1} \in A_{j+1}$ , the pair  $(\mathbb{B}_j, \mathbb{B}_{j+1})$  is suitable, which is formally defined as:

$$A'_j = \left\{ \mathbb{B}_j \in A_j : \mu_{j+1}(\{\mathbb{B}_{j+1} : R(\mathbb{B}_j, \mathbb{B}_{j+1}) = 1\}) \leq \frac{1}{3 \cdot 2^{2^{j+1}}} \right\}.$$

Next, we lower bound the measure  $\mu_j(A'_j)$ :

**Lemma 9.20.** *Assume that the high probability event from Lemma 9.19 holds. Then, for any  $j \geq j_0$ ,*

$$\mu_j(A'_j) \geq \frac{2}{3 \cdot 2^{2^j}}.$$

*Proof.* Fix some stream  $\bar{x}$  such that the high probability event of Lemma 9.19 holds; this fixes the values of  $Y_{\mathbb{B}}$  for all  $\mathbb{B}$ . If  $\mathbb{B}_j \sim \mu_j$  and  $\mathbb{B}_{j+1} \sim \mu_{j+1}$  are drawn independently, then

$$\begin{aligned} \mathbb{E}_{\mathbb{B}_j \sim \mu_j} \mathbb{E}_{\mathbb{B}_{j+1} \sim \mu_{j+1}} [R(\mathbb{B}_j, \mathbb{B}_{j+1})] &= \Pr_{\mathbb{B}_j \sim \mu_j, \mathbb{B}_{j+1} \sim \mu_{j+1}} [R(\mathbb{B}_j, \mathbb{B}_{j+1}) = 1] \\ &= \Pr_{\mathbb{B}_j, \mathbb{B}_{j+1}} \left[ \max(|Y_{\mathbb{B}_j \setminus \mathbb{B}_{j+1}}|, |Y_{\mathbb{B}_{j+1} \setminus \mathbb{B}_j}|) > C_0 \epsilon_j 2^{j/2} \sqrt{n}, \mathbb{B}_j \in A_j, \mathbb{B}_{j+1} \in A_{j+1} \right] \\ &\leq \Pr_{\mathbb{B}_j, \mathbb{B}_{j+1}} \left[ |Y_{\mathbb{B}_j \setminus \mathbb{B}_{j+1}}| > C_0 \epsilon_j 2^{j/2} \sqrt{n}, \mathbb{B}_j \in A_j, \mathbb{B}_{j+1} \in A_{j+1} \right] \end{aligned} \quad (31)$$

$$+ \Pr_{\mathbb{B}_j, \mathbb{B}_{j+1}} \left[ |Y_{\mathbb{B}_{j+1} \setminus \mathbb{B}_j}| > C_0 \epsilon_j 2^{j/2} \sqrt{n}, \mathbb{B}_j \in A_j, \mathbb{B}_{j+1} \in A_{j+1} \right], \quad (32)$$

We focus on bounding (31); (32) is bounded in a similar fashion. For any  $\mathbb{B}_j \in A_j$  and  $\mathbb{B}_{j+1} \in A_{j+1}$  we have that by definition of  $A_j$  and  $A_{j+1}$ ,

$$\begin{aligned} |(\mathbb{B}_j \setminus \mathbb{B}_{j+1})(\bar{x})| &= |\mathbb{B}_j(\bar{x}) \setminus \mathbb{B}_{j+1}(\bar{x})| \leq |\mathbb{B}_j(\bar{x}) \Delta \mathbb{B}_{j+1}(\bar{x})| \leq |\mathbb{B}_j(\bar{x}) \Delta E| + |E \Delta \mathbb{B}_{j+1}(\bar{x})| \\ &\leq \epsilon_j^2 n + \epsilon_{j+1}^2 n \leq 2\epsilon_j^2 n, \end{aligned}$$

This implies that

$$|(\mathbb{B}_j \setminus \mathbb{B}_{j+1})(\bar{x})| = |(\mathbb{B}_j \setminus \mathbb{B}_{j+1})_{\leq 2\epsilon_j^2 n}(\bar{x})|,$$

hence by definition of  $\mu_{j,0}$  and by Lemma 9.19,

$$\begin{aligned} &\Pr_{\mathbb{B}_j, \mathbb{B}_{j+1}} \left[ |Y_{\mathbb{B}_j \setminus \mathbb{B}_{j+1}}| > C_0 \epsilon_j 2^{j/2} \sqrt{n}, \mathbb{B}_j \in A_j, \mathbb{B}_{j+1} \in A_{j+1} \right] \\ &\leq \Pr_{\mathbb{B}_j, \mathbb{B}_{j+1}} \left[ \left| Y_{(\mathbb{B}_j \setminus \mathbb{B}_{j+1})_{\leq 2\epsilon_j^2 n}} \right| > C_0 \epsilon_j 2^{j/2} \sqrt{n} \right] = \Pr_{\mathbb{B} \sim \mu_{j,0}} \left[ |Y_{\mathbb{B}}| > C_0 \epsilon_j 2^{j/2} \sqrt{n} \right] \\ &\leq \frac{1}{18 \cdot 2^{2^j} \cdot 2^{2^{j+1}}}. \end{aligned}$$

Similarly, (32) is bounded by the same quantity, hence by (31) and (32) above,

$$\mathbb{E}_{\mathbb{B}_j \sim \mu_j} \mathbb{E}_{\mathbb{B}_{j+1} \sim \mu_{j+1}} [R(\mathbb{B}_j, \mathbb{B}_{j+1})] \leq \frac{1}{9 \cdot 2^{2^j} \cdot 2^{2^{j+1}}}.$$

By Markov's inequality,

$$\begin{aligned} \Pr_{\mathbb{B}_j \sim \mu_j} [\mathbb{B}_j \in A_j \setminus A'_j] &= \Pr_{\mathbb{B}_j \sim \mu_j} \left[ \mu_{j+1} (\{\mathbb{B}_{j+1} : R(\mathbb{B}_j, \mathbb{B}_{j+1}) = 1\}) > \frac{1}{3 \cdot 2^{2^{j+1}}} \right] \\ &= \Pr_{\mathbb{B}_j \sim \mu_j} \left[ \mathbb{E}_{\mathbb{B}_{j+1} \sim \mu_{j+1}} [R(\mathbb{B}_j, \mathbb{B}_{j+1})] > \frac{1}{3 \cdot 2^{2^{j+1}}} \right] \\ &\leq \frac{\mathbb{E}_{\mathbb{B}_j \sim \mu_j} \mathbb{E}_{\mathbb{B}_{j+1} \sim \mu_{j+1}} [R(\mathbb{B}_j, \mathbb{B}_{j+1})]}{1/(3 \cdot 2^{2^{j+1}})} \leq \frac{1}{3 \cdot 2^{2^j}}. \end{aligned}$$

By (30),

$$\mu_j(A'_j) \geq \mu_j(A_j) - \mu_j(A_j \setminus A'_j) \geq \frac{1}{2^{2^j}} - \frac{1}{3 \cdot 2^{2^j}} = \frac{2}{3 \cdot 2^{2^j}}.$$

□

To complete the proof, we show how to define  $\pi_{E,j}$  inductively on  $j$  such that  $R(\pi_{E,j}, \pi_{E,j+1}) = 0$  for all  $j \geq j_0$ . First, we select  $\pi_{E,j_0}$  to be any element of  $A'_{j_0}$  such that  $|Y_{\pi_{E,j_0}}| \leq C_0 \sqrt{\log(1/\delta)} \sqrt{n}$ ; such an element exists from Lemma 9.19 and Lemma 9.20. Next, assume for  $j \geq j_0$  that  $\pi_{E,j}$  was already selected and select  $\pi_{E,j+1}$  to be any element  $\mathbb{B}_{j+1} \in A'_{j+1}$  such that  $R(\pi_{E,j}, \mathbb{B}_{j+1}) = 0$ . Such an element exists since  $\pi_{E,j} \in A'_j$ , hence by definition of  $A'_j$ ,

$$\mu_{j+1} (\{\mathbb{B}_{j+1} : R(\pi_{E,j}, \mathbb{B}_{j+1}) = 1\}) \leq \frac{1}{3 \cdot 2^{2^{j+1}}},$$

while  $\mu_{j+1}(A'_{j+1}) \geq 2/(3 \cdot 2^{2^{j+1}})$  by Lemma 9.20. By (28) and since we defined  $\pi_{E,j}$  such that  $|Y_{\pi_{E,j_0}}| \leq C_0 \sqrt{\log(1/\delta)} \sqrt{n}$ ,  $R(\pi_{E,j}, \pi_{E,j+1}) = 0$  and  $\pi_{E,j} \in A_j$  for all  $j$ , we have that

$$|Y_E| \leq |Y_{\pi_{E,j_0}}| + \sum_{j=j_0}^{\infty} |Y_{\pi_{E,j+1} \setminus \pi_{E,j}}| + \sum_{j=j_0}^{\infty} |Y_{\pi_{E,j} \setminus \pi_{E,j+1}}| \leq C_0 \sqrt{n} \left( \sqrt{\log(1/\delta)} + 2 \sum_{j=j_0}^{\infty} \epsilon_j 2^{j/2} \right).$$

This proves (26) as required.

## 9.4 Bounds on $\epsilon$ -Nets via Fractional Covering Numbers

The goal of this section is to prove Lemma 9.8. First, we use the following Martingale bound (see Section A.1 for an introduction to Martingales):

**Lemma 9.21** (Freedman's inequality [Fre75]). *Let  $\mathbf{y}_0, \dots, \mathbf{y}_n$  be a Martingale adapted to the filtration  $F_0, \dots, F_n$ , such that  $\sum_{i=1}^n \mathbb{E}[(\mathbf{y}_i - \mathbf{y}_{i-1})^2 \mid F_{i-1}] \leq s$  holds almost surely for some  $s > 0$ . Further, assume that  $|\mathbf{y}_i - \mathbf{y}_{i-1}| \leq M$  for all  $i$ . Then, for any  $t > 0$ ,*

$$\Pr[\mathbf{y}_n - \mathbf{y}_0 \geq t] \leq \exp\left(-\frac{t^2}{2(s + Mt)}\right).$$

We derive the following concentration bound for a dynamic set:

**Lemma 9.22.** *Let  $\mathbb{B}$  be a dynamic set with  $|\mathbb{B}| \leq m$ . Let  $\mathbf{I} \sim \text{Ber}(n, p)$ . Let  $\bar{\mathbf{x}} = \bar{\mathbf{x}}(\mathcal{A}, \mathbf{I})$ . Then, for any  $t \geq 0$ ,*

$$\Pr[|\mathbb{B}(\bar{\mathbf{x}}) \cap \bar{\mathbf{x}}_{\mathbf{I}}| \leq p|\mathbb{B}(\bar{\mathbf{x}})| - t] \leq \exp(-ct^2/(mp + t)).$$



*Proof.* This follows from Lemma 9.21. We apply this lemma with  $\mathbf{y}_i = |\mathbb{B}(\bar{x}) \cap \bar{x}_{I \cap [i]}| - p|\mathbb{B}(\bar{x}) \cap \bar{x}_{[i]}|$  and  $F_i = \sigma(I \cap [i])$ . We can substitute  $s = c'mp$ , due to the following reason: conditioned on  $F_{i-1}$ , we know  $\bar{x}_{[i]}$ , which implies that we know whether  $x_i \in \mathbb{B}(\bar{x})$ . If this holds true, then,

$$\mathbf{y}_i - \mathbf{y}_{i-1} = \begin{cases} 1 - p & \text{if } i \in I \text{ (holds with probability } p) \\ -p & \text{otherwise (holds with probability } 1 - p) \end{cases}.$$

By simple calculations we have  $\mathbb{E}[(\mathbf{y}_i - \mathbf{y}_{i-1})^2 \mid F_{i-1}] \leq c'p$  in this case. If  $x_i \notin \mathbb{B}(\bar{x})$ , then  $\mathbf{y}_i = \mathbf{y}_{i-1}$  and we have  $\mathbb{E}[(\mathbf{y}_i - \mathbf{y}_{i-1})^2 \mid F_{i-1}] = 0$ . Since  $\bar{x}_i \in \mathbb{B}(\bar{x})$  can hold true for at most  $m$  values of  $i$ , then  $\sum_{i=1}^n \mathbb{E}[(\mathbf{y}_i - \mathbf{y}_{i-1})^2 \mid F_{i-1}] \leq c'mp$ . Further, we can substitute  $M = 1$  in Lemma 9.21, and the result follows.  $\square$

We are ready to prove Lemma 9.8:

*Proof of Lemma 9.8.* Let  $\mathcal{N}$  be a 0-net for  $\mathcal{E}$  with minimal cardinality. For each  $\mathbb{B} \in \mathcal{N}$ , let  $\mathbb{B}_{\leq m}$  be the dynamic set that simulates  $\mathbb{B}$  up to the point that it has retained  $m$  elements, and then it discards all the remaining elements. We apply Lemma 9.22 on  $\mathbb{B}_{\leq m}$  to obtain that

$$\begin{aligned} \Pr[|\mathbb{B}(\bar{x})| \geq m, |\mathbb{B}(\bar{x}) \cap \bar{x}_I| \leq mp/2] &\leq \Pr[|\mathbb{B}_{\leq m}(\bar{x})| = m, |\mathbb{B}_{\leq m}(\bar{x}) \cap \bar{x}_I| \leq mp/2] \\ &\leq \Pr[|\mathbb{B}_{\leq m}(\bar{x}) \cap \bar{x}_I| \leq p|\mathbb{B}_{\leq m}(\bar{x})| - mp/2] \leq \exp(-cmp). \end{aligned}$$

The proof follows by a union bound over  $\mathbb{B} \in \mathcal{N}$ .  $\square$

## 10 Reductions Between Different Sampling Schemes

This section establishes a framework that enables one to obtain bounds with respect to one sampler in terms of bounds with respect to a different sampler. In particular, this shows how, given bounds on  $\epsilon$ -nets and  $\epsilon$ -approximations for the uniform sampler, one can obtain bounds for the Bernoulli and the reservoir sampler. Section 10.1 presents an overview of an abstract method to reduce between two sampling schemes, that is formally presented in Section 10.2. Section 10.3 shows how to obtain bounds with respect to the reservoir sampler given bounds for the uniform sampler. Section 10.4 presents bounds on the Bernoulli sampler based on bounds for the uniform sampler. Finally, Section 10.5 provides bounds with respect to the uniform sampler based on the Bernoulli sampler, which proves the auxiliary Lemma 6.2 and Lemma 7.2. Both Section 10.3 and Section 10.5 use the abstract reduction method of Section 10.2 while Section 10.4 utilizes the fact that the Bernoulli sampler can be presented as a mixture of uniform samplers  $\text{Uni}(n, k)$  for different values of  $k$ .

### 10.1 Intuition for the Reduction Method

For convenience, we consider samplers with no deletions, that are characterized by some distribution over subsets of  $[n]$ , such as  $\text{Uni}(n, k)$ , that is the uniform distribution over subsets of  $[n]$  of size  $[k]$ . We use  $I$  and  $I'$  to denote such random variables over subsets of  $[n]$  (e.g., they can be distributed  $\text{Ber}(n, p)$  or  $\text{Uni}(n, k)$ ).

In these reductions, our goal is to show that one sampling scheme  $I$  is at least as good as a different scheme  $I'$ . For example, that  $I$  attains  $\epsilon$ -approximations for values of  $\epsilon$  smaller than

those attained by  $I'$ . In other words, we would like to say that  $I$  is resilient to the adversary at least as well as  $I'$ . The above is equivalent to saying that the *worst* adversary for  $I'$  is *at least as bad* as the worst adversary for  $I$ . The above can be shown by reduction: given an adversary  $\mathcal{A}$  that is *bad* for  $I$ , we will construct an adversary  $\mathcal{A}'$  that is bad for  $I'$ .

Here we define  $\mathcal{A}'$ , that plays against a sampler that samples  $I'$ . The general idea for  $\mathcal{A}'$  is to simulate  $\mathcal{A}$ . However,  $\mathcal{A}$  is known to be bad against  $I$  while  $\mathcal{A}'$  plays against  $I'$ . To tackle this issue,  $\mathcal{A}'$  will simulate a sample  $J$  that has the same distribution as  $I$  and then simulate the actions of  $\mathcal{A}$  against  $J$ . Then,  $\mathcal{A}'$  will output the same stream output by the simulated  $\mathcal{A}$ .

We would like to show that  $\mathcal{A}'$  is bad against  $I'$ . In order to show that, we will have to assume that the simulated sample  $J$  is very close to the true sample  $I'$  with high probability (say, in symmetric difference of sets). Since  $\mathcal{A}$  is bad against  $I$  and  $J \sim I$ , the simulated actions of  $\mathcal{A}$  are bad against the simulated sample  $J$ . Since further  $J$  is very close to  $I'$ , then the simulated  $\mathcal{A}$  is bad also against  $I'$ . Since  $\mathcal{A}'$  outputs the same stream as the simulated  $\mathcal{A}$ , this implies that  $\mathcal{A}'$  is bad against  $I'$ , as required.

Notice that since  $\mathcal{A}'$  would like the simulated sample  $J$  to be similar to the true sample  $I'$ , it has to construct  $J$  based on  $I'$  and this defines a joint probability distribution between  $I'$  and  $J$ . Such a joint distribution is called *coupling*. Further, the simulation has to be performed in an online fashion: once  $\mathcal{A}'$  receives the actions taken by the sampler  $I'$  that it plays against, it has to immediately simulate the actions of the simulated sample  $J$ . In particular, once  $\mathcal{A}'$  knows whether  $t \in I'$ , it has to decide whether  $t \in J$ . Since the coupling between  $J$  and  $I'$  is constructed in an online fashion, we denote it an *online coupling*.

Using the notation above, the goal of  $\mathcal{A}'$  is to construct an online coupling of  $J$  and  $I'$  such that  $J \sim I$  and such that with high probability, the symmetric set difference between  $J$  and  $I'$  is small. This can be done, for example, if  $J \sim \text{Uni}(2k, k)$  and  $I' \sim \text{Ber}(2k, 1/2)$ : there,  $\mathcal{A}'$  will have to omit or add a approximately  $O(\sqrt{k})$  elements to  $I'$  to create  $J$ .

## 10.2 Abstract Reduction Method

We refer to sampling schemes that are oblivious to the adversary, namely that the choice to retain or discard an element is independent of the stream. Formally, we denote by  $I_t$  the set of indices of elements retained by the algorithm after seeing  $\bar{x}_{[t]}$ , for  $t \in [n]$ . An *oblivious sampling scheme* is one where  $\bar{I} = (I_1, \dots, I_n)$  are random variables jointly distributed, that are not a function of the adversary  $\mathcal{A}$ . This section compares one oblivious sampling scheme  $\bar{I} := (I_1, \dots, I_n)$  with another,  $\bar{I}' = (I'_1, \dots, I'_n)$ . It will be shown that if the adversary, given an input stream  $\bar{I}'$  can simulate a stream that is distributed according to  $\bar{I}$  such that with high probability, the input stream is close to the output stream in some sense, then the sampling scheme  $\bar{I}$  is at least as resilient to the adversary as  $\bar{I}'$ . We begin with the following definition of online simulation:

**Definition 10.1.** *An online simulator  $\mathcal{S}$  is an algorithm that receives a stream  $I_1, I_2, \dots, I_n$  of subsets of  $[n]$  and an unlimited pool of independent random bits and outputs a stream  $I'_1, \dots, I'_n$ , such that  $I'_t$  has to be computed before seeing  $I_{t+1}$ , for  $t \in [n]$ . In other words,  $I'_t$  depends only on  $I_1, \dots, I_t, I'_1, \dots, I'_{t-1}$  and on the randomness of the simulator. The joint distribution of  $\bar{I}$  and  $\bar{I}'$  is called an *online coupling* of  $\bar{I}$  to  $\bar{I}'$ . Equivalent, we can say that  $\bar{I}$  is *online coupled* to  $\bar{I}'$ .*

Notice that an online coupling of  $\bar{I}$  and  $\bar{I}'$  is also a *coupling* of these two random variables, which is any joint distribution between them. Further, notice that an online coupling is not a

symmetric notion: an online coupling of  $\bar{I}$  to  $\bar{I}'$  is not necessarily an online coupling of  $\bar{I}'$  to  $\bar{I}$ .

In cases that the sampler cannot delete elements from its sample, notice that  $I_t = I_n \cap [t]$ . To simplify the notation, we can write  $I = I_n$  and  $I \cap [t] = I_t$ . Hence, we have the following definition of online coupling for no-deletion samplers:

**Definition 10.2.** Let  $I$  and  $I'$  be jointly distributed random variables over  $[n]$ . We say that  $I$  is online coupled to  $I'$  if  $(I \cap [1], I \cap [2], \dots, I \cap [n])$  is online coupled to  $(I' \cap [1], I' \cap [2], \dots, I' \cap [n])$ .

In order to show that  $\bar{I}$  is more resilient to the adversary than  $\bar{I}'$ , it suffices to find an online coupling of  $I$  to  $I'$  such that  $I_n$  is similar to  $I'_n$  in some sense. To be more formal, let  $f(I_n, \bar{x})$  be some  $\{0, 1\}$ -valued function that we view as an indicator denoting whether the sub-sample  $\bar{x}_{I_n}$  fails to represent the full stream  $\bar{x}$ . For example,  $f$  can be an indicator of whether  $\bar{x}_{I_n}$  is *not* an  $\epsilon$ -approximation of  $\bar{x}$ . Our goal is to bound the failure probability with the worst adversary. Namely, to bound  $\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{\bar{I}}[f(I_n, \bar{x}(\mathcal{A}, \bar{I})) = 1]$ . Say that we already know how to bound a similar quantity for a different sampling scheme  $\bar{I}'$ . If we can online couple  $\bar{I}$  to  $\bar{I}'$ , then we can reduce between these two bounds:

**Lemma 10.3.** Let  $\bar{I}$  be online coupled to  $\bar{I}'$ . Let  $f, g: \{0, 1\}^n \times X^n \rightarrow \{0, 1\}$ . Then,

$$\begin{aligned} & \max_{\mathcal{A} \in \text{Adv}_n} \Pr_{\bar{I}}[f(I_n, \bar{x}(\mathcal{A}, \bar{I})) = 1] \\ & \leq \max_{\mathcal{A}' \in \text{Adv}_n} \Pr_{\bar{I}'}[g(I'_n, \bar{x}(\mathcal{A}', \bar{I}')) = 1] + \Pr[\exists \bar{x} \in X^n \text{ s.t. } f(I_n, \bar{x}) = 1 \text{ and } g(I'_n, \bar{x}) = 0]. \end{aligned} \quad (33)$$

Notice the second term in the right hand side of (33): it equals zero if  $I_n = I'_n$ , and, it is expected to be small if  $I_n \approx I'_n$  with high probability.

*Proof.* Let  $\mathcal{A}_{\max}$  be the adversary that achieves the maximum on the left hand side of (33). Let  $\mathcal{S}$  denote the simulation adversary that given  $\bar{I}'$  and some additional random string  $r$ , outputs  $\bar{I} = \mathcal{S}(\bar{I}', r)$ . We will create the following adversary  $\mathcal{A}'_r$  that operates on the stream  $\bar{I}'$  and has additional randomness  $r$ : it creates the sample  $\bar{I} = \mathcal{S}(\bar{I}', r)$ , simulates  $\mathcal{A}_{\max}$  on this sample and outputs the same stream as the simulated  $\mathcal{A}_{\max}$ . In particular, we have

$$\bar{x}(\mathcal{A}_r, \bar{I}') = \bar{x}(\mathcal{A}_{\max}, \bar{I}). \quad (34)$$

We view  $\mathcal{A}'_r$  as a distribution over deterministic adversaries  $\{\mathcal{A}'_r\}_{r \in \text{support}(r)}$ . Further, notice that each  $\mathcal{A}'_r$  defines an appropriate adversary. By (34),

$$\begin{aligned} & \max_{\mathcal{A} \in \text{Adv}_n} \Pr_{\bar{I}} [f(I_n, \bar{x}(\mathcal{A}, \bar{I})) = 1] = \Pr_{\bar{I}} [f(I_n, \bar{x}(\mathcal{A}_{\max}, \bar{I})) = 1] \\ & = \Pr_{\bar{I}, \bar{I}', r} \left[ f(I_n, \bar{x}(\mathcal{A}_{\max}, \bar{I})) = 1 \text{ and } g(I'_n, \bar{x}(\mathcal{A}'_r, \bar{I}')) = 1 \right] \\ & \quad + \Pr_{\bar{I}, \bar{I}', r} \left[ f(I_n, \bar{x}(\mathcal{A}_{\max}, \bar{I})) = 1 \text{ and } g(I'_n, \bar{x}(\mathcal{A}'_r, \bar{I}')) = 0 \right] \\ & = \Pr_{\bar{I}, \bar{I}', r} \left[ f(I_n, \bar{x}(\mathcal{A}_{\max}, \bar{I})) = 1 \text{ and } g(I'_n, \bar{x}(\mathcal{A}'_r, \bar{I}')) = 1 \right] \\ & \quad + \Pr_{\bar{I}, \bar{I}', r} \left[ f(I_n, \bar{x}(\mathcal{A}_{\max}, \bar{I})) = 1 \text{ and } g(I'_n, \bar{x}(\mathcal{A}_{\max}, \bar{I})) = 0 \right] \\ & \leq \Pr_{\bar{I}', r} \left[ g(I'_n, \bar{x}(\mathcal{A}'_r, \bar{I}')) = 1 \right] + \Pr_{\bar{I}, \bar{I}'} \left[ \exists \bar{x} \in X^n \text{ s.t. } f(\bar{I}, \bar{x}) = 1 \text{ and } g(\bar{I}', \bar{x}) = 0 \right], \end{aligned}$$

as required. □

### 10.3 Bounds for Reservoir Sampling via Uniform Sampling

Next, we show how to obtain bounds for reservoir sampling based on uniform sampling. The intuition is that the reservoir sampler gives less information than the uniform sampler: indeed, when the reservoir sampler selects an element, the adversary does not know whether this element will remain for the final sample, while this is not the case for the uniform sampler. The following holds:

**Lemma 10.4.** *The reservoir sampler  $\bar{I} = (I_1, \dots, I_n) \sim \text{Res}(n, k)$  can be online coupled to the Uniform sampler  $I' \sim \text{Uni}(n, k)$  such that  $I_n = I'$  with probability 1.*

Notice that we describe the uniform sampler using one index-set as it is an insertion-only scheme, while the reservoir sample has deletions hence we describe it using  $n$  index-sets. Lemma 10.4, in combination with Lemma 10.3, immediately implies that any high probability bound obtained for the uniform sampling, also holds true for the reservoir sampler:

*Proof of Theorem 6.1, reservoir sampling.* Let  $f(I, \bar{x}) = g(I, \bar{x})$  denote an indicator of whether  $\bar{x}_I$  fails to be an  $\epsilon$ -approximation for  $\bar{x}$ . Let  $\bar{I}$  denote the reservoir sampler and let  $I'$  denote the uniform sampler. Then, by Lemma 10.4, we can online couple  $\bar{I}$  to  $I'$  such that  $I_n = I'$ . By Lemma 10.3, we have

$$\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{\bar{I}} [f(I_n, \bar{x}(\mathcal{A}, \bar{I})) = 1] \leq \max_{\mathcal{A}' \in \text{Adv}_n} \Pr_{I'} [g(I', \bar{x}(\mathcal{A}', I')) = 1].$$

By Theorem 6.5, the right hand side is bounded by  $\delta$ , for a suitable value of  $\delta$ . Hence, the left hand side is bounded by the same quantity.

Lastly, notice that the assumption  $n \geq 2k$  in Theorem 6.5 translates to  $n \geq 3k$  in the reduction Lemma 10.4.  $\square$

*Proof of Theorem 7.1, reservoir sampling.* The proof follows the same steps as the proof for Theorem 6.1, while replacing  $\epsilon$ -approximations with  $\epsilon$ -nets and using Theorem 7.5 for the bound on the uniform sampler.  $\square$

Finally, we prove Lemma 10.4. First, an auxiliary lemma:

**Lemma 10.5.** *Let  $\bar{I} \sim \text{Res}(n, k)$  and fix  $t \in [n]$ . Then, conditioned on  $I_1, \dots, I_{t-1}, I_n \cap [t]$ , it holds that  $I_t$  is independent of  $I_n$ .*

*Proof.* The proof follows from the following steps:

- First, the conditional distribution of  $I_n$  conditioned on  $I_1 = I_1, \dots, I_t = I_t$  is only a function of  $I_t$ . That is due to the fact that  $I_1 \rightarrow I_2 \rightarrow \dots \rightarrow I_n$  is a Markov chain.
- This implies that the conditional distribution of  $I_n$  conditioned on  $I_1 = I_1, \dots, I_t = I_t, I_n \cap [t] = S$  is only a function of  $S$  and  $I_t$ .
- Conditioned on  $I_1 = I_1, \dots, I_t = I_t, I_n \cap [t] = S$ , we can write  $I_t = S \cup (I_t \setminus S)$ . Notice that due to the symmetry of deletion, namely, that the deleted element is chosen uniformly at random, one derives that the conditional probability of  $I_n$  conditioned on  $I_1 = I_1, \dots, I_t = I_t, I_n \cap [t] = S$  is not dependent on  $I_t \setminus S$ , hence it is only a function of  $S$ .

- The above implies that conditioned on  $I_n \cap [t] = S$ , the random vector  $(I_1, \dots, I_t)$  is independent of  $I_n$ .
- This further implies that conditioned on  $I_1, \dots, I_{t-1}, I_n \cap [t]$ , it holds that  $I_t$  is independent of  $I_n$ .

□

The following is a well-known fact that can be proved by induction:

**Lemma 10.6.** *Let  $\bar{I} = (I_1, \dots, I_n) \sim \text{Res}(n, k)$ . Then,  $I_n \sim \text{Uni}(n, k)$ .*

Using only Lemma 10.5 and Lemma 10.6, we can prove Lemma 10.4:

*Proof of Lemma 10.4.* Let  $I' \sim \text{Uni}(n, k)$  and  $\bar{I} \sim \text{Res}(n, k)$ . We will define a random variable  $\bar{J}$  that is online coupled to  $I'$  and show that both  $J_n = I'$  with probability 1 and that  $\bar{J}$  has the same distribution as  $\bar{I}$ . The sample  $\bar{J}$  is created using the following inductive argument: for  $t = 1, \dots, n$ , assume that we have already set  $J_1 = J_1, \dots, J_{t-1} = J_{t-1}$ , and that  $I' = I'$ , and recall that by definition of online simulation, we can set  $J_t$  to be any randomized function of  $J_1, \dots, J_{t-1}, I' \cap [t]$ . Specifically,  $J_t$  is drawn from the following conditional distribution: for any  $J_t$ ,

$$\begin{aligned} \Pr [J_t = J_t \mid J_1 = J_1, \dots, J_{t-1} = J_{t-1}, I' \cap [t] = I' \cap [t]] \\ = \Pr [I_t = J_t \mid I_1 = J_1, \dots, I_{t-1} = J_{t-1}, I_n \cap [t] = I' \cap [t]] \end{aligned}$$

Using the fact that the random coins used by the algorithm are independent of the sample  $I'$ , we derive that conditioned on  $J_1, \dots, J_{t-1}, I' \cap [t]$ , it holds that  $J_t$  is independent of  $I'$ . In combination with Lemma 10.5, it follows that

$$\begin{aligned} \Pr [J_t = J_t \mid J_1 = J_1, \dots, J_{t-1} = J_{t-1}, I' = I'] \\ = \Pr [J_t = J_t \mid J_1 = J_1, \dots, J_{t-1} = J_{t-1}, I' \cap [t] = I' \cap [t]] \\ = \Pr [I_t = J_t \mid I_1 = J_1, \dots, I_{t-1} = J_{t-1}, I_n \cap [t] = I' \cap [t]] \\ = \Pr [I_t = J_t \mid I_1 = J_1, \dots, I_{t-1} = J_{t-1}, I_n = I']. \end{aligned} \tag{35}$$

The following inductive argument shows that for  $t \in \{0, \dots, n\}$ , the conditional distribution  $J_1, \dots, J_t$  conditioned on  $I' = I'$  equals the conditional distribution of  $I_1, \dots, I_t$  conditioned on  $I_n = I'$ . For the base of the induction,  $t = 0$ , there is nothing to prove. For the induction step, assume that the above holds for  $t - 1$  and we will prove for  $t$ . Indeed, by the chain rule, the induction hypothesis, and (35),

$$\Pr [J_1 = J_1, \dots, J_t = J_t \mid I' = I'] \tag{36}$$

$$= \Pr [J_1 = J_1, \dots, J_{t-1} = J_{t-1} \mid I' = I'] \Pr [J_t = J_t \mid J_1 = J_1, \dots, J_{t-1} = J_{t-1}, I' = I'] \tag{37}$$

$$= \Pr [I_1 = J_1, \dots, I_{t-1} = J_{t-1} \mid I_n = I'] \Pr [I_t = J_t \mid I_1 = J_1, \dots, I_{t-1} = J_{t-1}, I_n = I'] \tag{38}$$

$$= \Pr [I_1 = J_1, \dots, I_t = J_t \mid I_n = I']. \tag{39}$$

This concludes the induction. It follows that  $J_n = I'$  with probability 1, since the conditional distribution of  $J_n$  conditioned on  $I' = I'$  equals the conditional distribution of  $I_n$  conditioned on  $I_n = I'$ , which constantly equals  $I'$ . This proves one of the guarantees on  $\bar{J}$ . Further, it implies

that  $J_n$  has the same distribution as  $I'$ , which, by Lemma 10.6 implies that  $J_n$  is distributed as  $I_n$ . In combination with (39), we derive that

$$\begin{aligned} \Pr[J_1 = J_1, \dots, J_n = J_n] &= \Pr[J_n = J_n] \Pr[J_1 = J_1, \dots, J_{n-1} = J_{n-1} \mid J_n = J_n] \\ &= \Pr[J_n = J_n] \Pr[J_1 = J_1, \dots, J_{n-1} = J_{n-1} \mid I' = J_n] \\ &= \Pr[I_n = I_n] \Pr[I_1 = J_1, \dots, I_{n-1} = J_{n-1} \mid I_n = J_n] \\ &= \Pr[I_1 = J_1, \dots, I_n = J_n], \end{aligned}$$

as required.  $\square$

## 10.4 Bounds for Bernoulli Sampling via Uniform Sampling

Here, our goal is to show that concentration guarantees on uniform sampling imply guarantees on Bernoulli sampling. Notice that the latter can be viewed as a mixture of uniform sampling schemes for different values of  $k$ . To be more precise, a Bernoulli sample  $\text{Ber}(n, p)$  can be obtained by first drawing  $k \sim \text{Bin}(n, p)$  and then drawing a uniform sample  $\text{Uni}(n, k)$ , where  $\text{Bin}$  denotes the binomial distribution. For this reason, it is *harder* to be adversarial against a Bernoulli sampler, because the adversary there does not know in advance what value of  $k$  is drawn.

To formalize this notion, we use  $f(I, \bar{x})$  as some indicator of failure of the sample  $\bar{x}_I$  to represent  $\bar{x}$ , for example, an indicator of whether  $\bar{x}_I$  is not an  $\epsilon$ -approximation of  $\bar{x}$ . One would like to minimize the probability that  $f = 1$ , against any adversary. The following lemma compares the failure probability of the Bernoulli sampling with that of the uniform sampling:

**Lemma 10.7.** *Let  $n \in \mathbb{N}$  and  $p \in [0, 1]$ . Let  $f: \{0, 1\}^n \times X^n \rightarrow \{0, 1\}$ . Then,*

$$\begin{aligned} &\max_{\mathcal{A} \in \text{Adv}_n} \Pr_{I \sim \text{Ber}(n, p)} [f(I, \bar{x}(\mathcal{A}, I)) = 1] \\ &\leq \max_{\substack{k \in \mathbb{N}: \\ np/2 \leq k \leq 3np/2}} \max_{\mathcal{A}^k \in \text{Adv}_n} \Pr_{I^k \sim \text{Uni}(n, k)} [f(I^k, \bar{x}(\mathcal{A}^k, I^k)) = 1] + 2 \exp(-cnp). \end{aligned}$$

The proof relies on a variant of Bernstein's inequality, on the tail of the binomial random variable:

**Lemma 10.8** (Bernstein's inequality). *Let  $k \sim \text{Bin}(n, p)$ . Then, for any  $\epsilon \in [0, 1]$ ,*

$$\Pr[|k - np| \geq \epsilon np] \leq 2 \exp(-\epsilon^2 np / 3).$$

*Proof of Lemma 10.7.* Let  $\mathcal{A}$  denote the maximizer with respect to the Bernoulli sample. First, we decompose the Bernoulli sample into a mixture of uniform samples. Let  $k \sim \text{Bin}(n, p)$  drawn from a binomial distribution; we have

$$\Pr_{I \sim \text{Ber}(n, p)} [f(I, \bar{x}(\mathcal{A}, I)) = 1] = \sum_{k=0}^n \Pr[k = k] \Pr_{I^k \sim \text{Uni}(n, k)} [f(I^k, \bar{x}(\mathcal{A}, I^k)) = 1].$$

First, summing the terms corresponding to  $np/2 \leq k \leq 3np/2$ , we have

$$\begin{aligned} &\sum_{k=\lceil np/2 \rceil}^{\lfloor 3np/2 \rfloor} \Pr[k = k] \Pr_{I^k \sim \text{Uni}(n, k)} [f(I^k, \bar{x}(\mathcal{A}, I^k)) = 1] \leq \max_{k: np/2 \leq k \leq 3np/2} \Pr_{I^k \sim \text{Uni}(n, k)} [f(I^k, \bar{x}(\mathcal{A}, I^k)) = 1] \\ &\leq \max_{k: np/2 \leq k \leq 3np/2} \max_{\mathcal{A}^k \in \text{Adv}_n} \Pr_{I^k \sim \text{Uni}(n, k)} [f(I^k, \bar{x}(\mathcal{A}^k, I^k)) = 1]. \end{aligned}$$

Next, the sum in the remaining terms is bounded by the probability that  $k \notin [np/2, 3np/2]$ , or equivalently, the probability that  $|k - np| > np/2$ . From Lemma 10.8, this is bounded by  $2 \exp(-cnp)$ .  $\square$

As a direct application, we derive Theorem 6.1 and Theorem 7.1 for the Bernoulli sampling, given the bounds corresponding to the uniform sampling:

*Proof of Theorem 6.1, Bernoulli sampling.* Let  $\delta \in (0, 1/2)$ . Define by  $f(I, \bar{x})$  the indicator of whether  $\bar{x}_I$  fails to be an  $\epsilon$ -approximation for  $\bar{x}$ , where  $\epsilon = C_0 \sqrt{(d + \log(1/\delta))/(np)}$  and  $C_0 > 0$  is a sufficiently large constant. From Theorem 6.5, for any  $k$  such that  $np/2 \leq k \leq 3np/2$  and any  $\mathcal{A} \in \text{Adv}_n$ , it holds that

$$\Pr_{I \sim \text{Uni}(n,k)} [f(I, \bar{x}(\mathcal{A}, I)) = 1] \leq \delta/2.$$

By Lemma 10.7, it follows that for any  $\mathcal{A} \in \text{Adv}_n$ ,

$$\Pr_{I \sim \text{Uni}(n,k)} [f(I, \bar{x}(\mathcal{A}, I)) = 1] \leq \delta/2 + \exp(-c_0 np), \quad (40)$$

for some universal constant  $c_0 > 0$ . Notice that if  $\epsilon > 1$  then the result trivially follows and otherwise, we have that

$$np \geq C_0^2 \log(1/\delta).$$

Assuming that  $C_0$  is sufficiently large, we have that

$$c_0 np \geq \log(2/\delta),$$

which implies that

$$\exp(-c_0 np) \leq \delta/2.$$

In combination with (40), this concludes the proof. Notice that the condition that  $2k \leq n$  in Theorem 6.5 translates here to  $3np \leq n$ .  $\square$

*Proof of Theorem 7.1, Bernoulli sampling.* The proof follows similar steps as the proof for Theorem 6.1, while replacing  $\epsilon$ -approximations with  $\epsilon$ -nets and using Theorem 7.5 for the bound on the uniform sampler.  $\square$

## 10.5 Bounds for Uniform Sampling via Bernoulli Sampling

This section reduces bounds for the uniform sampler  $I \sim \text{Uni}(2k, k)$  to bounds for the Bernoulli sampler  $I' \sim \text{Ber}(2k, p)$ . This is done via the method of online coupling, presented in Section 10.2. First, we present some simple well known properties of binary random variables, and then, we describe an online coupling of  $I$  to  $I'$ :

**Lemma 10.9.** *If  $\mathbf{y}$  and  $\mathbf{z}$  are two random variables over  $\{0, 1\}$ , then there exists a coupling (i.e. joint distribution) of them such that:*

1.  $\Pr[\mathbf{y} \neq \mathbf{z}] = |\Pr[\mathbf{y} = 1] - \Pr[\mathbf{z} = 1]|$ .
2. If  $\Pr[\mathbf{y} = 1] \geq \Pr[\mathbf{z} = 1]$  then  $\mathbf{y} \geq \mathbf{z}$  with probability 1.

*Proof.* One can couple the following way: first, draw a random variable  $\xi$  uniformly in  $[0, 1]$  and set  $y = 0$  if  $\xi \leq \Pr[y = 0]$  and  $z = 0$  if  $\xi \leq \Pr[z = 0]$ . This satisfies the requirements of the lemma.  $\square$

To online couple  $I$  to  $I'$  we use the coupling guaranteed from the next lemma:

**Lemma 10.10.** Fix  $k \in \mathbb{N}$  and  $p \in (0, 1)$ , and let  $I \sim \text{Uni}(2k, k)$  and  $I' \sim \text{Ber}(2k, p)$ . Then,  $I$  can be online-coupled to  $I'$  such that for any  $t \in [n]$  and any  $I_{t-1}, I'_{t-1} \subseteq [t-1]$ , the following holds:

1.

$$\begin{aligned} & \Pr [t \in I \Delta I' \mid I \cap [t-1] = I_{t-1}, I' \cap [t-1] = I'_{t-1}] \\ &= |\Pr [t \in I \mid I \cap [t-1] = I_{t-1}] - \Pr [t \in I' \mid I' \cap [t-1] = I'_{t-1}]| \\ &= \left| \frac{k - |I_{t-1}|}{2k - (t-1)} - p \right|, \end{aligned} \tag{41}$$

where  $\Delta$  denotes the symmetric set difference  $A \Delta B = (A \setminus B) \cup (B \setminus A)$ .

2. For any  $t \in [n]$  such that  $p \leq (k - |I \cap [t-1]|)/(2k - t + 1)$ , it holds that  $I' \cap \{t\} \subseteq I \cap \{t\}$ .

*Proof of Lemma 10.10.* By Lemma 10.9, conditioned on  $I \cap [t-1] = I_{t-1}, I' \cap [t-1] = I'_{t-1}$ , there is a joint probability distribution between the indicators  $\mathbb{1}_{t \in I}$  and  $\mathbb{1}_{t \in I'}$  such that

$$\begin{aligned} & \Pr [t \in I \Delta I' \mid I \cap [t-1] = I_{t-1}, I' \cap [t-1] = I'_{t-1}] \\ &= \Pr [\mathbb{1}_{t \in I} \neq \mathbb{1}_{t \in I'} \mid I \cap [t-1] = I_{t-1}, I' \cap [t-1] = I'_{t-1}] \\ &= |\Pr [\mathbb{1}_{t \in I} = 1 \mid I \cap [t-1] = I_{t-1}] - \Pr [\mathbb{1}_{t \in I'} = 1 \mid I' \cap [t-1] = I'_{t-1}]| \\ &= |\Pr [t \in I \mid I \cap [t-1] = I_{t-1}] - \Pr [t \in I' \mid I' \cap [t-1] = I'_{t-1}]|. \end{aligned}$$

Define the online sampler to obey this property, namely, while receiving the value of  $\mathbb{1}_{t \in I'}$ , it can sample  $\mathbb{1}_{t \in I}$  from its conditional distribution, conditioned on the obtained value of  $\mathbb{1}_{t \in I'}$  in the above coupling. This proves property 1. Property 2 follows from property 2 in Lemma 10.9.  $\square$

Notice that there is a unique way to define a coupling that satisfies the properties above and we term it the *online monotone coupling*. We proceed with applying the monotone online coupling to prove Lemma 6.3 and Lemma 7.3.

### 10.5.1 Proof of Lemma 6.3

We will in fact prove a more general lemma. We let  $\varphi(I, \bar{x})$  be some real valued function, that we view as some loss corresponding to how the sample  $\bar{x}_I$  represents the complete stream  $\bar{x}$ . We have the following:

**Lemma 10.11.** Let  $\varphi: \{0, 1\}^{2k} \times X^{2k} \rightarrow \mathbb{R}$  be a function that is  $L$ -Lipschitz in each coordinate of  $I$ , namely, for all  $\bar{x} \in X^n$ ,

$$|\varphi(I, \bar{x}) - \varphi(I', \bar{x})| \leq L|I \Delta I'|.$$

Then, for any  $t \geq 0$  and  $\delta \in (0, 1/2)$ ,

$$\begin{aligned} & \sup_{\mathcal{A} \in \text{Adv}_{2k}} \Pr_{I \sim \text{Uni}(2k, k)} [\varphi(I, \bar{x}(\mathcal{A}, I)) \geq t] \\ & \leq \sup_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I' \sim \text{Ber}(2k, 1/2)} \left[ \varphi(I', \bar{x}(\mathcal{A}', I')) \geq t - CL\sqrt{k \log(1/\delta)} \right] + \delta. \end{aligned}$$



This directly implies Lemma 6.3:

*Proof of Lemma 6.3.* Apply Lemma 10.11 with

$$\varphi(I, \bar{x}) = \max_{E \in \mathcal{E}} \left| \frac{|E \cap \bar{x}_I| - |E \cap \bar{x}_{[2k] \setminus I}|}{k} \right|.$$

This function is  $L = 1/k$ -Lipschitz with respect to each coordinate of  $I$ , as the maximum of  $L$ -Lipschitz functions is  $L$ -Lipschitz itself. This suffices to conclude the proof.  $\square$

To prove Lemma 10.11, we start with the following auxiliary property:

**Lemma 10.12.** *Let  $I \sim \text{Uni}(2k, k)$  and  $I' \sim \text{Ber}(2k, 1/2)$  be coupled according to the monotone online coupling. Then, for every  $\delta \in (0, 1/2)$ ,*

$$\Pr \left[ |I' \Delta I| \geq C \sqrt{k \log(1/\delta)} \right] \leq \delta.$$

First, one can bound the expected symmetric difference between  $I$  and  $I'$ :

**Lemma 10.13.** *Let  $I \sim \text{Uni}(2k, k)$  and  $I' \sim \text{Ber}(2k, 1/2)$  be coupled according to the monotone online coupling. Then,*

$$\mathbb{E} [|I' \Delta I|] \leq C \sqrt{k}.$$

*Proof.* Summing (41) over all  $t$  and taking expectation:

$$\mathbb{E} [|I \Delta I'|] = \sum_{t=1}^{2k} \Pr[t \in I \Delta I'] = \sum_{t=1}^{2k} \mathbb{E} \left| \frac{1}{2} - \frac{k - |I \cap [t-1]|}{2k - (t-1)} \right|.$$

By using Jensen's inequality and applying Lemma 5.5 with  $U = \{t, \dots, 2k\}$ ,  $k = k$  and  $n = 2k$ , we have

$$\begin{aligned} \mathbb{E} \left| \frac{1}{2} - \frac{k - |I \cap [t-1]|}{2k - (t-1)} \right| &\leq \sqrt{\mathbb{E} \left[ \left( \frac{1}{2} - \frac{k - |I \cap [t-1]|}{2k - (t-1)} \right)^2 \right]} \\ &= \sqrt{\text{Var} \left[ \frac{|I \cap [t-1]|}{2k - (t-1)} \right]} = \sqrt{\text{Var} \left[ \frac{|I \cap \{t, \dots, 2k\}|}{2k - (t-1)} \right]} \\ &= \frac{1}{2k - (t-1)} \sqrt{\text{Var} [|I \cap \{t, \dots, 2k\}|]} \leq \frac{1}{2k - (t-1)} \sqrt{\frac{(2k - (t-1))k}{2k}} \\ &= \frac{1}{2\sqrt{2k - (t-1)}}. \end{aligned}$$

Summing over  $t = 1, \dots, 2k$ , we have

$$\mathbb{E} [|I \Delta I'|] \leq \sum_{t=1}^{2k} \frac{1}{2\sqrt{2k - (t-1)}} = \sum_{t=1}^{2k} \frac{1}{2\sqrt{t}} \leq \sqrt{2k}.$$

$\square$

The next step is to show that with high probability,  $|\mathbf{I} \cap \mathbf{I}'|$  is close to its expectation. For that, the notion of Martingales is used. We use a standard notation, that is presented in Section A.1. In particular, we use the following commonly used corollary of Azuma's inequality (Lemma A.1):

**Lemma 10.14.** *Let  $F_0 \subseteq F_1 \subseteq \dots \subseteq F_n$  be a filtration such that  $F_0$  is the trivial  $\sigma$ -algebra. Let  $\mathbf{y}$  be a random variable that is  $F_n$  measurable, and assume that  $a_1, \dots, a_n$  are numbers such that  $|\mathbb{E}[\mathbf{y} \mid F_i] - \mathbb{E}[\mathbf{y} \mid F_{i-1}]| \leq a_i$  holds for all  $i$  with probability 1. Then, for all  $t > 0$ ,*

$$\Pr[\mathbf{y} - \mathbb{E}[\mathbf{y}] > t] \leq \exp\left(\frac{-t^2}{2\sum_i a_i^2}\right).$$

Notice that, as described in Section A.1,  $\mathbb{E}[\mathbf{y} \mid F_i]$  is a random variable, and a bound of  $|\mathbb{E}[\mathbf{y} \mid F_i] - \mathbb{E}[\mathbf{y} \mid F_{i-1}]| \leq a_i$  states that the information that is present in  $F_i$  and not in  $F_{i-1}$  does not significantly affect the conditional expectation of  $\mathbf{y}$ .

*Proof of Lemma 10.14.* We define the following Martingale, which is known as *Doob's Martingale*:  $\mathbf{y}_i = \mathbb{E}[\mathbf{y} \mid F_i]$ , for  $i = 0, \dots, n$ . Then,  $\mathbf{y}_0 = \mathbb{E}\mathbf{y}$  and  $\mathbf{y}_n = \mathbf{y}$ , and the proof follows directly from Lemma A.1.  $\square$

We are ready to bound the deviation of  $|\mathbf{I} \Delta \mathbf{I}'|$ :

**Lemma 10.15.** *Let  $\mathbf{I} \sim \text{Uni}(2k, k)$  and  $\mathbf{I}' \sim \text{Ber}(2k, 1/2)$  be coupled according to the monotone online coupling. Then, for every  $\delta \in (0, 1/2)$ ,*

$$\Pr\left[|\mathbf{I}' \Delta \mathbf{I}| - \mathbb{E}[|\mathbf{I}' \Delta \mathbf{I}|] \geq C\sqrt{k \log(1/\delta)}\right] \leq \delta.$$

*Proof.* For any  $t = 0, \dots, n$ , let  $F_t$  denote the  $\sigma$ -field that contains all the information up to (and including) round  $t$ ,  $F_t = \sigma(\mathbf{I} \cap [t], \mathbf{I}' \cap [t])$ . In order to apply Lemma 10.14, it is desirable to bound the differences  $|\mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathcal{F}_t] - \mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathcal{F}_{t-1}]|$  for all  $t \in [2k]$ . Notice that

$$\begin{aligned} & |\mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathcal{F}_t] - \mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathcal{F}_{t-1}]| \\ &= |\mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathbf{I} \cap [t], \mathbf{I}' \cap [t]] - \mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathbf{I} \cap [t-1], \mathbf{I}' \cap [t-1]]|. \end{aligned} \quad (42)$$

We would like to bound (42) for any realization of  $\mathbf{I}$  and  $\mathbf{I}'$ , namely, bounding for any  $S, S' \subseteq [t]$  the quantity

$$|\mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathbf{I} \cap [t] = S, \mathbf{I}' \cap [t] = S'] \quad (43)$$

$$- \mathbb{E}[|\mathbf{I} \Delta \mathbf{I}'| \mid \mathbf{I} \cap [t-1] = S \cap [t-1], \mathbf{I}' \cap [t-1] = S' \cap [t-1]]|. \quad (44)$$

Define four random variables,  $J_t, J'_t, J_{t-1}, J'_{t-1}$  in a joint probability space such that  $(J_t, J'_t)$  is distributed according to the joint distribution of  $(\mathbf{I}, \mathbf{I}')$  conditioned on  $\mathbf{I} \cap [t] = S, \mathbf{I}' \cap [t] = S'$  and  $(J_{t-1}, J'_{t-1})$  is distributed according to the joint distribution of  $(\mathbf{I}, \mathbf{I}')$  conditioned on  $\mathbf{I} \cap [t-1] = S \cap [t-1], \mathbf{I}' \cap [t-1] = S' \cap [t-1]$ . Then, we derive that (43) equals

$$|\mathbb{E}[|J_t \Delta J'_t|] - \mathbb{E}[|J_{t-1} \Delta J'_{t-1}|]| \leq \mathbb{E}[|J_t \Delta J'_t| - |J_{t-1} \Delta J'_{t-1}|] \quad (45)$$

and our goal is to bound the right hand side of (45). The joint distribution is defined by an inductive argument, defining for  $j \geq t$  the intersection of the above four random variables with

$[j]$  given the intersection with  $j - 1$ . Begin with  $j = t$ : here,  $J_t$  and  $J - \mathbf{1}_t$  are fixed to  $S$  and  $S'$ , respectively. Further, the intersections of  $J_{t-1}$  and  $J'_{t-1}$  with  $[t - 1]$  are fixed and equal  $S \cap [t - 1]$ , and their intersections with  $[t]$  are random drawn according to the monotone online coupling. Further, for  $j > t$ , we start by drawing a random variable  $\xi$  uniformly in  $[0, 1]$ , and for any  $\mathbf{U} \in \{J_t, J'_t, J_{t-1}, J'_{t-1}\}$  we set  $j \in \mathbf{U}$  if and only if  $\xi \leq \Pr[j \in \mathbf{U} \mid \mathbf{U} \cap [j - 1]]$ . The above defined joint distribution satisfies the following properties:

- From the proofs of Lemma 10.9 and Lemma 10.10, it follows that  $(J_t, J'_t)$  is distributed according to the joint distribution of  $(I, I')$  conditioned on  $\{I \cap [t] = S, I' \cap [t] = S'\}$  and  $(J_{t-1}, J'_{t-1})$  is distributed according to the joint distribution of  $(I, I')$  conditioned on  $\{I \cap [t - 1] = S \cap [t - 1], I' \cap [t - 1] = S' \cap [t - 1]\}$ .
- For any  $\mathbf{U}, \mathbf{U}' \in \{J_t, J'_t, J_{t-1}, J'_{t-1}\}$  and any  $j > t$ , if  $\Pr[j \in \mathbf{U} \mid \mathbf{U} \cap [j - 1]] \leq \Pr[j \in \mathbf{U}' \mid \mathbf{U}' \cap [j - 1]]$  then  $j \in \mathbf{U}$  implies  $j \in \mathbf{U}'$ .
- Notice that for any for any  $j > t$  and any  $\mathbf{U} \in \{J'_t, J'_{t-1}\}$  it holds that  $\Pr[j \in \mathbf{U} \mid \mathbf{U} \cap [j - 1]] = 1/2$ , hence  $j \in J'_t$  if and only if  $j \in J'_{t-1}$ .
- For any  $j > t$  and any  $\mathbf{U} \in \{J_t, J_{t-1}\}$ , it holds that

$$\Pr[j \in \mathbf{U} \mid \mathbf{U} \cap [j - 1]] = \frac{k - |\mathbf{U} \cap [j - 1]|}{2k - (j - 1)},$$

which is a monotone decreasing function of  $|\mathbf{U} \cap [j - 1]|$ . This implies the following properties:

- For any  $j > t$  such that  $|J_t \cap [j - 1]| = |J_{t-1} \cap [j - 1]|$ , it holds that  $j \in J_t$  if and only of  $j \in J_{t-1}$ .
- For any  $j > t$  such that  $|J_t \cap [j - 1]| \geq |J_{t-1} \cap [j - 1]|$ , it holds that  $\Pr[j \in J_t \mid J_t \cap [j - 1]] \leq \Pr[j \in J_{t-1} \mid J_{t-1} \cap [j - 1]]$ , hence,  $j \in J_t$  implies  $j \in J_{t-1}$ .
- For any  $j > t$  such that  $|J_t \cap [j - 1]| \leq |J_{t-1} \cap [j - 1]|$ , due to a similar argument,  $j \in J_{t-1}$  implies  $j \in J_t$ .
- From the above arguments, for any  $j$ , if  $||J_t \cap [j - 1]| - |J_{t-1} \cap [j - 1]|| = 1$  then either  $j \notin J_t \Delta J_{t-1}$  or  $|J_t \cap [j]| = |J_{t-1} \cap [j]|$ .
- It holds that  $||J_t \cap [t]| - |J_{t-1} \cap [t]|| \leq 1$ . From the above arguments, at the first  $j > t$  such that  $j \in J_t \Delta J_{t-1}$ , it holds that  $|J_t \cap [j]| = |J_{t-1} \cap [j]|$  and from that point onward,  $j \notin J_t \Delta J_{t-1}$ . In particular, there is at most one  $j > t$  such that  $j \in J_t \Delta J_{t-1}$ .
- It follows that there are at most two values of  $j$  such that  $j \in J_t \Delta J_{t-1}$ : one (possibly) for  $j = t$  and one (possibly) for  $j > t$ . Since the only possible  $j$  where  $j \in J'_t \Delta J'_{t-1}$  is  $j = t$ , it follows that

$$||J_t \Delta J'_t| - |J_{t-1} \Delta J'_{t-1}|| \leq 2. \quad (46)$$

From (45) and (46) we derive that

$$|\mathbb{E}[|I \Delta I'| \mid \mathcal{F}_t] - \mathbb{E}[|I \Delta I'| \mid \mathcal{F}_{t-1}]| \leq 2.$$

Applying Lemma 10.14, the proof follows.  $\square$

Lemma 10.13 and Lemma 10.15 together imply Lemma 10.12.

*Proof of Lemma 10.11.* We apply Lemma 10.3 with

$$f(\bar{x}, I) = \mathbb{1}(\varphi(I, \bar{x}) \geq t)$$

and

$$g(\bar{x}, I') = \mathbb{1}\left(\varphi(I', \bar{x}) \geq t - C_0 L \sqrt{k \log(1/\delta)}\right).$$

Here  $C_0 > 0$  is the universal constant guaranteed from Lemma 10.12 such that with probability  $1 - \delta$ ,

$$|I \Delta I'| \leq C_0 \sqrt{k \log(1/\delta)}.$$

Notice that for any  $I$  and  $I'$  such that  $|I \Delta I'| \leq C_0 \sqrt{k \log(1/\delta)}$  and any  $\bar{x} \in X^n$ , it holds that

$$|\varphi(I, \bar{x}) - \varphi(I', \bar{x})| \leq L |I \Delta I'| \leq C_0 L \sqrt{k \log(1/\delta)}.$$

Hence, for any such  $I, I'$  and any  $\bar{x}$  such that  $f(I, \bar{x}) = 1$ , it holds that  $g(I', \bar{x}) = 1$ . Applying Lemma 10.3, we derive that for  $I \sim \text{Uni}(2k, k)$  and  $I' \sim \text{Ber}(2k, 1/2)$ ,

$$\begin{aligned} \max_{\mathcal{A} \in \text{Adv}_{2k}} \Pr[\varphi(I, \bar{x}(\mathcal{A}, I)) \geq t] &= \max_{\mathcal{A} \in \text{Adv}_{2k}} \Pr[f(I, \bar{x}(\mathcal{A}, I)) = 1] \\ &\leq \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr[g(I', \bar{x}(\mathcal{A}', I')) = 1] + \delta \\ &= \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr\left[\varphi(I', \bar{x}(\mathcal{A}', I')) \geq t - C_0 L \sqrt{k \log(1/\delta)}\right] + \delta. \end{aligned}$$

□

### 10.5.2 Proof of Lemma 7.3

We prove the following property of the online monotone coupling:

**Lemma 10.16.** *Let  $I \sim \text{Uni}(2k, k)$  be coupled to  $I' \sim \text{Ber}(2k, 1/8)$  according to the monotone online coupling and fix  $m \leq n$ . Then,*

$$\Pr[I' \subseteq I \cup \{n - m + 1, \dots, n\}] \geq 1 - 2 \exp(-cm).$$

*Proof.* From Lemma 10.10, it suffices to show that with probability  $1 - e^{-cm}$ , for all  $t \leq 2k - m$ ,

$$1/8 \leq \frac{k - |I \cap [t - 1]|}{2k - (t - 1)}. \quad (47)$$

Let  $m_0, m_1, \dots, m_r$ , such that  $m_i = m \cdot 2^i$  and  $k < m_r \leq 2k$ . For any  $i = 0, \dots, r$ , apply Lemma 5.6 with  $n = 2k$ ,  $k = k$ ,  $U = \{n - m_i + 1, \dots, 2k\}$ ,  $\alpha = 1/2$  and  $I = I$ , deriving

$$\Pr\left[|I' \cap (\{n - m_i + 1, \dots, 2k\})| \leq \frac{m_i}{4}\right] \leq \Pr\left[\left|\frac{|I \cap U|}{k} - \frac{|U|}{2k}\right| \geq \frac{m_i}{4k}\right] \leq 2 \exp(-cm_i). \quad (48)$$

Summing the failure probabilities over  $i = 0, \dots, r$ , the sum is dominated by the first summand, and we derive that with probability at least  $1 - 2 \exp(-cm)$ , (48) fails to hold for all  $i = 0, \dots, r$ . Fix a realization  $I$  of  $I$  such that (48) fails for all  $i$ , and we will prove (47), to complete the proof. Fix  $t \leq 2k - m$ , and let  $i_0$  be the maximal  $i$  such that  $t \leq 2k - m_i + 1$ . Then,  $2k - (t - 1) < m_{i_0+1} = 2m_{i_0}$ . Further,

$$k - |I \cap [t - 1]| = |I \cap \{t, \dots, n\}| \geq |I \cap \{n - m_i + 1, \dots, n\}| \geq m_i/4.$$

We derive that

$$\frac{k - |I \cap [t - 1]|}{2k - (t - 1)} \geq \frac{1}{8}$$

as required.  $\square$

*Proof of Lemma 7.3.* Let  $I$  be coupled to  $I'$  according to the monotone online coupling. We wish to apply Lemma 10.3 with

$$f(I, \bar{x}) = \begin{cases} 1 & \exists E \in \mathcal{E}, |\bar{x}_{I \cup J} \cap E| \geq \epsilon \cdot 2k, \bar{x}_I \cap E = \emptyset, \\ 0 & \text{otherwise} \end{cases},$$

and

$$g(I', \bar{x}) = \begin{cases} 1 & \exists E \in \mathcal{E}, |\bar{x} \cap E| \geq \epsilon \cdot 2k \text{ and } |\bar{x}_I \cap E| \leq \epsilon/16 \cdot 2k \\ 0 & \text{otherwise} \end{cases}$$

Applying Lemma 10.16 with  $m = \lfloor \epsilon/16 \cdot 2k \rfloor$ , it holds with probability  $1 - 2 \exp(-cek)$  that

$$I' \subseteq I \cup \{2k - \lfloor \epsilon/16 \cdot 2k \rfloor + 1, n\}. \quad (49)$$

For any values  $I, I'$  such that (49) holds and any  $\bar{x} \in X^n$  such that  $f(I, \bar{x}) = 1$ , it also holds that  $g(I', \bar{x}) = 1$ . Indeed, let  $E \in \mathcal{E}$  be a set such that  $|E \cap \bar{x}| \geq \epsilon \cdot 2k$  and  $\bar{x}_I \cap E = \emptyset$ . From (49),

$$|\bar{x}_{I'} \cap E| \leq |\bar{x}_I \cap E| + |\{2k - \lfloor \epsilon/16 \cdot 2k \rfloor + 1, n\} \cap E| \leq \epsilon/16 \cdot 2k,$$

which implies that  $g(I', \bar{x}) = 1$ . From Lemma 10.3,

$$\begin{aligned} & \max_{\mathcal{A} \in \text{Adv}_{2k}} \Pr \left[ \text{Net}_{\mathcal{A}, I, \epsilon \cdot 2k, 0}^{2k}(\mathcal{E}) = 1 \right] = \max_{\mathcal{A} \in \text{Adv}_{2k}} \Pr_I [f(I, \bar{x}(\mathcal{A}, I)) = 1] \\ & \leq \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr_{I'} [g(I', \bar{x}(\mathcal{A}', I')) = 1] + 2 \exp(-cek) \\ & = \max_{\mathcal{A}' \in \text{Adv}_{2k}} \Pr \left[ \text{Net}_{\mathcal{A}', I', \epsilon \cdot 2k, \epsilon/16 \cdot 2k}^{2k}(\mathcal{E}) = 1 \right] + 2 \exp(-cek). \end{aligned}$$

$\square$

## 11 Continuous $\epsilon$ -Approximation

In the adversarial model we discuss in this paper, the general goal is that in the end of the process (after all elements have been sent by the adversary), the obtained sample would be an  $\epsilon$ -approximation of the entire adversarial sequence. However, in many practical scenarios of

interest, one might want the sample obtained so-far to be an  $\epsilon$ -approximation of the current adversarial sequence *at any point along the sequence* (and not just at the end of the sequence). We call this condition a *continuous  $\epsilon$ -approximation*. Note that such a requirement only makes sense for sampling procedures that allow deletions, like reservoir sampling. (For insertion-only samplers, like Bernoulli and uniform sampling, one cannot hope for the sample to approximate the stream until there is sufficient “critical mass” collected in the sample; this is not an issue with reservoir sampling, which overcomes this by sampling the first elements in the sequence with higher probability, but may also delete them later.)

Obtaining upper bounds for continuous  $\epsilon$ -approximation can be done easily by plugging-in our upper bounds for reservoir sampling to a block-box argument by Ben-Eliezer and Yogev [BEY20, Section 6]. There, it is shown that if one ensures that the current sample approximates the current sequence at  $O(\log n)$  carefully located “checkpoints” along the stream (while setting the error parameter to be  $\delta' = \Theta(\delta/\log n)$ ), then with probability  $1 - \delta$ , the sample is a continuous  $\epsilon$ -approximation for the sequence. That is, we have the following.

**Theorem 11.1** (Adversarial ULLNs – Quantitative Characterization). *Let  $\mathcal{E}$  be a family with Littlestone dimension  $d$ . Then, the sample size  $k(\mathcal{E}, \epsilon, \delta)$ , which suffices to produce a continuous  $\epsilon$ -approximation w.r.t  $\mathcal{E}$  satisfies:*

$$k(\mathcal{E}, \epsilon, \delta) \leq O\left(\frac{d + \log(1/\delta) + \log \log n}{\epsilon^2}\right).$$

*This bound is attained by the reservoir sampler  $\text{Res}(n, k)$ .*

Compared with the standard setting (as summarized in Theorem 2.3), the bound here has an additional  $\log \log n$  term in the numerator.

## 12 Online Learning

In this section, we prove an optimal bound on the regret of online classification. We first provide the formal definitions and then proceed with the formal statement and the proof.

### 12.1 Formal Definitions

Consider the setting of online prediction with binary labels; a learning task in this setting can be described as a guessing game between a learner and an adversary. The game proceeds in rounds  $t = 1, \dots, T$ , each consisting of the following steps:

- The adversary selects  $(x_t, y_t) \in X \times \{0, 1\}$  and reveals  $x_t$  to the learner.
- The learner provides a prediction  $\hat{y}_t \in \{0, 1\}$  of  $y_t$  and announces it to the adversary.
- The adversary announces  $y_t$  to the learner.

Notice that both the learner and the adversary are allowed to use private randomness.

The goal of the learner is to minimize the number of mistakes,  $\sum_t \mathbb{1}(y_t \neq \hat{y}_t)$ . Given a class  $\mathcal{E}$ , a learner  $\mathcal{L}$  and an adversary  $\mathcal{A}$ , the *regret* of the learner w.r.t  $\mathcal{E}$  is defined as the expected difference between the number of mistakes made by the learner and the number of mistakes made by the best  $E \in \mathcal{E}$ :

$$R_T(\mathcal{E}, \mathcal{L}, \mathcal{A}) := \mathbb{E} \left[ \sum_t \mathbb{1}(y_t \neq \hat{y}_t) - \min_{E \in \mathcal{E}} \sum_t \mathbb{1}(y_t \neq \mathbb{1}(x_t \in E)) \right].$$

The *optimal regret* is defined as the value of the the regret achieved by the best sampler against its worst adversary:

$$R_T(\mathcal{E}) = \min_{\mathcal{L}} \max_{\mathcal{A}} R_T(\mathcal{E}, \mathcal{L}, \mathcal{A}).$$

## 12.2 Statement and Proof

We prove the following theorem:

**Theorem 12.1.** *Let  $\mathcal{E}$  denote a class of Littlestone dimension  $d$ . Then, the expected regret  $R_T(\mathcal{E})$  for a  $T$ -round online learner is bounded by*

$$R_T(\mathcal{E}) \leq C\sqrt{dT} ,$$

where  $C > 0$  is a universal constant.

We use a bound by [RST15a] on the regret based on the sequential Rademacher complexity:

**Theorem 12.2** ([RST15a], Theorem 7). *The expected regret satisfies*

$$R_T(\mathcal{E}) \leq 2\text{Rad}_T(\mathcal{E}) .$$

We combine this with the bound on the sequential Rademacher complexity from Lemma 6.4, to complete the proof:

*Proof of Theorem 12.1.* By Theorem 12.2, by definition of the sequential Rademacher complexity and by Lemma 6.4,

$$R_T(\mathcal{E}) \leq 2\text{Rad}_T(\mathcal{E}) = 2\mathbb{E}_{I \sim \text{Ber}(n, 1/2)}[\text{Disc}_{\mathcal{A}, I}(\mathcal{E})] \leq C\sqrt{dT} .$$

This concludes the proof. □

## 13 Lower Bounds

In this section we state and prove our lower bounds. Our first lower bound applies to *any* family  $\mathcal{E}$ , showing that the linear dependence of our upper bounds in the Littlestone dimension is universally tight.

**Theorem 13.1** (A universal lower bound). *Let  $\mathcal{E}$  be a family with Littlestone dimension  $d$ . Then, there exists a (deterministic) adversary such that the following holds. For any algorithm that retains at most  $k \leq d$  items (without deletions), the adversary presents  $d$  items  $x_1, \dots, x_d$  such that*

$$(\exists E \in \mathcal{E}) : \bar{s} \cap E = \emptyset \quad \text{and} \quad \frac{|\bar{x} \cap E|}{|\bar{x}|} = 1 - \frac{k}{d},$$

with probability 1 over the algorithm's randomness, where  $\bar{x}$  denotes the adversarial stream and  $\bar{s}$  is the sample. In particular, any subset of the sample of  $k$  items retained by the algorithm does not form an  $\epsilon$ -approximation with respect to  $x_1, \dots, x_n$  unless  $\epsilon \geq 1 - \frac{k}{d}$ .

Our second result in this section shows the existence of families  $\mathcal{E}$  of Littlestone dimension  $d$  in which all  $\epsilon$ -approximations are of size  $\Omega(d/\epsilon^2)$ , so long as  $d = \Omega(\log(1/\epsilon))$ . Interestingly, the requirement that  $d$  is large enough is necessary: classical results in discrepancy theory [MWW93, Mat95] imply that when  $d = o(\log 1/\epsilon)$ , smaller  $\epsilon$ -approximations exist.

**Theorem 13.2** ( $\epsilon$ -approximation: quadratic lower bound). *Let  $d \in \mathbb{N}$  and  $\epsilon > 0$  where  $d \geq C \log(1/\epsilon)$  for a large absolute constant  $C > 0$ . Then, there exists a family  $\mathcal{E}$  with Littlestone dimension at most  $d$  and a subset  $\{x_1, \dots, x_n\} \subset X$  for which no subset of size less than  $c \cdot \frac{\text{Ldim}(\mathcal{E})}{\epsilon^2}$  is an  $\epsilon$ -approximation, where  $c > 0$  is a small absolute constant.*

We also prove similar results for  $\epsilon$ -nets (without the requirement that  $d$  is large enough).

**Theorem 13.3** ( $\epsilon$ -net: a super linear lower bound). *Let  $d \in \mathbb{N}$  and  $\epsilon > 0$ . Then, there exists a family  $\mathcal{E}$  with Littlestone dimension  $\leq d$  and a subset  $\{x_1, \dots, x_n\} \subset X$  for which no subset of length less than  $c \cdot \frac{\text{Ldim}(\mathcal{E}) \log(1/\epsilon)}{\epsilon}$  is an  $\epsilon$ -net, where  $c > 0$  is a small absolute constant.*

### 13.1 Proofs

*Proof of Theorem 13.1.* The proof generalizes the construction from [BEY20], which provided a lower bound for the family of one-dimensional thresholds.<sup>9</sup> Let  $\mathcal{T}$  be a tree of depth  $d$  which is shattered by  $\mathcal{E}$ . The tree  $\mathcal{T}$  can be thought of as a strategy for the adversary as follows:

1. Set  $\mathcal{T}_1 = \mathcal{T}$  and  $i = 1$ .
2. For  $i = 1, \dots, d$ 
  - (i) Pick  $x_i$  to be the item labelling the root of  $\mathcal{T}_i$  and present it to the algorithm.
  - (ii) If  $x_i$  was retained by the algorithm then continue to the next iteration with  $\mathcal{T}_{i+1}$  being the left subtree of  $\mathcal{T}_i$  (corresponding to the sets in  $\mathcal{E}_{\not\ni x_i}$ ).
  - (iii) Else, continue to the next iteration with  $\mathcal{T}_{i+1}$  being the right subtree of  $\mathcal{T}_i$  (corresponding to the sets in  $\mathcal{E}_{\ni x_i}$ ).

Thus, the adversary picks the elements  $x_1, \dots, x_d$  according to a path on the tree such that whenever  $x_i$  is retained by the algorithm then a left turn is taken and whenever  $x_i$  is not retained by the algorithm then a right turn is taken. Thus, since the tree is shattered, there exists a set  $E \in \mathcal{E}$  such that

$$E \cap \{x_1, \dots, x_n\} = \{x_i : x_i \text{ was not sampled by the algorithm}\}.$$

In particular,  $\bar{s} \cap E = \emptyset$ , and if the algorithm samples  $m \leq d$  points then  $\frac{|\bar{x} \cap E|}{|\bar{x}|} = 1 - \frac{m}{d}$ , as required.  $\square$

*Proof of Theorem 13.2.* The proof follows from standard probabilistic arguments, and shows that most families in a certain setting have bounded Littlestone dimension yet do not admit a small  $\epsilon$ -approximation. Suppose that  $d \geq \log(1/\epsilon)$  and let  $n = d/6\epsilon^2$ . Let  $F$  be a family of  $2^d \cdot d/\epsilon^2$  subsets of  $[n]$  of size  $n/2$ , picked uniformly at random among all such families, and note that (by definition and since  $d \geq \log(1/\epsilon)$ ) the Littlestone dimension of  $F$  is at most  $\log |F| = O(d)$ .

<sup>9</sup>We note that the proof from [BEY20] would give a lower bound of  $\Omega(\log d)$  for any family of Littlestone dimension  $d$  (as compared to the  $\Omega(d)$  lower bound we prove here); this follows since, roughly speaking, any such family “contains” a class of thresholds of dimension logarithmic in  $d$ .



We now claim that with high probability, there is no  $\epsilon$ -approximation of size less than  $n/2$  for  $F$ . Indeed, fix any subset  $S$  of size  $m \leq n/2$ . By a simple counting argument, the number of sets  $A$  of size  $n/2$  for which  $|d_A(S) - d_A([n])| \geq \epsilon$  is at least

$$\binom{m}{(\frac{1}{2} - \epsilon)m} \binom{n-m}{\frac{n-m}{2} - \epsilon m} \geq \binom{m}{\frac{m}{2}} \binom{n-m}{\frac{n-m}{2}} \cdot (1 - 2\epsilon)^{2\epsilon m} \geq \frac{2^n}{2n} \cdot e^{-3\epsilon^2 n} = \frac{2^n}{2n} \cdot e^{-d/2},$$

where the second inequality holds for  $\epsilon < 1/10$ .

Plugging in the right hand side above, and noting the negative correlation between the events at hand, the probability that  $F$  does not contain any such  $A$  with  $|d_A(S) - d_A([n])| \geq \epsilon$  is bounded by

$$\left(1 - \frac{2^n \cdot e^{-d/2}}{2n \binom{n}{n/2}}\right)^{|F|} \leq e^{-2^d \cdot e^{-d/2}} \leq e^{-1.2^d \cdot d / \epsilon^2}.$$

Taking a union bound over all (less than  $2^n = 2^{d/6\epsilon^2}$ ) possible subsets  $S \subseteq [n]$  of size at most  $[n]/2$ , it follows that with high probability (as a function of  $d$ ), no  $\epsilon$ -approximation exists.  $\square$

*Proof of Theorem 13.3.* The proof extends a simple probabilistic construction in the projective plane, suggested by Alon, Kalai, Matoušek, and Meshulam [AKMM02].

Consider the projective plane of order  $p$ , where we pick  $p = C/\epsilon$  for a suitable constant  $C$ . Recall that this projective plane has  $p^2 + p + 1$  points and lines, where each line consists of exactly  $p + 1$  points, and every two points are contained in exactly one line. For each line  $L$ , pick uniformly at random (and independently from choices for other lines) a subset  $H_L$  containing exactly half the elements of  $L$ ; we call such a subset a *half line*. Consider the family consisting of all such half lines  $H_L$ . As was shown in [AKMM02], with high probability every  $\epsilon$ -net for this family has size  $\Omega(p \log p)$ , whilst the VC dimension is at most 2. We claim that the same bound also holds for the Littlestone dimension.

**Claim 13.4.** *The Littlestone dimension of the family consisting of all half lines as above is at most 2.*

*Proof.* Suppose to the contrary a depth-3 tree exists as in the definition of the Littlestone dimension, and consider the elements  $x, y, z$  appearing in the internal nodes of the all-1 branch in this tree. By definition, all three elements must belong to some half line  $H_L$  from the family. However, since any two lines  $L_1 \neq L_2$  in the projective plane intersect in exactly one point, we have  $|H_{L_1} \cap H_{L_2}| \leq 1$ . It follows that there does not exist  $L' \neq L$  where  $x, y \in L'$ , and thus, no half line corresponds to the  $(1, 1, 0)$ -branch of the tree, a contradiction.  $\square$

The proof that no small  $\epsilon$ -net exists is a straightforward probabilistic proof similar in spirit to that of Theorem 13.2, and is given in detail in [AKMM02]. The proof bounds from above the probability of any fixed set of size (say)  $0.1p \log p$  to intersect all half lines, and then takes a union bound over all such sets.

Next, we show how to generalize the above to get a lower bound with linear dependence in  $d$ . Let  $p$  be as above, consider  $d$  copies of the projective plane of order  $p$  and let  $\mathcal{C}_1, \dots, \mathcal{C}_d$  be collections of half lines generated as above, one in each plane. Now let  $\mathcal{C}$  be the collection of all unions of exactly  $d/2$  half lines coming from different planes, namely,  $\mathcal{C}$  contains all sets of the form

$$H = H_{L_{j_1}}^{i_1} \cup H_{L_{j_2}}^{i_2} \cup \dots \cup H_{L_{j_{d/2}}}^{i_{d/2}},$$

where  $i_1 < i_2 < \dots < i_{d/2} \in [d]$ ,  $H_{L_{j_t}}^{i_t}$  is a half line from the  $i_t$  copy corresponding to the line  $L_{j_t}$  in that copy of the plane.

Consider the family  $\mathcal{C}$  with the underlying universe with  $d(p^2 + p + 1)$  points, containing all points from all  $d$  planes.

**Claim 13.5.** *The Littlestone dimension of  $\mathcal{C}$  is at most  $d$ .*

*Proof.* The proof is a straightforward extension of the proof of Claim 13.4. Suppose to the contrary that the Littlestone dimension is  $t > d$ . Let  $T$  be a labeled tree of depth  $t$  as in the definition of Littlestone dimension and consider its all-1 branch. This branch corresponds to some set  $H = H_{L_{j_1}}^{i_1} \cup H_{L_{j_2}}^{i_2} \cup \dots \cup H_{L_{j_{d/2}}}^{i_{d/2}}$ . In particular, all elements labeling nodes along the branch are contained in  $H$ .

By the pigeonhole principle, there exist three elements  $x, y, z$  along the branch (in this order) contained in the same half line  $H_L$  from one of the plane copies. We claim that there is no set in  $\mathcal{C}$  that corresponds to any branch which is all-1 up until (and not including)  $z$ , and takes the value 0 at  $z$ . Indeed, such a set  $H$ , if exists, will contain  $x, y$  but not  $z$ . However, this is a contradiction as in Claim 13.4: any set  $H$  that contains  $x, y$  must also contain all elements in the half line  $H_L$  containing them both, and thus  $z \in H$ .  $\square$

It remains to prove that there is no  $\epsilon$ -net of size  $o(d\epsilon^{-1} \log \epsilon^{-1})$ . But this follows easily from the  $\Omega(\epsilon^{-1} \log \epsilon^{-1})$  lower bound for each of the planes separately: there exists some absolute constant  $C > 0$  so that for each of the planes at hand, no  $\epsilon$ -net of size  $C\epsilon^{-1} \log \epsilon^{-1}$  exists. Consider now any set  $S$  of less than  $Cd\epsilon^{-1} \log \epsilon^{-1}/2$  points in our universe, the union of all planes; since each point belongs to exactly one plane, there exist  $d/2$  planes with less than  $C\epsilon^{-1} \log \epsilon^{-1}$  points. Let  $i_1 < i_2 < \dots < i_{d/2}$  denote their indices. It follows that there exists some set  $H = H_{L_{j_1}}^{i_1} \cup H_{L_{j_2}}^{i_2} \cup \dots \cup H_{L_{j_{d/2}}}^{i_{d/2}} \in \mathcal{C}$  not intersecting  $S$ . This completes the proof.  $\square$

## References

- [AKMM02] Noga Alon, Gil Kalai, Jiří Matoušek, and Roy Meshulam. Transversal numbers for hypergraphs arising in geometry. *Advances in Applied Mathematics*, 29(1):79 – 101, 2002.
- [AMS99] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences*, 58(1):137–147, 1999.
- [BEHW89] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36:929–965, 1989.
- [BEY20] Omri Ben-Eliezer and Eylon Yogev. The adversarial robustness of sampling. In *Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS)*, page 49–62, 2020.

- [BJWY20] Omri Ben-Eliezer, Rajesh Jayaram, David P. Woodruff, and Eylon Yogev. A framework for adversarially robust streaming algorithms. In *Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS)*, pages 63–80, 2020.
- [Bla54] David Blackwell. Controlled random walks. In *Proceedings of the International Congress of Mathematicians*, volume 3, pages 336–338, 1954.
- [Bla56] David Blackwell. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics*, 6(1):1–8, 1956.
- [BM15] Rémi Bardenet and Odalric-Ambrym Maillard. Concentration inequalities for sampling without replacement. *Bernoulli*, 21(3):1361–1385, 2015.
- [Bou04] Olivier Bousquet. Introduction to Statistical Learning Theory. In *Advanced lectures on machine learning*, volume 3176, pages 169–207. Springer, 2004.
- [BPS09] Shai Ben-David, Dávid Pál, and Shai Shalev-Shwartz. Agnostic online learning. In *Proceedings of the 22nd Conference on Learning Theory (COLT)*, 2009.
- [CGP<sup>+</sup>18] Timothy Chu, Yu Gao, Richard Peng, Sushant Sachdeva, Saurabh Sawlani, and Junxing Wang. Graph sparsification, spectral sketches, and faster resistance computation, via short cycle decompositions. In *Proceedings of the 59th IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 361–372, 2018.
- [Cha05] Sourav Chatterjee. Concentration inequalities with exchangeable pairs (Ph.D. thesis). *arXiv preprint math/0507526*, 2005.
- [CN20] Yeshwanth Cherapanamjeri and Jelani Nelson. On adaptive distance estimation. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, 2020.
- [DFH<sup>+</sup>15] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349(6248):636–638, 2015.
- [dlPn99] Victor H. de la Peña. A general class of exponential inequalities for martingales and ratios. *The Annals of Probability*, 27(1):537–564, 1999.
- [Dud73] Richard M. Dudley. Sample functions of the Gaussian process. *The Annals of Probability*, 1(1):66–103, 1973.
- [Dud78] Richard M. Dudley. Central limit theorems for empirical measures. *The Annals of Probability*, 6(6):899–929, 1978.
- [Dud84] Richard M. Dudley. A course on empirical processes. In P. L. Hennequin, editor, *École d’Été de Probabilités de Saint-Flour XII - 1982*, pages 1–142. Springer Berlin Heidelberg, 1984.
- [Dud87] Richard M. Dudley. Universal Donsker classes and metric entropy. *The Annals of Probability*, 15(4):1306–1326, 1987.

- [Fre75] David A. Freedman. On tail probabilities for martingales. *the Annals of Probability*, pages 100–118, 1975.
- [GHR<sup>+</sup>12] Anna C. Gilbert, Brett Hemenway, Atri Rudra, Martin J. Strauss, and Mary Wootters. Recovering simple signals. In *Information Theory and Applications Workshop (ITA)*, pages 382–391, 2012.
- [GHS<sup>+</sup>12] Anna C. Gilbert, Brett Hemenway, Martin J. Strauss, David P. Woodruff, and Mary Wootters. Reusable low-error compressive sampling schemes through privacy. In *IEEE Statistical Signal Processing Workshop (SSP)*, pages 536–539, 2012.
- [Han57] James Hannan. Approximation to bayes risk in repeated play. *Contributions to the Theory of Games*, 3:97–139, 1957.
- [Hau92] David Haussler. Decision theoretic generalizations of the pac model for neural net and other learning applications. *Information and Computation*, 100(1):78 – 150, 1992.
- [HKM<sup>+</sup>20] Avinatan Hassidim, Haim Kaplan, Yishay Mansour, Yossi Matias, and Uri Stemmer. Adversarially robust streaming algorithms via differential privacy. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, 2020.
- [HRS20] Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis of online and differentially private learning. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, 2020.
- [HW13] Moritz Hardt and David P. Woodruff. How robust are linear sketches to adaptive inputs? In *Proceedings of the 45th ACM Symposium on Theory of Computing (STOC)*, pages 121–130, 2013.
- [Lit88] Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2(4):285–318, 1988.
- [LP17] David A. Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Society, 2017.
- [LW94] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Information and computation*, 108(2):212–261, 1994.
- [Mat95] Jiří Matoušek. Tight upper bounds for the discrepancy of half-spaces. *Discrete & Computational Geometry*, 13:593–601, 1995.
- [Mat09] Jiří Matoušek. *Geometric Discrepancy: An Illustrated Guide*. Springer-Verlag Berlin Heidelberg, 2009.
- [MNS11] Ilya Mironov, Moni Naor, and Gil Segev. Sketching in adversarial environments. *SIAM Journal on Computing*, 40(6):1845–1870, 2011.
- [MWW93] Jiří Matoušek, Emo Welzl, and Lorenz Wernisch. Discrepancy and approximations for bounded VC-dimension. *Combinatorica*, 13(4):455–466, 1993.

- [NY15] Moni Naor and Eylon Yogev. Bloom filters in adversarial environments. In *Advances in Cryptology - CRYPTO 2015 - 35th Annual Cryptology Conference*, pages 565–584, 2015.
- [Rob51] Herbert Robbins. Asymptotically subminimax solutions of compound statistical decision problems. In *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*. The Regents of the University of California, 1951.
- [RS14] Alexander Rakhlin and Karthik Sridharan. Statistical learning and sequential prediction. *Book Draft*, 2014.
- [RS15] Alexander Rakhlin and Karthik Sridharan. On Martingale Extensions of Vapnik–Chervonenkis Theory with Applications to Online Learning. In *Measures of Complexity*, pages 197–215. Springer, 2015.
- [RST10] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Random averages, combinatorial parameters, and learnability. In *Advances in Neural Information Processing Systems*, pages 1984–1992, 2010.
- [RST15a] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning via sequential complexities. *J. Mach. Learn. Res.*, 16(1):155–186, 2015.
- [RST15b] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Sequential complexities and uniform martingale laws of large numbers. *Probability Theory and Related Fields*, 161(1-2):111–153, 2015.
- [Sau72] Norbert Sauer. On the density of families of sets. *Journal of Combinatorial Theory, Series A*, 13(1):145–147, 1972.
- [SBM<sup>+</sup>18] Chawin Sitawarin, Arjun Nitin Bhagoji, Arsalan Mosenia, Mung Chiang, and Prateek Mittal. DARTS: deceiving autonomous cars with toxic signs. *CoRR*, abs/1802.06430, 2018.
- [SSBD14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning*. Cambridge university press, 2014.
- [SSSSS10] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *The Journal of Machine Learning Research*, 11:2635–2670, 2010.
- [Tal94] Michel Talagrand. Sharper bounds for Gaussian and empirical processes. *The Annals of Probability*, 22(1):28–76, 1994.
- [Vap98] Vladimir N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, Inc., 1998.
- [VC71] Vladimir N. Vapnik and Alexey Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280, 1971.
- [VC74] Vladimir N. Vapnik and Alexey Y. Chervonenkis. *Theory of Pattern Recognition*. Nauka, Moscow, 1974.

- [Vit85] Jeffrey S. Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software*, 11(1):37–57, 1985.
- [WFRS18] Blake E. Woodworth, Vitaly Feldman, Saharon Rosset, and Nati Srebro. The everlasting database: Statistical validity at a fair price. In *Advances in Neural Information Processing Systems 31 (NeurIPS)*, pages 6532–6541, 2018.
- [WZ20] David P. Woodruff and Samson Zhou. Tight bounds for adversarially robust streams and sliding windows via difference estimators. *arXiv preprint arXiv:2011.07471*, 2020.

## A Probabilistic Material

### A.1 Filtration and Martingales

In this section we give a brief probability background to Martingales, considering only finite probability spaces. Recall that a probability space consists of a sample space  $\Omega$ , a  $\sigma$ -field  $F \subseteq \{0, 1\}^\Omega$  that contains all measurable events and a probability measure  $\mu$  over  $\Omega$ . With finite probability spaces, it is possible for  $F$  to contain all subsets of  $\Omega$ , however, smaller sets can be considered as well. For instance, if  $\mathbf{y}_1, \dots, \mathbf{y}_n$  are random variables over the finite space  $Y$ , then  $\sigma(\mathbf{y}_1)$ , the  $\sigma$ -field generated by  $\mathbf{y}_1$ , contains all the events that depend only on  $\mathbf{y}_1$ . Formally, we have  $\Omega = Y^n$  and  $\sigma(\mathbf{y}_1) = \{\{\mathbf{y}_1 \in U\} : U \subseteq Y\}$ , where  $\{Y \in U\} = \{(y_1, \dots, y_n) : y_1 \in U\}$  is the event that  $\mathbf{y}_1 \in U$ . Similarly, we can have sigma fields generated by multiple random variables, for instance,  $\sigma(y_1, y_3, y_4)$ , that contains all the events that depend only on these three random variables. It is in fact also possible to consider the  $\sigma$ -algebra generated by zero random variables  $\sigma(\{\}) = \{0, \Omega\}$  which is called the *trivial  $\sigma$ -algebra*.

Conditioning on more random variables results in a larger  $\sigma$ -algebra, namely, if  $i \leq j$  then  $\sigma(\mathbf{y}_1, \dots, \mathbf{y}_i) \subseteq \sigma(\mathbf{y}_1, \dots, \mathbf{y}_j)$ . Intuitively, larger  $\sigma$ -algebras contain more information. We say that a  $\sigma$ -field  $F$  is  $\mathbf{y}$ -measurable if  $\sigma(\mathbf{y}) \subseteq F$ , which intuitively holds whenever  $F$  contains all the information on  $\mathbf{y}$ . Further, a *filtration* is a collection of nested  $\sigma$ -algebras  $F_0 \subseteq F_1 \cdots \subseteq F_n$ .

One can define *conditional expectation* with respect to a  $\sigma$ -algebra. In our application, each  $\sigma$ -algebra will be generated by a collection of random variables, and it holds that

$$\mathbb{E}[\cdot \mid \sigma(\mathbf{y}_1, \dots, \mathbf{y}_k)] = \mathbb{E}[\cdot \mid \mathbf{y}_1, \dots, \mathbf{y}_k]. \quad (50)$$

Notice that the quantity in (50) is a function of  $\mathbf{y}_1, \dots, \mathbf{y}_k$ , hence it is also a random variable. Additionally, if  $F$  is the trivial  $\sigma$ -algebra then

$$\mathbb{E}[\cdot \mid F] = \mathbb{E}[\cdot \mid \sigma(\{\})] = \mathbb{E}[\cdot].$$

And if  $\mathbf{y}$  is  $F$ -measurable, then  $\mathbb{E}[\mathbf{y} \mid F] = \mathbf{y}$ .

A collection of random variables  $z_0, \dots, z_n$  defines a *Martingale adapted to the filtration*  $F_0 \subseteq \cdots \subseteq F_n$  if  $z_i$  is  $F_i$ -measurable and if for any  $i < j$ ,  $\mathbb{E}[z_j \mid F_i] = z_i$ . The simplest case is when  $F_i = \sigma(z_1, \dots, z_i)$ , and there, the martingale condition translates to  $\mathbb{E}[z_j \mid z_1, \dots, z_i] = z_i$ . However, in the general case  $F_i$  can have additional information on other random variables.

Remarkably, Martingales obey high probability bounds. Perhaps the most well known bound is Azuma’s inequality, which is an adaptation of Chernoff’s bound for Martingales:

**Lemma A.1.** Let  $\mathbf{y}_0, \dots, \mathbf{y}_n$  be a Martingale adapted to the filtration  $F_0, \dots, F_n$ . Let  $a_1, \dots, a_n \geq 0$  be numbers such that almost surely,  $|\mathbf{y}_i - \mathbf{y}_{i-1}| \leq a_i$ . Then, for any  $t \geq 0$ ,

$$\Pr[\mathbf{y}_n - \mathbf{y}_0 > t] \leq \exp\left(\frac{-t^2}{2\sum_i a_i^2}\right).$$

## A.2 Sampling Without Replacement

Further, we have the following version of Chernoff without replacement:

**Lemma A.2** ([BM15]). Let  $a_1, \dots, a_N \in \mathbb{R}$  and let  $I$  denote a uniformly random subset of  $[N]$  of size  $n \in \mathbb{N}$ . Let  $R = \max_i a_i - \min_i a_i$ . Then, for any  $t > 0$ ,

$$\Pr\left[\frac{1}{n}\sum_{i \in I} a_i - \frac{1}{N}\sum_{i=1}^N a_i > t\right] \leq \exp\left(\frac{-2nt^2}{R^2}\right).$$

Another without-replacement lemma:

**Lemma A.3** ([Cha05], Proposition 3.10). Let  $\{a_{ij}\}_{i,j \in [n]}$  be a collection of numbers from  $[0, 1]$ . Let  $Y = \sum_{i=1}^n a_{i\pi(i)}$  where  $\pi$  is drawn from the uniform distribution over the set of permutations of  $\{1, \dots, n\}$ . Then for any  $t \geq 0$ ,

$$\Pr[|Y - \mathbb{E}Y| \geq t] \leq 2\exp(-t^2/(4\mathbb{E}Y + 2t)).$$

*Proof of Lemma 5.6.* First item follows directly from Lemma A.2. The second item follows from Lemma A.3 as described below. Define  $m = |U|$ . Let  $\pi: [n] \rightarrow [n]$  be a uniformly random permutation and let  $I = \{i: \pi(i) \leq k\}$ . Define  $\{a_{ij}\}_{i,j \in [n]}$  by  $a_{ij} = 1$  if  $i \in U$  and  $j \leq k$ . Notice that for all  $i \in U$ ,  $a_{i\pi(i)} = 1$  if  $i \in I$  and for all  $i \notin U$ ,  $a_{i\pi(i)} = 0$ . Hence,

$$Y := \sum_{i=1}^n a_{i\pi(i)}$$

equals  $|I \cap U|$  and  $\mathbb{E}Y = km/n$ . From Lemma A.3 we derive that for any  $t \geq 0$ ,

$$\Pr[|Y - \mathbb{E}Y| \geq t] \leq \exp\left(-\frac{t^2}{4\mathbb{E}Y + 2t}\right).$$

Substitute  $t = \alpha\mathbb{E}Y$  and we get that

$$\begin{aligned} \Pr\left[\left|\frac{Y}{\mathbb{E}Y} - 1\right| \geq \alpha\right] &= \Pr[|Y - \mathbb{E}Y| \geq t] \leq \exp\left(-\frac{t^2}{4\mathbb{E}Y + 2t}\right) = \exp\left(-\frac{\alpha^2\mathbb{E}[Y]^2}{4\mathbb{E}Y + 2\alpha\mathbb{E}Y}\right) \\ &\leq \exp\left(-\frac{\alpha^2\mathbb{E}[Y]^2}{6\mathbb{E}Y}\right) = \exp\left(-\frac{\alpha^2 km}{6n}\right). \end{aligned}$$

□

*Proof of Lemma 5.5.* Denote  $|U| = m$ . For any  $i \in U$ , let  $z_i$  denote the indicator of whether  $i \in I$ , and notice that  $|U \cap I| = \sum_{i \in U} z_i$ . Therefore, we have

$$\mathbb{E}[|U \cap I|] = \mathbb{E}\left[\sum_{i \in U} z_i\right] = \sum_{i \in U} \mathbb{E}[z_i] = \sum_{i \in U} k/n = mk/n.$$

Next,

$$\mathbb{E}[|U \cap \mathbf{I}|^2] = \sum_{i,j \in U} \mathbb{E}z_i z_j = \sum_{i \in U} \mathbb{E}z_i^2 + 2 \sum_{i,j \in U: i < j} \mathbb{E}z_i z_j.$$

For the first term, since  $z_i$  is an indicator, we have

$$\sum_{i \in U} \mathbb{E}z_i^2 = \sum_{i \in U} \mathbb{E}z_i = mk/n.$$

For the second term, fix  $i < j$ , and we have

$$\begin{aligned} \mathbb{E}z_i z_j &= \Pr[z_i = 1, z_j = 1] = \Pr[\{i, j\} \subseteq \mathbf{I}] = \frac{1}{\binom{n}{k}} |\{I \subseteq [n]: |I| = k, i, j \in I\}| \\ &= \frac{\binom{n-2}{k-2}}{\binom{n}{k}} = \frac{(n-2)!k!(n-k)!}{n!(k-2)!(n-k)!} = \frac{k(k-1)}{n(n-1)} \leq \frac{k^2}{n^2}. \end{aligned}$$

We derive that

$$\mathbb{E}[|U \cap \mathbf{I}|^2] \leq \frac{mk}{n} + \frac{m^2 k^2}{n^2}.$$

Hence,

$$\text{Var}(|U \cap \mathbf{I}|) = \mathbb{E}[|U \cap \mathbf{I}|^2] - \mathbb{E}[|U \cap \mathbf{I}|] \leq \frac{mk}{n} + \frac{m^2 k^2}{n^2} - \frac{m^2 k^2}{n^2} = \frac{mk}{n},$$

as required. □