

Abstract of “Classic and Modern Challenges in Statistical Estimation” by Chun Hin Jasper Lee, Ph.D., Brown University, May 2021.

The rise of the Internet has generated and has enabled the collection of massive amounts of data. However, this modern ubiquity and abundance of data is worth little unless we can efficiently extract from it useful information and insights. In this thesis, we focus on the question of data efficiency, that is, optimizing the amount of data needed to accomplish a statistical task to some given accuracy and confidence guarantees. Drawing on tools from across probability, statistics and theoretical computer science, we propose optimally data-efficient algorithms for two basic estimation problems. The problems and their solutions, defined in distinct data-access models, highlight and give techniques to overcome important data-collection and data-utilization challenges faced by algorithm designers in the modern era.

The first problem is a classic and fundamental problem in statistics: what is the best way to estimate the mean of a real-valued distribution from independent samples? Under the minimal and essentially necessary assumption that the distribution has finite (but unknown) variance, we settle the problem by presenting an estimator with convergence tight to within a $1 + o(1)$ multiplicative factor. This contrasts previous works that are either only tight up to multiplicative constants, or require strong additional assumptions such as knowledge of the variance, or a bounded 4th moment (kurtosis) assumption. Our estimator construction and analysis gives a generalizable framework, tightly analyzing a sum of dependent random variables by viewing the sum implicitly as a 2-parameter ψ -estimator, and constructing bounds using mathematical programming techniques.

The second problem is a coin-flipping problem motivated by crowdsourcing applications. Given a mixture between two populations of coins, “positive” coins that each have—unknown and potentially different—bias $\geq \frac{1}{2} + \Delta$ and “negative” coins with bias $\leq \frac{1}{2} - \Delta$, we consider the task of estimating the fraction of positive coins to within a given accuracy through drawing coins from the mixture and flipping them. We give an adaptive algorithm and a fully-adaptive lower bound with matching sample complexity, simultaneously tight in all relevant problem parameters, up to a multiplicative constant. The fine-grained adaptive flavor of both our algorithm and lower bound contrasts with much previous work in distributional testing and learning.

Classic and Modern Challenges in Statistical Estimation

by

Chun Hin Jasper Lee

MA, University of Cambridge, 2018

MPhil, University of Cambridge, 2015

BA, University of Cambridge, 2014

A dissertation submitted in partial fulfillment of the
requirements for the Degree of Doctor of Philosophy
in the Department of Computer Science at Brown University

Providence, Rhode Island

May 2021

© Copyright 2021 by Chun Hin Jasper Lee

This dissertation by Chun Hin Jasper Lee is accepted in its present form by the Department of Computer Science as satisfying the dissertation requirement for the degree of Doctor of Philosophy.

Date _____
_____ Paul Valiant, Director

Recommended to the Graduate Council

Date _____
_____ Ronitt Rubinfeld, Reader
Massachusetts Institute of Technology

Date _____
_____ Eliezer Upfal, Reader

Approved by the Graduate Council

Date _____
_____ Andrew G. Campbell
Dean of the Graduate School

Curriculum Vitae

Jasper Lee was born and raised in Hong Kong. He spent 4 years at Churchill College, University of Cambridge for both his undergraduate and masters studies, earning a BA in 2014, an MPhil in 2015 and an MA in 2018. In Fall 2015, he moved to Brown to start his doctoral studies.

During his time at Brown, Jasper was the graduate teaching assistant for CS157: Design and Analysis of Algorithms for 3 years. He also created and lectured CS1951-W: Sublinear Algorithms for Big Data in his final year.

Jasper co-authored the following works during his doctoral studies:

- Optimal Sub-Gaussian Mean Estimation in \mathbb{R} . Jasper C.H. Lee and Paul Valiant, *in submission*
- Uncertainty about Uncertainty: Optimal Adaptive Algorithms for Estimating Mixtures of Unknown Coins. Jasper C.H. Lee and Paul Valiant, *SODA 2021*
- Finding an Approximate Mode of a Kernel Density Estimate. Jasper C.H. Lee, Jerry Li, Christopher Musco, Jeff M. Phillips and Wai Ming Tai, *in submission*
- Fast Algorithms for Computing Interim Allocations in Single-Parameter Environments. Amy Greenwald, Jasper Lee, Takehiro Oyakawa, *PRIMA 2018*
- Augmenting Stream Constraint Programming with Eventuality Conditions. Jasper C.H. Lee, Jimmy H.M. Lee and Zhuowei Zhong, *CP 2018*
- Optimizing Star-Convex Functions. Jasper C.H. Lee and Paul Valiant, *FOCS 2016*

Acknowledgements

First and foremost, I thank my advisor Paul Valiant for guiding my formative years in academia. I arrived at Brown with very limited exposure to theoretical computer science—I did not even know what a Chernoff bound was. Over the past years, Paul has patiently taught me everything starting from basic mathematics to advanced techniques (not to mention our casual conversations about physics and beyond), while extracting useful research ideas from my often chaotic thought process as I continue to grow as a researcher. He has furthermore opened my eyes to what it means to truly understand a mathematical idea in depth. Beyond technical training, Paul has also instilled in me the courage to tackle highly challenging problems. Without this courage, we would not have been able to resolve the fundamental problem of mean estimation for real-valued distributions with finite variance, which now forms a major part of this thesis, and will always remain as one of the works I am most proud of. I am grateful for Paul’s mentoring, and I hope our future collaborations will continue to be fruitful.

On a non-academic note, I also want to thank Paul for introducing me to the sport of rock climbing. It is the first form of physical exercise that I have found myself to actually enjoy!

*

*

*

Next, I would like to thank other mentors I have had throughout my graduate school career.

My thesis committee members, Eli Upfal and Ronitt Rubinfeld, have been supportive and have given me much career advice, especially early on during my PhD. I thank Ronitt

also for hosting me for a summer visit at MIT, where I had the opportunity to talk to more researchers.

My collaborator, Christopher Musco, has been very kind in offering me various advice and discussing lots of technical ideas with me. Chris has exposed me to many ideas in streaming/sketching algorithms, numerical linear algebra and sparse recovery, which I would not have learnt about otherwise. I am also grateful to Chris for hosting me for a (virtual, due to the COVID-19 pandemic) visit to his research group at NYU this semester.

In addition, I have benefited from taking courses from and talking to many faculty members here at Brown, in particular, Amy Greenwald, Anna Lysyanskaya, Philip Klein and Shriram Krishnamurthi. Thank you for your insights and advice throughout the many points in my PhD career.

*

*

*

I would not have been able to make it through the past 6 years without my family and friends.

To my dad Jimmy, my mum Scarlet and my brother Julian, thank you for your love my whole life. I wouldn't be where I am today without you.

To Ben, thanks for your awesome friendship. Thanks for the numerous illuminating conversations, and for all your support, especially when I was briefly hospitalized.

To Arun, Clayton and Uthsav, thanks for the PhD student camaraderie and for understanding the struggle.

To Gabe, thanks for being a great housemate during the early years of our PhD life.

To the "Wheelin' and Dealin'" crew, Ethan, Jared, Oliver, Sacha, Sean, Tony and Zach, thanks for all the fun, and for making pandemic life more bearable.

To friends and officemates at the department, Alexandra, Apoorvaa, Archita, Esha, Evgenios, Ghous, Lucy, Marilyn, Megumi, Nedi and Thomas, thanks for giving me a sense of community.

To YCH, thanks for the 10+ year friendship, and for being my first algorithms teacher.

To my former students and UTAs in CS157 and CS1951-W, thanks for being engaging and curious, which made teaching you a joy.

Lastly, I would like to thank the departmental administrative staff for their help throughout the years. A special thanks to Lauren and Lori for dealing with my many questions and requests.

*

*

*

The work in this thesis was partially supported by NSF award IIS-1562657 and a Kannellakis fellowship.

Contents

List of Tables	xi
List of Figures	xii
1 Introduction	1
1.1 Mean Estimation in \mathbb{R}	2
1.2 Estimating the Fraction of “Positive” Coins	3
2 Background	6
2.1 Probability	6
2.2 Concentration Inequalities via the Chernoff Bound	7
2.3 Sample Complexity Lower Bounds	8
2.4 Duality in Mathematical Programming	10
3 Optimal Sub-Gaussian Mean Estimation in \mathbb{R}	13
3.1 Overview	13
3.1.1 The Model and Main Result	13
3.1.2 Our Approach	15
3.1.3 Motivation: 3 rd -order corrections of the empirical mean	16
3.1.4 Key Contributions in Our Construction and Analysis	19
3.2 Related Work	20
3.3 Our Mean Estimator	22
3.3.1 Meaning of the Estimator	22
3.3.2 Structural Properties of the Estimator	23

3.3.3	Representing a Special Case of Estimator 1 as a ψ -Estimator	26
3.3.4	The relation to the Catoni-Giulini estimator	27
3.4	Analyzing our estimator	29
3.5	Proofs of Lemmas 3.11 and 3.12	34
3.5.1	Mathematical Programming and Duality Analysis	35
3.5.2	Proof of Lemma 3.11	37
3.5.3	Proof of Lemma 3.12	41
4	Uncertainty about Uncertainty: Optimal Adaptive Algorithms for Estimating Mixtures of Unknown Coins	42
4.1	Overview	42
4.1.1	Our Approaches and Results	44
4.1.2	Related Work	52
4.2	The Triangular Walk Algorithm	55
4.2.1	Implementing Algorithm 3	61
4.3	The Main Algorithm	65
4.4	Characterizing Single-Coin Algorithms	72
4.5	Fully-Adaptive Lower Bounds	74
4.5.1	Reduction to Single-Coin Adaptive Algorithms	76
4.5.2	Upper Bounding the Squared Hellinger Distance for Single-Coin Adaptive Algorithms	86
4.6	Proof of Proposition 4.22	90
4.6.1	"Central" Region	91
4.6.2	"High Discrepancy" Region	95
4.6.3	The Last Row	101
4.7	Experimental Results	104
4.8	Algorithm for Known Conditional Distributions of Coins	106
4.8.1	A Quadratic+Linear Programming Approach	107
4.8.2	Optimality of such linear estimators	110
4.9	Non-Adaptive Lower Bound	113
4.10	Remaining Proofs/Calculations of Results	114

List of Tables

4.1	Sample Complexity Upper and Lower Bounds	50
-----	--	----

List of Figures

4.1	Experimental Results	105
4.2	A QP formulation for computing the output coefficients in terms of the stopping rule	108
4.3	An LP formulation for finding the best stopping rule given an expected sample complexity	109
4.4	An LP formulation for finding the best stopping rule independent of the expected sample complexity for a single coin	110

Chapter 1

Introduction

Over the past decades, the Internet has driven the generation and collection of a massive amount of data, for example, from everyday user activities (e.g. social media, online shopping), to users providing data under incentives (e.g. crowdsourcing). Buried in all this data is *information* and *insight* hidden in the statistical noise, for example, various statistics of social networks and internet traffic. Modern data science studies the efficient extraction of such information from data at a large scale. There are two important and distinct notions of efficiency: *Data efficiency* concerns the amount of data required to complete a statistical task to high accuracy and confidence, or equivalently, given a certain amount of data, optimize the accuracy and confidence guarantees we can derive from the data. *Computational efficiency*, on the other hand, concerns the computational resources (e.g. time and space) we need to compute the statistical estimators and tests of interest.

In this thesis, we focus on two problems—one classic and foundational problem and one motivated by modern crowdsourcing settings—where achieving optimal data efficiency is the main algorithmic challenge, and computational efficiency follows naturally from the algorithms we design. These two problems highlight the challenges that algorithm designers face in effectively collecting and utilizing data at a large scale.

The first problem is the fundamental problem of estimating the mean of a real-valued distribution, under the minimal (and essentially necessary, see later) assumption that the underlying distribution has a finite but unknown variance. The second problem is an

estimation problem inspired by practical data management applications under the crowd-sourcing paradigm. The problem is theoretically modelled by a population of coins with unknown biases, where the only way to gain information about the population is through picking out and flipping the coins. The goal is to estimate the fraction of coins in the population with bias $> \frac{1}{2}$. For both of these problems, the main challenge is in optimizing data efficiency: most of the construction and analysis focuses on the statistical properties of the algorithms we propose, where the algorithms can be computed in time linear in the number of samples. In future work, we plan to study variants of the problems where both data efficiency and computational efficiency are central challenges, for example, in high dimensional mean estimation.

1.1 Mean Estimation in \mathbb{R}

In Chapter 3, we revisit and settle a fundamental problems in statistics: estimating the mean of a real-valued distribution from independently and identically distributed (i.i.d.) samples, in the high probability regime, under the minimal (and essentially necessary) assumption that the distribution has finite but unknown variance. The sample mean (also known as the empirical mean), despite its ubiquitous use, is known to have sample complexity that is exponentially sub-optimal in the failure probability δ [15] for heavy-tailed distributions. We propose an estimator with convergence tight not only in the big-O sense (namely, up to multiplicative factors), but tight up to only a $1 + o(1)$ factor. In contrast to prior works, our estimator does not require prior knowledge of the variance [15], and works across the entire gamut of distributions with finite variance, including those without any higher moments [15, 23]. Our estimator is furthermore computable in time linear in the number of samples.

Parameterized by the sample size n , the failure probability δ , and the variance σ^2 , our estimator has additive accuracy within $\sigma \cdot (1 + o(1)) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$ (where the logarithm is the natural logarithm), which is optimal up to the $1 + o(1)$ term. This asymptotically matches the convergence of the sample mean for the Gaussian distribution with the same variance, as well as the information-theoretic lower bound for estimating the mean of Gaussians.

The main challenge in this work lies in analyzing our proposed estimator: we need to derive tight concentration properties, however, our estimator is a sum of *dependent* terms, making it difficult to directly and tightly control its moment generating function for the purposes of deriving a Chernoff bound. To circumvent this obstacle, we view our estimator as a 2-parameter ψ -estimator, from which we can derive proxies of the estimator that are sums of independent terms, and reduce the main theorem to proving Chernoff bounds on these proxies instead. A further obstacle lies in proving tight Chernoff bounds for our estimator over the entire space of distributions with finite variance. We resolve this second obstacle by viewing the Chernoff bound analysis as a convex-concave programming problem (over an infinite dimensional space of distributions), and use convex-concave programming duality and linear programming duality (see Section 2.4) to make the analysis finite-dimensional and tractable.

The above estimator construction and analysis approach gives a framework generalizable to other problems.

This work is currently in submission.

1.2 Estimating the Fraction of “Positive” Coins

In Chapter 4, we consider a natural statistical estimation task, motivated by a practical crowdsourcing application, with an intriguing adaptive flavor. In the problem setting, there is a universe of coins of two types: “positive” coins each have a (potentially different) probability of heads that lies in the interval $[\frac{1}{2} + \Delta, 1]$, while “negative” coins lie in the interval $[0, \frac{1}{2} - \Delta]$, where $\Delta \in (0, \frac{1}{2}]$ parameterizes the “quality” of the coins. Our only access to the coins is by choosing a coin and then flipping it, without access to the true biases of the coins. An algorithm in this setting may employ arbitrary adaptivity—for example, flipping three different coins in sequence and then flipping the first coin 5 more times if and only if the results of the first 3 flips were heads, tails, heads. The challenge is to estimate the *fraction* ρ of coins that are of positive type, to within a given additive error ϵ , failing with probability at most δ , using as few coin flips (samples) as possible.

This model arose from a collaboration with colleagues in data science and database systems, about harnessing paid crowdsourced workers to estimate the “quality” of a database. Our model is a direct theoretical analog of the following practical problem, where sample complexity linearly translates into the amount of money that must be paid to workers, and thus even multiplicative factors crucially affect the usefulness of an algorithm. Given a set of data and a predicate on the data, the task is to estimate what fraction of the data satisfies the predicate—for example, estimating the proportion of records in a large database that contain erroneous data. After automated tools have labeled whatever portion of the data they are capable of dealing with, the remaining data must be processed via *crowdsourcing*, an emerging setting that potentially offers sophisticated capabilities but at the cost of unreliability. Namely, for each data item, one may ask many human users/workers online whether they think the item satisfies the predicate, with the caveat that the answers returned could be noisy. Each coin in the theoretical model corresponds to an item in the database, modelling the noisy response we get when we ask a random human worker online to evaluate the predicate on the item.

We exhibit an adaptive algorithm and a fully-adaptive lower bound which have sample complexities that match each other simultaneously in *all* 4 of the problem parameters (fraction ρ of positive coins, estimation additive error ϵ , coin quality parameter Δ and failure probability δ), up to multiplicative constants.

A key feature that makes this estimation problem distinct from many others studied in the literature is the richness of adaptivity available to the algorithm. Achieving a tight lower bound in this setting requires considering and bounding all possible uses of adaptivity available to an algorithm; and achieving an optimal algorithm requires choosing the appropriate adaptive information flow between different parts of the algorithm. Much of the previous work in the area of statistical estimation is focused on non-adaptive algorithms and lower bounds; however see [10], and in particular, Sections 4.1 and 4.2 of that work, for a survey of several distribution testing models that allow for adaptivity. In our setting there are two distinct kinds of adaptivity that an algorithm can leverage: 1) single-coin adaptivity, deciding how many times a particular coin should be flipped—a per-coin stopping rule—in terms of the results of its previous flips, and 2) cross-coin adaptivity,

deciding which coin to flip next in terms of the results of previous flips across *all* coins. Our final optimal algorithm (Section 4.3) leverages both kinds of adaptivity. In our tight lower bound analysis (Section 4.5), we overcome the technical obstacles presented by the richness of adaptivity by giving a reduction (Section 4.5.1) from fully-adaptive algorithms that leverage both kinds of adaptivity to single-coin adaptive algorithms that process each coin independently, valid for our specific lower bound instance.

The main *algorithmic* challenge in this problem is what we call “uncertainty about uncertainty”: we make no assumptions about the quality of the coins beyond the existence of a gap 2Δ between biases of the coins of different types (centered at $\frac{1}{2}$). Our algorithm must return estimates with small bias, and be sample-efficient at the same time, regardless of the bias of the coins, whether they are all deterministic, or all maximally noisy as allowed by the Δ parameter, or some quality in between. While intuitively the hardest settings to distinguish information theoretically involve coins with biases as close to each other as possible (and indeed our lower bound relies on mixtures of only $\frac{1}{2} \pm \Delta$ coins), settings with biases near but not equal to $\frac{1}{2} \pm \Delta$ introduce “uncertainty about uncertainty” challenges. We overcome this challenge with an estimator designed using 1D random walk theory.

In addition, to illustrate the difficulty of the “uncertainty about uncertainty” paradigm, we consider a relaxation of the problem where we have some knowledge of the coin biases. Assuming (perhaps unrealistically) that we know 1) the conditional distribution of biases of positive coins, and 2) the same for negative coins, and 3) an initial estimate of the mixture parameter ρ between the two distributions, then we show that it is easy—using mathematical programming techniques in Section 4.8.1—to construct an estimation algorithm with sample complexity that is optimal *by construction* up to a multiplicative constant (see Section 4.8.2).

This work appeared in SODA 2021 [37]. We thank Tim Kraska and Yeounoh Chung for bringing these problems to our attention in the data analytics setting, and for contributing to the simulation results in this work. We also thank an anonymous reviewer for asking about the δ dependence in lower bounds, which led to the current tight results.

Chapter 2

Background

In this chapter, we give basic definitions and facts relevant to the results in this thesis.

2.1 Probability

All probability distributions considered in this thesis are defined over \mathbb{R} .

We will also use the following standard distributions in our modelling and analyses:

Bernoulli coins The Bernoulli distribution with *bias/parameter* p is defined such that it returns a 1 with probability p and returns a 0 otherwise. The Bernoulli distribution is also called a Bernoulli coin, and a sample from the distribution is also called a coin flip. A result of 1 corresponds to a heads flip and a 0 corresponds to a tails flip.

Binomials The Binomial distribution $\text{Bin}(n, p)$ with parameters n and p , is defined as the sum of n independent Bernoulli coins with probability p . That is, it is the number of heads observed by flipping a bias- p coin independently for n times. In the thesis, we will use the notation $\text{Bin}(n, k, p)$ to denote the probability that a Binomial random variable with parameters n and p returns the value k , namely, $\binom{n}{k} p^k (1-p)^{n-k}$.

Gaussians The standard Gaussian distribution, with mean 0 and variance 1, is defined with the density function $\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$. The Gaussian distribution with mean μ and variance σ^2 is defined as adding μ to σ times the standard Gaussian.

For a discrete distribution, we define its *support* to be the elements on which it has non-zero probability mass.

We will also be concerned with the *existence* of moments (e.g. mean and variance) of a distribution. Many distributions, such as the examples above, have well-defined/finite mean and variance, as well as higher moments. However, distributions on \mathbb{R} do not necessarily have well-defined moments. For example, the Cauchy distribution, whose probability density function decays at a rate of $\Theta(1/x^2)$ away from 0, does not have a well-defined mean. On the other hand, the distribution with density decaying as $1/x^3$ does have a well-defined mean, but not a well-defined variance. In our work on optimal mean estimation, we will make the essentially necessary assumption that the underlying distribution has a finite but unknown variance.

2.2 Concentration Inequalities via the Chernoff Bound

One of the most powerful techniques for analyzing randomness is to quantify *concentration behavior*. That is, distributions that are the sum of many independent distributions (e.g. Gaussians and Binomials) will have most of their probability mass near their expectation. Quantitatively, we would want to upper bound the following *tail* probabilities: $\mathbb{P}(X \geq \mathbb{E}(X) + x)$ and $\mathbb{P}(X \leq \mathbb{E}(X) - x)$ for $x > 0$.

The standard approach to showing these concentration inequalities is via the Chernoff bound.

Fact 2.1. *For any random variable X over \mathbb{R} and every x , we have the following inequalities:*

$$\mathbb{P}(X \geq x) \leq \inf_{t>0} \frac{\mathbb{E}[e^{tX}]}{e^{tx}}$$

$$\mathbb{P}(X \leq x) \leq \inf_{t>0} \frac{\mathbb{E}[e^{-tX}]}{e^{-tx}}$$

The quantity $M_X(t) = \mathbb{E}[e^{tX}]$ is also known as the moment generating function of the random variable X .

In the thesis, we will both derive custom Chernoff bounds from Fact 2.1, as well as use standard corollaries such as the following bound for Binomial distributions:

Fact 2.2. Suppose $X \leftarrow \text{Bin}(n, p)$, and denote $\mu = np$. Then for any $\kappa > 0$,

$$\mathbb{P}(X \geq (1 + \kappa)\mu) \leq \left(\frac{e^\kappa}{(1 + \kappa)^{1+\kappa}} \right)^\mu$$

$$\mathbb{P}(X \leq (1 - \kappa)\mu) \leq \left(\frac{e^{-\kappa}}{(1 - \kappa)^{1-\kappa}} \right)^\mu$$

2.3 Sample Complexity Lower Bounds

Proving complexity lower bounds is a main challenge in theoretical computer science. While many lower bounds for, for example, time complexity remain elusive open problems, there have been various techniques developed to show sample complexity lower bounds. In this section, we describe one of the most common, yet effective, ways to show a sample complexity lower bound: showing an *indistinguishability* result, which is related to Le Cam's two-point method in the statistics literature.

A typical indistinguishability argument for a problem consists of two parts:

1. A reduction, that if we have an algorithm/estimator/test solving the problem at hand, then the same algorithm can be used to distinguish between two given scenarios. This reduction is typically straightforward.
2. A sample complexity lower bound for the simpler problem of distinguishing the two scenarios, that no algorithm taking too few samples can distinguish the scenarios with high probability. This means that any algorithm succeeding at the original problem must take at least a certain number of samples.

For concreteness, consider the example problem of estimating the mean of an unknown Bernoulli coin to within additive error ϵ , succeeding with probability at least $1 - \delta$. We want to show a worst-case sample complexity lower bound of $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ many coin flips¹. The argument, as above, consists of two parts:

1. Reduction: If we had a mean estimator with additive accuracy ϵ , failing with probability at most δ , then we can use it to distinguish the bias $\frac{1}{2} - \epsilon$ coin from the $\frac{1}{2} + \epsilon$ coin, also failing with probability at most δ .

¹Although a more fine-grained view shows that we in fact need only $O(\frac{p(1-p)}{\epsilon^2} \log \frac{1}{\delta})$ many coin flips to estimate the mean of a Bernoulli coin with bias p . The above lower bound takes the worst case over p .

2. Lower bound: We need to show that any estimator needs at least $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ coin flips to distinguish between those two coins.

Before introducing techniques for showing the lower bound, we note that this lower bound is tight up to multiplicative constants. For example, the sample mean requires only $O(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$ many coin flips, per Hoeffding's inequality (another standard corollary of the Chernoff bound approach to tail bounds).

To prove the indistinguishability result, we want to upper bound the *total variation distance* between n independent flips of the $\frac{1}{2} - \epsilon$ coin and the $\frac{1}{2} + \epsilon$ coin, due to the following fact.

Definition 2.3 (Total variation distance). Given two discrete probability distributions \mathbf{p}, \mathbf{q} over a finite set S , the total variation distance $d_{\text{TV}}(\mathbf{p}, \mathbf{q})$ between \mathbf{p} and \mathbf{q} is defined as:

$$\begin{aligned} d_{\text{TV}}(\mathbf{p}, \mathbf{q}) &= \frac{1}{2} \sum_{i \in S} |p_i - q_i| = \frac{1}{2} \|\mathbf{p} - \mathbf{q}\|_1 \\ &= \sup_{A \subseteq S} \mathbf{p}(A) - \mathbf{q}(A) \end{aligned}$$

The definition naturally generalizes to arbitrary distributions.

Fact 2.4. Consider a game, where an adversary picks arbitrarily either distribution \mathbf{p} or distribution \mathbf{q} , and we want an algorithm which, on input n independent samples from the chosen distribution, decide whether the samples came from \mathbf{p} or \mathbf{q} , succeeding with probability at least $1 - \delta$. Then, there is no algorithm \mathcal{A} such that:

$$\mathbb{P}(\mathcal{A} \text{ returns } \mathbf{p} \mid \text{adversary picked } \mathbf{p}) - \mathbb{P}(\mathcal{A} \text{ returns } \mathbf{p} \mid \text{adversary picked } \mathbf{q}) > d_{\text{TV}}(\mathbf{p}^{\otimes n}, \mathbf{q}^{\otimes n})$$

where $\mathbf{p}^{\otimes n}$ denotes the n -fold product distribution of \mathbf{p} . In particular, this implies that there is no algorithm \mathcal{A} such that both of the following hold:

- $\mathbb{P}(\mathcal{A} \text{ returns } \mathbf{p} \mid \text{adversary picked } \mathbf{p}) > \frac{1}{2} + \frac{1}{2} d_{\text{TV}}(\mathbf{p}^{\otimes n}, \mathbf{q}^{\otimes n})$
- $\mathbb{P}(\mathcal{A} \text{ returns } \mathbf{q} \mid \text{adversary picked } \mathbf{q}) > \frac{1}{2} + \frac{1}{2} d_{\text{TV}}(\mathbf{p}^{\otimes n}, \mathbf{q}^{\otimes n})$

So if $d_{\text{TV}}(\mathbf{p}^{\otimes n}, \mathbf{q}^{\otimes n}) < 1 - 2\delta$, there is no algorithm that will succeed in distinguishing between two distributions with probability $\geq 1 - \delta$ using only n samples.

The total variation distance between n -fold product distributions can be difficult to directly and tightly upper bound. In particular, using the union bound typically leads to weak or even trivial bounds. In these cases, we often bound other statistical distances and divergences first, and in turn use the proxy to bound the total variation distance.

The *KL-divergence*, an important information-theoretic notion, enjoys two crucial properties that are useful for sample complexity lower bound purposes.

Definition 2.5 (KL-divergence). Given two discrete probability distributions \mathbf{p}, \mathbf{q} over a finite set S , the KL-divergence $D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q})$ between \mathbf{p} and \mathbf{q} is defined as:

$$D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = \sum_{i \in S} p_i \log \frac{p_i}{q_i}$$

The definition naturally generalizes to arbitrary distributions.

Fact 2.6. *The KL-divergence is additive for product distributions:*

$$D_{\text{KL}}(\mathbf{p}_1 \otimes \mathbf{p}_2 \parallel \mathbf{q}_1 \otimes \mathbf{q}_2) = D_{\text{KL}}(\mathbf{p}_1 \parallel \mathbf{q}_1) + D_{\text{KL}}(\mathbf{p}_2 \parallel \mathbf{q}_2)$$

Furthermore, the KL-divergence satisfies the high probability Pinsker inequality:

$$d_{\text{TV}}(\mathbf{p}, \mathbf{q}) \leq 1 - \frac{1}{2} e^{-D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q})}$$

With the above fact, we can now complete the example of showing a tight sample complexity lower bound for estimating the mean of a coin. Let \mathbf{p} be the $\frac{1}{2} - \epsilon$ bias coin, and \mathbf{q} be the $\frac{1}{2} + \epsilon$ bias coin. Then,

$$D_{\text{KL}}(\mathbf{p}^{\otimes n} \parallel \mathbf{q}^{\otimes n}) = nD_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = \Theta(n\epsilon^2)$$

where the last equality is straightforward calculation.

Thus, for $d_{\text{TV}}(\mathbf{p}^{\otimes n}, \mathbf{q}^{\otimes n})$ to be at least $1 - 2\delta$, n must be at least $\Omega(\frac{1}{\epsilon^2} \log \frac{1}{\delta})$, for otherwise the high probability Pinsker inequality would yield a contradiction.

2.4 Duality in Mathematical Programming

Duality theory in mathematical programming plays a central role in our work on optimal real-valued mean estimation. In particular, we leverage linear programming duality and

convex-concave programming (min-max optimization) duality. We state these duality principles in this section.

We start with linear programming duality. Consider the following maximization linear program P in standard form, over a finite dimensional vector of real-valued variables \mathbf{x} :

$$\begin{aligned} & \text{maximize} && \mathbf{c} \cdot \mathbf{x} \\ & \text{subject to} && A\mathbf{x} \leq \mathbf{b} \\ & \text{where} && \mathbf{x} \geq 0 \end{aligned}$$

where the matrix A and vectors \mathbf{b} and \mathbf{c} are parameters in the program.

The dual linear program D of P is defined as the following minimization problem over the vector \mathbf{y} , which has length equal to the number of constraints in \mathbf{x} , namely the number of rows in A :

$$\begin{aligned} & \text{minimize} && \mathbf{b} \cdot \mathbf{y} \\ & \text{subject to} && A^T \mathbf{y} \geq \mathbf{c} \\ & \text{where} && \mathbf{y} \geq 0 \end{aligned}$$

Similarly, the dual of a minimization problem is defined by the construction of P from D , and it is trivial to see that duality is an *involution*, that is, the dual of a dual is the original *primal* problem.

By construction, the optimum of D (the minimization problem), $\text{opt}(D)$ is always at least the optimum of P (the maximization problem), $\text{opt}(P)$, if both optima exist. This is known as the *weak duality* of linear programs, which can be proved as follows. Consider an arbitrary pair of feasible solutions \mathbf{x} and \mathbf{y} for P and D respectively, then

$$\mathbf{c} \cdot \mathbf{x} = \mathbf{x}^T \mathbf{c} \leq \mathbf{x}^T A \mathbf{y} \leq \mathbf{b}^T \mathbf{y} = \mathbf{b} \cdot \mathbf{y}$$

where the first inequality is by the dual constraints and that $\mathbf{x} \geq 0$, and the second inequality is by the primal constraints and that $\mathbf{y} \geq 0$. Typically, for other kinds of optimization problems, the definition of the dual problem is going to imply weak duality also by construction, as in the linear programming case.

Much less trivial is the fact that linear programs enjoy *strong duality*. That is, if the primal and dual programs both have feasible solutions, then their optima are equal to each other, captured by the following fact.

Fact 2.7 (Strong duality for linear programs). *Given a primal linear program P that is feasible, let D be its LP dual. Then D is feasible and $\text{opt}(P) = \text{opt}(D)$, where $\text{opt}()$ maps a linear program to its optimal objective value.*

Therefore, taking the dual of a linear program does not quantitatively change the optimization problem. However, qualitatively, the dual is an alternative formulation of the problem, and can give additional insights. It can also enable and simplify analysis, by reasoning about a different form of the same optimization problem, which we use crucially in our mean estimation work.

The other duality principle we require for our mean estimation result is on max-min optimization problems. Consider the following general max-min optimization problem:

$$\max_{x \in S_x} \min_{y \in S_y} f(x, y)$$

for domains S_x, S_y and objective $f : S_x \times S_y \rightarrow \mathbb{R}$.

The duality notion we consider is simple: swapping the order of maximization and minimization. Thus, weak duality is just the standard, straightforward max-min inequality:

$$\max_{x \in S_x} \min_{y \in S_y} f(x, y) \leq \min_{y \in S_y} \max_{x \in S_x} f(x, y)$$

It is easy to check that strong duality, namely when the weak duality inequality is in fact an equality, does *not* hold for arbitrary objectives f and domains S_x and S_y . We are therefore interested in conditions under which strong duality holds. A strong duality result in this context is known as a *minimax theorem*. Perhaps the most well-known minimax theorem is the original: von Neumann's minimax theorem. In this thesis, we will require a generalization of von Neumann's theorem, which follows from Sion's minimax theorem [59].

Fact 2.8. *Suppose S_x and S_y are convex sets, at least one of which is compact. Also suppose the objective function $f : S_x \times S_y \rightarrow \mathbb{R}$ is continuous, convex in the first argument (i.e. for all $y \in S_y$, $f(\cdot, y)$ is convex) and concave in the second argument. Then*

$$\sup_{x \in S_x} \inf_{y \in S_y} f(x, y) = \inf_{y \in S_y} \sup_{x \in S_x} f(x, y)$$

For our purposes, it suffices to assume that the maxima and minima exist, that the sup and inf in the above can be replaced by max and min.

Chapter 3

Optimal Sub-Gaussian Mean Estimation in \mathbb{R}

3.1 Overview

We revisit one of the most fundamental problems in statistics: estimating the mean of a real-valued distribution, using as few independent samples from it as possible. Our proposed estimator has convergence that is optimal not only in a big-O sense (i.e. “up to multiplicative constants”), but tight to a $1 + o(1)$ factor, under the minimal (and essentially necessary, see below) assumption of the finiteness of the variance. Previous works, discussed further in Section 3.2, are either only big-O tight [34, 50, 2], or require additional strong assumptions such as the variance being known to the estimator [15] or assumptions that allow for accurate estimates of the variance, such as the kurtosis (fourth moment) being finite [15, 23].

3.1.1 The Model and Main Result

Given a set of i.i.d. samples from a real-valued distribution, the goal is to return, with extremely high probability, an accurate estimate of the distribution’s mean. Specifically, given a sample set X of size n consisting of independent draws from a real-value distribution D , an (ϵ, δ) -estimator of the mean is a function $\hat{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}$ such that, except with failure

probability $\leq \delta$, the estimate $\hat{\mu}(X)$ is within ϵ of the true mean $\mu(D)$. Namely,

$$\mathbb{P}(|\hat{\mu}(X) - \mu(D)| \leq \epsilon) \geq 1 - \delta \quad (3.1)$$

The goal is to find the optimal tradeoff between the sample size n , and the error parameters ϵ and δ , for the distribution D . Fixing any two of the three parameters and minimizing the third yields essentially equivalent reformulations of the problem: we can fix ϵ, δ and minimize the *sample complexity* n ; we can fix δ, n and minimize *error* ϵ ; or we can fix ϵ, n and minimize the *failure probability* δ (maximizing the *robustness* $1 - \delta$).

Perhaps the most standard and well-behaved setting for mean estimation is when the distribution D is a Gaussian. The sample mean (the empirical mean) is a provably optimal estimator in our sense when D is Gaussian: for any $\epsilon, \delta > 0$, the sample mean $\mu(X)$ is an (ϵ, δ) -estimator when given a sample set of size $n = (2 + o(1)) \frac{\sigma^2(D) \cdot \log \frac{1}{\delta}}{\epsilon^2}$ (all logarithms will be base e); and there is *no* (ϵ, δ) -estimator for Gaussians if the constant 2 in the previous expression for the sample size is changed to any smaller number.

The main result of this paper is an estimator that performs as well, on *any* distribution with finite variance, as the sample mean does on a Gaussian, without knowledge of the distribution or its variance:

Theorem 3.1. *Estimator 1, given $\delta, n > 0$, defines a function $\hat{\mu}$ such that with probability at least $1 - \delta$, given a sample set X of size n , yields an estimate with error*

$$|\hat{\mu}(X) - \mu(D)| \leq \sigma(D) \cdot (1 + o(1)) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$$

Here, the $o(1)$ term tends to 0 as $(\frac{\log \frac{1}{\delta}}{n}, \delta) \rightarrow (0, 0)$. Furthermore, as evidenced by the Gaussian case, there is no estimator which, under the same settings, produces an error that improves on our guarantees by more than a $1 + o(1)$ multiplicative factor.

We have parameterized the above theorem in terms of fixing the sample size n and the robustness parameter δ and asking for the minimum error ϵ ; however, because of the simple functional form of the bounds of Theorem 3.1, we can equivalently rephrase it as saying that, for any ϵ, δ , (a reparameterized) Estimator 1 is an (ϵ, δ) estimator using $(2 + o(1)) \frac{\sigma(D)^2}{\epsilon^2} \log \frac{1}{\delta}$ samples; or for any n, ϵ , Estimator 1 gives an estimate that is $\delta = \exp(-\frac{n\epsilon^2}{(2+o(1)) \cdot \sigma^2(D)})$ -robust.

For each of these formulations, the performance is optimal up to the $1 + o(1)$ factor, as evidenced by the well-known Gaussian case, as explained above.

We make the following observations regarding the main (minimal) assumption in the theorem, namely the finiteness of the variance of the unknown distribution. First, imposing further assumptions about the finiteness of higher moments will not yield any improvements to the result, since matching lower bounds are provided by Gaussians, for which all moments are finite. Second, as shown by Devroye et al. [23], relaxing the finite variance assumption by only assuming, say, the finiteness of the $(1 + \beta)^{\text{th}}$ moment for some $\beta < 1$ will yield strictly worse sample complexity. In particular, the sample complexity will have an ϵ -dependence that is $\omega(1/\epsilon^2)$. Thus, our result shows that mean estimation can be performed at a sub-Gaussian rate, with the optimal multiplicative constant of 2 in the sample complexity, if and only if the variance of the underlying distribution is finite.

We also contrast with previous works that attain optimal sub-Gaussian convergence but make additional assumptions such as the finiteness of the kurtosis (4th moment) [15, 23]. The gap in assumptions between those works and this work is not only theoretical, but also of practical consequence: Pareto distributions (power law distributions) are known to be good models of certain real-world phenomena, and for a shape parameter (i.e. exponent or Pareto index α) in the range $(2, 4]$, the variance exists, but not the kurtosis.

3.1.2 Our Approach

We briefly describe the main features of our estimator, as a setting for what follows, and to distinguish it from prior work. At the highest level: in order to return a δ -robust estimate of the mean, our estimator “throws out the $\frac{1}{3} \log \frac{1}{\delta}$ most extreme points in the sample set”, and returns the mean of what remains. More specifically, outliers are thrown out in a *weighted* manner, where we throw out a *fraction* of each data point, with the fraction proportional to the square of its distance from a median-of-means initial guess for the mean, where the fraction is capped at 1, and the proportionality constant is chosen so that the total weight thrown out equals exactly $\frac{1}{3} \log \frac{1}{\delta}$. See Estimator 1 for full details, but we stress here that the estimator is simple to implement—it may be computed in linear time—and therefore applicable in practice.

The above description is rather different from the typical M-estimator/ ψ -estimator approach of Catoni [15] and other works in this area. However, as we see in Section 3.3, our estimator can be reinterpreted as a 2-parameter ψ -estimator, and the proof of our main result will crucially rely on this reformulation. In Section 3.3.4 we examine the similarities and differences between our estimator and the Catoni-Giulini estimator [16], which is a particular instantiation of the approach in [15].

3.1.3 Motivation: 3rd-order corrections of the empirical mean

Perhaps the most non-obvious part of our estimator is throwing out exactly $\frac{1}{3} \log \frac{1}{\delta}$ many samples. We motivate this quantity in this section, by considering the special case of estimating the mean of asymmetric—very biased—Bernoulli distributions, which is in some sense an extremal case for our setting.

Example 3.2. Consider the mean estimation problem, given n samples from a Bernoulli distribution supported on 0 and 1, where the probability of drawing 1 equals some parameter p . Thus the number of 1s observed is distributed as the Binomial distribution $Bin(n, p)$, of mean np and variance $np(1 - p)$. The interesting regime for us is when p is very small, and thus $1 - p \approx 1$, and the Binomial distribution is essentially the Poisson distribution $Poi(np)$ of mean and variance $\lambda = np$. In this setting, the mean estimation problem becomes: given a sample k from $Poi(np)$, and the parameters n and δ , return an estimate that, except with failure probability δ , is as close as possible to p (or equivalently np). Given a Poisson sample $k \leftarrow Poi(np)$, returning simply k is a natural estimate of np ; however, since Poisson distributions are slightly skewed, it turns out that one should instead return the correction $k - \frac{1}{3} \log \frac{1}{\delta}$.

Explicitly, the Poisson distribution has pmf $poi(\lambda, k) = \frac{\lambda^k e^{-\lambda}}{k!}$, whose logarithm, using Stirling's approximation for the factorial, expanding to 3rd order in k , and dropping lower-order terms in λ is $-\frac{(k-\lambda)^2}{2\lambda} + \frac{(k-\lambda)^3}{6\lambda^2}$. The 2nd-order term here corresponds to a Gaussian centered at $k = \lambda$ of variance λ , which is a standard approximation for the Poisson distribution. However, crucially, the 3rd order term, corresponding to the positive skewness of the Poisson distribution, increases the pmf to the right of $k = \lambda$ and decreases it by an essentially symmetric factor to the left.

Seeking a δ -robust estimation of λ from a single sample of k , we are concerned, essentially, with the interval where the Poisson pmf is greater than δ , or equivalently, where the log pmf is greater than $\log \delta$. The quadratic $-\frac{(k-\lambda)^2}{2\lambda}$ in the first term of the above approximation equals $\log \delta$ when $k = \lambda \pm \sqrt{2\lambda \log \frac{1}{\delta}}$, and this interval is centered at λ . However, crucially, when we take into account the 3rd-order term, the interval where $\text{poi}(\lambda, k) \geq \delta$ essentially shifts to become $k = \frac{1}{3} \log \frac{1}{\delta} + \lambda \pm \sqrt{2\lambda \log \frac{1}{\delta}}$. Thus, given a single sample, one can δ -robustly estimate the mean of a Poisson distribution similarly well as the Gaussian of same mean and variance, but only if one returns the sample minus $\frac{1}{3} \log \frac{1}{\delta}$.

Thus, the $\frac{1}{3} \log \frac{1}{\delta}$ term in our estimator arises essentially from a 3rd order correction to the sample mean, at least in the special case of Bernoulli distributions.

We give another example illustrating our estimator as being a "3rd order correction". Suppose, for this section only, as in [15], that one knows the variance $\sigma^2(D)$ of the distribution in question, or has a good estimate of it.

Example 3.3. Given samples x_1, \dots, x_n from a distribution of mean 0 and variance 1 and bounded higher moments, suppose our goal is to construct a slight variant of the empirical mean that will robustly return an estimate that is close to 0, the true mean; we consider estimates of the form $\frac{1}{n} \sum_{i=1}^n (x_i + c(x_i))$ for some function $c : \mathbb{R} \rightarrow \mathbb{R}$. Explicitly, given a bound b , we want our estimate to be between $\pm b$, with as high probability as possible. For simplicity we will consider the positive case, namely, bounding $\mathbb{P}_{x_1, \dots, x_n}(\frac{1}{n} \sum_i (x_i + c(x_i)) \geq b)$. With a view towards deriving a Chernoff bound, we rearrange, multiply by an arbitrary positive constant α , and exponentiate inside the probability to yield that this probability equals $\mathbb{P}_{x_1, \dots, x_n}(\exp(\alpha \sum_i (x_i + c(x_i) - b)) \geq 1)$; by Markov's inequality, this probability is at most $\mathbb{E}_{x_1, \dots, x_n}(\exp(\alpha \sum_i (x_i + c(x_i) - b)))$, for our choice of $\alpha > 0$. We set $\alpha = b$. Since each x_i is independent, this probability becomes $\mathbb{E}_x(\exp(b(x + c(x) - b)))^n$.

Considering the empirical estimator, where $c(x_i) = 0$, we thus have that the probability the empirical mean estimate exceeds b is at most the n^{th} power of $\mathbb{E}_x(\exp(bx - b^2))$, where this expression can be expanded to 3rd order as

$$e^{-b^2} \left(1 + b \mathbb{E}(x) + \frac{1}{2} b^2 \mathbb{E}(x^2) + \frac{1}{6} b^3 \mathbb{E}(x^3) + O(x^4) \right)$$

As we assumed the data distribution has mean 0 and variance 1, we can simplify the above

expression to

$$e^{-b^2} \left(1 + \frac{1}{2}b^2 + \frac{1}{6}b^3 \mathbb{E}(x^3) + O(b^4) \right)$$

Ignoring, for the moment, the 3rd or higher-order terms, this expression is $e^{-b^2} (1 + \frac{1}{2}b^2) \approx e^{-b^2/2}$, whose n^{th} power equals $e^{-b^2 n/2}$, which is exactly the bound one would expect for the standard *Gaussian* case, of the probability that the empirical mean of n samples is more than b from the true value. However, the 3rd order term is a crucial obstacle here, as the third moment $\mathbb{E}(x^3)$ could be of either sign, skewing either the left tail or right tail to have substantially more mass than in our benchmark of the Gaussian case.

We thus choose a correction function $c(x_i)$ so as to cancel out this 3rd-order term and improve the estimate in this regime: to cancel out the term $\frac{1}{6}b^3 \mathbb{E}(x^3)$ in the 3rd-order expansion of our Chernoff bound $\mathbb{E}_x(\exp(b(x - b)))$, we replace x by $x - \frac{1}{6}x^3 b^2$, yielding a bound on the failure probability of the n^{th} power of

$$e^{-b^2} \left(1 + \frac{1}{2}b^2 + O(b^4) \right) = e^{-b^2/2 + O(b^4)}$$

as desired.

For the sake of clarity, we can change variables, letting the leading term of our probability bound $e^{-b^2 n/2}$ equal δ , and thus the correction $-\frac{1}{6}x^3 b^2$ becomes $c(x) \equiv -\frac{1}{n}x^3 \frac{1}{3} \log \frac{1}{\delta}$, meaning the correction amounts essentially to a 3rd moment correction, split n ways and scaled by the same $\frac{1}{3} \log \frac{1}{\delta}$ of our main estimator.

We explicitly relate this estimator to Estimator 1 by pointing out that, when none of the samples x_i are “truncated” by Estimator 1 (namely, $\hat{\alpha} x_i^2 \leq 1$ always), then $\hat{\alpha} \equiv \frac{\log(1/\delta)}{3n\hat{\sigma}}$ may be expressed in terms of the empirical variance $\hat{\sigma}$; taking $\kappa = 0$ for simplicity, the returned estimate will be $\frac{1}{n} \sum_i x_i - \alpha x_i^3 = \frac{1}{n} \sum_i x_i - \frac{1}{\hat{\sigma} n} x_i^3 \frac{1}{3} \log \frac{1}{\delta}$, which equals the above-derived “3rd-order corrected estimator” when the empirical variance is the true variance, 1.

In the above example we showed that Chernoff bounds for the empirical mean deteriorate for distributions with large 3rd moments (skew), and that adding a 3rd-order correction to the empirical mean corrects for this, leaving essentially “Gaussian-like” performance. These calculations motivate several features of Estimator 1—including the $\frac{1}{3} \log \frac{1}{\delta}$ parameter, and the 3rd-order terms in the expression for $\hat{\mu}$ —even though the overall form of

Estimator 1 is rather different, as it must work in all regimes and not just in the cartoon asymptotic regime considered in this example.

3.1.4 Key Contributions in Our Construction and Analysis

In addition to settling the fundamental sample complexity question of mean estimation, we point out that the estimator construction and analysis may also be of independent interest. In particular, the analysis framework—as described below—is generalizable to other problem settings and estimator constructions.

Our overall analysis framework may be viewed as a Chernoff bound—showing exponentially small probability of estimation error via bounds on a moment generation function (expectation of an exponentiated real-valued random variable). However, since we seek to analyze our estimator to sub-constant accuracy, many standard approaches fail to yield the required resolution. We point out three crucial components of our approach.

First, our estimator (Estimator 1) is *not* a sum of independent terms, which is fundamental to standard Chernoff bound approaches, and thus we instead reformulate our estimator as a 2-parameter ψ -estimator (see Definition 3.7). This technique rewrites our estimate $\hat{\mu}$ as the first coordinate of the root $(\hat{\mu}, \hat{\alpha})$ of a system of 2 equations $\psi_{\mu}(\hat{\mu}, \hat{\alpha}) = 0$ and $\psi_{\alpha}(\hat{\mu}, \hat{\alpha}) = 0$, where the functions $\psi_{\mu}(\hat{\mu}, \hat{\alpha}) = \sum_i \psi_{\mu}(x_i, \hat{\mu}, \hat{\alpha})$ and $\psi_{\alpha}(\hat{\mu}, \hat{\alpha}) = \sum_i \psi_{\alpha}(x_i, \hat{\mu}, \hat{\alpha})$ are explicitly sums of a corresponding function applied to each of the n independent data points in the sample set. Thus we have bought independence at the price of making the estimator an implicit function, introducing two new variables. One-dimensional estimators of this form are standard: for example, Catoni’s [15] mean estimator in the case of known variance is a (1 parameter) ψ -estimator for which he proves finite sample concentration. However, adding another dimension— $\hat{\alpha}$, a new implicit variable whose value the estimator will ultimately discard—is less standard, without standard analysis techniques, yet significantly increases the expressive power of such estimators [60]. Our high-level approach is to find carefully chosen linear combinations of the functions ψ_{μ} and ψ_{α} , each of which is now a sum of independent terms, and prove Chernoff bounds about these linear combinations.

Second, even after identifying these linear combinations of ψ functions, the corresponding Chernoff bound analysis is difficult to directly tackle. The Chernoff bound analysis,

as it turns out, is essentially equivalent to bounding a max-min optimization problem where the maximization is over the set of real-valued probability measures with mean 0 and variance 1. In other words, the max-min optimization problem can be interpreted as having uncountably infinitely many variables. In order to drastically simplify the problem and make it amenable to analysis, we use convex-concave programming and linear programming duality techniques to reduce the problem to a pure minimization problem with a small finite number of variables, which we can analyze tightly.

We believe that the above two ideas—1) reformulating an estimator as a multi-parameter ψ -estimator, so as to find a proxy of the estimator that is a sum of independent variables, and 2) viewing the corresponding Chernoff bound analysis as an optimization problems and applying relevant duality techniques—form a general analysis framework which expands the space of possible estimators that are amenable to *tight* analysis.

3.2 Related Work

There is a long history of work on real-valued mean estimation in a variety of models. In the problem setting we adopt, where the sole assumption is on the finiteness of the second moment, the median-of-means algorithm [34, 50, 2] has long been known to have sample complexity tight to within constant multiplicative factors, albeit with a sub-optimal constant. In seminal work, Catoni [15] improved this sample complexity to essentially optimal (tight up to a $1 + o(1)$ factor), by focusing on the special cases where the variance of the underlying distribution is known or the 4th moment is finite and bounded (in which case the second moment can be accurately estimated). We stress however that the finiteness of the 4th moment is nonetheless a much stronger assumption than our minimal assumption on the finiteness of the variance (see the discussion at the end of Section 3.1.1).

Moving beyond the original problem formulation, Devroye et al. [23] drew the distinction between a *single- δ estimator*, which takes in the robustness parameter δ as input, versus a *multiple- δ estimator*, which does not take any δ as input, but still provides guarantees across a wide range of δ values. In their work, making the same finite kurtosis assumption

as Catoni, they achieved a multiple- δ estimator with essentially optimal sample complexity, for a wide range of δ values. It is thus natural and prudent to ask whether a multiple- δ estimator can exist for the entire class of distributions with finite variance, for a meaningful range of δ values. Unfortunately, Devroye et al. [23] showed strong lower bounds answering the question in the negative. Hence, in this work, our proposed estimator is (and must be) a single- δ estimator, taking δ as input.

Many applications have arisen from the success of sub-Gaussian mean estimation, showing how to leverage or extend Catoni-style estimators to new settings, achieving sub-Gaussian performance on problems such as regression, empirical risk minimization, and online learning (bandit settings): for example see [48, 8, 16, 9].

A separate but closely related line of work is on *high dimensional* mean estimation. While estimators generalizing the “median-of-means” construction were found to have statistical convergence tight to multiplicative constants, until recently, such estimators took super-polynomial time to compute [45]. A recent line of work [31, 20, 38], started by Hopkins [31], thus focuses on the computational aspect, and brought the computation time first down to polynomial time, with subsequent work bringing it further down to quadratic time using spectral methods.

A recent comprehensive survey by Lugosi and Mendelson [43] explains much of the above works in greater detail.

Mean estimation is also well studied in various more restrictive settings, and we highlight a recent line of work seeking to find optimal convergence under *differential privacy* constraints. Kamath et al. [36] studies the differentially private mean estimation problem in the constant probability regime, and shows strong sample complexity separations from our unrestricted setting. Duchi, Jordan and Wainwright [25, 26] study the problem under the stricter constraint of *local* differential privacy. See Kamath et al. [36] for a more comprehensive literature review on differentially private mean estimation.

Similar in style to our approach of “throwing out the most extreme $c \log \frac{1}{\delta}$ samples and returning the mean of the rest”, Oliveira and Orenstein [51] show that this “trimmed mean” estimator enjoys similar tight-up-to-constants guarantees as the median-of-means estimator. This approach was significantly expanded by Lugosi and Mendelson to encompass the

high dimensional case [44], which has the additional feature of being optimally *robust* (up to constants) to adversarial contamination in the samples. Recent work by Diakonikolas et al. [24] improves on this result by giving the first polynomial-time computable estimator that achieves the same optimal robustness guarantees. Furthermore, in the absence of adversarial contamination, Diakonikolas et al. [24] also simplify the arguments in [31, 20, 38] for showing a computationally efficient sub-Gaussian mean estimator in high dimensions.

Part of our tight analysis relies on insights from mathematical programming and duality; see [54] for a detailed discussion of prior works that use such mathematical programming and duality tools to either design or analyze statistical estimators [53, 49, 66, 67, 35, 63].

3.3 Our Mean Estimator

In this section, we present our estimator (Estimator 1), as well as its reformulation as a 2-parameter ψ -estimator. We then present some perspective and basic structural properties of the estimator that will serve as a foundation for the analysis to follow.

Estimator 1 The Main Estimator

Inputs:

- n independent samples $\{x_i\}$ from the unknown underlying distribution D (guaranteed to have finite variance)
 - Confidence parameter δ
1. Compute the median-of-means estimate κ : evenly partition the data into $\log \frac{1}{\delta}$ groups and let κ be the median of the set of means of the groups.
 2. Find the solution $\hat{\alpha}$ to the monotonic, piecewise-linear equation $\sum_i \min(\hat{\alpha}(x_i - \kappa)^2, 1) = \frac{1}{3} \log \frac{1}{\delta}$
 3. Output: $\hat{\mu} = \kappa + \frac{1}{n} \sum_i (x_i - \kappa)(1 - \min(\hat{\alpha}(x_i - \kappa)^2, 1))$
-

3.3.1 Meaning of the Estimator

Consider the expression in Step 3 for the final returned value of the estimator, $\hat{\mu} = \kappa + \frac{1}{n} \sum_i (x_i - \kappa)(1 - \min(\hat{\alpha}(x_i - \kappa)^2, 1))$. Without the final min expression, the expression $\kappa +$

$\frac{1}{n} \sum_i (x_i - \kappa) \cdot 1$ computes exactly the sample mean. The factor $(1 - \min(\hat{\alpha}(x_i - \kappa)^2, 1))$ may be thought of as a weight on the i^{th} element, between 0 and 1, where a weight of 1 leaves that element as is, but a weight towards 0 essentially throws out part of the sample x_i and instead defaults to the median-of-means estimate κ . Thus, rather than either keeping or discarding each entry, the weight $\min(\hat{\alpha}(x_i - \kappa)^2, 1)$ specifies what *fraction* of the i^{th} sample to discard.

The condition in Step 2 of Estimator 1 picks α so that the total, weighted, number of discarded samples equals $\frac{1}{3} \log \frac{1}{\delta}$. The expression $\min(\hat{\alpha}(x_i - \kappa)^2, 1)$ specifying what fraction of each x_i to discard says, essentially, that this fraction should be proportional to the square of the deviation of x_i from the mean estimate κ , capped at 1 so that we do not discard “more than 100% of” any sample x_i .

3.3.2 Structural Properties of the Estimator

We point out three basic structural properties of Estimator 1 that both shed light on the estimator itself, and will be crucial to its analysis.

First, the estimator is “affine invariant” in the sense that, if its input samples $\{x_i\}$ undergo an affine map $x \rightarrow ax + b$ then its output will be mapped correspondingly.

Lemma 3.4. *Suppose X is a set of samples in \mathbb{R} . Then for any $\delta > 0$ and any scale $a > 0$ and shift b ,*

$$\hat{\mu}(aX + b, \delta) = a \hat{\mu}(X, \delta) + b$$

where $\hat{\mu}$ denotes the output of Estimator 1.

The above lemma follows trivially from the fact that the median-of-means estimate also respects shift and scale in the input samples, and that α is chosen in Step 2 of Estimator 1 so that $\min(\alpha(x_i - \kappa)^2, 1)$ does not depend on the affine parameters a, b .

Second, as is well known, the median-of-means estimate κ of Step 1, while not as accurate as what we will eventually return, is robust in the sense that, with probability at least $1 - \delta/2$, the median-of-means estimate has additive error from the true mean that is at most $O(\sigma \cdot \sqrt{\frac{\log \frac{1}{\delta}}{n}})$ —proportional to the eventual guarantees of our estimator, but with somewhat worse proportionality constant.

Fact 3.5 ([32]). *For any distribution D with mean μ and standard deviation σ , the median-of-means estimate κ , on input n samples, satisfies*

$$\mathbb{P} \left(|\kappa - \mu| > O \left(\sigma \sqrt{\frac{\log \frac{1}{\delta}}{n}} \right) \right) \leq \delta$$

Third, if we temporarily ignore Step 1, treating κ as a free parameter, we show that the final output of the estimator, $\hat{\mu}$, varies very little with κ . Combined with the accuracy guarantees of the median-of-means estimate, the difference in the final estimate between using the median-of-means as κ versus using the *true* mean as κ is inconsequential (a $o(1)$ factor) compared to the total additive error we aim for. Therefore, for the purposes of *analysis*, it suffices to assume that κ takes the value of the true mean (though an estimator could not do this in practice, as the true mean is unknown).

Lemma 3.6. *Consider a fixed sample set X of size n , and a confidence parameter δ . Let $e(X, \delta, \kappa)$ denote Estimator 1 but where Step 1 is omitted and κ is instead considered as an input. Then,*

$$\left| \frac{d e(X, \delta, \kappa)}{d \kappa} \right| = O \left(\sqrt{\frac{\log \frac{1}{\delta}}{n}} \right)$$

Fact 3.5 shows that, except with δ probability, the median-of-means estimate is within $O \left(\sigma \sqrt{\frac{\log \frac{1}{\delta}}{n}} \right)$ of the true mean, and multiplying this by the Lipschitz constant $O \left(\sqrt{\frac{\log \frac{1}{\delta}}{n}} \right)$ from Lemma 3.6 shows that the change in output of Estimator 1, between using the median-of-means versus setting $\kappa = 0$, has magnitude $O \left(\sigma \sqrt{\frac{\log \frac{1}{\delta} 2}{n}} \right) = o \left(\sigma \sqrt{\frac{\log \frac{1}{\delta}}{n}} \right)$. This discrepancy is therefore a $o(1)$ fraction of the additive error guaranteed by Theorem 3.1.

We now prove Lemma 3.6.

Proof. We compute the derivatives with respect to $\hat{\alpha}$ and κ of the $\hat{\mu}$ (computed in Step 3 of Estimator 1), and the expression on the left hand side of Step 2, which we denote $v \equiv \sum_i \min(\hat{\alpha}(x_i - \kappa)^2, 1)$. We note that for terms where $\min(\hat{\alpha}(x_i - \kappa)^2, 1) = 1$, all derivatives are 0, so we adopt the notation “ $\Sigma_{<}$ ” to denote summing only over those indices i for which $\hat{\alpha}(x_i - \kappa)^2 < 1$. Thus we have

$$\frac{dv}{d\kappa} = 2 \sum_{<} \hat{\alpha}(x_i - \kappa)$$

$$\begin{aligned}\frac{dv}{d\hat{\alpha}} &= \sum_{<} (x_i - \kappa)^2 \\ \frac{d\hat{\mu}}{d\kappa} &= 1 + \frac{1}{n} \sum_{<} (-1 + 3\hat{\alpha}(x_i - \kappa)^2) \\ \frac{d\hat{\mu}}{d\hat{\alpha}} &= -\frac{1}{n} \sum_{<} (x_i - \kappa)^3\end{aligned}$$

Recall that $\hat{\alpha}$ is defined implicitly so as to make the expression $v = \frac{1}{3} \log \frac{1}{\delta}$; thus in Estimator 1, if we change κ at a rate of 1, then $\hat{\alpha}$ also changes at rate $-\frac{dv}{d\kappa} / \frac{dv}{d\hat{\alpha}}$ to keep v unchanged. Thus, the overall derivative of the estimate with respect to changing κ equals $\frac{d\hat{\mu}}{d\kappa} - \frac{d\hat{\mu}}{d\hat{\alpha}} \frac{dv}{d\kappa} / \frac{dv}{d\hat{\alpha}}$. We bound this from the derivatives computed above.

To bound $\frac{d\hat{\mu}}{d\kappa}$, we note that the number of indices *not* in the sum " $\sum_{<}$ " is at most $\frac{1}{3} \log \frac{1}{\delta}$ because each such i contributes 1 to the left hand side of the condition in Step 2 of Estimator 1 and the right hand side equals $\frac{1}{3} \log \frac{1}{\delta}$. Thus the initial terms of $\frac{d\hat{\mu}}{d\kappa}$ are bounded as $1 + \frac{1}{n} \sum_{<} (-1) \leq \frac{1}{3n} \log \frac{1}{\delta}$. The remaining part of $\frac{d\hat{\mu}}{d\kappa}$, namely $\frac{1}{n} \sum_{<} 3\hat{\alpha}(x_i - \kappa)^2$ is $\frac{3}{n}$ times the corresponding terms in $v \leq \frac{1}{3} \log \frac{1}{\delta}$ itself, and thus is at most $\frac{1}{n} \log \frac{1}{\delta}$. Thus $\frac{d\hat{\mu}}{d\kappa} = O(\frac{1}{n} \log \frac{1}{\delta})$.

We now bound the remaining term $-\frac{d\hat{\mu}}{d\hat{\alpha}} \frac{dv}{d\kappa} / \frac{dv}{d\hat{\alpha}}$. Since for each index i in " $\sum_{<}$ " we have $|x_i - \kappa| \leq \frac{1}{\sqrt{\hat{\alpha}}}$, we may bound $\frac{d\hat{\mu}}{d\hat{\alpha}}$, involving a 3rd moment term, by the simpler $|\frac{d\hat{\mu}}{d\hat{\alpha}}| \leq \frac{1}{n} \sum_{<} |x_i - \kappa|^2 / \sqrt{\hat{\alpha}}$. Combining this, with the other derivatives and the bound $\hat{\alpha} \leq \frac{1}{3} \log \frac{1}{\delta} / \sum_{<} (x_i - \kappa)^2$ from the previous paragraph yields:

$$\left| \frac{d\hat{\mu}}{d\hat{\alpha}} \frac{dv}{d\kappa} / \frac{dv}{d\hat{\alpha}} \right| \leq \frac{2\sqrt{\hat{\alpha}}}{n} \left| \frac{\sum_{<} (x_i - \kappa) \sum_{<} (x_i - \kappa)^2}{\sum_{<} (x_i - \kappa)^2} \right| \leq \frac{2\sqrt{\frac{1}{3} \log \frac{1}{\delta}}}{n} \left| \frac{\sum_{<} (x_i - \kappa)}{\sqrt{\sum_{<} (x_i - \kappa)^2}} \right| \leq \sqrt{\frac{4 \log \frac{1}{\delta}}{3n}}$$

where the last inequality is Cauchy-Schwarz applied to the sequence $(x_i - \kappa)$ and the all-1s sequence. \square

The above three structural properties allow us to drastically simplify the analysis: the affine invariance means it is sufficient to show our estimator works for the special case of distributions with mean 0 and variance 1; the second and third properties mean that errors in κ effectively do not matter, and, for distributions with mean 0, it is sufficient to omit Step 1 and instead just analyze the case where $\kappa = 0$.

We point out that Estimator 1 when modified to set $\kappa = 0$ (independently of the samples) is *no longer* affine invariant, nor is its reformulation as a ψ -estimator in Section 3.3.3. The

structural properties in this section show that, instead of analyzing the actual estimator (Estimator 1 which is affine invariant), it suffices to analyze this artificially simplified, although no longer affine invariant, estimator which sets $\kappa = 0$, on distributions with mean 0 and variance 1. Explicitly, in the rest of the paper we will show Proposition 3.10 (Section 3.4), which analyzes the mean-0 variance-1 case of the ψ -estimator defined below in Definition 3.7; the discussion of this section shows that this proposition implies our main result, Theorem 3.1.

3.3.3 Representing a Special Case of Estimator 1 as a ψ -Estimator

As discussed in Section 3.1.4, our estimator, even its simplified version with $\kappa = 0$, is not a sum of independent terms, making it difficult to tightly bound its moment generating function, and hence also difficult to prove its concentration around the true mean using a Chernoff-style bound. Our solution is to reformulate Estimator 1, with the simplifying assumption that $\kappa = 0$, as a 2-parameter ψ -estimator, as defined in Definition 3.7. This reformulation defines our estimate $\hat{\mu}$ implicitly in terms of two new functions ψ_μ and ψ_α that are indeed sums of n independent terms, each term depending on a single x_i . We will use this representation crucially for the concentration analysis of the estimator.

Definition 3.7. Consider Estimator 1 but with Step 1 replaced with “ $\kappa = 0$ ”. The estimator can be equivalently expressed as follows:

1. Input: n independent samples $X = x_1, \dots, x_n$
2. Solve for the (unique) pair $(\hat{\mu}, \hat{\alpha})$ satisfying $\psi_\mu = 0$ and $\psi_\alpha = 0$, where the functions are defined as follows:

$$\psi_\mu(X, \hat{\mu}, \hat{\alpha}) = \sum_{i=1}^n \left(\hat{\mu} - x_i \left(1 - \min(\hat{\alpha} x_i^2, 1) \right) \right)$$

$$\psi_\alpha(X, \hat{\mu}, \hat{\alpha}) = \sum_{i=1}^n \left(\min(\hat{\alpha} x_i^2, 1) - \frac{1}{3n} \log \frac{1}{\delta} \right)$$

(Note that $\hat{\alpha} > 0$ always)

3. Output: $\hat{\mu}$ from the previous step

We will sometimes omit $\hat{\mu}$ from the arguments of ψ_α since $\hat{\mu}$ is not used in the definition of the function. We will often refer to the pair (ψ_μ, ψ_α) as a 2-element vector ψ .

For convenience in the rest of the paper, we define $\hat{v} \equiv \frac{\log(1/\delta)}{3n\hat{\alpha}}$, which we refer to as the “truncated empirical variance”; this is because, if we modify the $\psi_\alpha = 0$ condition by removing the “truncation” of taking the min with 1, then the resulting condition, when expressed in terms of $\hat{v} = \frac{\log(1/\delta)}{3n\hat{\alpha}}$ and rearranged, is exactly the condition that \hat{v} is the empirical variance: $\frac{1}{n} \sum_{i=1}^n x_i^2$. Thus $\hat{\alpha}$ may be thought of as a proxy for the empirical variance, as $\hat{v} = \frac{\log(1/\delta)}{3n\hat{\alpha}}$ equals the empirical variance, except in cases when samples are far enough from 0 that they are “truncated” by the “min”.

Interestingly, in the case that none of the samples are “truncated”, (and $\kappa = 0$), the overall output of the estimator becomes $\frac{1}{n} \sum_i x_i - \hat{\alpha} x_i^3 = \frac{1}{n} \sum_i x_i - \frac{\log(1/\delta)}{3n\hat{v}} x_i^3$, namely, $\hat{\mu}$ is “the empirical mean, corrected by subtracting $\frac{1}{3n} \log \frac{1}{\delta}$ times the ratio of the empirical 3rd moment over the empirical 2nd moment.”

Proof that Definition 3.7 is equivalent to Estimator 1 when κ is set to 0. Fix a set of samples $X = \{x_i\}$. We observe that Estimator 1, with the additional simplifying assumption that $\kappa = 0$, can be represented by the following 2 equations.

$$\begin{aligned} \sum_i \min(\hat{\alpha} x_i^2, 1) &= \frac{1}{3} \log \frac{1}{\delta} \\ \hat{\mu} &= \frac{1}{n} \sum_i x_i (1 - \min(\hat{\alpha} x_i^2, 1)) \end{aligned} \tag{3.2}$$

Estimator 1 solves for $\hat{\alpha}$ in the first line, and uses this $\hat{\alpha}$ value to compute the estimate $\hat{\mu}$ in the second line. The two conditions of Equation 3.2 are equivalent to the two conditions $\psi_\alpha = 0$, $\psi_\mu = 0$ respectively, and thus the two estimators are equivalent. \square

3.3.4 The relation to the Catoni-Giulini estimator

As a side note, in this section we discuss the similarities and significant differences between our estimator and the Catoni-Giulini estimator [15, 16], which has optimal convergence assuming knowledge of the variance.

Definition 3.8. Define (as will be used in this section only) the function T that “truncates” its input t to a specified range $[-r, r]$, as $T_r(t) = \min(r, \max(-r, t))$; and define the influence function $\psi(t) = t - \frac{1}{6}t^3$.

We use Definition 3.8 to re-express the final estimate $\hat{\mu}$ returned in Step 3 of Estimator 1—where κ is the initial median-of-means estimate and $\hat{\lambda} = \sqrt{6\hat{\alpha}} = \sqrt{\frac{2\log(1/\delta)}{n\hat{\sigma}}}$ comes from the $\hat{\alpha}$ computed in Step 2 of Estimator 1, or equivalently from the truncated empirical variance $\hat{\sigma}$ defined from $\hat{\alpha}$ in Section 3.3.3—as

$$\hat{\mu} = \kappa + \frac{1}{n\hat{\lambda}} \sum_i \psi(T_{\sqrt{6}}(\hat{\lambda}(x_i - \kappa)))$$

On the other hand, the Catoni-Giulini estimator chooses $\hat{\lambda} \approx \sqrt{\frac{2\log(1/\delta)}{n\sigma^2}}$ using the *true* variance σ^2 (see Proposition 2.4 of [15]), and solves for $\hat{\mu}$ in the very similar equation

$$0 = \frac{1}{n} \sum_i \psi(T_{\sqrt{2}}(\hat{\lambda}(x_i - \hat{\mu})))$$

With this view, the two estimators are similar, albeit with the crucial differences that 1) our $\hat{\alpha}$ is computed from the data while the Catoni-Giulini $\hat{\alpha}$ is computed from the true variance, and 2) we use the truncation constant $\sqrt{6}$ instead of the $\sqrt{2}$ used by Catoni and Giulini. As Catoni and Giulini noted in [15, 16], the constant $\sqrt{2}$ is the largest “truncation constant” r for which the function $\psi(T_r(t))$ is monotonic. Monotonicity of the “influence function” ψ is a crucial proof technique in [15], responsible for some of the generality of that paper. Further, the non-monotonicity of our analog of ψ leads our overall estimator to be non-monotonic in its data: there is an input vector to Estimator 1, that, when some of its entries are increased, actually leads to a *smaller* final estimate $\hat{\mu}$, which is counterintuitive. One may thus ask if the $\sqrt{6}$ in our estimator can be replaced with $\sqrt{2}$, to make the estimator monotonic. Intriguingly, this ruins its performance, though subtler modifications of our estimator to make it monotonic may be possible.

We also point out that a *variant* of our estimator may be expressed as a (2-parameter) ψ -estimator—with no separate median-of-means preprocessing step needed, in line with the ψ -estimators of [15, 16]. This ψ -estimator formulation might be of independent interest given the huge body of work analyzing such estimators, though this formulation is not used in our analysis (which instead uses the setup of Section 3.3.3).

Definition 3.9 (A ψ -estimator variant of Estimator 1). Given n independent samples $X = x_1, \dots, x_n$, solve for $(\hat{\mu}, \hat{\lambda})$ satisfying the following equations, and return $\hat{\mu}$:

$$\sum_{i=1}^n T_{\sqrt{\delta}}(\hat{\lambda}(x_i - \hat{\mu}))^2 = 2 \log \frac{1}{\delta} \quad \text{and} \quad \sum_{i=1}^n \psi(T_{\sqrt{\delta}}(\hat{\lambda}(x_i - \hat{\mu}))) = 0 \quad (3.3)$$

From this perspective, Steps 2 and 3 of Estimator 1 may be viewed as a single update of an iterative algorithm to solve Equation 3.3, starting with the guess $\hat{\mu} = \kappa$ found by the median-of-means algorithm. The Lipschitz analysis in Section 3.3.2 (the third property in Section 3.3.2, see Lemma 3.6) essentially shows that these updates converge extremely rapidly, which can further be used to show that this concise estimator also satisfies the guarantees of Theorem 3.1.

3.4 Analyzing our estimator

In this section, we outline the proof of our main theorem, restated as follows.

Theorem 3.1. *Estimator 1, given $\delta, n > 0$, and a sample set X of n independent samples from distribution D , will, with probability at least $1 - \delta$ over the sampling process, yield an estimate $\hat{\mu}$ with error at most $|\hat{\mu}(X) - \mu(D)| \leq \sigma(D) \cdot (1 + o(1)) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$. Here, the $o(1)$ term tends to 0 as $(\frac{\log \frac{1}{\delta}}{n}, \delta) \rightarrow (0, 0)$.*

The discussion of the structural properties of Estimator 1 in Section 3.3.2 shows that it is sufficient to instead show that, for any distribution of mean 0 and variance 1, the ψ -estimator of Definition 3.7 will return an estimate $\hat{\mu}$ that is close to 0, except with tiny probability. Recall also that, since the ψ -estimator solves for $(\hat{\mu}, \hat{\alpha})$ such that $\psi(X, \hat{\mu}, \hat{\alpha}) = 0$ (where X is the sample set) and returns $\hat{\mu}$, the claim that the returned estimate will be close to 0 is equivalent to saying that *every* $(\hat{\mu}, \hat{\alpha})$ pair with $\hat{\mu}$ far from 0 must violate the equation, namely $\psi(X, \hat{\mu}, \hat{\alpha}) \neq 0$. We thus prove the following proposition (Proposition 3.10), to yield Theorem 3.1. Note that the failure probability in Proposition 3.10 is $\delta/2$ (instead of δ , as in Theorem 3.1), accounting for an additional $\delta/2$ probability that the median-of-means estimate in Step 1 of Estimator 1 fails.

Proposition 3.10. *There exists a universal constant $c > 0$ such that, fixing $\epsilon' = \left(1 + \frac{c \log \log \frac{1}{\delta}}{\log \frac{1}{\delta}}\right) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$,*

we have that for all distributions D with mean 0 and variance 1, with probability at least $1 - \frac{\delta}{2}$ over the set of samples X , for all $\hat{\mu}, \hat{\alpha}$ where $|\hat{\mu}| > \epsilon'$ and $\hat{\alpha} > 0$, the vector $\psi(X, \hat{\mu}, \hat{\alpha}) \neq 0$.

Proposition 3.10 asks us to show that, with high probability, $\psi(X, \hat{\mu}, \hat{\alpha})$ is not at the origin for *any* choice of $|\hat{\mu}| > \epsilon', \hat{\alpha}$; instead, as a proof strategy, we choose a finite bounded mesh of $\hat{\mu}, \hat{\alpha}$ and show that the function $\psi(X, \hat{\mu}, \hat{\alpha})$ is 1) not just nonzero, but far from the origin on this set, 2) Lipschitz in between mesh elements, and 3) monotonic (in an appropriate sense) outside the mesh bounds. Step 1), discussed below, contains the most noteworthy part of the proof, a mathematical programming-inspired bound to help complete a delicate Chernoff bound argument.

For simplicity, we reparameterize to work with $\hat{\nu} \equiv \frac{\log(1/\delta)}{3n\hat{\alpha}}$ (the “truncated empirical variance”) instead of $\hat{\alpha}$: the mesh we analyze, covering the most delicate region for analysis, will span the interval $\hat{\nu} \in [0.05, 55.5]$, namely, where the truncated empirical variance $\hat{\nu}$ is within a constant factor of the true variance of 1. Note that this should *not* be taken to imply that $\hat{\nu} \in [0.05, 55.5]$ with high probability—the truncated empirical variance is not designed to be a good estimate of the variance, merely as a step in robustly estimating the mean; and further, accurate estimates of the variance are simply impossible in general without further assumptions such as bounds on the distribution’s 3rd or 4th moments. We also want to distinguish our estimator from Catoni’s [15]: Catoni’s estimator relies on having a high-precision estimate of the variance (to within a $1 + o(1)$ factor) in order to achieve the desired performance. By contrast, our estimator is robust against wild inaccuracies of the (truncated) empirical variance $\hat{\nu}$ compared to the true variance of 1. In short, the approach of our estimator should be viewed as distinct from Catoni’s, since, while Catoni’s estimator relies on an initial good guess at the variance, ours thrives in the inevitable situations where $\hat{\nu}$ is far from 1.

We return to describing our strategy for analyzing the performance of our estimator. For each $\hat{\mu}, \hat{\nu} = \frac{\log(1/\delta)}{3n\hat{\alpha}}$ that we analyze (from the finite mesh): instead of directly showing that, with $\geq 1 - \frac{\delta}{2}$ probability, $\psi(X, \hat{\mu}, \hat{\alpha})$ is far from the origin in some direction, we instead *linearize* this claim; we prove the stronger claim that there exists a specific direction $\mathbf{d}(\hat{\nu})$ such that with $\geq 1 - \frac{\delta}{2}$ probability, $\psi(X, \hat{\mu}, \hat{\alpha})$ is more than $\frac{1}{\log(1/\delta)}$ distance from the origin in direction \mathbf{d} (specifically we lower bound the dot product $\mathbf{d}(\hat{\nu}) \cdot \psi(\hat{\mu}, \hat{\alpha})$, while we upper

bound each coordinate of \mathbf{d} inversely with the Lipschitz coefficients of ψ). The crucial advantage of this reformulation is that, since each of ψ_μ, ψ_α is a sum of n terms, that are each a function of an independent sample x_i from D , the dot product $\mathbf{d}(\hat{\nu}) \cdot \psi(X, \hat{\mu}, \hat{\alpha})$ is thus also a sum of n independent terms, and thus we finish the proof with a Chernoff bound, Lemma 3.11. The Chernoff bound argument itself is standard; however, to bound the resulting expression requires an extremely delicate analysis that we pull out into a separate 4-variable inequality expressed as Lemma 3.14—see the discussion around the lemma for more details and for motivation of the analysis from a mathematical programming perspective.

We state the crucial Chernoff bound (Lemma 3.11) and the Lipschitz bounds (Lemma 3.12), and then use them to prove Proposition 3.10. We prove Lemmas 3.11 and 3.12 in the next section, along with the statement and proof of the delicate component that is Lemma 3.14.

Lemma 3.11. *Consider an arbitrary distribution D with mean 0 and variance 1. There exists a universal constant c where the following claim is true. Fixing $\hat{\mu} = \epsilon' = \left(1 + \frac{c \log \log \frac{1}{\delta}}{\log \frac{1}{\delta}}\right) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$, then for all δ smaller than some universal constant, and for all $\hat{\nu} \in [0.05, 55.5]$, there exists a vector $\mathbf{d}(\hat{\nu})$ where $d_\mu \geq 0$, and both $\sqrt{\frac{n}{\log(1/\delta)}} |d_\mu|, |d_\alpha|$ are bounded by a universal constant, such that*

$$\mathbb{P}_{X \leftarrow D^n} \left(\mathbf{d}(\hat{\nu}) \cdot \psi \left(X, \hat{\mu} = \epsilon', \hat{\alpha} = \frac{\log(1/\delta)}{3n\hat{\nu}} \right) > \frac{1}{\log \frac{1}{\delta}} \right) \geq 1 - \frac{\delta}{\log^4 \frac{1}{\delta}}$$

Furthermore, for $\hat{\nu} = 0.05$ we have $d_\mu = \sqrt{3.75 \frac{\log(1/\delta)}{n}}$, $d_\alpha = \sqrt{3}$; and for $\hat{\nu} = 55.5$ we have $d_\mu = 0$, $d_\alpha < 0$.

Lemma 3.12. *Consider an arbitrary set of n samples X . Consider the expressions $\psi_\mu(X, \hat{\mu}, \hat{\alpha}), \psi_\alpha(X, \hat{\alpha})$, reparameterized in terms of $\hat{\nu} \equiv \frac{\log(1/\delta)}{3n\hat{\alpha}}$ in place of $\hat{\alpha}$. Suppose the equation $\psi_\alpha(X, \hat{\alpha}) = 0$ has a solution in the range $\hat{\nu} \in [0.05, 55.5]$. Then the functions $\sqrt{\frac{\log(1/\delta)}{n}} \psi_\mu(X, \hat{\mu}, \hat{\alpha})$ and $\psi_\alpha(X, \hat{\alpha})$ are Lipschitz with respect to $\hat{\nu}$ on the entire interval $\hat{\nu} \in [0.05, 55.5]$, with Lipschitz constant $c \log \frac{1}{\delta}$ for some universal constant c .*

We now prove Proposition 3.10, which per our previous discussion, implies our main result, Theorem 3.1.

Proof of Proposition 3.10. As in Lemma 3.11, we fix $\epsilon' = \left(1 + \frac{c \log \log \frac{1}{\delta}}{\log \frac{1}{\delta}}\right) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$, where c is some universal constant.

By symmetry, instead of considering positive and negative $\hat{\mu}$, it suffices to consider the case $\hat{\mu} > \epsilon'$ (as opposed to $\hat{\mu} < -\epsilon'$) and show that this case succeeds with probability at least $1 - \frac{\delta}{4}$.

To prove the claim, we first prove a stronger statement on a restricted domain, that with probability at least $1 - \frac{\delta}{4}$ over the randomness of the sample set X , for each $\hat{\nu} \in [0.05, 55.5]$ there exists a vector $\mathbf{d} = (d_\mu, d_\alpha)$ such that $\mathbf{d} \cdot \psi(X, \epsilon', \hat{\alpha} \equiv \frac{\log(1/\delta)}{3n\hat{\nu}}) > 0$, with $d_\mu \geq 0$ throughout, and, for $\hat{\nu} = 0.05$ we have $d_\mu = \sqrt{3.75 \frac{\log(1/\delta)}{n}}$, $d_\alpha = \sqrt{3}$; and for $\hat{\nu} = 55.5$ we have $d_\mu = 0$, $d_\alpha < 0$.

We will first apply Lemma 3.11 to each $\hat{\nu}$ in a discrete mesh: let M consist of evenly spaced points between 0.05 and 55.5 with spacing $1/\log^3 \frac{1}{\delta}$ (thus with $\Theta(\log^3 \frac{1}{\delta})$ many points).

By Lemma 3.11 and a union bound over these $\Theta(\log^3 \frac{1}{\delta})$ points, we have that with probability at least $1 - \frac{\delta}{\Theta(\log \frac{1}{\delta})}$ (which is at least $1 - \frac{\delta}{4}$ for δ smaller than some universal constant) over the set of n samples X , for all $\hat{\nu} \in M$, there exists a vector $\mathbf{d}(\hat{\nu})$ such that $\mathbf{d}(\hat{\nu}) \cdot \psi(X, \hat{\mu} = \epsilon', \hat{\alpha} \equiv \frac{\log(1/\delta)}{3n\hat{\nu}}) > 1/\log \frac{1}{\delta}$, where \mathbf{d} further satisfies the desired positivity and boundary conditions, and where both $\sqrt{\frac{n}{\log(1/\delta)}} |d_\mu|, |d_\alpha|$ are bounded by a universal constant. For the rest of the proof, we will only consider sets of samples X satisfying the above condition.

Now consider an arbitrary $\hat{\nu}' \in [0.05, 55.5] \setminus M$ and consider the vector ψ evaluated at $\hat{\alpha}' = \frac{\log(1/\delta)}{3n\hat{\nu}'}$. We wish to extend the dot product inequality to hold also for $\hat{\nu}'$. If $\psi_\alpha \neq 0$ then there is nothing to prove: set $d_\mu = 0$ and $d_\alpha = \text{sign}(\psi_\alpha)$; otherwise, $\psi_\alpha = 0$ means we may apply Lemma 3.12 to conclude that both $\sqrt{\frac{n}{\log(1/\delta)}} \psi_\mu(X, \hat{\mu}, \hat{\alpha}')$ and $\psi_\alpha(X, \hat{\mu}, \hat{\alpha}')$ are Lipschitz with respect to $\hat{\nu}'$ on the interval $\hat{\nu}' \in [0.05, 55.5]$, with Lipschitz constant $c \log \frac{1}{\delta}$ for some universal constant c .

Consider the closest $\hat{\nu} \in M$ to $\hat{\nu}'$, which by definition of M is at most $1/\log^3 \frac{1}{\delta}$ away. By assumption on X , there exists a vector \mathbf{d} such that $\mathbf{d} \cdot \psi(X, \hat{\mu} = \epsilon', \hat{\alpha} = \frac{\log(1/\delta)}{3n\hat{\nu}}) > 1/\log \frac{1}{\delta}$, with $d_\mu \geq 0$ and both $\sqrt{\frac{n}{\log(1/\delta)}} |d_\mu|, |d_\alpha|$ are bounded by a universal constant. Because of the Lipschitz bounds on ψ , combined with the bounds on the size of the d_μ, d_α , we conclude that the Lipschitz constant of the dot product (treating the vector \mathbf{d} as fixed) is $O(\log \frac{1}{\delta})$. Thus, the large positive dot product at $\hat{\nu}$ implies at least a positive dot product nearby at

$\hat{\nu}'$: $\mathbf{d} \cdot \psi(X, \hat{\mu} = \epsilon', \hat{\nu}') > \frac{1}{\log \frac{1}{\delta}} - O(\log \frac{1}{\delta}) \frac{1}{\log^3 \frac{1}{\delta}} > 0$, for sufficiently small δ as given in the proposition statement.

Having shown the stronger version of the claim for the restriction $\hat{\mu} = \epsilon'$ and $\hat{\nu} \in [0.05, 55.5]$ we now extend to the entire domain via three monotonicity arguments. Explicitly, assume the set of samples X satisfies the dot product inequality above with the vector function $\mathbf{d}(\hat{\nu})$, where $\mathbf{d}(\hat{\nu})$ satisfies the boundary conditions at $\hat{\nu} = 0.05$ and 55.5 specified in Lemma 3.11. From this assumption, we will show that $\psi \neq 0$ for *any* positive $\hat{\nu} = \frac{\log(1/\delta)}{3n\hat{\alpha}}$, and for *any* $\hat{\mu} \geq \epsilon'$.

First consider $\hat{\nu} > 55.5$ (still fixing $\hat{\mu} = \epsilon'$). The function $\psi_\alpha = \sum_{i=1}^n (\min(\hat{\alpha}x_i^2, 1) - \frac{1}{3n} \log \frac{1}{\delta})$ is an increasing function of $\hat{\alpha}$, and thus a decreasing function of $\hat{\nu} \equiv \frac{\log(1/\delta)}{3n\hat{\alpha}}$. Since for $\hat{\nu} = 55.5$, the dot product $\mathbf{d} \cdot \psi > 0$ with $d_\mu = 0, d_\alpha < 0$, the dot product will thus remain positive for this same choice of \mathbf{d} as we increase $\hat{\nu}$ from 55.5 .

Next, for $\hat{\nu} < 0.05$ (again still fixing $\hat{\mu} = \epsilon'$), we analogously show that the dot product of $\psi(X, \epsilon', \hat{\alpha} \equiv \frac{\log(1/\delta)}{3n\hat{\nu}})$ with the fixed vector $\mathbf{d}(0.05)$ will increase as we decrease $\hat{\nu}$. The i^{th} term in the sums defining ψ_μ or ψ_α depends on $\hat{\alpha}$ (and thus $\hat{\nu}$) only in the factor $\min(\hat{\alpha}x_i^2, 1)$. Further, there is no dependence unless the first term attains the min, namely $|x_i| \leq \sqrt{1/\hat{\alpha}}$, which in turn is upper bounded by $\sqrt{0.15 \frac{n}{\log(1/\delta)}}$ because of our assumption that $\hat{\nu} < 0.05$. Thus, the only i^{th} terms in the dot product which have $\hat{\alpha}$ dependent are simply equal to $d_\mu \hat{\alpha}x_i^3 + d_\alpha \hat{\alpha}x_i^2 = \hat{\alpha}x_i^2(d_\alpha + x_i d_\mu)$. By our choice of $d_\mu(0.05) = \sqrt{3.75 \frac{\log(1/\delta)}{n}}$ and $d_\alpha(0.05) = \sqrt{3}$ from Lemma 3.11, the expression $(d_\alpha + x_i d_\mu) \geq \sqrt{3} - \sqrt{0.15} \sqrt{3.75}$ is thus always non-negative, and thus the overall dot product cannot decrease as we send $\hat{\alpha}$ to ∞ —equivalently, sending $\hat{\nu}$ to 0 —as desired.

We have thus shown that, for all non-negative $\hat{\alpha} = \frac{\log(1/\delta)}{3n\hat{\nu}}$, there is a vector \mathbf{d} with $d_\mu \geq 0$ whose dot product with $\psi(X, \epsilon', \hat{\alpha})$ is greater than 0 . We complete the proof by noting that the only dependence on $\hat{\mu}$ in ψ is that ψ_μ is (trivially) increasing in $\hat{\mu}$. Since $d_\mu \geq 0$, increasing $\hat{\mu}$ from ϵ' will only increase the dot product, and thus the dot product remains strictly greater than 0 , implying that $\psi(X, \hat{\mu}, \hat{\alpha}) \neq 0$ as desired. \square

3.5 Proofs of Lemmas 3.11 and 3.12

The main purpose of this section is to present and motivate the proof of Lemma 3.11—since our results are tight across such a wide parameter space, the resulting inequalities are somewhat subtle. After, we also present the short proof of Lemma 3.12.

Lemma 3.11. *Consider an arbitrary distribution D with mean 0 and variance 1. There exists a universal constant c where the following claim is true. Fixing $\hat{\mu} = \epsilon' = \left(1 + \frac{c \log \log \frac{1}{\delta}}{\log \frac{1}{\delta}}\right) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$, then for all δ smaller than some universal constant, and for all $\hat{\nu} \in [0.05, 55.5]$, there exists a vector $\mathbf{d}(\hat{\nu})$ where $d_\mu \geq 0$, and both $\sqrt{\frac{n}{\log(1/\delta)}} |d_\mu|, |d_\alpha|$ are bounded by a universal constant, such that*

$$\mathbb{P}_{X \leftarrow D^n} \left(\mathbf{d}(\hat{\nu}) \cdot \psi \left(X, \hat{\mu} = \epsilon', \hat{\alpha} = \frac{\log(1/\delta)}{3n\hat{\nu}} \right) > \frac{1}{\log \frac{1}{\delta}} \right) \geq 1 - \frac{\delta}{\log^4 \frac{1}{\delta}} \quad (3.4)$$

Furthermore, for $\hat{\nu} = 0.05$ we have $d_\mu = \sqrt{3.75 \frac{\log(1/\delta)}{n}}$, $d_\alpha = \sqrt{3}$; and for $\hat{\nu} = 55.5$ we have $d_\mu = 0$, $d_\alpha < 0$.

We start the analysis via standard Chernoff bounds on the complement of the probability in Equation 3.4 via Lemma 3.13, before pausing to discuss how mathematical programming and duality insights lead to the formulation of the crucial Lemma 3.14; we then complete the proof.

Lemma 3.13. *Consider an arbitrary distribution D with mean 0 and variance 1. For all sufficiently small δ , for any $\hat{\mu}, \hat{\alpha}$ and vector $\mathbf{d} = (d_\mu, d_\alpha)$, we have*

$$\mathbb{P}_{X \leftarrow D^n} \left(\mathbf{d} \cdot \psi(X, \hat{\mu}, \hat{\alpha}) \leq \frac{1}{\log \frac{1}{\delta}} \right) \leq 2 \left(e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \mathbb{E}_{x \leftarrow D} (e^{d_\mu x(1 - \min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)} \right)^n$$

Proof. We upper-bound the probability by exponentiating the negation of both sides of the expression inside the probability, and then using Markov's inequality:

$$\begin{aligned} & \mathbb{P}_{X \leftarrow D^n} \left(\mathbf{d}(\hat{\nu}) \cdot \psi(X, \hat{\mu}, \hat{\alpha}) \leq \frac{1}{\log \frac{1}{\delta}} \right) \\ &= \mathbb{P}_{X \leftarrow D^n} \left(e^{-\mathbf{d}(\hat{\nu}) \cdot \psi(X, \hat{\mu}, \hat{\alpha})} \geq e^{-\frac{1}{\log \frac{1}{\delta}}} \right) \\ &\leq 2 \mathbb{E}_{X \leftarrow D^n} \left(e^{-\mathbf{d}(\hat{\nu}) \cdot \psi(X, \hat{\mu}, \hat{\alpha})} \right) \quad \text{by Markov's inequality; and } e^{\frac{1}{\log(1/\delta)}} \leq 2 \text{ for sufficiently small } \delta \\ &= 2 \mathbb{E}_{x \leftarrow D} (e^{-\mathbf{d}(\hat{\nu}) \cdot \psi(x, \hat{\mu}, \hat{\alpha})})^n \quad \text{by independence} \end{aligned}$$

$$= 2 \left(e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \mathbb{E}_{x \leftarrow D} (e^{d_\mu x(1 - \min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)}) \right)^n \quad \text{substituting the def. of } \psi \quad (3.5)$$

□

3.5.1 Mathematical Programming and Duality Analysis

In order to show Lemma 3.11, we aim to find bounds on the failure probability that are as strong as possible. Appealing to Lemma 3.13 that we have just proven, recall that, as in the standard Chernoff bound methodology, we are still free to choose the parameters d_μ, d_α , which we do so as to minimize the resulting bound on the failure probability. Phrased abstractly, the goal is, for the $\hat{\mu}, \hat{\alpha}$ of Lemma 3.11, to show that, for any distribution D of mean 0 and variance 1, there is a choice $\mathbf{d} = (d_\mu, d_\alpha)$ that makes Equation 3.5 sufficiently small. Phrased as an optimization problem, our goal is to evaluate (or tightly bound):

$$\max_D \min_{\mathbf{d}=(d_\mu, d_\alpha)} e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \mathbb{E}_{x \leftarrow D} (e^{d_\mu x(1 - \min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)}) \quad (3.6)$$

where D ranges over distributions of mean 0 and variance 1.

We will use convex-concave programming and linear programming duality to significantly simplify the max-min program in Equation 3.6 before we dive into the part of analysis that is ad hoc for this problem. We wish to emphasize here again that the steps of 1) writing an estimator as a multi-parameter ψ -estimator and finding an analogous lemma to our Lemma 3.11, then 2) using mathematical programming duality to simplify the Chernoff bound analysis, are a framework generalizable for tightly analyzing other estimators.

For simplicity of exposition, assume that we restrict the support of D to some sufficiently fine-grained finite set, meaning that the maximization in Equation 3.6 is now finite-dimensional, albeit an arbitrarily large finite number. For each support element x , let D_x be a variable representing the probability of choosing x under distribution D . The expectation component of Equation 3.6 may now be expressed as sum that is a linear function in the variables D_x :

$$\max_D \min_{\mathbf{d}=(d_\mu, d_\alpha)} e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \sum_x D_x \cdot e^{d_\mu x(1 - \min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)} \quad (3.7)$$

Using the standard max-min inequality (a form of weak duality in optimization), we have that Equation 3.7 is upper bounded by swapping the maximization and minimization (Equation 3.8), meaning that the vector \mathbf{d} no longer depends on the distribution D .

$$\min_{\mathbf{d}=(d_\mu, d_\alpha)} \max_D e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \sum_x D_x \cdot e^{d_\mu x(1-\min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)} \quad (3.8)$$

Crucially, however, Equation 3.8 is not just an upper bound on Equation 3.7, but is in fact *equal* to it, due to Sion's minimax theorem [59]. To apply Sion's minimax theorem, it suffices to check that 1) both \mathbf{d} and D are constrained to be in convex sets, at least one of which is compact, 2) the objective is convex in \mathbf{d} and 3) concave in the variables D_x . For the first condition, we note that the set of distributions on a finite domain is compact. The objective is convex in \mathbf{d} since the objective is the sum of exponentials that are each linear in \mathbf{d} . And the objective is concave in D_x because it is in fact a linear function of D .

The guarantee of Sion's minimax theorem means that we may work with Equation 3.8 instead of Equation 3.7 without sacrificing tightness in our analysis. This justifies why we are free to choose $\mathbf{d} = (d_\mu, d_\alpha)$ in Lemma 3.11 that does not depend on the distribution D .

To further simplify the problem in Equation 3.8, we note again that both the objective and the constraints on D are linear in the variables D_x , meaning that the inner maximization is in fact a linear program. We can then apply linear programming (strong) duality to yield the following equivalent optimization (Equation 3.9). We note that, as above, for the purposes of upper bounding Equation 3.6, it suffices to only use weak duality. Strong duality however guarantees that this step does not introduce slack into the analysis.

The three variables V, M, S in the inner minimization below are the dual variables corresponding to the three constraints on distribution D originally: that D has variance 1, mean 0, and total probability mass 1.

$$\min_{\mathbf{d}=(d_\mu, d_\alpha)} \min_{V, M, S} V + S \quad (3.9)$$

for all x : $Vx^2 + Mx + S \geq e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta} + d_\mu x(1-\min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)}$

We have thus reduced the infinite-dimensional optimization problem of Equation 3.6 to the five-dimensional problem of Equation 3.9 (or six dimensions, if we include the universal quantification for $x \in \mathbb{R}$), a significant simplification. We bound Equation 3.9 by explicitly choosing values for $\mathbf{d} = (d_\mu, d_\alpha), V, M, S$ as functions of $\hat{\alpha}, n, \log \frac{1}{\delta}$, and showing

that they jointly satisfy the constraint of Equation 3.9, for all x . We factor out the terms in the exponential that do not depend on x ; we make the variable substitutions $y \equiv \sqrt{\hat{\alpha}}x$ and $\hat{\nu} \equiv \frac{\log(1/\delta)}{3n\hat{\alpha}}$ to replace dependence on $\hat{\alpha}, n, \log \frac{1}{\delta}$ with dependence on the single variable $\hat{\nu}$; taking the log of both sides (and swapping sides) yields an expression that is recognizable in the following lemma, where the multipliers of $1, y, y^2$ respectively on the right hand side are essentially our choices of S, M, V :

Lemma 3.14. *For all $\hat{\nu} \in [0.05, 55.5]$, there exist $a > 0$ and b such that*

$$\forall y \in \mathbb{R} : ay(1 - \min(y^2, 1)) - b \cdot \min(y^2, 1) \leq \log \left(1 + ay + y^2 \hat{\nu} \left(-3 + \frac{a\sqrt{6}}{\sqrt{\hat{\nu}}} - b \right) \right)$$

where $a \in [C, C']$ and $b \in [-C', C']$ for positive constants C, C' . Further, for $\hat{\nu} = 0.05$, the pair $a = 0.75, b = \sqrt{3}$ works.

We emphasize that the application of Lemma 3.14 in the proof of Lemma 3.11 below is straightforward, though finding the particular form of Lemma 3.14 is not. Further, one would not seek a result of the form of Lemma 3.14 without the guarantees of this section, derived via duality and mathematical programming, showing that “results of the form of Lemma 3.14 encompass the full power of the Chernoff bounds of Equation 3.5.” See the end of Section 3.5.2 for the proof of Lemma 3.14.

3.5.2 Proof of Lemma 3.11

We now prove Lemma 3.11 by combining the Chernoff bound analysis of Lemma 3.13 with the inequality from Lemma 3.14. We point out that the proof below is direct, without any reference to duality or mathematical programming; however, the discussion of Section 3.5.1 was crucial to discovering the right formulation for Lemma 3.14. We prove Lemma 3.14 at the end of the section.

Lemma 3.11. *Consider an arbitrary distribution D with mean 0 and variance 1. There exists a universal constant c where the following claim is true. Fixing $\hat{\mu} = \epsilon' = \left(1 + \frac{c \log \log \frac{1}{\delta}}{\log \frac{1}{\delta}} \right) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}$, then for all δ smaller than some universal constant, and for all $\hat{\nu} \in [0.05, 55.5]$, there exists a vector $\mathbf{d}(\hat{\nu})$ where $d_{\mu} \geq 0$, and both $\sqrt{\frac{n}{\log(1/\delta)}} |d_{\mu}|, |d_{\alpha}|$ are bounded by a universal constant, such that*

$$\mathbb{P}_{X \leftarrow D^n} \left(\mathbf{d}(\hat{\nu}) \cdot \psi \left(X, \hat{\mu} = \epsilon', \hat{\alpha} = \frac{\log(1/\delta)}{3n\hat{\nu}} \right) > \frac{1}{\log \frac{1}{\delta}} \right) \geq 1 - \frac{\delta}{\log^4 \frac{1}{\delta}}$$

Furthermore, for $\hat{\nu} = 0.05$ we have $d_\mu = \sqrt{3.75 \frac{\log(1/\delta)}{n}}$, $d_\alpha = \sqrt{3}$; and for $\hat{\nu} = 55.5$ we have $d_\mu = 0$, $d_\alpha < 0$.

Proof. Start with the bound on the probability of failure given by Lemma 3.13:

$$2 \left(e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \mathbb{E}_{x \leftarrow D} \left(e^{d_\mu x(1 - \min(\hat{\alpha}x^2, 1)) - d_\alpha \min(\hat{\alpha}x^2, 1)} \right) \right)^n$$

For $\hat{\nu} \in [0.05, 55.5)$ we bound the exponential inside the expectation via the exponential of Lemma 3.14; we also use Lemma 3.14 to choose d_μ, d_α for us (the $\hat{\nu} = 55.5$ case is covered at the end). Namely, in Lemma 3.14 use $\hat{\nu}$ as given, substitute $x \sqrt{\hat{\alpha}} \equiv y$ (where $\hat{\alpha} \equiv \frac{\log(1/\delta)}{3n\hat{\nu}}$ as always), and choose $d_\mu \equiv a \sqrt{\hat{\alpha}}$, and $d_\alpha \equiv b$ —in particular, for $\hat{\nu} = 0.05$ this gives $d_\mu(0.05) = 0.75 \sqrt{\hat{\alpha}} = 0.75 \sqrt{\frac{\log(1/\delta)}{3n\hat{\nu}}} = \sqrt{3.75 \frac{\log(1/\delta)}{n}}$. Thus the failure probability is bounded by

$$\begin{aligned} & 2 \left(e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \mathbb{E}_{\substack{x \leftarrow D \\ y = x \sqrt{\hat{\alpha}}}} \left(1 + ay + y^2 \hat{\nu} \left(-3 + \frac{a \sqrt{6}}{\sqrt{\hat{\nu}}} - b \right) \right) \right)^n \\ &= 2 \left(e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta}} \left(1 + \frac{\log \frac{1}{\delta}}{3n} \left(-3 + 3d_\mu \sqrt{\frac{2n}{\log(1/\delta)}} - d_\alpha \right) \right) \right)^n \quad (\text{mean 0, variance 1}) \\ &\leq 2 \left(e^{-d_\mu \hat{\mu} + d_\alpha \frac{1}{3n} \log \frac{1}{\delta} + \frac{\log \frac{1}{\delta}}{3n} \left(-3 + 3d_\mu \sqrt{\frac{2n}{\log(1/\delta)}} - d_\alpha \right)} \right)^n \quad \text{since } 1 + z \leq e^z \text{ for any } z \\ &\leq 2e^{-d_\mu \sqrt{\frac{2n}{\log(1/\delta)}} c \log \log \frac{1}{\delta} - \log \frac{1}{\delta}} \quad \text{substituting } \hat{\mu} = \epsilon' = \left(1 + \frac{c \log \log \frac{1}{\delta}}{\log \frac{1}{\delta}} \right) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}} \\ &\leq \frac{\delta}{\log^4 \frac{1}{\delta}} \end{aligned}$$

where the last inequality holds for large enough c , since $d_\mu \sqrt{\frac{n}{\log(1/\delta)}} = \frac{a}{\sqrt{3\hat{\nu}}}$ is greater than some positive constant.

We prove the $\hat{\nu} = 55.5$ case now. We choose $d_\mu = 0$ and $d_\alpha = -4$, substituting into the bound of Equation 3.5 to yield

$$\begin{aligned} 2 \left(e^{-\frac{4}{3n} \log \frac{1}{\delta}} \mathbb{E}_{x \leftarrow D} \left(e^{4 \min(\hat{\alpha}x^2, 1)} \right) \right)^n &\leq 2\delta^{4/3} \mathbb{E}_{x \leftarrow D} \left(1 + 54\hat{\alpha}x^2 \right)^n \quad \text{for } y \in [0, 1], e^{4y} \leq 1 + 54y \\ &= 2\delta^{4/3} (1 + 55.5\hat{\alpha})^n \quad \text{since } D \text{ has variance 1} \end{aligned}$$

$$\begin{aligned} &\leq 2\delta^{4/3} e^{n \cdot 54 \frac{\log(1/\delta)}{3.55.5n}} \quad 1 + z \leq e^z; \text{ substituting def. of } \hat{a} \\ &= 2\delta^{4/3} \delta^{-\frac{54}{3.55.5}} \leq 2\delta^{1.009} \end{aligned}$$

which is bounded as desired for small enough δ . \square

We now prove Lemma 3.14.

Lemma 3.14. *For all $\hat{v} \in [0.05, 55.5]$, there exist $a > 0$ and b such that*

$$\forall y \in \mathbb{R} : ay(1 - \min(y^2, 1)) - b \cdot \min(y^2, 1) \leq \log\left(1 + ay + y^2 \hat{v} \left(-3 + \frac{a\sqrt{6}}{\sqrt{\hat{v}}} - b\right)\right) \quad (3.10)$$

where $a \in [C, C']$ and $b \in [-C', C']$ for positive constants C, C' . Further, for $\hat{v} = 0.05$, the pair $a = 0.75, b = \sqrt{3}$ works.

Proof. We first prove the special case of 1) $\hat{v} = 0.05$, before moving to the general case of 2) $\hat{v} \in (0.05, 55.5]$. We note that our choice of $a(\hat{v}), b(\hat{v})$ is *not* continuous in \hat{v} at 0.05, but the usage of the lemma does not require any continuity. We choose a, b at the edge case $\hat{v} = 0.05$ for convenience.

1) For $\hat{v} = 0.05$, we choose $a = 0.75, b = \sqrt{3}$. This special case of Equation 3.10 simplifies to:

$$\forall y \in \mathbb{R} : 0.75y(1 - \min(y^2, 1)) - \sqrt{3} \cdot \min(y^2, 1) \leq \log(1 + 0.75y + 0.174y^2)$$

(where 0.174 is a lower bound on $\hat{v}(-3 + \frac{a\sqrt{6}}{\sqrt{\hat{v}}} - b)$). This is a 1-dimensional bound and can be easily analyzed in many ways. For the range $y \in [-1, 1]$: the right hand side is at least $\log(1 + 0.75y)$, which in this range is at least $.75y - .75y^2$, which is easily shown to be greater than the polynomial expression that the left hand side reduces to in this range, $0.75y - \sqrt{3}y^2 - 0.75y^3$. For the remaining range, $y \notin [-1, 1]$, the left hand side is the constant $-\sqrt{3}$, and it is easy to check that the quadratic in the argument of the right hand side, $1 + 0.75y + 0.174y^2$, always exceeds $e^{-\sqrt{3}}$.

2) To show Equation 3.10 for the rest of the range of $\hat{v} \in (0.05, 55.5]$, we choose a to be the positive root of the quadratic equation $\sqrt{\hat{v}}(a^2 - 12) + \sqrt{6}a = 0$ and let $b = 3 - a^2/2$ —we will see the motivation for this choice shortly. For now, note that the definition of a implies $a \leq \sqrt{12}$, for otherwise $\sqrt{\hat{v}}(a^2 - 12) + \sqrt{6}a$ would be greater than 0.

Our proof will analyze the sign of the derivative with respect to y of the difference between the right and left hand sides of Equation 3.10. For the critical region $|y| \leq 1$ this derivative equals:

$$\frac{a + 2y\hat{\nu}(-3 + \frac{a\sqrt{6}}{\sqrt{\hat{\nu}}} - b)}{1 + ay + y^2\hat{\nu}(-3 + \frac{a\sqrt{6}}{\sqrt{\hat{\nu}}} - b)} - a + 3ay^2 + 2by \quad (3.11)$$

The crucial step is to choose a to be the positive root of the quadratic equation $\sqrt{\hat{\nu}}(a^2 - 12) + \sqrt{6}a = 0$ and let $b = 3 - a^2/2$, after which Equation 3.11 miraculously factors as

$$\frac{1}{3a^2} \cdot \frac{y(y + \frac{2}{a})(y + \frac{2}{a} - \frac{a}{3})^2}{y^2 + (\frac{4}{a} - \frac{a}{3})y + (\frac{4}{a^2} - \frac{1}{3})}$$

From this expression for the derivative, it is straightforward to read off its sign. The discriminant of the quadratic in the denominator is $\frac{1}{9}(a^2 - 12) > 0$, meaning the denominator is always positive. The squared term in the numerator cannot affect the overall sign. Thus the sign of the derivative equals the sign of $y(y + \frac{2}{a})$, meaning that the difference between the right and left side of Equation 3.10 is monotonically increasing for $y > 0$, and unimodal for $y < 0$, having non-positive derivative for $y \in [\frac{2}{a}, 0]$ and nonnegative derivative for smaller y . Thus to show the inequality holds for all $y \in [-1, 1]$ it suffices to check it at $y = 0$ and $y = -1$.

The $y = 0$ case is trivial as both sides of Equation 3.10 equal 0.

For $y = -1$, Equation 3.10, after expressing both $\sqrt{\hat{\nu}}$ and b in terms of a becomes

$$\frac{a^2}{2} - 3 \leq \log\left(-2 + a - \frac{36}{a^2 - 12}\right) \quad (3.12)$$

For $a \in [0, \sqrt{12})$, the inverse of the rational expression inside the log is bounded by its linear approximation, $\frac{\sqrt{12-a}}{\sqrt{12}}$. Calling this a new variable $z = \frac{\sqrt{12-a}}{\sqrt{12}}$, which is between 0 and 1, Equation 3.12 becomes the claim that $6(1-z)^2 - 3 \leq -\log z$, which is easily verified for $z \in (0, 1]$.

Lastly, we show Equation 3.10 for $|y| > 1$. Reexpressing b and $\sqrt{\hat{\nu}}$ in terms of a , the left hand side of the inequality is the constant value $-b = -3 + \frac{a^2}{2}$ (independent of y), while the right hand side is $\log(1 + ay + \frac{3a^2}{12-a^2}y^2)$. Analyzing the quadratic inside the log shows that the right hand side has a minimum of $\frac{a^2}{12}$, attained at $y = -\frac{12-a^2}{6a}$.

When the location of this minimum, $y = -\frac{12-a^2}{6a}$, is inside the interval $[-1, 1]$, then because this quadratic is monotonic to either side of the minimum, the fact that we have already proven Equation 3.10 for $y = \pm 1$ implies the inequality holds for all y further from 0.

The remaining case is when the minimum is not in $[-1, 1]$, namely, $-\frac{12-a^2}{6a} < -1$, meaning $a < 1.59$; since a is monotonic in \hat{v} , a is at least its value when $\hat{v} = 0.05$, namely $a \geq 1.003$. Equation 3.10 thus reduces to showing that, for $a \in [1.003, 1.59]$ we have $\frac{a^2}{2} - 3 \leq \log \frac{a^2}{12}$, which is trivially implied, substituting $z = \frac{a^2}{12}$, by the inequality $6z - 3 \leq \log z$ for $z \in [0.083, 0.22]$, yielding the claim. \square

3.5.3 Proof of Lemma 3.12

Lemma 3.12. *Consider an arbitrary set of n samples X . Consider the expressions $\psi_\mu(X, \hat{\mu}, \hat{\alpha})$, $\psi_\alpha(X, \hat{\alpha})$, reparameterized in terms of $\hat{v} \equiv \frac{\log(1/\delta)}{3n\hat{\alpha}}$ in place of $\hat{\alpha}$. Suppose the equation $\psi_\alpha(X, \hat{\alpha}) = 0$ has a solution in the range $\hat{v} \in [0.05, 55.5]$. Then the functions $\sqrt{\frac{\log(1/\delta)}{n}} \psi_\mu(X, \hat{\mu}, \hat{\alpha})$ and $\psi_\alpha(X, \hat{\alpha})$ are Lipschitz with respect to \hat{v} on the entire interval $\hat{v} \in [0.05, 55.5]$, with Lipschitz constant $c \log \frac{1}{\delta}$ for some universal constant c .*

Proof. Consider the \hat{v} derivative of $\psi_\alpha(X, \hat{\mu}, \hat{\alpha} \equiv \frac{\log(1/\delta)}{3n\hat{v}}) = \sum_{i=1}^n \left(\min\left(\frac{\log(1/\delta)}{3n\hat{v}} x_i^2, 1\right) - \frac{1}{3n} \log \frac{1}{\delta} \right)$. The \hat{v} derivative of $\min\left(\frac{\log(1/\delta)}{3n\hat{v}} x_i^2, 1\right)$ is either $-\frac{\log(1/\delta)}{3n\hat{v}^2} x_i^2 = -\frac{1}{\hat{v}} \hat{\alpha} x_i^2$ or 0, depending on which term in the min is the smallest, and in either case has magnitude at most $\frac{1}{\hat{v}} \min(\hat{\alpha} x_i^2, 1)$. Thus the overall \hat{v} derivative of $\psi_\alpha(X, \hat{\mu}, \hat{\alpha})$ has magnitude at most $\frac{1}{\hat{v}} \sum_i \min(\hat{\alpha} x_i^2, 1)$. Since, we are guaranteed that $\sum_{i=1}^n \min(\hat{\alpha} x_i^2, 1) = \frac{1}{3} \log \frac{1}{\delta}$ for some $\hat{v} \in [0.05, 55.5]$, we thus have that the derivative is within a constant factor of this across the entire range, as desired.

Similarly, consider the \hat{v} derivative of $\psi_\mu(X, \hat{\mu}, \hat{\alpha}) = \sum_{i=1}^n \left(\hat{\mu} - x_i \left(1 - \min(\hat{\alpha} x_i^2, 1) \right) \right)$. The i^{th} term of this is the \hat{v} derivative of $\min(\hat{\alpha} x_i^2, 1)$, which is either $-\frac{1}{\hat{v}} \hat{\alpha} x_i^2$ or 0 depending on whether $x_i \leq \sqrt{1/\hat{\alpha}}$, and thus the magnitude of this derivative may be bounded by $\frac{1}{\hat{v} \sqrt{\hat{\alpha}}} \sum_{i=1}^n \min(\hat{\alpha} x_i^2, 1)$. Since $\sum_{i=1}^n \min(\hat{\alpha} x_i^2, 1)$ is bounded by a constant times $\log \frac{1}{\delta}$ (as in the last paragraph), and $\frac{1}{\hat{v} \sqrt{\hat{\alpha}}}$ is bounded by a constant times $\frac{1}{\sqrt{\hat{v} \hat{\alpha}}} = \sqrt{\frac{3n}{\log(1/\delta)}}$, the magnitude of the derivative of $\sqrt{\frac{\log(1/\delta)}{n}} \psi_\mu(X, \hat{\mu}, \hat{\alpha})$ is bounded by a constant times $\log \frac{1}{\delta}$, as desired. \square

Chapter 4

Uncertainty about Uncertainty: Optimal Adaptive Algorithms for Estimating Mixtures of Unknown Coins

4.1 Overview

We consider a natural statistical estimation task, motivated by a practical setting, with an intriguing adaptive flavor. We provide a new adaptive algorithm and a matching fully adaptive lower bound, tight up to multiplicative constants.

In our problem setting, there is a universe of coins of two types: positive coins each have a (potentially different) probability of heads that lies in the interval $[\frac{1}{2} + \Delta, 1]$, while negative coins lie in the interval $[0, \frac{1}{2} - \Delta]$, where $\Delta \in (0, \frac{1}{2}]$ parameterizes the “quality” of the coins. Our only access to the coins is by choosing a coin and then flipping it, without access to the true biases of the coins. An algorithm in this setting may employ arbitrary adaptivity—for example, flipping three different coins in sequence and then flipping the first coin 5 more times if and only if the results of the first 3 flips were heads, tails, heads. The challenge is to estimate the *fraction* ρ of coins that are of positive type, to within a given

additive error ϵ , using as few coin flips (samples) as possible. We assume because of the symmetry of the problem (between positive and negative coins) that $\rho \leq \frac{1}{2}$.

This model arose from a collaboration with colleagues in data science and database systems, about harnessing paid crowdsourced workers to estimate the “quality” of a database. Our model is a direct theoretical analog of the following practical problem, where sample complexity linearly translates into the amount of money that must be paid to workers, and thus even multiplicative factors crucially affect the usefulness of an algorithm. Given a set of data and a predicate on the data, the task is to estimate what fraction of the data satisfies the predicate—for example, estimating the proportion of records in a large database that contain erroneous data. After automated tools have labeled whatever portion of the data they are capable of dealing with, the remaining data must be processed via *crowdsourcing*, an emerging setting that potentially offers sophisticated capabilities but at the cost of unreliability. Namely, for each data item, one may ask many human users/workers online whether they think the item satisfies the predicate, with the caveat that the answers returned could be noisy. In the case that the workers have no ability to distinguish the predicate, we cannot hope to succeed; however, if the histograms of detection probabilities for positive versus negative data have a gap between them (the gap is 2Δ in the model above), then the challenge is to estimate ρ as accurately as possible, from a limited budget of queries to workers [21].

A key feature that makes this estimation problem distinct from many others studied in the literature is the richness of adaptivity available to the algorithm. Achieving a tight lower bound in this setting requires considering and bounding all possible uses of adaptivity available to an algorithm; and achieving an optimal algorithm requires choosing the appropriate adaptive information flow between different parts of the algorithm. Much of the previous work in the area of statistical estimation is focused on non-adaptive algorithms and lower bounds; however see [10], and in particular, Sections 4.1 and 4.2 of that work, for a survey of several distribution testing models that allow for adaptivity. In our setting there are two distinct kinds of adaptivity that an algorithm can leverage: 1) single-coin adaptivity, deciding how many times a particular coin should be flipped—a per-coin stopping rule—in terms of the results of its previous flips, and 2) cross-coin adaptivity,

deciding which coin to flip next in terms of the results of previous flips across *all* coins. Our final optimal algorithm (Section 4.3) leverages both kinds of adaptivity. In our tight lower bound analysis (Section 4.5), we overcome the technical obstacles presented by the richness of adaptivity by giving a reduction (Section 4.5.1) from fully-adaptive algorithms that leverage both kinds of adaptivity to single-coin adaptive algorithms that process each coin independently, valid for our specific lower bound instance. We discuss the approaches and challenges of our lower bound in more detail in Section 4.1.1.

The main *algorithmic* challenge in this problem is what we call “uncertainty about uncertainty”: we make no assumptions about the quality of the coins beyond the existence of a gap 2Δ between biases of the coins of different types (centered at $\frac{1}{2}$). If we relaxed the problem, and assumed (perhaps unrealistically) that we know 1) the conditional distribution of biases of positive coins, and 2) the same for negative coins, and 3) an initial estimate of the mixture parameter ρ between the two distributions, then we show that it is easy—using mathematical programming techniques in Section 4.8.1—to construct an estimation algorithm with sample complexity that is optimal *by construction* up to a multiplicative constant (see Section 4.8.2). On the other hand, our algorithm for the original setting has to return estimates with small bias, and be sample efficient at the same time, regardless of the bias of the coins, be they all deterministic, or all maximally noisy as allowed by the Δ parameter, or some quality in between. While intuitively the hardest settings to distinguish information theoretically involve coins with biases as close to each other as possible (and indeed our lower bound relies on mixtures of only $\frac{1}{2} \pm \Delta$ coins), settings with biases near but not equal to $\frac{1}{2} \pm \Delta$ introduce “uncertainty about uncertainty” challenges. The two kinds of adaptivity available to the algorithm allow us to meet these challenges by trading off, optimally, between 1) investigating a single coin to reduce uncertainty about its bias, and 2) apportioning resources between different coins to reduce uncertainty about the ground truth fraction ρ , which is the objective of the problem.

4.1.1 Our Approaches and Results

To motivate the new algorithms of this paper, we start by describing the straightforward analysis of perhaps the most natural approach to the problem, which is non-adaptive,

based on subsampling.

Example 4.1. Recall that it takes $\Omega(\frac{1}{\Delta^2})$ samples to distinguish a coin of bias $\frac{1}{2} - \Delta$ from a coin of bias $\frac{1}{2} + \Delta$. We can therefore imagine an algorithm that chooses a random subset of the coins, and flips each coin $\Omega(\frac{1}{\Delta^2})$ many times. Asking for $\Theta(\frac{1}{\Delta^2} \log \frac{1}{\epsilon})$ flips from each coin guarantees that all but ϵ fraction of the coins in the subset will be accurately classified. Given an accurate classification of m randomly chosen coins, we use the fraction of these that appear positive as an estimate on the overall mixture parameter ρ . Estimating ρ to within error ϵ requires $m = O(\frac{\rho}{\epsilon^2})$ randomly chosen coins. Overall, taking $\Theta(\frac{1}{\Delta^2} \log \frac{1}{\epsilon})$ samples from each of $m = \Theta(\frac{\rho}{\epsilon^2})$ coins uses $\Theta(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\epsilon})$ samples.

As we will see, the above straightforward algorithm is potentially wasteful in samples by up to a $\log \frac{1}{\epsilon}$ factor, since it makes $\Theta(\frac{1}{\Delta^2} \log \frac{1}{\epsilon})$ flips for every single coin, yet—since $\Omega(\frac{1}{\Delta^2})$ samples suffices to label a coin with constant accuracy—each sample beyond the first $\Theta(\frac{1}{\Delta^2})$ samples from a single coin gives increasing certainty yet diminishing information-per-coin. If we can save on this $\log \frac{1}{\epsilon}$ factor without sacrificing impractical constants, then our approach leads to significant practical savings in samples, and thus monetary cost—in regimes, such as crowdsourcing, where gathering data is by far the most expensive part of the estimation process.

Algorithmic Construction

We give two algorithmic constructions. Algorithm 3, which we call the Triangular Walk algorithm, is single-coin adaptive, and is theoretically almost-tight in sample complexity. Second, Algorithm 6 has the optimal sample complexity, by combining the Triangular Walk algorithm with a new (and surprisingly) non-adaptive component (Algorithm 4).

The Triangular Walk algorithm (Algorithm 3) is designed for the specific *practical* parameter regime where ρ is small: in our earlier crowdsourcing example, practitioners typically preprocess data items by using automated techniques and heuristics to classify a majority of the items, before leaving to crowdsourced workers a small number of items that cannot be automatically classified. These automated filtering techniques usually flag significantly more “negative” items than “positive” items as “unclassifiable automatically”, resulting in a

small fraction ρ of positive items among the ones selected for crowdsourced human classification. The intuition behind our approach, then, is to try to abandon sampling (frequent) negative coins as soon as possible, after $\Theta(\frac{1}{\Delta^2})$ samples, while being willing to investigate (infrequent) positive coins up to depth $\Theta(\frac{1}{\Delta^2} \log \frac{1}{\epsilon})$. Thus we disproportionately bias our investment of resources towards the rare and valuable regime. Using techniques from random walk theory, we design a linear estimator based on this behavior (Algorithm 2), whose expectation across many coins yields a robust estimator, Algorithm 3, as shown in Theorem 4.2 (restated and proved in Section 4.2).

Theorem 4.2. *Given coins where a ρ fraction of the coins have bias $\geq \frac{1}{2} + \Delta$, and $1 - \rho$ fraction have bias $\leq \frac{1}{2} - \Delta$, then running Algorithm 3 on $t = \Theta(\frac{\rho}{\epsilon^2} \log \frac{1}{\delta})$ randomly chosen coins will estimate ρ to within an additive error of $\pm\epsilon$, with probability at least $1 - \delta$, with an expected sample complexity of $O(\frac{\rho}{\epsilon^2 \Delta^2} (1 + \rho \log \frac{1}{\epsilon}) \log \frac{1}{\delta})$.*

The analysis of Algorithm 3 uses only standard concentration inequalities, and thus the big-O notation for the sample complexity does not hide large constants. As further evidence of the good practical performance of Algorithm 3, Section 4.7 shows simulation-based experimental results, run on settings with practical problem parameters for crowdsourcing applications. These results demonstrate the advantages of our algorithm as compared with the straightforward majority vote algorithm as well as the state-of-the-art algorithm [21] (which does not enjoy any theoretical guarantees).

As for our second, optimal, algorithmic construction (Algorithm 6 in Section 4.3), we combine the adaptive techniques from the Triangular Walk algorithm with a non-adaptive estimation component. More concretely, in the regimes where Algorithm 3 is not optimal, Algorithm 6 uses Algorithm 3 to first give a 2-approximation of ρ , before using this information to non-adaptively estimate ρ much more accurately, while keeping the variance of the estimate small, to control the sample complexity. The theoretical guarantees of Algorithm 6 are shown in Theorem 4.3 (restated and proved in Section 4.3).

Theorem 4.3 (Informal). *Given coins where a ρ fraction of the coins have bias $\geq \frac{1}{2} + \Delta$, and $1 - \rho$ fraction have bias $\leq \frac{1}{2} - \Delta$, then for large enough constant c , running Algorithm 6 on a budget of $B \geq c \frac{\rho}{\Delta^2 \epsilon^2}$ coin flips will estimate ρ to within an additive error of $\pm\epsilon$, with probability at least $2/3$. If*

the algorithm is repeated $\Theta(\log \frac{1}{\delta})$ times, and the median estimate is returned, then the probability of failure is at most δ .

Lower Bounds and Discussion

Complementary to our algorithm, we show a matching lower bound of $\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ samples for a success probability of $1 - \delta$ for the problem. Crucially, our bounds match across choices of all four parameters, $\rho, \epsilon, \Delta, \delta$. To show the lower bound, we use the following setup: consider a scenario where all positive coins have bias *exactly* $\frac{1}{2} + \Delta$ and all negative coins have bias *exactly* $\frac{1}{2} - \Delta$.

The overall intuition for our lower bound is that, for each coin, even flipping it enough to learn whether it is a positive or negative coin will tell us little about whether the true fraction of positive coins is ρ versus $\rho + \epsilon$, and thus the flow of information to our algorithm is at most a slow trickle. To capture this intuition, we aim to decompose the analysis into a sum of coin-by-coin bounds; however, the key challenge is the *cross-coin adaptivity* that is available to the algorithm.

To demonstrate the challenge of tightly analyzing cross-coin adaptivity, consider the following natural attempt at a lower bound.

1. Consider flipping a fair coin S to choose between a universe with ρ fraction of positive coins, versus $\rho + \epsilon$ fraction.
2. The aim is to bound the amount of mutual information that the entire transcript of an adaptive coin-flipping algorithm can have with the coin S .
3. Suppose this mutual information can be bounded by the mutual information of the sub-transcript of the i^{th} coin with S , summed over all i .
4. Thus consider and bound the amount of mutual information between the sub-transcript of just coin i alone, with S ; and sum these bounds over all coins at the end.

While one would intuitively expect the bounds of Step 4 to be small for each single coin, cross-coin adaptivity allows for each single-coin sub-transcript to encode a lot of mutual

information via its *length*, which may be adaptively chosen by the algorithm in light of information gathered across all other coins. The amount of mutual information about S in a sub-transcript may be linear in the number of times *other* coins have been flipped, implying that summing up such mutual information across all coins would yield a bound that uselessly grows quadratically with the number of flips, instead of linearly.

Our approach: We show that no fully-adaptive algorithm can distinguish the following two scenarios with probability at least $1 - \delta$, using $o(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ samples: 1) when a ρ fraction of the coins are positive, and 2) when a $\rho + \epsilon$ fraction of the coins are positive. This is formalized as the following theorem (Theorem 4.4), and proved in Section 4.5.

Theorem 4.4. *For $\rho \in [0, \frac{1}{2})$ and $\epsilon \in (0, 1 - 2\rho]$, the following two situations are impossible to distinguish with at least $1 - \delta$ probability using an expected $o(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ samples: A) ρ fraction of the coins have probability $\frac{1}{2} + \Delta$ of landing heads and $1 - \rho$ fraction of the coins have probability $\frac{1}{2} - \Delta$ of landing heads, versus B) $\rho + \epsilon$ fraction of the coins have probability $\frac{1}{2} + \Delta$ of landing heads and $1 - (\rho + \epsilon)$ fraction of the coins have probability $\frac{1}{2} - \Delta$ of landing heads. This impossibility crucially includes fully-adaptive algorithms.*

In Section 4.5.1, we capture rather generally via Lemmas 4.19 and 4.20 the above intuitive decomposition of a many-coin adaptive algorithm into its single-coin contributions, but via a careful simulation argument that precludes the kind of information leakage between coins that we described above. More explicitly, instead of decomposing a single transcript into many (possibly correlated) sub-transcripts, we relate an n -coin transcript to n *separate* runs of the algorithm (each on freshly drawn random coins), where in the i^{th} run, coin i is authentically sampled (from either the ρ scenario or the $\rho + \epsilon$ scenario), while all the remaining coins are simulated by the algorithm. Crucially, since the remaining simulated coins do not depend on the “real” scenario, no cross-coin adaptivity can leak any information about the real world to coin i , beyond the information gained from flipping coin i itself.

Furthermore, Lemmas 4.19 and 4.20 apply to a broad variety of problem settings, where the population of random variables can be arbitrary and not necessarily Bernoulli coins. We believe these lemmas are of independent interest beyond this work, and can be a useful

tool for proving lower bounds for other problem settings, for example a Gaussian variant of the current problem, where instead of being input a noisy yes/no answer on the positivity of an item, we instead receive a numerical Gaussian-distributed score with mean, say, > 1 for positive items and < 0 for negative items.

Given the decomposition lemmas (Lemmas 4.19 and 4.20), completing the lower bound analysis for the current problem requires upper bounding the squared Hellinger distance between running any single-coin adaptive algorithm on the two coin populations described earlier, with slightly different positive-to-negative mixture ratios. This forms the bulk (and technical parts) of the proof of Theorem 4.4.

Non-adaptive bounds: As motivation for the algorithmic results of this paper, it is reasonable to ask, given Theorem 4.4's lower bound of $\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ on the number of samples for our problem, is it possible that a *non-adaptive* algorithm can approach this performance, or is the adaptive flavor of Algorithms 3 or 6 required? We briefly describe how the framework of the "natural attempt" (the numbered list above) in fact yields a lower bound for *non-adaptive* algorithms that is a $\log \frac{1}{\rho}$ factor higher than that of Theorem 4.4, namely $\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\rho})$, when $\rho \geq \epsilon^2$

Given a random variable S that uniformly chooses between scenarios " ρ " and " $\rho + \epsilon$ " respectively, and a sample of size n from a coin that has bias $\frac{1}{2} + \Delta$ with probability ρ or $\rho + \epsilon$ respectively, and bias $\frac{1}{2} - \Delta$ otherwise, what is the mutual information between the n observed flips (from a single coin) and the scenario variable S ? A non-adaptive algorithm must fix the number of queries n independent of the observed outcomes from the coins, where the information about S is the sum received from sampling each coin. Thus the optimal such algorithm chooses n that maximizes the mutual information per sample. Estimates of this mutual information in the relevant cases are not too difficult, as this is the mutual information between a univariate distribution that is the mixture of two binomials, with a fair coin that determines the mixture probabilities. In terms of Δ , and $\rho \geq \epsilon^2$, some calculation shows that the optimal value of n is $\Theta(\frac{1}{\Delta^2} \log \frac{1}{\rho})$, which yields the above-claimed non-adaptive lower bound on sample complexity of $\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\rho})$. See Section 4.9 for the complete calculations.

For the constant- ρ (and constant probability) regime, this lower bound is in fact tight.

	Upper Bound	Lower Bound
Adaptive	$O(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ (Algorithm 6)	$\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ (Section 4.5)
Non-adaptive	$O(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\epsilon} \log \frac{1}{\delta})$ (Trivial, Example 4.1) $O(\frac{1}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ (Algorithm 4 for $\rho = \Theta(1)$)	$\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\rho})$ (For $\rho \geq \epsilon^2$ and constant δ)

Table 4.1: Sample Complexity Upper and Lower Bounds

A major component of our final algorithm, Algorithm 4, when run on a single constant quality ($\Delta = \Omega(1)$) coin with the parameter $f = f_0(p)$ as defined in Definition 4.12, is a non-adaptive unbiased estimator for the indicator function of the positivity of the coin, with small variance and constant sample complexity. For a low quality coin ($\Delta = o(1)$), we can simulate a flip of a constant quality coin by taking the majority result of $\Theta(1/\Delta^2)$ low quality coin flips. Returning the mean of $O(\frac{1}{\epsilon^2})$ repetitions of Algorithm 4 on different coins yields an ϵ -accurate estimate of ρ . The total sample complexity is $O(\frac{1}{\epsilon^2 \Delta^2})$, which matches the non-adaptive lower bound in the constant- ρ regime.

In summary, we have the adaptive and non-adaptive bounds in Table 4.1. As shown in Table 4.1, the non-adaptive bounds match each other and the adaptive bounds only in the regime where $\rho = \Theta(1)$ (and in the trivial $\epsilon = \Theta(1)$ regime). In the non-constant ρ regime, the non-adaptive lower bound is asymptotically larger than the adaptive lower bound, demonstrating the need for adaptivity in the design of our final optimal algorithm.

Practical Considerations

The keen-eyed reader might notice that the algorithmic results in Theorems 4.2 and 4.3 both depend on the unknown ground truth ρ , so thus these bounds are not immediately invocable by a user. We present two approaches to address this issue.

The first approach is to note that Algorithm 3 can be interpreted as an *anytime* algorithm: it can produce an estimate at any point in its execution. As more coins are used in Algorithm 3, the estimate simply gains accuracy. Section 4.2.1 discusses this approach in more detail, and our experiments in Section 4.7 are also run using the same approach. Because of its simplicity, we recommend this method in practice.

A complication arising from this approach is the fact the sample complexity bound of Theorem 4.2 is an *expected* sample complexity bound. Thus there are potential issues introduced by abruptly stopping the algorithm after a fixed budget of samples, which might inadvertently introduce bias to the estimate. Section 4.2.1 also shows how to analyze and address this issue.

The second, theoretically more interesting approach is to fix a total budget of allowable coin flips, and have the algorithm “discover” the optimal achievable accuracy ϵ just from interacting with the different coins. Our presentation and analysis of Algorithm 6, in Section 4.3, follows this approach. We point out that Algorithm 3 can also be made to have this theoretical guarantee, as demonstrated by the invocation of Algorithm 3 in Algorithm 6.

Designing Optimal Estimators when Coin Biases are Known

By contrast with the above results that analyze the “uncertainty about uncertainty” regime with unknown populations of coins, we shed light on the algorithmic challenges of that regime by providing a tight analysis in the case where knowledge about the populations of coins can be leveraged by the algorithm. In particular, we give a bootstrapping approach which takes some initial guess of ρ along with knowledge of the coin population, and produces an optimal-by-construction estimator that can be used improve on the initial estimate. Explicitly, consider the regime where we know 1) the distribution of coin biases conditioned on being a positive coin, 2) analogously for negative coins and 3) also ρ itself, for bootstrapping purposes even though in practice we would only have a guess. Suppose further that we are given 4) the constraint that we will invest at most n_{\max} flips on a single coin, controlling both the sample complexity but also the computational complexity we can afford to compute the optimal estimator. In Section 4.8.1, we use quadratic and linear programming techniques to find a single-coin adaptive algorithm, taking the form of a linear estimator, with the minimum variance possible subject to the constraint that, even if our knowledge of ρ is wrong, the estimator is still unbiased. This construction yields the following theorem.

Theorem 4.5. *Suppose we are given 1) the distribution of coin biases conditioned on being a positive coin, 2) the analogous distribution for negative coins and 3) the mixture parameter ρ (which, again,*

is a circular assumption but useful for a bootstrapping approach). Suppose further that we are given 4) the parameter n_{\max} , which controls the maximum depth of the triangular walk.

Then, following the method described in Section 4.8.1, we can find the linear estimator for ρ that minimizes variance, subject to a) the expected output of the estimator on input a random positive coin is 1 and b) the analogous expected output for a random negative coin is 0.

Moreover, if the objective of the linear program in Figure 4.4 is U , then the expected sample complexity of the constructed linear estimator is $O(\frac{1}{U\epsilon^2} \log \frac{1}{\delta})$, which will estimate ρ to within an additive error of ϵ with probability at least $1 - \delta$.

We further show in Section 4.8.2 that this linear estimator construction is in fact optimal in sample complexity, up to constant multiplicative factors, in the regime of constant probability success and subject to the same constraint that each coin can only be flipped at most n_{\max} times. The following theorem captures the exact guarantees.

Theorem 4.6. *Suppose we are given the 4 pieces of data as in Theorem 4.5 above.*

The linear estimator produced from solving the linear program in Figure 4.4, as described in Theorem 4.5, has total expected sample complexity that is within a constant factor of any optimal fully-adaptive algorithm with $\geq \frac{2}{3}$ probability of success, subject to the same constraint that no coin is flipped more than n_{\max} many times.

The proof of this theorem—like our main lower bound of Theorem 4.4—also relies on Lemma 4.19 to relate fully-adaptive algorithms to single-coin-adaptive algorithms; and constant-factor tightness comes from the fact that the linear programs minimizing the variance of a linear estimator versus maximizing squared Hellinger distance are within a constant factor of each other.

4.1.2 Related Work

A related line of work considers the scenario where all positive coins have identical bias (not necessarily greater than $1/2$), and negative coins also have identical bias (strictly less than the positive coins' bias), with the ultimate goal of *identifying* any single coin that is positive (or “heavy” in the terminology of these works). The problem has been studied and solved optimally in the context where the biases and positive-negative proportions are

known [18], and also when none of this information is known [46, 33]. Such problems may be seen as a special case of bandit problems.

Another related line of work concerns the *learning* of distributions of (e.g. coin) parameters over a population, which arises in various scientific domains [41, 42, 47, 52, 22, 4]. In particular, the works of Lord [41], and Kong et al. [61, 64] consider a model similar to ours, with the crucial difference that each coin is sampled a fixed number t many times—instead of allowing adaptive sampling as in the current work—with the objective of learning the distribution of biases of the coins in the universe.

Since an earlier version of this paper was posted on arXiv, more recent work by Brennan et al. [7] considers a generalization of our setting, but because of different motivation and parameterization, both their upper and lower bounds are not directly comparable with ours.

Our problem also sits in the context of estimation and learning tasks with *noisy* or *uncalibrated* queries. The noiseless version of our problem would be when $\Delta = \frac{1}{2}$ and thus $\frac{1}{2} \pm \Delta$ equals either 0 or 1. That is, all coins are either deterministically heads or deterministically tails, and thus estimating the mixture parameter ρ is equivalent to estimating the parameter of a *single* coin with bias ρ , which has a standard analysis. Prior works have considered noisy versions of well-studied computational problems, such as (approximate) sorting and maximum selection under noisy access to pairwise comparisons [29, 27] and maximum selection under access to uncalibrated numerical scores that are consistent with some global ranking [65].

Furthermore, our problem can be interpreted as a special case of the “testing collections of distributions” model introduced by Levi, Ron and Rubinfeld [39, 40], modulo the distinction between testing and parameter estimation. In their model, a collection of m distributions (D_1, \dots, D_m) (over the same domain) is given to the tester, and the task is to test whether the collection satisfies a particular *property*, where a property in this case is defined as a subset of m -tuples of distributions. In the *query* access model, one is allowed to name an index $i \in \{1, \dots, m\}$ and get a fresh sample from the distribution D_i . Our problem can be analogously phrased in this model, where the distributions are over the domain $\{0, 1\}$, and the property in question is whether the fraction ρ of distributions in the collection

having bias $\geq 1/2$ is greater than some threshold τ .

We highlight other distribution testing models that allow for adaptive sampling access. For example, in testing contexts, conditional sampling oracles have been considered [17, 11, 14, 13, 28, 1], where a subset of the domain is given as input to the oracle, which in turn outputs a sample from the underlying unknown distribution conditioned on the subset. Evaluation oracles have also been considered [55, 3, 30, 12], where the testing algorithm has access to an oracle that evaluates the probability mass function or the cumulative mass function of the underlying distribution. See the survey by Canonne [10] for detailed comparisons between the different standard and specialized access models, along with a discussion of recent results.

Adaptive lower bounds of problems related to testing monotonicity of high-dimensional Boolean functions have a somewhat similar setup to ours, where binary decisions adaptively descend a decision tree according to probabilities that depend both on the algorithm and its (unknown) input that it seeks to categorize [5, 19]. Lower bounds in these works rely on showing that the probabilities of reaching any leaf in the decision tree under the two scenarios that they seek to distinguish are either exponentially small or within a constant factor of each other. This proof technique is powerful yet does not work in our setting, as many adaptive algorithms have high-probability outcomes that yield non-negligible insight into which of the two scenarios we are in. By contrast, our proof technique involves showing that, while such “insightful” outcomes may be realized with high probability, in these cases we must pay a correspondingly high sample complexity cost somewhere *else* in the adaptive tree.

A crucial part of our lower bound proof, Lemma 4.20, involves carefully “decomposing” fully-adaptive (multi-coin) algorithms into their single-coin components. Work by Braverman et al. [6] gives a data processing inequality in the context of communication lower bounds, whose proof uses similar ideas to how we prove Lemma 4.20.

As described at the beginning of the introduction, results of this work have practical applications in *crowdsourcing* algorithms in the context of data science and beyond. Theoretical studies with similar aims to our own have been undertaken on handling potentially noisy answers from crowdsourced workers due to lack of expertise [58, 56], (including this

work); in practice it is also crucial to understand how to incentivize workers to answer truthfully [57]. Our work also addresses directly the practical problem proposed by Chung et al. [21], to issue queries to potentially unreliable crowdsourced workers in order to estimate the fraction of records containing “wrong” data within a database; here adaptive queries are a natural capability of the model.

4.2 The Triangular Walk Algorithm

In this section, we present the *Triangular Walk* algorithm (Algorithm 3) for the problem, in the regime where both ρ and the coin biases are unknown. This is an important subroutine of our main, optimal algorithm; and the Triangular Walk algorithm itself can be used as an estimator in its own right. We demonstrate later in Section 4.7, with simulation results, that this algorithm offers practical advantages over the straightforward majority vote estimator mentioned in the introduction, as well as the state-of-the-art method used in practice.

The Triangular Walk algorithm leverages only single-coin adaptivity, and makes no use of cross-coin adaptivity. At the heart of our algorithm is an estimator (Algorithm 2) that works coin-by-coin, in the regime $\Delta \geq \frac{1}{4}$; subsequently we show how to use this estimator to solve the general problem, with an arbitrary (but known) Δ .

We describe an asymmetric estimator (Algorithm 2) that, given sampling access to a single coin of bias p , returns a real number whose expectation is in $[1 \pm \frac{\epsilon}{2}]$ if $p \geq \frac{3}{4}$, and whose expectation is in $[\pm \frac{\epsilon}{2}]$ if $p \leq \frac{1}{4}$. The estimator is asymmetric in the sense that it will quickly “give up on” coins with $p \leq \frac{1}{4}$, taking only a constant number of samples from them in expectation, while it will more deeply investigate the rare and interesting case of $p \geq \frac{3}{4}$. Below, c will be a constant that emerges from the analysis, where $c \log \frac{1}{\epsilon}$ coin flips suffice to yield an empirical fraction of heads within $poly(\epsilon)$ of the ground truth, p .

Our overall algorithm robustly combines estimates from running Algorithm 2 on many coins via the standard median-of-means technique. To deal with the general case when Δ might be much smaller than $\frac{1}{4}$, we “simulate a $\frac{1}{4}$ -quality coin” by running Algorithm 2 not on individual flips, but rather on the majority vote of blocks of $\Theta(\frac{1}{\Delta^2})$ flips; this majority vote will convert a coin of bias $\leq \frac{1}{2} - \Delta$ to a simulated coin of bias $\leq \frac{1}{4}$, and symmetrically,

Algorithm 2 Single-coin estimate

Given: a coin of bias p , error parameter ϵ

1. Let $n \leftarrow 0$ *(representing the total number of coin flips so far)*
 2. Let $k \leftarrow 0$ *(representing the total number of observed heads so far)*
 3. Repeat:
 - (a) Flip the coin, and increment $n \leftarrow n + 1$
 - (b) If heads, increment $k \leftarrow k + 1$
 - (c) If $2k \leq n$, **return** 0 and **halt** *(majority of flips are tails, evidence that p is small)*
 - (d) If $n = c \log \frac{1}{\epsilon}$, **return** $\min(4, \frac{n}{2k-n})$ and **halt** *(enough flips for concentration)*
-

Algorithm 3 Triangular walk algorithm

Given: t coins, quality parameter Δ , error parameter ϵ , and failure probability δ

1. For each coin: simulate a new “virtual” coin by computing the majority of $\Theta(\frac{1}{\Delta^2})$ flips each time a “virtual” flip is requested; run Algorithm 2 on each virtual coin, using, inputting ϵ unchanged, and record the returned estimates.
 2. Partition the returned estimates into $\Theta(\log \frac{1}{\delta})$ groups and compute the mean of each group.
 3. **Return** the median of the $\Theta(\log \frac{1}{\delta})$ means, or 0 if any of the groups in step 2 are empty.
-

convert a coin of bias $\geq \frac{1}{2} + \Delta$ to a simulated coin of bias $\geq \frac{3}{4}$.

Theorem 4.2. *Given coins where a ρ fraction of the coins have bias $\geq \frac{1}{2} + \Delta$, and $1 - \rho$ fraction have bias $\leq \frac{1}{2} - \Delta$, then running Algorithm 3 on $t = \Theta(\frac{\rho}{\epsilon^2} \log \frac{1}{\delta})$ randomly chosen coins will estimate ρ to within an additive error of $\pm\epsilon$, with probability at least $1 - \delta$, with an expected sample complexity of $O(\frac{\rho}{\epsilon^2 \Delta^2} (1 + \rho \log \frac{1}{\epsilon}) \log \frac{1}{\delta})$.*

The rest of this section concerns the (relatively straightforward) proof of Theorem 4.2, via an analysis of Algorithms 2 and 3; Section 4.4 instead formulates a more general algorithmic framework that adds some perspective to Algorithm 2, and whose abstractions will be crucial to the lower bound analysis in Section 4.5.

Intuition and analysis of Algorithm 2: Recall that Algorithm 2 is designed to work for coins of constant noise-quality Δ , namely, coins have bias either $\leq \frac{1}{4}$ or $\geq \frac{3}{4}$, and nothing in between. Algorithm 2 halts under two conditions: either the majority of observed flips have been tails—Step 3(c)—or our budget of coin flips (for that coin) is exhausted—Step 3(d). The first stopping condition is designed to make it more likely to halt early for

negative coins (coins with bias $p \leq \frac{1}{4}$), even though *all* coins may have a significant chance of halting early. Importantly, the chance of Algorithm 2 halting early depends on the coin's bias p , which is a priori unknown. The output coefficients in Step 3(d) are designed so that the expected output, given any negative coin (of bias $\leq \frac{1}{4}$), is close to 0, and similarly close to 1 given a positive coin (of bias $\geq \frac{3}{4}$). Furthermore, the output coefficients are all bounded by a constant, which gives a constant bound on the variance of the estimate.

Lemma 4.7 captures the guarantees we need from Algorithm 2 in order to analyze the triangular walk algorithm, Algorithm 3.

Lemma 4.7. *If Algorithm 2 is run with a sufficiently large universal constant c , then the following statements hold.*

1. *Given an arbitrary negative coin (having bias $p \leq \frac{1}{4}$), the output of Algorithm 2 has expectation in $[\pm \frac{\epsilon}{2}]$ and variance upper bounded by ϵ^2 . Furthermore, the expected sample complexity in this case is upper bounded by a constant.*
2. *Given an arbitrary positive coin (having bias $p \geq \frac{3}{4}$), the output of Algorithm 2 has expectation in $[1 \pm \frac{\epsilon}{2}]$ and variance upper bounded by a constant. The expected sample complexity in this case is (trivially) upper bounded by $c \log \frac{1}{\epsilon}$.*

The overall expected sample complexity, when the fraction of positive coins is ρ , is $O(1 + \rho \log \frac{1}{\epsilon})$.

Proof. Consider running Algorithm 2 on a coin of bias p , and let $n_{\max} = c \log \frac{1}{\epsilon}$ be the number of coin flips after which the algorithm always halts in Step 3(d). Consider running Algorithm 2 on a sequence of n_{\max} flips of the coin (even if the algorithm may halt early before exhausting the sequence of flips). If the sequence is majority-tails, then the algorithm must halt early via Step 3(c) at some point, and thus return 0.

For a negative coin, namely with bias $p \leq \frac{1}{4}$, the chance of observing majority-heads after $c \log \frac{1}{\epsilon}$ coin flips is $\epsilon^{O(c)}$, and we choose c so that this probability is $O(\epsilon^2)$, so that (given that estimates returned by Algorithm 2 are always bounded by 4), for negative coins, the expected estimate of Algorithm 2 is in $[0, \frac{\epsilon}{2}]$ and the variance is at most ϵ^2 .

By contrast, for a positive coin, with bias $p \geq \frac{3}{4}$, the fraction of observed heads after $c \log \frac{1}{\epsilon}$ flips will concentrate around $p \geq \frac{3}{4}$. The challenge is to choose nonzero output

coefficients in Step 3(d) of Algorithm 2 that will average out to 1 in expectation, despite the fact that many of these sequences of coin flips will lead Algorithm 2 to terminate early in Step 3(c) and output 0. Moreover, as mentioned earlier, the proportion of sequences that will halt early depends on p itself, which is a priori unknown.

The key idea, from standard results on random walks, is that, *conditioned on* k out of n_{\max} flips landing heads, the probability of reaching n_{\max} flips—without ever halting in Step 3(c) by having a temporary majority of tails—is *independent* of p , and is in fact expressed by the formula $\frac{2k-n_{\max}}{n_{\max}}$. Conditioned on k out of n_{\max} flips being heads, whether Algorithm 2 halts early depends only on the permutation of the coin flips, and each such permutation of k heads out of n_{\max} flips is *equally* likely. We thus apply the following standard random walk result to derive the aforementioned formula—where heads is interpreted as a +1 step in a 1-D random walk, tails is interpreted as a -1 step, and observing k out of n heads is analogous to reaching position $v = 2k - n$ in the random walk.

Fact 4.8 (The Ballot Theorem). *Consider a 1-D walk that starts at the origin, and moves one step in either the positive or negative direction at each time. The number of paths from the origin that end at v at time n_{\max} , which do not revisit the origin, is a $\frac{|v|}{n_{\max}}$ fraction of the total number of paths from the origin to v at time n_{\max} .*

Algorithm 2, in Step 3(d), returns $\min(4, \frac{n_{\max}}{2k-n_{\max}})$, which equals $\frac{n_{\max}}{2k-n_{\max}}$ when $k \geq \frac{5}{8}n_{\max}$. In light of Fact 4.8, in Algorithm 2, conditioned on $k \geq \frac{5}{8}n_{\max}$ out of n_{\max} coin flips being heads, the nonzero coefficient $\frac{n_{\max}}{2k-n_{\max}}$ will be output in Step 3(d) with probability $\frac{2k-n_{\max}}{n_{\max}}$, and thus the conditional expected output is exactly 1. For $p \geq \frac{3}{4}$, the probability of $k \geq \frac{5}{8}n_{\max}$ flips being heads is $1 - \epsilon^{O(c)}$ by our choice of n_{\max} . The expected output of Algorithm 2 will therefore be within $\epsilon/2$ of 1 for a sufficiently large choice of the constant c .

The variance of Algorithm 2 given a positive coin is clearly upper bounded by a constant, simply because the output coefficients are bounded by 4.

We lastly analyze the expected sample complexity of Algorithm 2, run on negative and positive coins. For positive coins, we can simply upper bound the sample complexity by $n_{\max} = \Theta(\log \frac{1}{\epsilon})$, which is tight if the coin has bias $p = 1$. For negative coins, even ignoring the halting conditions, the probability of getting $> \frac{n}{2}$ heads after n coin flips decreases

exponentially in n . Since the algorithm halts if it ever observes majority-tails, proceeding for many flips becomes exponentially unlikely. Thus the expected number of flips of the algorithm is bounded by a constant in the case of a negative coin. \square

Analyzing Algorithm 3: We conclude by proving Theorem 4.2, which analyzes Algorithm 3.

Proof of Theorem 4.2. For this proof, we assume that $\rho = \Omega(\epsilon^2)$. Otherwise, the case is handled in Step 3 of Algorithm 3, which returns the valid estimate of 0.

At a high-level, Algorithm 3 runs Algorithm 2 repeatedly on independently chosen coins.

Observe that in Step 1 of Algorithm 3, for each coin we simulate a new “virtual” coin, by using the majority vote of $\Theta(1/\Delta^2)$ coin flips to compute each requested coin flip. By Chernoff bounds, if each given coin has bias either $p \leq \frac{1}{2} - \Delta$ or $p \geq \frac{1}{2} + \Delta$, then the corresponding virtual coin will have bias $p \leq \frac{1}{4}$ and $p \geq \frac{3}{4}$ respectively. Therefore, by Lemma 4.7, the output of Step 1 for each coin is a random variable with expectation in $[\rho \pm \frac{\epsilon}{2}]$. As for the variance of the output, we do the following calculation. Let X_0 denote the random variable that is the output of Algorithm 2 when given a random *negative* coin, and similarly X_1 for a random positive coin. The output of Algorithm 2, which we call Y , is thus distributed as X_1 with ρ probability and as X_0 with $1 - \rho$ probability. The variance of Y is

$$\begin{aligned} \text{Var}[Y] &= \rho \text{Var}[X_1] + (1 - \rho) \text{Var}[X_0] + \underset{i \leftarrow \text{Bernoulli}(\rho)}{\text{Var}} [\mathbb{E}[X_i]] \\ &\leq O(\rho) + (1 - \rho)\epsilon^2 + \rho(\mathbb{E}[X_1])^2 + (1 - \rho)(\mathbb{E}[X_0])^2 \\ &\leq O(\rho) + \epsilon^2 + O(\rho) + O(\epsilon^2) \\ &= O(\rho) \end{aligned}$$

Steps 2 and 3 of Algorithm 3 are the median-of-means method for estimating the mean of a (real-valued) random variable. Using $t = \Theta(\frac{\rho}{\epsilon^2} \log \frac{1}{\delta})$ coins, each of the $\Theta(\log \frac{1}{\delta})$ groups will have $\Theta(\frac{\rho}{\epsilon^2})$ coins and hence outputs from Algorithm 2. By Chebyshev’s inequality, with constant probability, the sample mean of each group’s estimates will be within $O(\sqrt{\frac{\epsilon^2}{\rho}})$

standard deviations of the expected output of Algorithm 2. The estimation error is therefore equal to $O(\epsilon)$, with a multiplicative constant that can be made arbitrarily small by adjusting the constant in the choice of the number of coins t . Step 3 computes the median of $\Theta(\log \frac{1}{\delta})$ such sample means, which boosts the success probability from constant to $1 - \delta$, via standard uses of Chernoff bounds.

Lastly, the total expected sample complexity is the product of 1) the choice of t in the theorem statement, 2) $\Theta(1/\Delta^2)$ which is the number of coin flips used for each majority vote in Step 1, and 3) the sample complexity of Algorithm 2 as stated in Lemma 4.7, yielding $O(\frac{\rho}{\epsilon^2 \Delta^2} (1 + \rho \log \frac{1}{\epsilon}) \log \frac{1}{\delta})$. \square

While Theorem 4.2 gives ϵ as input to Algorithm 3 and then asks how many coins are needed to achieve this ϵ error, it will be useful as a preliminary step of our optimal Algorithm 6 to consider the performance of Algorithm 3 where these two roles for ϵ are decoupled. Explicitly, how many coins or samples does it take for Algorithm 3 to achieve error ϵ_1 , when Algorithm 3 is given ϵ_2 as input? We will use this result in the regime where the failure probability for Algorithm 3 should be a constant, and thus for simplicity we omit δ from the following statement.

Corollary 4.9. *Given coins where a ρ fraction of the coins have bias $\geq \frac{1}{2} + \Delta$, and $1 - \rho$ fraction have bias $\leq \frac{1}{2} - \Delta$, then, for parameters $\epsilon_1, \epsilon_2 > 0$, running Algorithm 3 on $t = \Theta(\frac{\rho}{\epsilon_1^2})$ randomly chosen coins with parameter $\epsilon = \epsilon_2$ will estimate ρ to within an additive error of $\pm \epsilon_1$, with failure probability at most $0.1 + O(t \cdot \text{poly}(\epsilon_2))$, with an expected sample complexity of $O(\frac{\rho}{\epsilon_1^2 \Delta^2} (1 + \rho \log \frac{1}{\epsilon_2}))$. Note that the degree of the polynomial term (in ϵ_2) in the failure probability can be made arbitrarily high, by choosing a large constant c in Step 3(d) of Algorithm 2.*

Proof. (Sketch) The proof is essentially the same as that of Theorem 4.2.

The crucial difference is that, instead of interpreting the ϵ parameter of Algorithm 3 as an (additive) error parameter for the produced estimate, we interpret ϵ (which we parameterize as ϵ_2 in the corollary statement) as a “coin misclassification probability”. We explain this in more detail.

As shown in the proof of Theorem 4.2, with $n_{\max} = \Theta(\log(1/\epsilon_2))$, the probability that Algorithm 2 produces a non-0 estimate for a negative coin is $\text{poly}(\epsilon_2)$. As for a positive

coin, based on the analysis using Fact 4.8, the probability that n_{\max} flips of a positive coin resulting in fewer than $k = \frac{5}{8}n_{\max}$ heads is also $\text{poly}(\epsilon_2)$. Conditioned on such failure not happening, the expected value of Algorithm 2 on a positive coin is *exactly* 1.

Therefore, we can interpret Algorithm 3 as follows. Taking a union bound over the probabilities of the aforementioned failure modes, there is at most $O(t \cdot \text{poly}(\epsilon_2))$ probability that any of the t coins are “misclassified”. Conditioned on that not happening, Algorithm 2 is just a (meta-)Bernoulli coin that flips positive with probability ρ and negative with probability $1 - \rho$, in expectation. Explicitly, assuming (as happens with probability $1 - O(t \cdot \text{poly}(\epsilon_2))$ that none of the t coins are “misclassified”), each negative coin will yield an output of exactly 0, and each positive coin will yield an output of expectation exactly 1 and constant variance. Thus, given that a ρ fraction of coins from the underlying distribution are positive, the output will be exactly ρ in expectation (except with $O(t \cdot \text{poly}(\epsilon_2))$ misclassification probability) and has $O(\rho)$ variance. By Chebyshev’s inequality, the mean output over $t = \Theta(\frac{\rho}{\epsilon_1^2})$ coins will be within $\pm\epsilon_1$, except with probability 0.1 (choosing the multiplicative constant in the definition of t appropriately), as desired.

The variance of Algorithm 2 has already been bounded by $O(\rho)$ as in Theorem 4.2, and so the performance of Algorithm 3 can be analyzed by a straightforward application of Chebyshev’s inequality, yielding the accuracy part of the corollary statement, as well as the 0.1 probability term in the failure probability. \square

4.2.1 Implementing Algorithm 3

Later in Section 4.7, we give experimental results to demonstrate the performance of Algorithm 3 in practice. Here we address other concerns regarding the practical implementation and use of the algorithm.

The first concern is the fact that ρ , the ground truth that we are trying to estimate, appears in the sample complexity bound. We note that, once we fix an error parameter ϵ and the constant c in Algorithm 2, the overall algorithm of Algorithm 3 is an *anytime* algorithm: it can produce an estimate given any number of samples/coin flips, and the estimate simply gets more accurate (until it is as small as ϵ) as it gets a larger sample size. Crucially, the algorithm execution does not depend on the value of ρ itself. Thus, the fact

that ρ appears in the sample complexity has no bearing on the execution on the algorithm. In Section 4.3, we present our final algorithm (Algorithm 6) in a form where, given a budget B of coin flips, the algorithm “discovers” the correct ϵ based on the unknown answer ρ and the tight sample complexity formula. Algorithm 3 can enjoy the same guarantee following its invocation in Algorithm 6. The key insight is that, instead of fixing an ϵ for Algorithm 2, we use the budget B to derive a cutoff for the maximum number of flips we invest in a single coin.

The second issue is on the practical parameter regime of the noise parameter Δ . In this work, we study the asymptotics of the sample complexity as $\Delta \rightarrow 0$, but in practice, the quality of yes/no questions being asked will have at least constant correlation with the truth. To run our algorithms, then, we would ignore Step 1 of Algorithm 3, namely simulating “virtual” coins from real coins, and use real coin flips directly in Algorithm 2.

The third practical concern is on the non-zero output coefficients in Algorithm 2, in Step 3(d). In the algorithm and its subsequent analysis, we gave coefficients with a simple form of $\min(4, \frac{n}{2^{k-n}})$, which allowed for a straightforward analysis with Chernoff bounds. However, these output coefficients may not be the best possible, recalling that the objective of these coefficients is to make sure that the expected output of any underlying coin of bias $p \geq \frac{3}{4}$ is as close to 1 as possible. A simple observation is that the expected output of Algorithm 2 is in fact a smooth polynomial in p with coefficients determined by the output coefficients in the algorithm. Therefore, in practice, once we fix the constant c (or the entire quantity $c \log \frac{1}{\epsilon}$) in the algorithm description (Step 3(d) again), we can run a local search/gradient-based method to find output coefficients that make the expected output polynomial as close to 1 as possible, with initial coefficients being the ones given in Algorithm 2. These new output coefficients are reusable in practice as long as the noise parameter Δ in practice is not lower than the Δ for which the output coefficients were generated.

The fourth and last concern we address here is related to the first concern and the “any-time algorithm” implementation of Algorithm 3. The concentration results were phrased in terms of the number of coins that need to be sampled, namely the number of times Algorithm 2 is called by Algorithm 3. Each run of Algorithm 2 then flips each coin an *a priori*

unknown number of times to produce estimates. Since the sample complexity of a single triangular walk (Algorithm 2) is random, the concentration results only give an expected overall sample complexity for the algorithm. On the other hand, in practice one may wish to impose a *fixed* budget for sample complexity and simply use the entire budget. Such an approach introduces the issue that the triangular walk started last will probably not have finished by the time the budget is exhausted. How then can we aggregate the estimates obtained from the completed triangular walks without introducing bias in subtle ways?

Here we show the surprising result that the most natural algorithm does in fact work as is: that is, we take the average estimate of all the completed walks, ignoring the incomplete walk in progress. As an example, to demonstrate that this success is unintuitive and nontrivial, if we instead separately run two executions of Algorithm 3 with *separate* budgets, and averaged the estimates of all completed walks across both executions, this average *would* be biased; but if we computed the average of those walks completed under each budget, separately, then—by the main result of this section—each average would be unbiased, and we could average these averages together to yield an unbiased estimator.

To show that this “budgeted estimator” is unbiased, we view it as the following two-stage estimator: 1) We estimate without bias the distribution over states (n, k) that the triangular walk (Algorithm 2) terminates at, when given a randomly chosen coin from the universe, namely the numbers $\{\alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)\}$ (using notation defined in Definition 4.14). 2) We simply take the dot product of this distribution with the corresponding output values $\{v_{n,k}\}$ (as defined in Algorithm 2).

In order to perform step 1), that is to estimate the distribution of termination over the states, we use the estimator $i_{(n,k)} / \sum_{m,j} i_{(m,j)}$ where $i_{(n,k)}$ is the number of observed walks that terminated at (n, k) , ignoring incomplete walks. The following proposition shows that the estimation in step 1) is unbiased, from which it follows that the entire estimator is indeed also unbiased.

Proposition 4.10. *Given a budget $T > 0$, and suppose we repeatedly run an adaptive algorithm A on a single coin until we have flipped the coin T times in total. Given a set of outcomes for the algorithm A , indexed by $k \in \{1, \dots, K\}$, let p_k be a probability distribution over outcomes, and let t_k denote the number of coin flips taken to reach this outcome. When an outcome using t coin flips*

is drawn, if t is less than or equal to the remaining budget, then t is subtracted from the remaining budget; and otherwise the most recent outcome is discarded as “over budget” and the algorithm terminates. Let i_k be the number of times that outcome k is drawn. Then $i_k / \sum_j i_j$ is an unbiased estimator of p_k .

Proof. Given the coin budget T , the possible sequences of samples can be classified into the following cases. Either 1) the sequence ends exactly at time T , or 2) the sequence ends with a time interval of length t_m for some m , which in turn ends *after* time T . For a vector \mathbf{i} , whose k^{th} index denotes the number of times outcome k occurs, the dot product with vector \mathbf{t} counts the total number of coin flips used by this sequence. Thus, if $\mathbf{i} \cdot \mathbf{t} = T$, then the probability of \mathbf{i} occurring equals

$$\binom{i_1 + \dots + i_K}{i_1; \dots; i_K} p_1^{i_1} \dots p_K^{i_K}$$

This expression captures all cases where we use *exactly* our budget T . In the remaining cases, there is a final (discarded) outcome m that goes “over budget”. In this case, $\mathbf{i} \cdot \mathbf{t} \in [T - t_m + 1, T - 1]$, and the probability of observing \mathbf{i} and discarding m equals

$$\binom{i_1 + \dots + i_K}{i_1; \dots; i_K} p_1^{i_1} \dots p_K^{i_K} \cdot p_m$$

Therefore, the expectation of $i_k / \sum_j i_j$ can be written as

$$\begin{aligned} & \sum_{\substack{\text{vector } \mathbf{i} \\ \text{s.t. } \mathbf{i} \cdot \mathbf{t} = T}} \binom{i_1 + \dots + i_K}{i_1; \dots; i_K} p_1^{i_1} \dots p_K^{i_K} \cdot \frac{i_k}{i_1 + \dots + i_K} \\ & + \sum_m \sum_{\substack{\text{vector } \mathbf{i} \\ \text{s.t. } \mathbf{i} \cdot \mathbf{t} \in [T - t_m + 1, T - 1]}} \binom{i_1 + \dots + i_K}{i_1; \dots; i_K} p_1^{i_1} \dots p_K^{i_K} \cdot p_m \cdot \frac{i_k}{i_1 + \dots + i_K} \end{aligned}$$

Now observe that

$$\binom{i_1 + \dots + i_K}{i_1; \dots; i_K} \frac{i_k}{i_1 + \dots + i_K} = \binom{i_1 + \dots + (i_k - 1) + \dots + i_K}{i_1; \dots; i_k - 1; \dots; i_K}$$

meaning that, letting the vector \mathbf{i}' equal the vector \mathbf{i} with its k^{th} entry decreased by 1, the expectation can be rewritten and simplified as

$$p_k \left[\sum_{\substack{\mathbf{i}' \\ \text{s.t. } \mathbf{i}' \cdot \mathbf{t} = T'}} \binom{i'_1 + \dots + i'_K}{i'_1; \dots; i'_K} p_1^{i'_1} \dots p_K^{i'_K} + \sum_m \sum_{\substack{\mathbf{i}' \\ \text{s.t. } \mathbf{i}' \cdot \mathbf{t} \in [T' - t_m + 1, T' - 1]}} \binom{i'_1 + \dots + i'_K}{i'_1; \dots; i'_K} p_1^{i'_1} \dots p_K^{i'_K} \cdot p_m \right]$$

where $T' = T - t_k$. The term inside the square brackets sums to 1, as we observed at the beginning of the proof, but substituting T' for T . Thus the expectation is p_k , as desired. \square

4.3 The Main Algorithm

Here we present our main algorithm, Algorithm 6, analyzed in Theorem 4.3, which uses $O(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ samples, matching the fully-adaptive lower bound we prove in Section 4.5.

Algorithm 6 uses the Triangular Walk estimator as a subroutine and has a hybrid flavor, combining both (single-coin) adaptive and non-adaptive techniques, where the algorithm is increasingly adaptive for smaller values of ρ . Crucially, in the adaptive component of Algorithm 6, we use the Triangular Walk estimator to provide a 2-approximation to ρ , and a variant of the algorithm to “filter” out most negative coins such that we get a constant ratio of positive vs negative coins, to reduce variance. The coins “surviving” the filter are then fed into a new, non-adaptive algorithm (Algorithm 4) that we call “refined sampling”, which like Algorithm 2 flips different coins a *different* number of times, yet the number of flips is chosen *non-adaptively*; the information from different coins is combined in a subtle way.

As a general motivation, consider taking t coins, flipping them n times each, and trying to estimate the fraction of positive coins. For a slightly different setting that may have cleaner intuition, consider having sample access to many univariate Gaussian distributions of bounded variance, some of which have mean ≤ 0 and some of which have mean ≥ 1 , where the goal is to estimate the fraction of “positive” Gaussians with as few samples as possible. If we take n samples from a given distribution, then testing whether the sample mean is $> \frac{1}{2}$ lets us correctly determine its identity with probability $1 - \exp(-n)$, incentivizing us to choose a large n . However, for a fixed budget on the total number of samples across all distributions, choosing many samples per distribution means we can only sample from a limited number of distributions, introducing sampling errors across distributions (as opposed to within each distribution), and thus introducing a variance into our estimate inversely proportional to the number of coins sampled, and thus $O(n/T)$ for a total budget of T . This is the classic bias-variance tradeoff, where larger n induces a better

bias but worse variance.

While in many settings, one might try to find an optimal n that balances these two concerns, the right answer here is instead to combine the two approaches: sample some distributions many times, to get a low-bias signal, and also sample many distributions a few times, to get a low-variance signal; and combine these two signals with care. Explicitly, the coefficients in Step 3 of Algorithm 4 are carefully chosen so that their contributions “telescope” in expectation between distributions sampled different numbers of times, allowing, essentially, all the high-variance terms to cancel out without worsening the bias.

We first present the non-adaptive component (Algorithm 4) of Algorithm 6 for estimating smooth functions f on the underlying coin bias p , which has constant expected sample complexity, with *zero bias*, at the cost of $O(1)$ variance instead of $O(\rho)$ variance as in Algorithm 3. This will be combined with a single-coin adaptive “filtering” component such that only an $O(\rho)$ fraction of coins will be used in running Algorithm 4, giving an overall $O(\rho)$ variance in Algorithm 6.

Think of the function f as being analogous to the output coefficient $\frac{n}{2^{k-n}}$ of Algorithm 2—correcting for a probabilistic filtering mechanism, such that the expected output of f for those coins that survive filtering will be essentially 0 for negative coins ($p \leq \frac{1}{4}$), 1 for positive coins ($p \geq \frac{3}{4}$), and smoothly transitions between 0 and 1 in between. See later in Definition 4.12 for the precise instantiation of $f(p)$ we need.

Let $\text{Bin}(n, p, k)$ denote the probability that a Binomial distribution with n trials and bias p outputs k .

Algorithm 4 Refined Sampling

Input: sample access to a coin of bias p ; target function $f : [0, 1] \rightarrow \mathbb{R}$

1. Choose a number of coin flips n that is a power of 2, choosing 2^i with probability $\frac{\sqrt{8}-1}{\sqrt{8}}(2^i)^{-1.5}$, where $\frac{\sqrt{8}-1}{\sqrt{8}}$ is the normalizing constant so that the probabilities sum to 1.
 2. Flip the coin n times, and let k be the number of observed heads.
 3. Return $\frac{n^{1.5}\sqrt{8}}{\sqrt{8}-1} \left(f\left(\frac{k}{n}\right) - \sum_{i=0}^{n/2} f\left(\frac{i}{n/2}\right) \cdot \binom{n/2}{i} \binom{n/2}{k-i} / \binom{n}{k} \right)$
-

The sum in Step 3 of the algorithm is omitted if the power of 2 chosen for the number of coin flips is $n = 1$, in which case $\binom{n/2}{i}$ would be undefined. We now describe the properties

of Algorithm 4 in Lemma 4.11.

Lemma 4.11. *Given a coin of bias p , and given a function $f : [0, 1] \rightarrow \mathbb{R}$ that is bounded by a universal constant, and has 2nd derivative bounded by a universal constant, then Algorithm 4 will return an estimate of $f(p)$ that has bias 0, variance $O(1)$, and uses $O(1)$ samples in expectation.*

Proof. The expected number of coin flips taken by Algorithm 4 is the sum of a fixed geometric series, and is thus $O(1)$ as desired.

We bound the variance of the algorithm by showing that, for each depth n , the values returned in Step 3 will have magnitude $O(n^{0.5})$. Consider the sum in the second term of the expression of Step 3. The expression $\binom{n/2}{i}\binom{n/2}{k-i}/\binom{n}{k}$ can be interpreted as: given a sequence of n coin tosses of which k were heads, if a random subsequence of length $n/2$ is chosen, what is the probability that i heads are chosen. This distribution has expectation $\frac{k}{2}$, and variance $< n$. Since f has second derivative bounded by a constant, the difference of f from $f(\frac{k}{n})$ is upper and lower bounded by quadratics centered at $\frac{k}{n}$. Thus the difference between $f(\frac{k}{n})$ and the expected value of $f(\frac{i}{n/2})$ when i is drawn from the distribution with pmf $\binom{n/2}{i}\binom{n/2}{k-i}/\binom{n}{k}$ is bounded by a constant times the variance of the random variable $\frac{i}{n/2}$, namely $O(\frac{1}{n})$. Therefore, when multiplied by $\frac{n^{1.5}}{\sqrt{8}-1}$, the output of Step 3 will be bounded by $O(n^{0.5})$ as desired. Since in Step 1, n is chosen with probability $\frac{\sqrt{8}-1}{n^{1.5}}$, the contribution to the variance from a particular n is at most $\frac{\sqrt{8}-1}{n^{1.5}}O(n^{0.5})^2 = O(n^{-0.5})$; summing this bound over all n that are powers of 2 yields a constant, $O(1)$, variance, since geometric series converge.

To analyze the expectation of the values returned in Step 3 of Algorithm 4, we show that it telescopes across the different depths n . Namely, consider the expected contribution just of the second (sum) term at level n , $-\sum_{k=0}^n \text{Bin}(n, p, k) \sum_{i=0}^{n/2} f\left(\frac{i}{n/2}\right) \cdot \binom{n/2}{i}\binom{n/2}{k-i}/\binom{n}{k}$. The coefficient in this expression of a given $f(\frac{i}{n/2})$ equals $-\sum_{k=0}^n \text{Bin}(n, p, k) \binom{n/2}{i}\binom{n/2}{k-i}/\binom{n}{k}$; from the discussion at the start of the proof, the k^{th} term of this sum can be reinterpreted as the probability that, in n tosses of a coin of bias p , we have k heads total, and i heads among the first $n/2$ tosses; summed over all k this is clearly just the probability that i heads will be observed among $n/2$ tosses, namely $\text{Bin}(\frac{n}{2}, p, i)$. Thus the expected value of the sum term of Step 3 at level n is $-\sum_{i=0}^{n/2} \text{Bin}(\frac{n}{2}, p, i) f(\frac{i}{n/2})$, which is exactly the negation of the expectation of the first term of Step 3, at level $n/2$. (The multiplier $\frac{n^{1.5}\sqrt{8}}{\sqrt{8}-1}$ in Step 3 is exactly canceled out by the probability of choosing n in Step 1.)

Thus the expected output of the algorithm, considering only contributions up to some depth $n = 2^i$, collapses to just the expectation of the first term of Step 3 at the deepest level, n . This expected output is thus $\sum_{k=0}^n f(\frac{k}{n}) \cdot \text{Bin}(n, p, k)$, namely the expected value of $f(\frac{k}{n})$ when k is drawn from a binomial distribution with n trials and bias p . Since the binomial distribution $\text{Bin}(n, p, \cdot)$ has expectation pn and variance $< n$, and since f has 2nd derivative bounded by a constant, we have that this expectation converges to $f(p)$ for large n ; namely, $|f(p) - \sum_{k=0}^n f(\frac{k}{n}) \cdot \text{Bin}(n, p, k)| = O(\frac{1}{n})$. Thus, as n goes to infinity, we see that the expected output of Algorithm 4 converges to $f(p)$, as claimed. \square

We now give a new non-adaptive algorithm, Algorithm 5, in order to motivate the choice of $f(p)$ that we use for Algorithm 4 within Algorithm 5. Algorithm 5 will be a major component of our final algorithm, Algorithm 6.

Algorithm 5 Optimal Algorithm given an estimate $\hat{\rho}$

Given: A total budget B of coin flips, quality parameter Δ , and an estimate $\hat{\rho}$ that is within a factor of 2 of ρ

1. Run the following on $t = \Theta(\Delta^2 B)$ randomly drawn coins. For each coin: simulate a new “virtual” coin by computing the majority of $\Theta(\frac{1}{\Delta^2})$ flips each time a “virtual” flip is requested, so that each virtual coin will have probability either $p \leq \frac{1}{4}$ or $p \geq \frac{3}{4}$.
 - (a) For each virtual coin, flip it at most $d = \Theta(\log \frac{1}{\hat{\rho}})$ times but stop if at any point the majority of flips are tails.
 - (b) If the previous step did not stop early, then run Algorithm 4 for the function $f_d(p)$ of Definition 4.12.
 2. Return $\frac{1}{t}$ times the sum of all the values output by Algorithm 4 in Step 3(b).
-

As mentioned above, the choice of $f(p)$ is a correction for the filtering mechanism. Concretely, in Algorithm 5, Step 2(a) will stop early on negative coins with probability that is increasingly high for smaller ρ , significantly reducing the number of coin flips; and in Step 2(b) we exactly compensate for this (a priori) unknown early stopping probability by running the unbiased Algorithm 4 on an appropriately chosen function $f_d(p)$ that is exactly the inverse of this early stopping probability, for positive coins, and 0 for negative coins:

Definition 4.12. Given a depth d , let $f_d(p) : [0, 1] \rightarrow \mathbb{R}$ be defined to equal 0 for $p \leq \frac{1}{4}$; and for $p \geq \frac{3}{4}$, let $f_d(p)$ equal 1 divided by the probability that a sequence of d flips of a coin of

bias p never has a majority-tails initial sequence; for $\frac{1}{4} < p < \frac{3}{4}$, let $f_d(p)$ be chosen so as to smoothly connect the regions $p \leq \frac{1}{4}$ and $p \geq \frac{3}{4}$ so that $f_d(p)$ has second derivative bounded by a universal constant (independent of d).

With this choice of $f(p)$, we state and prove Proposition 4.13, which gives the soundness and sample complexity bounds for Algorithm 5.

Proposition 4.13. *On input 1) a budget $B = O(\frac{\rho}{\epsilon^2 \Delta^2})$ of coin flips, 2) the quality parameter Δ and 3) a 2-approximation $\hat{\rho}$ of ρ , Algorithm 5 returns an estimate of ρ that has additive error at most ϵ with probability at least 0.99, using at most B coin flips.*

Proof. We first show the expected output of Algorithm 5 equals ρ . For each positive coin, Step 1 transforms it into a “virtual” coin of probability $p \geq \frac{3}{4}$; this coin will “survive” Step 1(a) with probability exactly $1/f_d(p)$, by definition of $f_d(p)$ in Definition 4.12. Thus Algorithm 4 will return an estimate of $f_d(p)$, with bias 0. Multiplying through by the survival probability $1/f_d(p)$, and by the probability ρ that a positive coin will be drawn, we see that, over t coins, the expected contribution to the estimate from Step 2 of the *positive* coins will be ρ . For each negative coin, by definition $f_d(p) = 0$, so the expected contribution from these coins, added over all $\leq t$ of them, and scaled by $\frac{1}{t}$ in Step 2, will be 0.

To bound the variance of the output of Step 2, we note that at most a 2ρ fraction of the coins reach Step 1(b): a ρ fraction of the coins are positive; meanwhile, negative coins, where $p \leq \frac{1}{4}$, have $\exp(-d)$ probability of surviving Step 1(a), which can be made $\leq \rho$ since $d = \Theta(\log \frac{1}{\rho})$. Thus the output returned in Step 2 is $\frac{1}{t}$ times the sum of t independent trials of a process that, with probability $\leq 2\rho$ outputs a random variable whose expected squared magnitude is bounded by a constant (by Lemma 4.11). For $t = \Theta(\frac{\rho}{\epsilon^2})$, the expected squared magnitude—and hence the variance—of the output of the algorithm is thus bounded by $O(\frac{t\rho}{t^2}) = O(\epsilon^2)$. Thus by Chebyshev’s inequality, Step 2 will return an estimate accurate to within $O(\epsilon)$, with constant probability.

Lastly, we need to verify that Step 1 will exceed the coin flip budget only with small constant probability. It suffices, using a Markov’s inequality argument, to bound the expected number of coin flips used in the steps. We consider the number of (“virtual”) flips from Step 1(a), and also Step 1(b), and then multiply by $\Theta(\frac{1}{\Delta^2})$ as described in Step 1.

For a negative coin, the expected number of flips until a majority-tails initial sequence is observed in Step 1(a) is constant by standard random walk analysis, leading to an $O(\Delta^2 B)$ term; for positive coins, there are on average $O(\rho \Delta^2 B)$ of them, so we could afford to flip each $O(\frac{1}{\rho})$ times, but Step 1(a) uses only at most $d = \Theta(\log \frac{1}{\rho})$ flips. Step 1(b) is run on an expected $\leq 2\rho$ fraction of the coins, as explained above; and by Lemma 4.11, Algorithm 4 takes $O(1)$ expected samples, for a total bound of $O(\rho \Delta^2 B)$ virtual flips from Step 1(b). (Algorithm 4 could thus afford to take up to $O(\frac{1}{\rho})$ samples on average, so, interestingly, there is a lot of slack here.)

Thus in total we use $O(\Delta^2 B) = O(\frac{\rho}{\epsilon^2})$ virtual flips, each requiring $\Theta(1/\Delta^2)$ real flips, corresponding to expected sample complexity of $O(B) = O(\frac{\rho}{\epsilon^2 \Delta^2})$. \square

Having analyzed Algorithm 5, we can now present our final optimal algorithm, stated as Algorithm 6. The theoretical guarantees are given in Theorem 4.3, restated and proved below.

We stress again that, in our presentation of Algorithm 6, the error parameter ϵ (of Theorem 4.3) is not known, since it depends on the budget B and the unknown ground truth ρ , yet the returned estimate will have this optimal ϵ accuracy regardless. This is achieved by Algorithm 6's calls to Algorithm 3 and Algorithm 5, which collectively cover all regimes of how ρ and ϵ relate to each other, yielding optimal error guarantees in each case.

Theorem 4.3. *Given coins where a ρ fraction of the coins have bias $\geq \frac{1}{2} + \Delta$, and $1 - \rho$ fraction have bias $\leq \frac{1}{2} - \Delta$, then running Algorithm 6 on a budget of B coin flips will estimate ρ to within an additive error of $\pm\epsilon$, with probability at least $2/3$, where ϵ is implicitly defined by the relation $B = \Theta(\frac{\rho}{\Delta^2 \epsilon^2})$ based on the unknown ground truth ρ . If the algorithm is repeated $\Theta(\log \frac{1}{\delta})$ times, and the median estimate is returned, then the probability of failure is at most δ .*

Proof. Given the fixed total sample complexity budget of B coin flips, and fixing the unknown ground truth ρ , the target additive error parameter ϵ is defined by the sample complexity equation $B = \Theta(\frac{\rho}{\epsilon^2 \Delta^2})$. There are two cases, either $\log \frac{1}{\epsilon} \leq c/\rho$ for some sufficiently small universal constant c (in which case we show that, with high probability, Step 1 will output a correct answer and then halt), or the inequality is in the opposite direction

Algorithm 6 Optimal Algorithm

Given: A total budget B of coin flips and quality parameter Δ

1. Use Algorithm 3 in Section 4.2 on $O(\Delta^2 B)$ many coins (a small fraction of B), using an “ ϵ ” that is $\Theta(1/(\Delta^2 B))$, and a constant δ . Let $\hat{\rho}$ be the returned estimate of ρ .
 - (a) If Algorithm 3 ever tries to use more than $B/4$ coin flips total, then terminate Algorithm 3 and move onto the next step.
 - (b) Otherwise, return the estimate produced by Algorithm 3.
 2. Use Algorithm 3 on $\Theta(\sqrt{\Delta^2 B})$ freshly drawn coins, using again an “ ϵ ” that is $\Theta(1/(\Delta^2 B))$, and a constant δ . The returned estimate $\hat{\rho}$ will be a 2-approximation to ρ . If in this step, Algorithm 3 tries to use more than $B/4$ coin flips, terminate and fail, which happens only with small constant probability.
 3. Run Algorithm 5 on input $B/2$, Δ , and $\hat{\rho}$, and return its answer.
 4. (If a sub-constant failure probability δ is desired, then repeat the entire algorithm $\Theta(\log \frac{1}{\delta})$ times and return the median of the outputs, ignoring invocations that failed.)
-

(in which case, with high probability, either Step 1 still produces a correct answer and halts, or Steps 2 and 3 will output a correct answer).

In the case where $\log \frac{1}{\epsilon} \leq c/\rho$, we use Corollary 4.9 with parameters $\epsilon_1 = \epsilon$ and $\epsilon_2 = \Theta(\frac{1}{\Delta^2 B}) = \Theta(\frac{1}{t})$: Algorithm 3 will have error $\pm\epsilon$, except with failure probability $0.1 + O(t \cdot \text{poly}(\epsilon_2)) = 0.1 + O(t \cdot \text{poly}(\frac{1}{t}))$, where, as noted in Corollary 4.9, we may make the polynomial superlinear to make this failure probability $0.1 + o(1)$. Further, the expected sample complexity is $O(\frac{\rho}{\epsilon^2 \Delta^2}) = O(B)$ in the case where $\log \frac{1}{\epsilon} \leq c/\rho$, so by Markov’s inequality, for appropriate constants we can ensure that Algorithm 3 uses $\leq B/4$ samples with high constant probability. Thus in this case, the algorithm will correctly terminate in Step 1(b) with high probability.

Next, we analyze the case where $\log \frac{1}{\epsilon} \geq c/\rho$. By Corollary 4.9, as above, if Step 1(b) is reached then its answer will be ϵ -accurate except with some small constant probability. Otherwise, since Steps 2 and 3 are statistically independent of Step 1, we can just analyze these steps for the case $\log \frac{1}{\epsilon} \geq c/\rho$, ignoring what happened in Step 1.

We first claim that Step 2 will return a 2-approximation $\hat{\rho}$ of ρ with high constant probability. As before, we use Corollary 4.9 with $\epsilon_2 = \epsilon$; since (from the algorithm and the parameters of the theorem) this step uses $t = \Theta(\sqrt{\Delta^2 B}) = \Theta(\frac{\sqrt{\rho}}{\epsilon})$ coins, solving the

equation $t = \Theta(\frac{\rho}{\epsilon_1})$ of the Corollary yields $\epsilon_1 = \Theta(\sqrt{\epsilon \sqrt{\rho}}) = O(\sqrt{\epsilon})$. Since we are in the regime where $\log \frac{1}{\epsilon} \geq c/\rho$, we have that $\epsilon_1 = O(\sqrt{\epsilon}) \leq O(e^{-c/\rho}) \ll \rho/2$ for sufficiently small ρ , meaning that we will approximate ρ to within $\pm \rho/2$, giving us a 2-approximation. The failure probability is $0.1 + o(1)$ as above. From Corollary 4.9, the expected sample complexity, in our case $\log \frac{1}{\epsilon} \geq c/\rho$ will be $O(\frac{\rho}{\epsilon_1^2 \Delta^2} \rho \log \frac{1}{\epsilon_2})$; substituting in the definitions of ϵ_1, ϵ_2 yields $O(\frac{\rho^{3/2}}{\epsilon \Delta^2} \log \frac{1}{\epsilon})$. Since $\rho = O(1)$ and $\log \frac{1}{\epsilon} = o(\frac{1}{\epsilon})$ this expected sample complexity is thus $O(\frac{\rho}{\epsilon^2 \Delta^2}) = O(B)$ and Markov's inequality implies the algorithm exceeds its sample bound in Step 2 with an arbitrarily small constant probability.

We conclude by invoking Proposition 4.13 to show that the estimate returned in Step 3 by Algorithm 5 is accurate to within additive error ϵ except with small constant probability. \square

4.4 Characterizing Single-Coin Algorithms

As a crucial first step towards the lower bounds of Section 4.5 that analyze how information from many different coins may interact, in this section we describe a unified framework for characterizing (adaptive) algorithms that flip only a single coin. Section 4.5.1 will then show a general structural result describing how any adaptive multi-coin algorithm may be broken into single-coin subroutines that may then be analyzed in light of the characterization of this section.

The most general form of an adaptive single-coin algorithm is a decision tree, where each node is a coin flip, and has two outgoing edges denoting the outcome of the coin flip, heads or tails; the current node captures the outcome of the entire sequence of coin flips so far, and thus for each node, a generic algorithm specifies a probability of halting, versus continuing from that node.

Via a (standard) symmetrization argument, instead of considering the state of the algorithm to be an arbitrary sequence of coin flips, we instead aggregate this information into a pair (n, k) representing the number of coin flips, and the number of heads observed so far. In outline, one may prove by induction on the number of coin flips n that any such decision tree may be "symmetrized" so that its stopping probability at each node $(n' \leq n, k)$ depends only on n' and k , while preserving, for any underlying coin bias p ,

the total probability of hitting the set of decision tree nodes that represent observing k total heads out of n' flips. The inductive step relies on the fundamental property that, conditioned on observing exactly k heads out of n' coin flips, the distribution over all such sequences of coin flips is independent of the coin bias p , and depends only on the stopping probabilities along each of the $\binom{n}{k}$ paths in the decision tree. This is a direct generalization of the analogous observation in the triangular walk algorithm section (Section 4.2), and is analyzed in slightly different form in Equation 4.1 below.

We thus consider single-coin algorithms as random walks (Algorithm 7) on the structure of the Pascal Triangle, in which the states are represented by pairs (n, k) , where n is the total number of flips of the coin so far, and $k \leq n$ is the number of "heads" responses. At each state (n, k) , the algorithm terminates with some probability $\gamma_{n,k}$, else the algorithm may request a further coin flip and continue the walk. The collection of parameters $\gamma_{n,k}$ we call a *stopping rule*, and specifies that algorithm's behavior.

Algorithm 7 Triangular Walk

Input: a coin of bias p

1. Initialize state (n, k) to $(0, 0)$.
 2. Repeat until termination:
 - (a) With probability $\gamma_{n,k}$, terminate and output (n, k) .
 - (b) Otherwise, sample one more coin flip. Increment n , and increment k by the result of the flip (0 or 1).
-

This formulation of single-coin algorithms, which we call a *triangular walk*, reveals structure that will be useful to the rest of the analysis of this paper. In particular, since the overall objective of running an adaptive coin-flipping algorithm is to recover information about the bias p of the coin (while minimizing expected sample complexity), it is fortuitous (as we will see) that the outcome of such an algorithm depends on p in an unexpectedly transparent way. This is given in Definition 4.14.

Definition 4.14. Given a stopping rule $\{\gamma_{n,k}\}$, we define coefficients $\{\alpha_{n,k}\}$, $\{\beta_{n,k}\}$, and $\{\eta_{n,k}\}$, so that, for any $p \in [0, 1]$, the triangular walk with stopping rule $\{\gamma_{n,k}\}$ on a coin of bias p , the coefficients have the semantics: $\alpha_{n,k}p^k(1-p)^{n-k}$ represents the probability that the walk terminates at (n, k) , with all such probabilities summing to 1; $\beta_{n,k}p^k(1-p)^{n-k}$ represents the

probability that the triangular walk encounters (n, k) , whether or not it terminates there, and $\eta_{n,k} p^k (1-p)^{n-k}$ is the probability that the triangular walk encounters (n, k) but does *not* terminate there. Each of these reparameterizations of the stopping rule may be derived from $\{\gamma_{n,k}\}$ using the following relations.

$$\begin{aligned} \beta_{0,0} &= 1 & (4.1) \\ \beta_{n+1,k+1} &= \beta_{n,k+1} \cdot (1 - \gamma_{n,k+1}) + \beta_{n,k} \cdot (1 - \gamma_{n,k}) \\ \alpha_{n,k} &= \beta_{n,k} \cdot \gamma_{n,k} \\ \eta_{n,k} &= \beta_{n,k} - \alpha_{n,k} \quad (= \beta_{n,k} \cdot (1 - \gamma_{n,k})). \end{aligned}$$

Consider the original setting, where one has a universe of (different) coins; one might repeatedly run a single-coin algorithm on coins drawn from the universe, and somehow combine their outputs into a final answer. There are many conceivable ways of aggregating the outputs of single-coin algorithms into an estimate, and the lower bounds of Section 4.5 consider them all. However, a particularly natural and powerful approach is to construct a linear estimator, namely to have the single-coin algorithm output a real number coefficient $v_{n,k}$ at each termination node, with the overall algorithm estimating the expected output of the single-coin algorithm, across the coins in the universe. Algorithm 3 works this way, using the median-of-means method (instead of taking the sample mean) to estimate the expected output of Algorithm 2. Such linear estimators are surprisingly flexible, and are known to be optimal in certain classes of estimation tasks [63].

4.5 Fully-Adaptive Lower Bounds

We show in this section that Algorithm 6 is optimal in all four problem parameters ρ, ϵ, Δ and δ , even when compared to all fully-adaptive algorithms that are adaptive across different coins. In particular, we show the following indistinguishability result (Theorem 4.4).

Theorem 4.4. *For $\rho \in [0, \frac{1}{2}]$ and $\epsilon \in (0, 1 - 2\rho]$, the following two situations are impossible to distinguish with at least $1 - \delta$ probability using an expected $o(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta})$ samples: A) ρ fraction of the coins have probability $\frac{1}{2} + \Delta$ of landing heads and $1 - \rho$ fraction of the coins have probability $\frac{1}{2} - \Delta$ of landing heads, versus B) $\rho + \epsilon$ fraction of the coins have probability $\frac{1}{2} + \Delta$ of landing heads*

and $1 - (\rho + \epsilon)$ fraction of the coins have probability $\frac{1}{2} - \Delta$ of landing heads. This impossibility crucially includes fully-adaptive algorithms.

With the algorithmic result of Theorem 4.3, this lower bound is therefore tight to within a constant factor. We note that the restrictions $\rho < \frac{1}{2}$ and $\epsilon \leq 1 - 2\rho$ reflect the symmetry of the problem, where the pair $\rho, \rho + \epsilon$ is exactly as hard to distinguish as the pair $1 - \rho - \epsilon, 1 - \rho$, yielding analogous results for the symmetric parameter regime.

Example 4.15. Even in the constant failure probability regime, the $\Omega(\frac{\rho}{\epsilon^2 \Delta^2})$ lower bound requires significant analysis, forming the bulk of the remainder of this paper, but two special cases have direct proofs. When $\Delta = \Theta(1)$ we can prove a $\Omega(\frac{\rho}{\epsilon^2})$ lower bound without the Δ dependence: consider the case where all coins are unbiased and perfect, meaning that the only source of randomness is from the mixture of coins, which is itself a Bernoulli distribution of bias either ρ or $\rho + \epsilon$. We quote the standard fact that, in order to estimate a Bernoulli coin flip of bias ρ to up to additive ϵ , we need $\Omega(\frac{\rho}{\epsilon^2})$ samples to succeed with constant probability; this can be proven by a standard (squared) Hellinger distance argument. On the other hand, it is also straightforward to prove a $\frac{1}{\Delta^2}$ lower bound (covering the regime where ρ and ϵ are constant): consider the easiest regime for ρ and ϵ , where $\rho = 0$ and $\epsilon = 1$; thus coins either all have $\frac{1}{2} + \Delta$ bias or all have $\frac{1}{2} - \Delta$ bias. To distinguish whether we have access to positive coins or negative coins requires $\Omega(\frac{1}{\Delta^2})$ samples.

In order to show Theorem 4.4, we use the Hellinger distance and KL-divergence between probability distributions as proxies for bounding the total variation distance.

Definition 4.16 (Hellinger Distance). Given two discrete distributions P and Q , the *Hellinger distance* $H(P, Q)$ between them is

$$\frac{1}{\sqrt{2}} \sqrt{\sum_i (\sqrt{p_i} - \sqrt{q_i})^2} = \sqrt{1 - \sum_i \sqrt{p_i q_i}}$$

Definition 4.17 (KL-divergence). Given two discrete distributions P and Q , the *KL-divergence* $D_{\text{KL}}(P||Q)$ between them is

$$\sum_i p_i \log \frac{p_i}{q_i}$$

The following facts capture how the Hellinger distance and KL-divergence can be used to show sample complexity lower bounds.

Fact 4.18 (Chapter 2.4, [62]). *For any two distributions P and Q over the same domain, we have*

$$\ell_1(P, Q) \leq \sqrt{2}H(P, Q)$$

and furthermore, for any event E ,

$$P(E) + Q(\bar{E}) \geq \frac{1}{2}e^{-D_{\text{KL}}(P\|Q)}$$

The second inequality is also known as the high-probability Pinsker inequality.

Recall from the introduction that, the main challenge in proving a general lower bound for our problem lies in analyzing the two kinds of adaptivity that algorithms may employ that were both absent in the special cases of Example 4.15. Explicitly, when taking samples from a given coin, we can choose whether to ask for another sample based on A) previous results of this coin, and also B) previous results of all the other coins. This first kind of adaptivity, “single-coin adaptivity”, is crucially used in the algorithms presented in the rest of the paper (e.g. the “shape” of the stopping rule for our triangular-walk algorithms); in Proposition 4.22 we analyze the best possible performance of such triangular stopping rules. The most interesting part of the proof of Theorem 4.4 consists of showing that the second kind of adaptivity (cross-coin adaptivity) cannot help in the lower bound setting, which we analyze via general Hellinger distance/KL-divergence inequalities (Lemmas 4.19 and 4.20) in Section 4.5.1.

4.5.1 Reduction to Single-Coin Adaptive Algorithms

In this section, we give two related but distinct reductions to single-coin adaptive algorithms. The first is a general decomposition (“direct sum”) inequality that decomposes the squared Hellinger distance of running a fully-adaptive algorithm on two different coin populations into the sum of squared Hellinger distances of running single-coin adaptive algorithms on the two coin populations. This inequality will lead to a constant probability sample complexity lower bound. The second inequality instead decomposes the KL

divergence into (a constant times) a sum of squared Hellinger distances, however with an additional slight restriction that the two coin populations being considered must be very close to each other. The upside to using this second inequality is that, an upper bound on the KL divergence combined with the high probability Pinsker inequality allows us to obtain a *high probability* sample complexity lower bound, which in particular is tight in *all* parameters of the problem, up to a multiplicative constant.

Both of the following inequalities are applicable to populations of variables *beyond* Bernoulli coins. We believe that the general inequalities are of independent interest to the community, since they would be applicable and useful for proving lower bounds on a variety of scenarios involving, for example, a Gaussian variant of the current problem, where instead of getting yes/no answers on the positivity of an item, one gets a real-valued score which correlates with the positivity of the item.

We phrase both lemmas as upper bounds on distances between distributions of the *transcript* of an algorithm, which when combined with the data processing inequality immediately yields upper bounds on distances between distributions of the algorithm's *output*. See, for example, the very end of the proof of Theorem 4.4.

Lemma 4.19. *Consider a problem setting where there is a collection of random variables, and an adaptive algorithm can draw variables from the collection and draw independent samples from the drawn variables. Now consider an arbitrary algorithm that iteratively samples from random variables drawn from the collection, choosing each subsequent variable to sample in an arbitrary adaptive manner based on the results of previous sample outcomes. Suppose the algorithm terminates almost surely. Consider two arbitrary collections of random variables, denoted by distributions \mathcal{A} and \mathcal{B} over the set of possible random variables. Let H_{full}^2 be the squared Hellinger distance between the transcript of a single run of the algorithm where 1) the random variables are drawn from \mathcal{A} versus where 2) the random variables are drawn from \mathcal{B} . Furthermore, let H_i^2 be the squared Hellinger distance between the two scenarios, but instead of running the algorithm as is, we only use random variable i (as drawn either from \mathcal{A} or \mathcal{B} depending on the scenario) and simulate all other random variables as independent random variables that are themselves drawn from the mixture distribution*

$\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}$. Then

$$H_{\text{full}}^2 \leq \sum_{\text{variable } i} H_i^2$$

Proof. It suffices to prove the result for deterministic algorithms, since squared Hellinger distance is linear with respect to mixtures of distributions *with distinct outcomes*, and a randomized algorithm is simply a mixture of deterministic algorithms which also records which of the algorithms the random coins picked. Furthermore, the following proof is phrased in terms of the special case where the collection of random variables are Bernoulli coins (which is the setting considered in this paper). Barring measure-theoretic formalization issues that we do not discuss, the proof generalizes directly to populations of arbitrary random variables.

A deterministic fully-adaptive algorithm is a decision tree, where each node is labeled by the identity of the coin the algorithm chooses to flip next conditioned on reaching this node, and each edge out of a node is labeled by a heads or tails result for this coin. We can view a run of the algorithm as follows: 1) first draw all the random coins from either \mathcal{A} or \mathcal{B} depending on the scenario, and then 2) flip these coins according to this fully-adaptive algorithm—we view choosing the coins from \mathcal{A} or \mathcal{B} as happening at the beginning since all these samples are free and only the coin flips themselves are counted. After step 1, fixing the bias of each coin, the probability of ending up at the i^{th} leaf of the decision tree is simply the probability (over coin flips) that every edge along the path from the root to that leaf is followed. Note that each edge is a probabilistic event depending on only one coin. Therefore, this probability can be factored into a product of probabilities, one term for each of the coins. For example, suppose the path to leaf i involves coin j returning 5 heads in a row, then getting some particular sequence from flipping some *other* coins, then coin j returning another 2 heads followed by 3 tails. Then, if coin j has bias p_j , it contributes $p_j^{5+2}(1-p_j)^3$ to the probability product.

We denote by $q_{j,i}^{\mathcal{A}}$ the *expected* contribution of coin j to the probability product for leaf i , over the randomness of \mathcal{A} on the bias of coin j . In the previous example, $q_{j,i}^{\mathcal{A}}$ would be equal to $\mathbb{E}_{p \leftarrow \mathcal{A}}[p^7(1-p)^3]$. We similarly define $q_{j,i}^{\mathcal{B}}$. Explicitly, for leaf i and coin j , $q_{j,i}^{\mathcal{A}}$ is the expectation (over p drawn from \mathcal{A}) of p to the exponent of the number of “heads” edges on the path from the root to node i in the decision tree, times $(1-p)$ to the exponent of the

number of “tails” edges on this path.

Using this notation, the probability of the algorithm reaching leaf i , when the coins are sampled from distribution \mathcal{A} , would be $\prod_{\text{coin } j} q_{j,i}^{\mathcal{A}}$, since each coin is sampled from \mathcal{A} independently; let $\prod_{\text{coin } j} q_{j,i}^{\mathcal{B}}$ be the respective probability for sampling from \mathcal{B} .

Since the total probability of reaching all leaves i must equal 1, this expression yields the immediate corollary, that for any distribution \mathcal{A} over $[0, 1]$,

$$\sum_{\text{leaf } i} \prod_{\text{coin } j} q_{j,i}^{\mathcal{A}} = 1 \quad (4.2)$$

We can now express the squared Hellinger distance with this notation. For any two distributions \mathbf{a} and \mathbf{b} , 1 minus their squared Hellinger distance can be rewritten as $\sum_i \sqrt{a_i b_i}$. In our context, the summation is over leaves i , and thus the squared Hellinger distance between the two scenarios in question is

$$H_{\text{full}}^2 = 1 - \sum_{\text{leaf } i} \sqrt{\prod_{\text{coin } j} q_{j,i}^{\mathcal{A}} \prod_{\text{coin } j} q_{j,i}^{\mathcal{B}}} \quad (4.3)$$

Since $q_{j,i}^{\mathcal{A}}$ and $q_{j,i}^{\mathcal{B}}$ are both non-negative, we simplify the summand as

$$\begin{aligned} & \sqrt{\prod_{\text{coin } j} q_{j,i}^{\mathcal{A}} \prod_{\text{coin } j} q_{j,i}^{\mathcal{B}}} \\ &= \left(\prod_{\text{coin } j} \frac{q_{j,i}^{\mathcal{A}} + q_{j,i}^{\mathcal{B}}}{2} \right) \left(\prod_{\text{coin } j} \frac{2\sqrt{q_{j,i}^{\mathcal{A}} q_{j,i}^{\mathcal{B}}}}{q_{j,i}^{\mathcal{A}} + q_{j,i}^{\mathcal{B}}} \right) \\ &\geq \left(\prod_{\text{coin } j} \frac{q_{j,i}^{\mathcal{A}} + q_{j,i}^{\mathcal{B}}}{2} \right) \left[1 - \sum_{\text{coin } j} \left(1 - \frac{2\sqrt{q_{j,i}^{\mathcal{A}} q_{j,i}^{\mathcal{B}}}}{q_{j,i}^{\mathcal{A}} + q_{j,i}^{\mathcal{B}}} \right) \right] \end{aligned} \quad (4.4)$$

where the inequality holds because each $\frac{2\sqrt{q_{j,i}^{\mathcal{A}} q_{j,i}^{\mathcal{B}}}}{q_{j,i}^{\mathcal{A}} + q_{j,i}^{\mathcal{B}}}$ is less than or equal to 1 by the AM-GM inequality (and at least 0), and therefore we can apply the union bound by treating each term as a probability—namely, for any $x_j \in [0, 1]$ we have $\prod_j x_j \geq 1 - \sum_j (1 - x_j)$.

Observe that our definition of q , being an expectation, is thus linear in the distribution in its superscript, and thus $\frac{1}{2}(q_{j,i}^{\mathcal{A}} + q_{j,i}^{\mathcal{B}}) = q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}}$, and therefore the right hand side of the inequality can be rewritten as

$$\left(\prod_{\text{coin } j} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) \left[1 - \sum_{\text{coin } j} \left(1 - \frac{\sqrt{q_{j,i}^{\mathcal{A}} q_{j,i}^{\mathcal{B}}}}{q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}}} \right) \right] \quad (4.5)$$

Thus the sum of Equation 4.5 over all leaves is at most $1 - H_{\text{full}}^2$. We simplify the summation by changing the summation variable in Equation 4.5 from j to k , and distributing the initial product so as to form three additive terms (the " $j \neq k$ " in the bounds of the last product below is because the $j = k$ term gets canceled by the denominator from the last term in Equation 4.5):

$$\begin{aligned} & \left(\sum_{\text{leaf } i} \prod_{\text{coin } j} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) - \sum_{\text{coin } k} \left(\sum_{\text{leaf } i} \prod_{\text{coin } j} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) \\ & - \sum_{\text{coin } k} \left(\sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) \sqrt{q_{k,i}^{\mathcal{A}} q_{k,i}^{\mathcal{B}}} \right) \end{aligned}$$

We know by Equation 4.2 that $\left(\sum_{\text{leaf } i} \prod_{\text{coin } j} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) = 1$, and so the sum can be written as

$$1 - \sum_{\text{coin } k} \left(1 - \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) \sqrt{q_{k,i}^{\mathcal{A}} q_{k,i}^{\mathcal{B}}} \right)$$

which by definition of H_k is equal to $1 - \sum_{\text{coin } k} H_k^2$: by Equation 4.3, 1 minus the squared Hellinger distance between the view of the algorithm when the k^{th} coin is from \mathcal{A} versus from \mathcal{B} , where all remaining coins are drawn from the mixture $\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}$ equals $\sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^{\frac{\mathcal{A}}{2} + \frac{\mathcal{B}}{2}} \right) \sqrt{q_{k,i}^{\mathcal{A}} q_{k,i}^{\mathcal{B}}}$.

Summarizing, we have shown that $1 - H_{\text{full}}^2 \geq 1 - \sum_{\text{coin } k} H_k^2$, from which the lemma statement follows. \square

We now give the KL-divergence decomposition lemma (Lemma 4.20) which will yield a tight *high probability* sample complexity lower bound, but makes a further assumption than Lemma 4.19 that the two coin populations are close to each other. As a note, the definition of H_i^2 is slightly different in this lemma from the definition in Lemma 4.19, and is not a typographical mistake.

Lemma 4.20. *Consider an arbitrary algorithm that iteratively flips coins from a collection of coins, choosing each subsequent coin to flip in an arbitrary adaptive manner based on the results of previous flips. Suppose the algorithm terminates almost surely. Consider two arbitrary mixtures of coins,*

denoted by distributions \mathcal{A} and \mathcal{B} over the coin bias $[0, 1]$. Let D_{full} be the KL-divergence between the transcript of a single run of the algorithm where 1) the coins are drawn from the mixture $\rho\mathcal{A} + (1 - \rho)\mathcal{B}$ versus where 2) the coins are drawn from $(\rho + \epsilon)\mathcal{A} + (1 - \rho - \epsilon)\mathcal{B}$, where $\rho \in [0, \frac{1}{2})$, $\epsilon \in (0, 1 - 2\rho]$ and $\epsilon < \rho$. Furthermore, let H_i^2 be the squared Hellinger distance between the two scenarios, but instead of running the algorithm as is, we only use coin i (as drawn either from the ρ -fraction mixture or the $(\rho + \epsilon)$ -fraction mixture depending on the scenario) and simulate all other coins as independent coins drawn from the ρ -fraction mixture. Then

$$D_{\text{full}} = O\left(\sum_{\text{coin } i} H_i^2\right)$$

The proof of Lemma 4.20 is similar to that of Lemma 4.19 by viewing algorithms as decision trees, with the crucial difference that, rather than using the AM-GM inequality, Lemma 4.20 instead bounds the KL-divergence via a quadratic bound $\log(1 + x) \geq x - x^2$, valid for $x \in [-\frac{1}{2}, 1]$.

Proof. (For the following proof, the set-up up to and including Equation 4.6 is essentially the same as that in the proof of Lemma 4.19, analogously, up to and including Equation 4.2. For completeness, we include the context for the specific notation we use in this proof.)

It suffices to prove the result for deterministic algorithms, since both squared Hellinger distance and KL-divergence are linear with respect to mixtures of distributions *with distinct outcomes*, and a randomized algorithm is simply a mixture of deterministic algorithms which also records which of the algorithms the random coins picked.

A deterministic fully-adaptive algorithm is a decision tree, where each node is labeled by the identity of the coin the algorithm chooses to flip next conditioned on reaching this node, and each edge out of a node is labeled by a heads or tails result for this coin. We can view a run of the algorithm as follows: 1) first draw all the random coins from either $\rho\mathcal{A} + (1 - \rho)\mathcal{B}$ or $(\rho + \epsilon)\mathcal{A} + (1 - \rho - \epsilon)\mathcal{B}$ depending on the scenario, and then 2) flip these coins according to this fully-adaptive algorithm. After step 1, fixing the bias of each coin, the probability of ending up at the i^{th} leaf of the decision tree is simply the probability (over coin flips) that every edge along the path from the root to that leaf is followed. Note that each edge is a probabilistic event depending on only one coin. Therefore, this probability can be factored into a product of probabilities, one term for each of the coins. For example,

suppose the path to leaf i involves coin j returning 5 heads in a row, then getting some particular sequence from flipping some *other* coins, then coin j returning another 2 heads followed by 3 tails. Then, if coin j has bias p_j , it contributes $p_j^{5+2}(1-p_j)^3$ to the probability product.

We denote by $q_{j,i}^{\mathcal{A}}$ the *expected* contribution of coin j to the probability product for leaf i , over the randomness of \mathcal{A} on the bias of coin j . In the previous example, $q_{j,i}^{\mathcal{A}}$ would be equal to $\mathbb{E}_{p \leftarrow \mathcal{A}}[p^7(1-p)^3]$. We similarly define $q_{j,i}^{\mathcal{B}}$. Explicitly, for leaf i and coin j , $q_{j,i}^{\mathcal{A}}$ is the expectation (over p drawn from \mathcal{A}) of p to the exponent of the number of "heads" edges on the path from the root to node i in the decision tree, times $(1-p)$ to the exponent of the number of "tails" edges on this path.

To simplify notation, we also denote by $q_{j,i}^{\rho}$ as the above probability product for coin j and leaf i when coin j is drawn from the ρ -mixture, namely $\rho\mathcal{A} + (1-\rho)\mathcal{B}$. Note that, by definition, $q_{j,i}^{\rho} = \rho q_{j,i}^{\mathcal{A}} + (1-\rho)q_{j,i}^{\mathcal{B}}$. We use analogous notation for the $(\rho + \epsilon)$ -mixture.

Using this notation, the probability of the algorithm reaching leaf i , when the coins are sampled from distribution \mathcal{A} , would be $\prod_{\text{coin } j} q_{j,i}^{\mathcal{A}}$, since each coin is sampled from \mathcal{A} independently, with $\prod_{\text{coin } j} q_{j,i}^{\mathcal{B}}$ being the respective probability for sampling from \mathcal{B} . The observation holds similarly for coin distribution that are the ρ -mixture or $(\rho + \epsilon)$ -mixture of \mathcal{A} and \mathcal{B} .

Since the total probability of reaching all leaves i must equal 1, this expression yields the immediate corollary, that for any collection of distributions C_j over $[0, 1]$ indexed by coin j (imagine C_j each being one of \mathcal{A} , \mathcal{B} , the ρ -mixture of the two or the $(\rho + \epsilon)$ -mixture of the two)

$$\sum_{\text{leaf } i} \prod_{\text{coin } j} q_{j,i}^{C_j} = 1 \quad (4.6)$$

We can now express the KL-divergence with the above notation.

$$\begin{aligned} D_{\text{full}} &= - \sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^{\rho} \right) \log \left(\frac{\prod_{\text{coin } k} q_{k,i}^{\rho+\epsilon}}{\prod_{\text{coin } k} q_{k,i}^{\rho}} \right) \\ &= - \sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^{\rho} \right) \sum_{\text{coin } k} \log \left(1 + \epsilon \frac{q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}}}{\rho q_{k,i}^{\mathcal{A}} + (1-\rho)q_{k,i}^{\mathcal{B}}} \right) \end{aligned}$$

where the second line follows from the definition of $q_{k,i}^{\rho} = \rho q_{k,i}^{\mathcal{A}} + (1-\rho)q_{k,i}^{\mathcal{B}}$. Further observe that the multiplier to ϵ in the second line is upper bounded by $1/\rho$ in magnitude. Since

Taylor's theorem gives that $\log(1+x) = x - \Theta(x^2)$ for $x \leq 1$, we have when $\epsilon/\rho \leq 1$, that

$$D_{\text{full}} = - \sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^\rho \right) \sum_{\text{coin } k} \left(\epsilon \frac{q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}}}{q_{k,i}^\rho} - \Theta(\epsilon^2) \frac{(q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}})^2}{(q_{k,i}^\rho)^2} \right)$$

We can further simplify the expression by observing that for any fixed coin k ,

$$\sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^\rho \right) \frac{q_{k,i}^{\mathcal{A}}}{q_{k,i}^\rho} = \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^\rho \right) q_{k,i}^{\mathcal{A}} = 1$$

where the second equality is by Equation 4.6. This observation holds also when we replace the mixture \mathcal{A} with the mixture \mathcal{B} . Therefore, we have

$$\sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^\rho \right) \sum_{\text{coin } k} \epsilon \frac{q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}}}{q_{k,i}^\rho} = 0$$

meaning that

$$D_{\text{full}} = \Theta(\epsilon^2) \sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^\rho \right) \sum_{\text{coin } k} \frac{(q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}})^2}{(q_{k,i}^\rho)^2} = \sum_{\text{coin } k} \Theta(\epsilon^2) \sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^\rho \right) \frac{(q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}})^2}{(q_{k,i}^\rho)^2}$$

It remains to show that the right hand side is bounded by $\sum_{\text{coin } k} H_k^2$, where, as in the lemma statement, H_k^2 is the squared Hellinger distance between a single run of the algorithm when coin k is drawn either from the ρ -mixture of \mathcal{A} and \mathcal{B} or the $(\rho + \epsilon)$ -mixture, and all other coins are simulated and simply drawn from the ρ -mixture. To see this, we write out what H_k^2 is, using the definition of squared Hellinger distance:

$$\begin{aligned} H_k^2 &= \sum_{\text{leaf } i} \left(\sqrt{q_{k,i}^\rho \prod_{\text{coin } j \neq k} q_{j,i}^\rho} - \sqrt{q_{k,i}^{\rho+\epsilon} \prod_{\text{coin } j \neq k} q_{j,i}^\rho} \right)^2 \\ &= \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^\rho \right) \left(\sqrt{q_{k,i}^\rho} - \sqrt{q_{k,i}^{\rho+\epsilon}} \right)^2 \\ &= \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^\rho \right) q_{k,i}^\rho \left(1 - \sqrt{\frac{q_{k,i}^{\rho+\epsilon}}{q_{k,i}^\rho}} \right)^2 \\ &= \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^\rho \right) q_{k,i}^\rho \left(1 - \sqrt{1 + \epsilon \frac{q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}}}{\rho q_{k,i}^{\mathcal{A}} + (1-\rho)q_{k,i}^{\mathcal{B}}}} \right)^2 \end{aligned}$$

where the last line is again by definition that $q_{k,i}^\rho = \rho q_{k,i}^{\mathcal{A}} + (1 - \rho)q_{k,i}^{\mathcal{B}}$. By reasoning we used above, as long as $\epsilon < \rho$, we have this expression being equal to

$$\begin{aligned} H_k^2 &= \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^\rho \right) q_{k,i}^\rho \left(\Theta(\epsilon) \frac{q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}}}{\rho q_{k,i}^{\mathcal{A}} + (1 - \rho)q_{k,i}^{\mathcal{B}}} \right)^2 \\ &= \Theta(\epsilon^2) \sum_{\text{leaf } i} \left(\prod_{\text{coin } j \neq k} q_{j,i}^\rho \right) q_{k,i}^\rho \frac{(q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}})^2}{(q_{k,i}^\rho)^2} \\ &= \Theta(\epsilon^2) \sum_{\text{leaf } i} \left(\prod_{\text{coin } j} q_{j,i}^\rho \right) \frac{(q_{k,i}^{\mathcal{A}} - q_{k,i}^{\mathcal{B}})^2}{(q_{k,i}^\rho)^2} \end{aligned}$$

which is exactly the term in the sum over coin k for D_{full} , showing the lemma. \square

For the lower bound proof at hand, we show Corollary 4.21 in the next subsection, which upper bounds the squared Hellinger distance for single-coin adaptive algorithms by a quantity that is proportional to the expected number of samples taken by the algorithm.

Corollary 4.21. *Consider an arbitrary single-coin adaptive algorithm. Let H^2 be the squared Hellinger distance between a single run of the algorithm where 1) a coin with bias $\frac{1}{2} + \Delta$ is used with probability ρ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise, versus a run of the algorithm where 2) a coin with bias $\frac{1}{2} + \Delta$ is used with probability $\rho + \epsilon$ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise. Furthermore, let $\mathbb{E}_\rho[n]$ and $\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$ be the expected number of coin flips during a run of the algorithm where we use a $\frac{1}{2} + \Delta$ coin with probability ρ and $\rho + \frac{\epsilon}{2}$ respectively, and a $\frac{1}{2} - \Delta$ coin otherwise. If all of ρ , ϵ , Δ and ϵ/ρ are smaller than some universal absolute constant, then*

$$\max \left[\frac{H^2}{\mathbb{E}_\rho[n]}, \frac{H^2}{\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]} \right] = O \left(\frac{\epsilon^2 \Delta^2}{\rho} \right)$$

Using Corollary 4.21 and Lemma 4.20, we now complete the proof of the main high probability indistinguishability result (Theorem 4.4) for fully-adaptive algorithms. We note again that Lemma 4.19, which is applicable to more general coin populations with fewer restrictions than Lemma 4.20, can be used to derive a constant probability sample complexity lower bound with essentially the same proof as follows, with the exception that we would use the Hellinger distance inequality in Fact 4.18 instead of the high-probability Pinsker inequality.

Proof of Theorem 4.4. Letting \mathcal{A} be a population of coins that all have $\frac{1}{2} + \Delta$ probability, with \mathcal{B} a population of coins that all have $\frac{1}{2} - \Delta$ probability, our goal is to show the indistinguishability of $\rho\mathcal{A} + (1 - \rho)\mathcal{B}$ from $(\rho + \epsilon)\mathcal{A} + (1 - \rho - \epsilon)\mathcal{B}$. We apply Lemma 4.20 and use the lemma's conclusion, that $D_{\text{full}} = O\left(\sum_{\text{coin } i} H_i^2\right)$.

Next, for each i , the quantity H_i^2 of Lemma 4.20 describes the squared Hellinger distance between an induced *single-coin* algorithm run on a single coin from scenario A versus B respectively (with the remaining coins being simulated, from scenario A with a ρ -fraction mixture). We thus bound H_i^2 from Corollary 4.21. As in the corollary, let $\mathbb{E}_{i,\rho}[n]$ denote the expected number of samples from coin i when running the induced algorithm (for coin i) on a mixture that uses a $\frac{1}{2} + \Delta$ coin with probability ρ and a $\frac{1}{2} - \Delta$ coin otherwise. Thus Corollary 4.21 yields that $H_i^2 = O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right) \cdot \mathbb{E}_{i,\rho}[n]$. Summing, combined with the result from Lemma 4.20 above, yields

$$D_{\text{full}} \leq O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right) \cdot \sum_{\text{coin } i} \mathbb{E}_{i,\rho}[n]$$

Crucially, the sum (over choice of coin i) of the expected number of flips $\mathbb{E}_{i,\rho}[n]$ (when running the algorithm induced for coin i) can be viewed in a different way: the i^{th} term is exactly the expected number of times that coin i is flipped when running the overall algorithm where every coin is drawn from the ρ -fraction mixture in scenario A . Namely, this sum counts the total expected number of coin flips (across all coins i), for the algorithm run in the setting where all coins are drawn from the ρ -fraction mixture. Thus, for an algorithm that uses $o\left(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta}\right)$ flips in expectation, we conclude that

$$D_{\text{full}} \leq O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right) \cdot o\left(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\delta}\right) = o\left(\log \frac{1}{\delta}\right)$$

We conclude by using the high-probability Pinsker inequality. In the notation of the inequality, given an algorithm that attempts to classify whether it is in scenario A or B , let P, Q respectively be the distributions of its output in scenarios A, B respectively; let E be the event that the algorithm outputs "scenario B ". Then the probability that the algorithm is wrong is $P(E) + Q(\bar{E})$. By the high-probability Pinsker inequality this failure probability is at least $\frac{1}{2}e^{-D_{\text{KL}}(P\|Q)} \geq \frac{1}{2}e^{-D_{\text{full}}} \geq \frac{1}{2}e^{-o(\log \frac{1}{\delta})} = \frac{1}{2}\delta^{o(1)} \gg \delta$ as desired, where the first inequality is the data processing inequality for KL-divergence. \square

4.5.2 Upper Bounding the Squared Hellinger Distance for Single-Coin Adaptive Algorithms

In this section, we prove Corollary 4.21, though significant technical details are deferred to Section 4.6. Explicitly, we analyze a simplified scenario in Proposition 4.22, after discussing why each of the simplifying assumptions does not give up generality, and cannot affect the key “squared Hellinger distance per sample” quantity by more than a constant factor.

1. Consider a single-coin algorithm. We restrict our attention to algorithms that only stop once they have seen a number of coin flips that is exactly a power of 2. Any stopping rule S that potentially stops in between powers of 2 could be converted into an almost-equivalent rule S' by collecting coin flips up to the next power of 2 and discarding them as necessary so as to simulate S : this will sacrifice at most a factor of 2 in sample complexity, and can only increase our Hellinger distance (since discarding data is a form of “data processing” and thus we may apply the data processing inequality). Thus the new S' will have “squared Hellinger distance per sample” at least half that of S .
2. By standard symmetrization arguments, a single-coin algorithm can always be implemented such that decisions only depend on the *number* of flips for a coin as well as the *number of observed “heads”*, as opposed to the explicit *sequence* of heads/tails observations. Thus we restrict our attention to stopping rules in the sense of Algorithm 7, specified in full generality by a triangle of stopping coefficients $\{\gamma_{n,k}\}$.
3. There is in some sense a “phase change” once an algorithm has received $\Omega(\frac{1}{\Delta^2})$ samples from a single coin: after this point, the algorithm might have good information about whether the coin is of type $\frac{1}{2} + \Delta$ versus type $\frac{1}{2} - \Delta$, and might productively make subtle adaptive decisions after this point. We restrict our analysis to the regime where *no coin is flipped more than $10^{-8}/\Delta^2$ times*: formally, we show an impossibility result in the following stronger setting, where we assume that whenever a single coin is flipped $10^{-8}/\Delta^2$ times, then the coin’s true bias (either $\frac{1}{2} + \Delta$ or $\frac{1}{2} - \Delta$) is *immediately* revealed to the algorithm. Thus any coin flips beyond $10^{-8}/\Delta^2$ that an algorithm desires can instead be simulated at no cost.

Formally, an impossibility result in this setting with “advice” (Proposition 4.22) implies the analogous result in the original setting (Corollary 4.21) by the data processing inequality for Hellinger distance (since Hellinger distance is an f -divergence): simulating additional coin flips in terms of “advice” is itself “data processing”, and thus can only decrease the Hellinger distance. Thus the setting without advice has smaller-or-equal Hellinger distance, and uses greater-or-equal number of samples, and hence the bound on their ratio in Proposition 4.22 implies the corresponding bound in Corollary 4.21.

Proposition 4.22. *Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that 1) is non-zero only for n that are powers of 2, and 2) $\gamma_{10^{-8}/\Delta^2,k} = 1$ for all k , that is the random walk always stops if it reaches $10^{-8}/\Delta^2$ coin flips. Suppose that given a coin, after a random walk on the Pascal triangle according to the stopping rule, the position (n,k) that the walk ended at is always revealed, and furthermore, if $n = 10^{-8}/\Delta^2$, then the bias of the coin is also revealed. Let H^2 be the squared Hellinger distance between a single run of the above process where 1) a coin with bias $\frac{1}{2} + \Delta$ is used with probability ρ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise versus 2) a coin with bias $\frac{1}{2} + \Delta$ is used with probability $\rho + \epsilon$ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise. Furthermore, let $\mathbb{E}_\rho[n]$ and $\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$ be the expected number of coin flips during a run of the algorithm where we use a $\frac{1}{2} + \Delta$ coin with probability ρ and $\rho + \frac{\epsilon}{2}$ respectively, and a $\frac{1}{2} - \Delta$ coin otherwise. If all of ρ , ϵ , Δ and ϵ/ρ are smaller than some universal absolute constant, then*

$$\max \left[\frac{H^2}{\mathbb{E}_\rho[n]}, \frac{H^2}{\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]} \right] = O \left(\frac{\epsilon^2 \Delta^2}{\rho} \right)$$

It remains to prove Proposition 4.22. For the rest of the section, we shall use the notation $h_{n,k}^+ = (\frac{1}{2} + \Delta)^k (\frac{1}{2} - \Delta)^{n-k}$ and $h_{n,k}^- = (\frac{1}{2} - \Delta)^k (\frac{1}{2} + \Delta)^{n-k}$ for convenience. The proofs for upper bounding $\frac{H^2}{\mathbb{E}_\rho[n]}$ and $\frac{H^2}{\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]}$ are essentially the same, and here we give the high-level outline of the proof for bounding the latter, with calculations and details in Section 4.6.

The first step in the proof is the following lemma that writes out the squared Hellinger distance induced by a given stopping rule $\{\gamma_{n,k}\}$, whose proof can be found in Section 4.10. The expression in Lemma 4.23 avoids square roots and in other ways simplifies aspects of the squared Hellinger distance by estimating terms to within a constant factor, which is folded into a multiplicative “big- Θ ” term at the start of the expression. The two lines

in the expression below capture the different forms of the Hellinger distance for stopping *before* the last row versus *at* the last row—recall that we prove impossibility under the stronger model where, upon reaching the last row the algorithm receives the true bias of the coin (as “advice”). Thus the squared Hellinger distance coefficients from elements of the last row are typically much larger than for other rows, capturing the cases when this advice is valuable. Recall from Definition 4.14 that $\{\alpha_{n,k}\}$ is defined from the stopping rule $\{\gamma_{n,k}\}$, so that when multiplied by $h_{n,k}^+$ or $h_{n,k}^-$ respectively, it equals the probability of encountering (n,k) without necessarily stopping there, in the cases of positive and negative bias respectively.

Lemma 4.23. *Consider the two probability distributions in Proposition 4.22 over locations (n,k) in the Pascal triangle of depth $10^{-8}/\Delta^2$ and bias $p \in \{\frac{1}{2} \pm \Delta\}$, generated by the given stopping rule $\{\gamma_{n,k}\}$ in the two cases 1) a coin with bias $\frac{1}{2} + \Delta$ is used with probability ρ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise versus 2) a coin with bias $\frac{1}{2} + \Delta$ is used with probability $\rho + \epsilon$ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise. If ϵ/ρ is smaller than some universal constant, then the squared Hellinger distance between these two distributions can be written as*

$$\Theta(\epsilon^2) \left[\sum_{n < \frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)^2} \right. \\ \left. + \sum_{n = \frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \frac{\frac{h_{n,k}^+}{\rho} + h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} \right]$$

Intuitively, Lemma 4.23 breaks up the squared Hellinger distance into its contributions from each location (n,k) in the triangle, with the coefficient $\alpha_{n,k}$ depending on the stopping rule (proportional to the algorithm’s probability of stopping at location (n,k)), and the remaining portion of the expression depending only on n,k,Δ,ρ , with the ϵ dependence already factored out in the initial $\Theta(\epsilon^2)$ term.

The rest of the analysis uses the above tools to upper bound the squared Hellinger distance per sample. We defer the concrete details and calculations of the proof of Proposition 4.22 to the next section, Section 4.6. The high level idea of the analysis is to split the expression of Lemma 4.23 for the total squared Hellinger distance per sample into three

components, with the contribution from each location (n, k) assigned to either 1) the last row $n = \frac{10^{-8}}{\Delta^2}$, 2) a “high discrepancy region” where $h_{n,k}^+/h_{n,k}^- \geq 1/\rho^{0.1}$ which is towards the right of the triangle, potentially contributing large amounts to the squared Hellinger distance and 3) a “central” region that is the rest of the triangle. The last row, because of the nature of “advice”, clearly needs its own analysis. As for the rest of the triangle, we divide it into the “central” and “high discrepancy” regions, and bound their contributions to the squared Hellinger distance per sample using different strategies. For the central region, the key insight is that the squared Hellinger distance term is bounded by a well-behaved quadratic function in that region. On the other hand, for the high discrepancy region, the key observation is that the region is defined such that it is a large number of standard deviations away from where a non-stopping random walk on the Pascal triangle should concentrate, and thus it is very unlikely for the algorithm to enter that region. We take additional care to show that, for any stopping rule used by any algorithm, it *cannot* sufficiently skew the distribution of where the walk ends up—for example, while the distribution might skew to the right if the algorithm stops whenever it enters the “left” side of the triangle, we show that this cannot significantly save on expected sample complexity nor substantially increase the squared Hellinger distance per sample. The analysis for the high discrepancy region makes crucial use of our simplifying assumption that the stopping rule *only* stops at powers of 2 coin flips, letting us analyze large sequences of coin flips at a time, where we may take advantage of the tight concentration of the Binomial distribution over sufficiently many coin flips to bound the effect of any skewing-towards-the-right that can be introduced by the stopping rule.

Propositions 4.33, 4.26 and 4.24 in Section 4.6 assert that for each of the respective regions, their contribution to the squared Hellinger distance, divided by the expected sample complexity, is at most $O(\epsilon^2 \Delta^2 / \rho)$. Summing up the three terms is an upper bound on the *total* squared Hellinger distance per expected sample of $O(\epsilon^2 \Delta^2 / \rho)$, completing the proof of Proposition 4.22.

4.6 Proof of Proposition 4.22

This section proves Proposition 4.22, which upper bounds the squared Hellinger distance per sample for any single-coin algorithm of (without loss of generality) a particular form stated in the proposition. We state the proposition again for the reader's convenience, as well as Lemma 4.23 that is introduced in Section 4.5.2, which simplifies the expression of the squared Hellinger distance by sacrificing a constant factor.

Proposition 4.22. *Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that 1) is non-zero only for n that are powers of 2, and 2) $\gamma_{10^{-8}/\Delta^2,k} = 1$ for all k , that is the random walk always stops if it reaches $10^{-8}/\Delta^2$ coin flips. Suppose that given a coin, after a random walk on the Pascal triangle according to the stopping rule, the position (n,k) that the walk ended at is always revealed, and furthermore, if $n = 10^{-8}/\Delta^2$, then the bias of the coin is also revealed. Let H^2 be the squared Hellinger distance between a single run of the above process where 1) a coin with bias $\frac{1}{2} + \Delta$ is used with probability ρ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise versus 2) a coin with bias $\frac{1}{2} + \Delta$ is used with probability $\rho + \epsilon$ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise. Furthermore, let $\mathbb{E}_\rho[n]$ and $\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$ be the expected number of coin flips during a run of the algorithm where we use a $\frac{1}{2} + \Delta$ coin with probability ρ and $\rho + \frac{\epsilon}{2}$ respectively, and a $\frac{1}{2} - \Delta$ coin otherwise. If all of ρ, ϵ, Δ and ϵ/ρ are smaller than some universal absolute constant, then*

$$\max \left[\frac{H^2}{\mathbb{E}_\rho[n]}, \frac{H^2}{\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]} \right] = O \left(\frac{\epsilon^2 \Delta^2}{\rho} \right)$$

As mentioned in Section 4.5.2, the proofs for bounding $\frac{H^2}{\mathbb{E}_\rho[n]}$ and $\frac{H^2}{\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]}$ are essentially identical, so here we present the proof only for the latter.

Lemma 4.23. *Consider the two probability distributions in Proposition 4.22 over locations (n,k) in the Pascal triangle of depth $10^{-8}/\Delta^2$ and bias $p \in \{\frac{1}{2} \pm \Delta\}$, generated by the given stopping rule $\{\gamma_{n,k}\}$ in the two cases of 1) a coin with bias $\frac{1}{2} + \Delta$ is used with probability ρ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise versus 2) a coin with bias $\frac{1}{2} + \Delta$ is used with probability $\rho + \epsilon$ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise. If ϵ/ρ is smaller than some universal constant, then the squared*

Hellinger distance between these two distributions can be written as

$$\Theta(\epsilon^2) \left[\sum_{n < \frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left((\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^- \right) \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)^2} \right. \\ \left. + \sum_{n = \frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left((\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^- \right) \frac{\frac{h_{n,k}^+}{\rho} + h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \right]$$

We perform separate analyses on three regions of the Pascal triangle: 1) the last row $n = \frac{10^{-8}}{\Delta^2}$, 2) a "high discrepancy region" where $h_{n,k}^+/h_{n,k}^- \geq 1/\rho^{0.1}$ which is towards the right of the triangle, potentially contributing large amounts to the squared Hellinger distance and 3) a "central" region that is the rest of the triangle. We shall show that each region contributes small squared Hellinger distance per sample, and thus their sum bounds the total squared Hellinger distance per sample, completing the proof of Proposition 4.22.

We present the three analyses in the order of central region (Section 4.6.1), high discrepancy region (Section 4.6.2) and the last row (Section 4.6.3).

4.6.1 "Central" Region

For the purposes of this section, define $b_{n,k,\rho+\frac{\epsilon}{2}}$ to equal $((\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^-)$, so that $\alpha_{n,k}b_{n,k,\rho+\frac{\epsilon}{2}}$ is the probability of reaching and stopping at location (n, k) under a $\rho + \frac{\epsilon}{2}$ mixture of the two coin types. Further, let $R_{n,k,\rho}$ be defined to equal $\frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)^2}$, which is the contribution of location (n, k) to the squared Hellinger distance *per unit of probability mass that stops there*.

By Lemma 4.23, the contribution to the squared Hellinger distance from the central region of the triangle is bounded by the sum, over this region, of $\epsilon^2 \alpha_{n,k} b_{n,k,\rho+\frac{\epsilon}{2}} R_{n,k,\rho}$.

Proposition 4.24. *For an arbitrary stopping rule, the contribution of the central region to the squared Hellinger distance, divided by the (total) expected sample complexity $\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$ of the walk using a $\rho + \frac{\epsilon}{2}$ mixture of $\frac{1}{2} \pm \Delta$ coins, is at most $O(\epsilon^2 \Delta^2 / \rho)$. Explicitly, with notation for b and R defined in the previous paragraphs, we have*

$$\epsilon^2 \sum_{n < \frac{10^{-8}}{\Delta^2}, k \text{ s.t. } \frac{h_{n,k}^+}{h_{n,k}^-} < \frac{1}{\rho^{0.1}}} \alpha_{n,k} b_{n,k,\rho+\frac{\epsilon}{2}} R_{n,k,\rho} = O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right) \mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$$

Proof. We upper bound this quantity here by instead 1) replacing $R_{n,k,\rho}$ by a similar quantity $\widehat{R}_{n,k,\rho}$ that is an upper bound on R in the central region, and 2) summing over the entire triangle instead of just the central region. Let $\widehat{R}_{n,k,\rho} = 2 \left(\min \left(\frac{h_{n,k}^+}{h_{n,k}^-}, \frac{1}{\rho^{0.1}} \right) - 1 \right)^2$. This bounds R in the region where $h_{n,k}^+ / h_{n,k}^- \leq 1/\rho^{0.1}$: in this regime, $\widehat{R} = \frac{2(h_{n,k}^+ - h_{n,k}^-)^2}{(h_{n,k}^-)^2}$. The numerator of R is at most $\frac{1}{2}$ of the numerator of \widehat{R} , and the denominator of R is at least $\frac{1}{2}$ of the denominator of \widehat{R} .

Thus we instead prove the related fact that

$$\epsilon^2 \sum_{n \leq \frac{10-8}{\Delta^2}, k \in [0..n]} \alpha_{n,k} b_{n,k,\rho+\frac{\epsilon}{2}} \widehat{R}_{n,k,\rho} = O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right) \mathbb{E}_{\rho+\frac{\epsilon}{2}}[n] \quad (4.7)$$

We prove this by induction on a row i , where we define $A_{n,k}^i$ to be the stopping probabilities (corresponding to the product $\alpha_{n,k} b_{n,k,\rho+\frac{\epsilon}{2}}$) for the variant of the given stopping rule where we *force* the rule to stop at row i if it reaches this row; analogously define $\mathbb{E}_{\rho+\frac{\epsilon}{2}}^i[n]$ to be the expected number of samples taken by this rule. We consider how both the left hand side and $\mathbb{E}_{\rho+\frac{\epsilon}{2}}^i[n]$ change as we increase i by 1, and show that the ratio of their change is $O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right)$.

See Lemma 4.25 for a proof of this fact. The proof of the lemma rely on the concrete definitions of $b_{n,k,\rho+\frac{\epsilon}{2}}$ and $R_{n,k,\rho}$, and so both the lemma statement and the proof write out the expressions for purposes of calculations.

As a proof sketch of the ground covered by Lemma 4.25: if for some location (i, k) some amount of probability mass m continues down to row $i + 1$ instead of stopping here, then the expected number of samples increases by exactly m . Meanwhile, this probability mass m will end up split between locations $(i + 1, k)$ and $(i + 1, k + 1)$, where for a coin of bias p (that will be $\frac{1}{2} \pm \Delta$), we will have $m(1 - p)$ mass going left and mp mass going right, contributing to $A_{i+1,k}^{i+1}$ and $A_{i+1,k+1}^{i+1}$ entries respectively. The change in the left hand side of Equation 4.7 induced by sending mass m down to level $i + 1$ is thus expressed as a linear combination of 3 evaluations of the function $\widehat{R}_{n,k,\rho}$. Since $\widehat{R}_{n,k,\rho}$ is essentially a quadratic function of the ratio $\frac{h_{n,k}^+}{h_{n,k}^-}$, this linear combination evaluates to the difference between a quadratic evaluated at 1 point, versus the weighted average of the quadratic at 2 surrounding points, and is bounded by $m \cdot O\left(\frac{\Delta^2}{\rho^{0.2}}\right)$ essentially because of the second derivative of the quadratic in the central region.

□

Lemma 4.25. For any (n, k) ,

$$\begin{aligned} & \epsilon^2 \eta_{n,k} \left[\left(\left(\rho + \frac{\epsilon}{2} \right) h_{n+1,k+1}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n+1,k+1}^- \right) \times 2 \left(\min \left(\frac{h_{n+1,k+1}^+}{h_{n+1,k+1}^-}, \frac{1}{\rho^{0.1}} \right) - 1 \right)^2 \right. \\ & \quad \left. + \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n+1,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n+1,k}^- \right) \times 2 \left(\min \left(\frac{h_{n+1,k}^+}{h_{n+1,k}^-}, \frac{1}{\rho^{0.1}} \right) - 1 \right)^2 \right] \\ & \leq \epsilon^2 \eta_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \left[2 \left(\min \left(\frac{h_{n,k}^+}{h_{n,k}^-}, \frac{1}{\rho^{0.1}} \right) - 1 \right)^2 + O \left(\frac{\Delta^2}{\rho^{0.2}} \right) \right] \end{aligned}$$

Proof. It suffices to show that the left hand side of the inequality is upper bounded by the right hand side, substituting in both options for the minimum. For the $1/\rho^{0.1}$ case, since both summands on the left hand side are upper bounded by the $1/\rho^{0.1}$ case of their expressions, the inequality follows trivially and in fact without the excess term of $O(\Delta^2/\rho^{0.2})$.

We now prove the other case, for which it is sufficient to show that

$$\begin{aligned} & \epsilon^2 \eta_{n,k} \left[\left(\left(\rho + \frac{\epsilon}{2} \right) h_{n+1,k+1}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n+1,k+1}^- \right) \times 2 \left(\frac{h_{n+1,k+1}^+}{h_{n+1,k+1}^-} - 1 \right)^2 \right. \\ & \quad \left. + \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n+1,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n+1,k}^- \right) \times 2 \left(\frac{h_{n+1,k}^+}{h_{n+1,k}^-} - 1 \right)^2 \right] \\ & \leq \epsilon^2 \eta_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \left[2 \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1 \right)^2 + O \left(\frac{\Delta^2}{\rho^{0.2}} \right) \right] \end{aligned}$$

when $h_{n,k}^+/h_{n,k}^- \leq 1/\rho^{0.1}$.

In turn, we can break this inequality into a conjunction of two inequalities, that

$$h_{n+1,k+1}^+ \left(\frac{h_{n+1,k+1}^+}{h_{n+1,k+1}^-} - 1 \right)^2 + h_{n+1,k}^+ \left(\frac{h_{n+1,k}^+}{h_{n+1,k}^-} - 1 \right)^2 \leq h_{n,k}^+ \left[\left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1 \right)^2 + O \left(\frac{\Delta^2}{\rho^{0.2}} \right) \right]$$

and

$$h_{n+1,k+1}^- \left(\frac{h_{n+1,k+1}^+}{h_{n+1,k+1}^-} - 1 \right)^2 + h_{n+1,k}^- \left(\frac{h_{n+1,k}^+}{h_{n+1,k}^-} - 1 \right)^2 \leq h_{n,k}^- \left[\left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1 \right)^2 + O \left(\frac{\Delta^2}{\rho^{0.2}} \right) \right]$$

again assuming that $h_{n,k}^+/h_{n,k}^- \leq 1/\rho^{0.1}$.

For the first inequality, observe that

$$\frac{h_{n+1,k+1}^+}{h_{n+1,k+1}^-} = \frac{\frac{1}{2} + \Delta h_{n,k}^+}{\frac{1}{2} - \Delta h_{n,k}^-} \quad \text{and} \quad \frac{h_{n+1,k}^+}{h_{n+1,k}^-} = \frac{\frac{1}{2} - \Delta h_{n,k}^+}{\frac{1}{2} + \Delta h_{n,k}^-}$$

and also $h_{n+1,k+1}^+ = h_{n,k}^+(\frac{1}{2} + \Delta)$ and $h_{n+1,k}^+ = h_{n,k}^+(\frac{1}{2} - \Delta)$. We therefore factor out and drop the $h_{n,k}^+$ on both sides, simplify, and reduce to showing that

$$\left(\frac{1}{2} + \Delta\right) \left(\frac{\frac{1}{2} + \Delta h_{n,k}^+}{\frac{1}{2} - \Delta h_{n,k}^-} - 1\right)^2 + \left(\frac{1}{2} - \Delta\right) \left(\frac{\frac{1}{2} - \Delta h_{n,k}^+}{\frac{1}{2} + \Delta h_{n,k}^-} - 1\right)^2 \leq \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1\right)^2 + O\left(\frac{\Delta^2}{\rho^{0.2}}\right)$$

The left hand side is

$$\begin{aligned} & \left(\frac{1}{2} + \Delta\right) \left(\frac{\frac{1}{2} + \Delta h_{n,k}^+}{\frac{1}{2} - \Delta h_{n,k}^-} - 1\right)^2 + \left(\frac{1}{2} - \Delta\right) \left(\frac{\frac{1}{2} - \Delta h_{n,k}^+}{\frac{1}{2} + \Delta h_{n,k}^-} - 1\right)^2 \\ &= \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 \left(\frac{(\frac{1}{2} + \Delta)^3}{(\frac{1}{2} - \Delta)^2} + \frac{(\frac{1}{2} - \Delta)^3}{(\frac{1}{2} + \Delta)^2}\right) - 2 \frac{h_{n,k}^+}{h_{n,k}^-} \left(\frac{(\frac{1}{2} + \Delta)^2}{\frac{1}{2} - \Delta} + \frac{(\frac{1}{2} - \Delta)^2}{\frac{1}{2} + \Delta}\right) + 1 \\ &= \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 (1 + O(\Delta^2)) - 2 \frac{h_{n,k}^+}{h_{n,k}^-} \left(\frac{(\frac{1}{2} + \Delta)^2}{\frac{1}{2} - \Delta} + \frac{(\frac{1}{2} - \Delta)^2}{\frac{1}{2} + \Delta}\right) + 1 \\ &\leq \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 (1 + O(\Delta^2)) - 2 \frac{h_{n,k}^+}{h_{n,k}^-} + 1 \\ &= \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1\right)^2 + O(\Delta^2) \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 \\ &\leq \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1\right)^2 + O\left(\frac{\Delta^2}{\rho^{0.2}}\right) \end{aligned}$$

where the last inequality holds again because we have $h_{n,k}^+/h_{n,k}^- \leq 1/\rho^{0.1}$ by our case analysis.

For the second inequality, via similar reasoning as above, we only need to show that

$$\left(\frac{1}{2} - \Delta\right) \left(\frac{\frac{1}{2} + \Delta h_{n,k}^+}{\frac{1}{2} - \Delta h_{n,k}^-} - 1\right)^2 + \left(\frac{1}{2} + \Delta\right) \left(\frac{\frac{1}{2} - \Delta h_{n,k}^+}{\frac{1}{2} + \Delta h_{n,k}^-} - 1\right)^2 \leq \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1\right)^2 + O\left(\frac{\Delta^2}{\rho^{0.2}}\right)$$

The left hand side is

$$\begin{aligned} & \left(\frac{1}{2} - \Delta\right) \left(\frac{\frac{1}{2} + \Delta h_{n,k}^+}{\frac{1}{2} - \Delta h_{n,k}^-} - 1\right)^2 + \left(\frac{1}{2} + \Delta\right) \left(\frac{\frac{1}{2} - \Delta h_{n,k}^+}{\frac{1}{2} + \Delta h_{n,k}^-} - 1\right)^2 \\ &= \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 \left(\frac{(\frac{1}{2} + \Delta)^2}{\frac{1}{2} - \Delta} + \frac{(\frac{1}{2} - \Delta)^2}{\frac{1}{2} + \Delta}\right) - 2 \frac{h_{n,k}^+}{h_{n,k}^-} \left(\frac{1}{2} + \Delta + \frac{1}{2} - \Delta\right) + 1 \\ &\leq \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 (1 + O(\Delta^2)) - 2 \frac{h_{n,k}^+}{h_{n,k}^-} + 1 \\ &= \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1\right)^2 + O(\Delta^2) \left(\frac{h_{n,k}^+}{h_{n,k}^-}\right)^2 \end{aligned}$$

$$\leq \left(\frac{h_{n,k}^+}{h_{n,k}^-} - 1 \right)^2 + O\left(\frac{\Delta^2}{\rho^{0.2}} \right)$$

with reasoning as in the previous inequality, thus completing the proof of the lemma. \square

4.6.2 "High Discrepancy" Region

Proposition 4.26. *Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that 1) is non-zero only for n that are powers of 2, and 2) $\gamma_{10^{-8}/\Delta^2,k} = 1$ for all k , that is the random walk always stops after $10^{-8}/\Delta^2$ coin flips. Let*

$$H_{disc}^2 = \Theta(\epsilon^2) \sum_{n < \frac{10^{-8}}{\Delta^2}, k \text{ s.t. } \frac{h_{n,k}^+}{h_{n,k}^-} \geq \frac{1}{\rho^{0.1}}} \alpha_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)^2}$$

be the contribution to the squared Hellinger distance by the "high discrepancy" region. Furthermore, again let $\mathbb{E}_{\rho + \frac{\epsilon}{2}}[n]$ be the expected number of coin flips on this random walk, where we use a $\frac{1}{2} + \Delta$ coin with probability $\rho + \frac{\epsilon}{2}$ (instead of ρ or $\rho + \epsilon$), and a $\frac{1}{2} - \Delta$ coin otherwise. If all of ρ , ϵ , Δ and ϵ/ρ are smaller than some universal absolute constant, then

$$\frac{H_{disc}^2}{\mathbb{E}_{\rho + \frac{\epsilon}{2}}[n]} = O\left(\frac{\epsilon^2 \Delta^2}{\rho} \right)$$

The key observation for this section is that the "high discrepancy" region is in fact at least $\Omega(\log \frac{1}{\rho})$ standard deviations away from where a random walk on the triangle (without a stopping rule) would concentrate; and thus it is very unlikely for the random walk to enter the region. However, the existence of a stopping rule could potentially skew the distribution of the random walk on each row towards the "high discrepancy" side of the triangle, while saving on sample complexity by stopping early whenever the walk enters the other side of the triangle. In this section, we essentially show that this cannot happen.

The analysis in this section relies on our assumption that the stopping rule only stops at rows that are powers of 2 (unlike the analysis of the previous section). Intuitively, if the random walk ends up very far to the right, then there must be a single region of rows $[2^i .. 2^{i+1}]$ where, without any stopping rule on intermediate rows to guide it, the walk still somehow makes unlikely progress to the right. More explicitly, if the distribution of

reaching-and-not-stopping-at row 2^{i+1} is skewed significantly far to the right of the distribution of reaching-and-not-stopping-at row 2^i (despite the intervening process being strictly a binomially distributed random walk), then the only way this could have occurred is if an overwhelming fraction of the probability mass reaching row 2^{i+1} stops there. Namely, if probability mass m emerges below row 2^{i+1} and skewed far to the right, the potential Hellinger distance gains this induces will be more than counterbalanced by the huge addition to sample complexity induced by the overwhelming (relative to m) probability of stopping at row 2^{i+1} .

We utilize the following fact, essentially a consequence of a Binomial distribution being upper bounded by a corresponding Gaussian.

Fact 4.27. *Let $\text{Bin}(n, p, k)$ denote the probability that a Binomial distribution with n trials and bias p has value k . If Δ is sufficiently small, then there exists some absolute constant C such that for all $n \geq 1$, and for both $\frac{1}{2} + \Delta$ and $\frac{1}{2} - \Delta$ substituted in the expression “ $\frac{1}{2} \pm \Delta$ ” below,*

$$\sum_{k \in [0..n]} e^{\frac{(k - (\frac{1}{2} \pm \Delta)n)^2}{n}} \text{Bin}(n, \frac{1}{2} \pm \Delta, k) \leq C$$

The sum of the pointwise products of the Binomial pmf and the inverse Gaussian can instead be re-expressed as the evaluation of a convolution between corresponding functions evaluated at a single point. We express this straightforward corollary below, and use it crucially in this section and the next.

Fact 4.28. *Consider the sequences $f_{n,k}^+(m) = e^{\frac{(k - (\frac{1}{2} + \Delta)n - m)^2}{n}}$ for $m \in \mathbb{Z}$, and $f_{n,k}^-(m) = e^{\frac{(k - (\frac{1}{2} - \Delta)n - m)^2}{n}}$ for $m \in \mathbb{Z}$. Let $\text{Bin}(n, p)$ be the pmf of the Binomial distribution with n trials and bias p . If Δ is sufficiently small, then there exists some absolute constant C such that for all $n \geq 1$ and all k ,*

$$(f_{n,k}^+ * \text{Bin}(n, \frac{1}{2} + \Delta))(k) \leq C$$

and

$$(f_{n,k}^- * \text{Bin}(n, \frac{1}{2} - \Delta))(k) \leq C$$

To start lower bounding the expected sample complexity of the random walk, we start with the following two lemmas stating that if there is probability c of reaching a right tail on a particular power-of-2 row, then there must be a tail on the previous power-of-2 row

that the walk has high probability reaching. These are formalized as Lemma 4.29 and 4.30 for $\frac{1}{2} + \Delta$ coins and $\frac{1}{2} - \Delta$ coins respectively. The crux of the arguments are (weighted) averaging arguments based on Fact 4.28.

Lemma 4.29. *Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that is non-zero only for n that are powers of 2. For a coin with bias $\frac{1}{2} + \Delta$, suppose at row 2^j there is some position $k \in [(\frac{1}{2} + \Delta)2^j..2^j]$ such that the total probability mass of the random walk reaching positions $\geq k$ at row 2^j is at least c . Then, there must be some position $k' \in [0..2^{j-1}]$ at row 2^{j-1} such that the probability of reaching positions $\geq k'$ at that row is at least $\frac{c}{C} \cdot f_{2^{j-1},k}^+(k') = \frac{c}{C} \cdot e^{\frac{(k - (\frac{1}{2} + \Delta)2^{j-1} - k')^2}{2^{j-1}}}$, where the constant C and the function $f_{n,k}^+$ are defined in Fact 4.28.*

Proof. Let us denote by D_n^\downarrow the vector (over $k \in [0..n]$) of probabilities that the random walk using a coin of bias $\frac{1}{2} + \Delta$ reaches but does not stop at the location (n, k) . Similarly, let us denote by D_n the vector (over $k' \in [0..n]$) of probabilities that the random walk using a $\frac{1}{2} + \Delta$ coin reaches the location (n, k) (and can either stop at or leave the location).

Consider the vector I that is 1 for all coordinates ≤ 0 , and 0 otherwise. Then for any vector v , $(v * I)(k) = \sum_{i \leq k} v(i)$, using "*" to denote convolution.

Assume for the sake of contradiction that the statement is false, namely that for all $k' \in [0..2^{j-1}]$, $(D_{2^{j-1}} * I)(k') < \frac{c}{C} \cdot f_{2^{j-1},k}^+(k')$. Then, since $D_{2^{j-1}}^\downarrow \leq D_{2^{j-1}}$ pointwise, we have for all $k' \in [0..2^{j-1}]$, $(D_{2^{j-1}}^\downarrow * I)(k') < \frac{c}{C} \cdot f_{2^{j-1},k}^+(k')$. Observe that $D_{2^{j-1}}^\downarrow * I$ is constant for all coordinates ≤ 0 , and that $f_{2^{j-1},k}^+$ is a decreasing function in the same region if $k \in [(\frac{1}{2} + \Delta)2^j..2^j]$ (as in the lemma assumption), and therefore $D_{2^{j-1}}^\downarrow * I < f_{2^{j-1},k}^+$ also for that region since the inequality holds at coordinate 0. As for coordinates $> 2^j$, $D_{2^{j-1}}^\downarrow * I$ is 0, whilst $f_{2^{j-1},k}^+$ is strictly positive. It follows that the inequality also holds for coordinates $> 2^j$, and thus it holds everywhere.

From this, using the commutativity of convolution, we have

$$\begin{aligned} D_{2^j} * I &= \left(D_{2^{j-1}}^\downarrow * \text{Bin}(2^{j-1}, \frac{1}{2} + \Delta) \right) * I \\ &= (D_{2^{j-1}}^\downarrow * I) * \text{Bin}(2^{j-1}, \frac{1}{2} + \Delta) \\ &< f_{2^{j-1},k}^+ * \text{Bin}(2^{j-1}, \frac{1}{2} + \Delta) \end{aligned}$$

which holds pointwise, in particular at coordinate k . However, $(D_{2^j} * I)(k) = c$ by assumption, but $f_{2^{j-1},k}^+ * \text{Bin}(2^{j-1}, \frac{1}{2} + \Delta)(k) \leq c$ by Fact 4.28, which is a contradiction. \square

Lemma 4.30. Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that is non-zero only for n that are powers of 2. For a coin with bias $\frac{1}{2} - \Delta$, suppose at row 2^j there is some position $k \in [(\frac{1}{2} - \Delta)2^j .. 2^j]$ such that the total probability mass of the random walk reaching positions $\geq k$ at row 2^j is at least c . Then, there must be some position $k' \in [0..2^{j-1}]$ at row 2^{j-1} such that the probability of reaching positions $\geq k'$ at that row is at least $\frac{c}{C} \cdot f_{2^{j-1},k}^-(k') = \frac{c}{C} \cdot e^{\frac{(k - (\frac{1}{2} - \Delta)2^{j-1} - k')^2}{2^{j-1}}}$, where the constant C and the function $f_{n,k}^-$ are defined in Fact 4.28.

Proof. The proof is completely analogous to that of Lemma 4.29. \square

In order to conclude the sample complexity lower bound corresponding to a particular row, we need the following lemma saying that, if we repeatedly apply Lemma 4.29 (or Lemma 4.30), then some row 2^j will have a large probability of stopping at that row, which will contribute a large amount to the overall sample complexity. Further, when 2^j is smaller (corresponding to fewer samples taken before stopping), the probability bound induced by the following lemma will be correspondingly higher, so that the product of the row and its stopping probability (i.e., a lower bound on total sample complexity) will be high for the j produced by the lemma.

Lemma 4.31. Consider an arbitrary sequence of numbers $\{g_j\}_{j \in [0..J]}$ such that $\sum_j g_j = K$. Let r_j be chosen arbitrarily such that $r_j \geq \frac{1}{C} e^{g_j^2/2^j}$, where C is the constant in Fact 4.28 and let $\pi_j = \prod_{i=j}^J r_i$. Furthermore suppose that $K^2 \geq 100 \log(2C) \cdot 2^J$. Then there exists $j \in [0..J]$ such that $\pi_j 2^{j-J} \geq e^{0.01 K^2/2^j}$.

Proof. Taking logarithms and rearranging, we see that it suffices to show the existence of j such that $(j - J - 1) \log(2C) + \sum_{i=j}^J \frac{g_i^2}{2^i} \geq 0.01 \frac{K^2}{2^j}$.

The sequence $\frac{K}{5} \cdot 0.8^{J-j}$ for $j \in [0..J]$ sums up to less than K . Since $\sum_j g_j = K$, there must exist a j such that $g_j \geq \frac{K}{5} 0.8^{J-j}$. Therefore, $\frac{g_j^2}{2^j} \geq \frac{K^2}{25} \frac{0.64^{J-j}}{2^j} = \frac{K^2}{25} \frac{1.28^{J-j}}{2^j}$.

It suffices to show that $\frac{K^2}{25} \frac{1.28^{J-j}}{2^j} \geq 0.01 \frac{K^2}{2^j} + (J - j + 1) \log(2C)$. It is easy to check that a sufficient condition is $K^2/2^j \geq 100 \log(2C)$, as assumed in the lemma statement; thus we conclude the above inequality for all $j \in [0..J]$. \square

Now we use Lemmas 4.29, 4.30 and 4.31 to prove the sample complexity lower bound corresponding to a particular row (Lemma 4.32). Afterwards we shall combine these bounds across all possible power-of-2 rows to prove Proposition 4.26.

Lemma 4.32. Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that is non-zero only for n that are powers of 2. For a mixture coin that has bias $\frac{1}{2} + \Delta$ with probability $\rho + \frac{\epsilon}{2}$ and bias $\frac{1}{2} - \Delta$ otherwise, suppose at row 2^{J+1} there is some position $k \in [(\frac{1}{2} + \Delta)2^{J+1}..2^{J+1}]$ such that the probability mass of the random walk reaching positions $\geq k$ at row 2^{J+1} is c . If $k \geq (\frac{1}{2} + \Delta)2^{J+1} + \sqrt{100 \log(2C)}2^{\frac{J}{2}} + 1$, then the expected sample complexity of a single random walk using the above mixture coin is at least $2^{J-1} \cdot c \cdot e^{0.01(k - (\frac{1}{2} + \Delta)2^{J+1})^2/2^J}$.

We point out that the restriction on k (that it lies at least a constant number of standard deviations to the right of its mean) includes the entire high discrepancy region, as analyzed in this section, and further includes all of the larger yet analogous region for the analysis of the last row in the next section.

Proof of Lemma 4.32. The probability of the random walk reaching positions $\geq k$ at row 2^{J+1} using a mixture coin is the sum of $\rho + \frac{\epsilon}{2}$ times such probability of the random walk using a $\frac{1}{2} + \Delta$ coin and $1 - \rho - \frac{\epsilon}{2}$ times such probability of the random walk using a $\frac{1}{2} - \Delta$ coin. Since the total probability of this walk reaching positions $\geq k$ equals c , at least half this probability must come from one of the two coin types. Explicitly, at least one of the following two statements has to be true: 1) the probability that the random walk using a coin with $\frac{1}{2} + \Delta$ bias reaches positions $\geq k$ at row 2^{J+1} is at least $c/(2\rho + \epsilon)$, or 2) the same probability but using a $\frac{1}{2} - \Delta$ coin instead is at least $c/(2 - 2\rho - \epsilon)$.

For case 1, we repeatedly apply Lemma 4.29 to generate a sequence of $\{k_j\}$ from $j = J$ backwards (and $k_{J+1} = k$), until $k_{j^*} < (\frac{1}{2} + \Delta)2^{j^*}$ or $j^* = 0$. By induction, the probability of reaching positions $\geq k_j$ at row 2^j is at least $\frac{c}{2\rho + \epsilon} \cdot \prod_{i=j}^J \frac{1}{C} e^{\frac{(k_{i+1} - (\frac{1}{2} + \Delta)2^i) - k_i)^2}{2^i}}$. We would now apply Lemma 4.31 with $g_i = k_{i+1} - k_i - (\frac{1}{2} + \Delta)2^i$ for $i \geq j^*$, and $g_i = 0$ for $i < j^*$, noting that K in that lemma that we get is $K = \sum_{i=j^*}^J k_{i+1} - k_i - (\frac{1}{2} + \Delta)2^i \geq k_{J+1} - k_{j^*} - \sum_{i=j^*}^J (\frac{1}{2} + \Delta)2^i > k_{J+1} (= k) - (\frac{1}{2} + \Delta)2^{J+1} - 1$ since $k_{j^*} < (\frac{1}{2} + \Delta)2^{j^*}$ if $j^* > 0$ and $k_0 \leq 1$ when $j^* = 0$. Since we assumed in the lemma statement that $k \geq (\frac{1}{2} + \Delta)2^{J+1} + \sqrt{100 \log(2C)}2^{\frac{J}{2}} + 1$, we have $K^2/2^J \geq 100 \log(2C)$.

Therefore, as a result of applying Lemma 4.31, we know that there exists j such that

$$2^{j-J} \prod_{i=j}^J \frac{1}{C} e^{\frac{(k_{i+1} - (\frac{1}{2} + \Delta)2^i) - k_i)^2}{2^i}} \geq e^{0.01(k - (\frac{1}{2} + \Delta)2^{J+1})^2/2^J}$$

Thus in case 1, we multiply the left hand side by $c/(2\rho + \epsilon)2^J$ to give a lower bound

on the expected sample complexity of the random walk, using a $\frac{1}{2} + \Delta$ coin. We thus use the above inequality to conclude a lower bound of $2^J \frac{c}{2\rho+\epsilon} e^{0.01(k - (\frac{1}{2} + \Delta)2^{J+1})^2/2^J}$ for the expected sample complexity conditioned on a $\frac{1}{2} + \Delta$ coin. Since the mixture coin has probability $\rho + \frac{\epsilon}{2}$ of being a $\frac{1}{2} + \Delta$ coin, the lemma statement follows.

The proof for case 2 is completely analogous, using Lemma 4.30 instead of Lemma 4.29, and noting that $k - (\frac{1}{2} - \Delta)2^{J+1} \geq k - (\frac{1}{2} + \Delta)2^{J+1} \geq 0$. \square

Equipped with Lemma 4.32, we prove Proposition 4.26.

Proof of Proposition 4.26. The general strategy is to show using Lemma 4.32 that, for each row (from 1 to $10^{-8}/\Delta^2$), if there is some probability c_J for the random walk to reaching the high discrepancy region, then: 1) the total expected sample complexity must be large, and 2) by Lemma 4.23, if there is probability c_J of reaching the high discrepancy region at row 2^J , then the contribution to the squared Hellinger distance by the high discrepancy region at row 2^J is upper bounded by $\Theta(c_J \epsilon^2 / \rho^2)$. Thus the squared Hellinger distance per sample complexity for the high discrepancy region of each row is small, and our bounds are in fact strong enough for us to simply take a union bound over the rows and lose by no more than a constant factor. We now formalize the above argument.

Consider the rows 2^{J+1} for $J \in [-1, (\log_2 \frac{10^{-8}}{\Delta^2}) - 1]$. Recall that the high discrepancy region consists of coordinates $k \in [0, 2^{J+1}]$ such that $h_{2^{J+1},k}^+ / h_{2^{J+1},k}^- \geq 1/\rho^{0.1}$. Observe that

$$\frac{h_{2^{J+1},k}^+}{h_{2^{J+1},k}^-} = \left(\frac{1 + 2\Delta}{1 - 2\Delta} \right)^{2k - 2^{J+1}}$$

and therefore the high discrepancy region consists of k such that $2k - 2^{J+1} \geq \frac{.1 \log \frac{1}{\rho}}{\log \frac{1+2\Delta}{1-2\Delta}}$, implying that

$$k \geq 2^J + \frac{.1 \log \frac{1}{\rho}}{\log \frac{1+2\Delta}{1-2\Delta}} \geq \frac{1}{2} 2^{J+1} + \frac{.099 \log \frac{1}{\rho}}{4\Delta}$$

Furthermore, since $J \leq (\log_2 \frac{10^{-8}}{\Delta^2}) - 1$, we have $2^J \leq \frac{0.01}{2\Delta^2}$, which for sufficiently small ρ and Δ (both smaller than some absolute constant, with no requirements on how they depend on each other) means that $\frac{.099 \log \frac{1}{\rho}}{4\Delta} \geq \Delta 2^{J+1} + \sqrt{100 \log(2C) 2^J} + 1$. Thus the coordinates k in the high discrepancy region always satisfy the precondition of Lemma 4.32.

Now note that for sufficiently small ρ (smaller than some absolute constant),

$$\left(k - \left(\frac{1}{2} + \Delta\right)2^{J+1}\right)^2 \geq \left(\frac{.098 \log \frac{1}{\rho}}{4\Delta}\right)^2 \geq \frac{10^{-8}(\log \frac{1}{\rho})^2}{\Delta^2}$$

Therefore, if the probability of the random walk using a random coin reaches the high discrepancy region at row 2^{J+1} is c_{J+1} , then by Lemma 4.32, the total expected sample complexity of the random walk must be at least $2^{J-1} \cdot c_{J+1} \cdot e^{0.01 \frac{10^{-8}(\log \frac{1}{\rho})^2}{\Delta^2 \cdot 2^J}}$.

We can now upper bound the ratio between the high discrepancy region contribution to the squared Hellinger distance and the total expected sample complexity of the random walk by

$$\begin{aligned} & \frac{\sum_{J \in [-1, (\log_2 \frac{10^{-8}}{\Delta^2}) - 1]} \Theta\left(\frac{\epsilon^2}{\rho^2}\right) c_{J+1}}{\mathbb{E}_{\rho + \frac{\epsilon}{2}}[n]} \\ &= \Theta\left(\frac{\epsilon^2}{\rho^2}\right) \sum_{J \in [-1, (\log_2 \frac{10^{-8}}{\Delta^2}) - 1]} \frac{c_{J+1}}{\mathbb{E}_{\rho + \frac{\epsilon}{2}}[n]} \\ &\leq \Theta\left(\frac{\epsilon^2}{\rho^2}\right) \sum_{J \in [-1, (\log_2 \frac{10^{-8}}{\Delta^2}) - 1]} \frac{c_{J+1}}{2^{J-1} \cdot c_{J+1} \cdot e^{0.01 \frac{10^{-8}(\log \frac{1}{\rho})^2}{\Delta^2 \cdot 2^J}}} \\ &= \Theta\left(\frac{\epsilon^2}{\rho^2}\right) \sum_{J \in [-1, (\log_2 \frac{10^{-8}}{\Delta^2}) - 1]} \frac{\rho^{0.02 \cdot \log \frac{1}{\rho} \cdot \frac{10^{-8}}{2\Delta^2 \cdot 2^J}}}{2^{J-1}} \\ &\leq \Theta\left(\frac{\epsilon^2}{\rho^2}\right) \sum_{J \in [-1, (\log_2 \frac{10^{-8}}{\Delta^2}) - 1]} \frac{2^{-\frac{10^{-8}}{2\Delta^2 \cdot 2^J}} \rho^{0.02 \cdot \log \frac{1}{\rho}}}{2^{J-1}} \quad \text{since for sufficiently small } \rho, \text{ we have } \rho^{0.02 \cdot \log \frac{1}{\rho}} < \frac{1}{2} \\ &= \Theta\left(\frac{\epsilon^2 \Delta^2 \rho^{0.02 \log \frac{1}{\rho}}}{\rho^2}\right) \quad \text{as the sum is bounded by } O(\rho^{0.02 \cdot \log \frac{1}{\rho}} \Delta^2) \\ &= O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right) \end{aligned}$$

□

4.6.3 The Last Row

We lastly analyze the squared Hellinger distance contribution from the last row of the triangle.

Proposition 4.33. Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ that 1) is non-zero only for n that are powers of 2, and 2) $\gamma_{10^{-8}/\Delta^2,k} = 1$ for all k , that is the random walk always stops after $10^{-8}/\Delta^2$ coin flips. Let

$$H_{last}^2 = \Theta(\epsilon^2) \sum_{n=\frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left((\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^- \right) \frac{\frac{h_{n,k}^+}{\rho} + h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-}$$

be the contribution of the squared Hellinger distance from the last row of the triangle, namely row $10^{-8}/\Delta^2$. Furthermore, again let $\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$ be the expected number of coin flips on this random walk, where we use a $\frac{1}{2} + \Delta$ coin with probability $\rho + \frac{\epsilon}{2}$ (instead of ρ or $\rho + \epsilon$), and a $\frac{1}{2} - \Delta$ coin otherwise. If all of ρ , ϵ , Δ and ϵ/ρ are smaller than some universal absolute constant, then

$$\frac{H_{last}^2}{\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]} = O\left(\frac{\epsilon^2 \Delta^2}{\rho}\right)$$

The squared Hellinger distance contribution from the last row has a different form from the rest of the triangle, and can be large even outside the previously “high discrepancy” region. While the term $(\frac{h_{n,k}^+}{\rho} + h_{n,k}^-)/(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)$ is still upper bounded by $1/\rho^2$ everywhere, it may be as large as $\Theta(1/\rho)$ even in when $h_{n,k}^+/h_{n,k}^- = \Theta(1)$. The intuition for this section is again that despite having a stopping rule that may have subtle effects on the distribution, it is impossible to skew the distribution of the random walk so much that it appears mostly in the “high discrepancy” side of the triangle. We shall use Lemma 4.32 again along with a case analysis and a weighted averaging argument to show sample complexity lower bounds, which lets us upper bound the squared Hellinger distance contribution per expected sample, as required.

Proof. We separate the last row again into a “high discrepancy” region and a “central” region, but with a different criterion: whether

$$\frac{\frac{h_{n,k}^+}{\rho} + h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \geq \frac{C}{\rho}$$

where C is the constant specified in Fact 4.28. The criterion can be equivalently stated as whether $h_{n,k}^+/h_{n,k}^- \geq r$ for some $r = \Theta(1)$.

For the “central” region, suppose there is probability $c_{n,k}$ of reaching position (n, k) in that region for the random walk that uses a $\rho + \frac{\epsilon}{2}$ mixture random coin. Consider an

alternate form of the squared Hellinger distance contribution that is within a constant factor of that presented in the proposition statement, assuming that ϵ/ρ is small:

$$\Theta(\epsilon^2) \sum_{n=\frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left(\frac{h_{n,k}^+}{\rho} + h_{n,k}^- \right)$$

This approximation holds since $\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-$ and $(\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^-$ are within constant factors of each other. Note that $c_{n,k} = \alpha_{n,k}((\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^-)$, and so when $h_{n,k}^+/h_{n,k}^- \leq r = \Theta(1)$, we have both $\alpha_{n,k}h_{n,k}^+ = O(c)$ and $\alpha_{n,k}h_{n,k}^- = O(c)$. Thus the squared Hellinger contribution of this location is upper bounded by $O(c\epsilon^2/\rho)$, yet the total sample complexity is lower bounded by $\Omega(cn) = \Omega(c/\Delta^2)$, giving a fraction that is $O(\epsilon^2\Delta^2/\rho)$.

For the ‘‘central’’ region, recall that it consists of the locations where

$$\frac{\frac{h_{n,k}^+}{\rho} + h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \quad (4.8)$$

ranges from $\frac{c}{\rho}$ to $\frac{1}{\rho^2}$. We separate this region into $O(\log \frac{1}{\rho})$ buckets delimited by consecutive powers of 2.

Suppose there is probability c_{disc} of the $\rho + \frac{\epsilon}{2}$ mixture coin random walk entering the ‘‘high discrepancy’’ region in the last row. Note that the geometric sequence $1, 0.8, 0.64, \dots$ converges to 5, and therefore the sequence $\{\frac{c_{\text{disc}}}{5}0.8^i\}_i$ sums to c_{disc} . If we took only the first $O(\log \frac{1}{\rho})$ many terms, they sum to strictly less than c_{disc} . By a standard averaging argument, there must exist a bucket such that the probability of reaching locations in that bucket i is greater than $\frac{c_{\text{disc}}}{5}0.8^i$. Note that the values (Equation 4.8) inside bucket i range from $C \cdot 2^i/\rho$ to $2 \cdot C \cdot 2^i/\rho$. It is possible to calculate that the locations k within bucket i satisfy $k \geq \frac{n}{2} + \frac{1}{16} \frac{i \log C}{\Delta}$, where n again is $10^{-8}/\Delta^2$. For a sufficiently small Δ , this lower bound in location is at least $(\frac{1}{2} + \Delta)n + \sqrt{100 \log(2C) \cdot n} + 1$, and therefore we can apply Lemma 4.32.

Furthermore, the locations are also at least $(i \cdot 10^{-2} \log C)/\Delta$ away from $(\frac{1}{2} + \Delta)n$, and so Lemma 4.32 guarantees a sample complexity of at least $\Theta(1/\Delta^2) \times \frac{c_{\text{disc}}}{5}0.8^i \times e^{0.01((i \cdot 10^{-2} \log C)/\Delta)^2 \cdot \Delta^2 \times 10^8} \geq \Theta(1/\Delta^2) \times \frac{c_{\text{disc}}}{5}0.8^i \times e^{100(i \log C)^2} \geq \Omega(2^i c_{\text{disc}}/\Delta^2)$, where the last inequality is true because the large exponential term has a logarithm that is quadratic in i and with a base that is a lot greater than $1/0.8$.

The squared Hellinger distance contribution from the ‘‘high discrepancy’’ region in the last row is upper bounded by $O(c_{\text{disc}}\epsilon^2 2^i/\rho)$, and we have shown a sample complexity

lower bound of $\Omega(2^i c_{\text{disc}}/\Delta^2)$. We therefore conclude a fraction of $O(\epsilon^2 \Delta^2/\rho)$ for the squared Hellinger distance per expected sample, for contributions from the “high discrepancy” region in the last row.

Summarizing, both the “high discrepancy” and “central” region contribute no more than $O(\epsilon^2 \Delta^2/\rho)$ times $\mathbb{E}_{\rho+\frac{\epsilon}{2}}[n]$ to H_{last}^2 . Therefore, the proposition follows from summing the two contributions. \square

4.7 Experimental Results

We give simulation results to demonstrate the practical efficacy of our proposed algorithm. In our experimental setups, we compare the convergence rates of 1) our algorithm (“T-WALK (15)” on the plots), 2) the natural majority vote method mentioned in the Introduction (“VOTING” on the plots) and 3) the “SWITCH” method proposed in previous work by Chung et al. [21] which has been observed to perform well in practice, but does not have a theoretical analysis. For our algorithm, we choose the maximum number of flips for a single coin to be 15 ($= c \log \frac{1}{\epsilon}$) in Algorithm 2. We also make the assumption that the noise parameter satisfies $\Delta \geq 0.3$, meaning that we can use Algorithm 2 directly instead of using Algorithm 3 to simulate virtual coins before feeding them into Algorithm 2. To further improve the practical performance of Algorithm 2, we ran a local search method to improve on the non-zero output coefficients in Step 3(d) of Algorithm 2, using the assumption that $\Delta \geq 0.3$. Concretely, recall that we output a non-zero coefficient when the maximum number of coin flips (15) has occurred and the majority of coin flips has been heads. Thus for $k \in \{8, \dots, 15\}$, we output 8: 0, 9: 6.913, 10: 5.032, 11: 2.101, 12: 0.636, 13: 1.965, 14: 1.016, 15: 1.009. We note again that these coefficients are *reusable* in practice, as long as the $\Delta \geq 0.3$ assumption can be made.

Figure 4.1 presents the experimental results, for “coin quality” $\Delta = 0.3$ or 0.4 , and ground truth fraction of positive coins ρ taking representative values $0.005, 0.01, 0.03$, or 0.1 . For each plot, the x -axis corresponds to the number of coin flips, with all algorithms eventually converging to the ground truth for enough coin flips. Standard deviation bars are computed over 10 runs of each different setting. The estimates, given a strict budget of

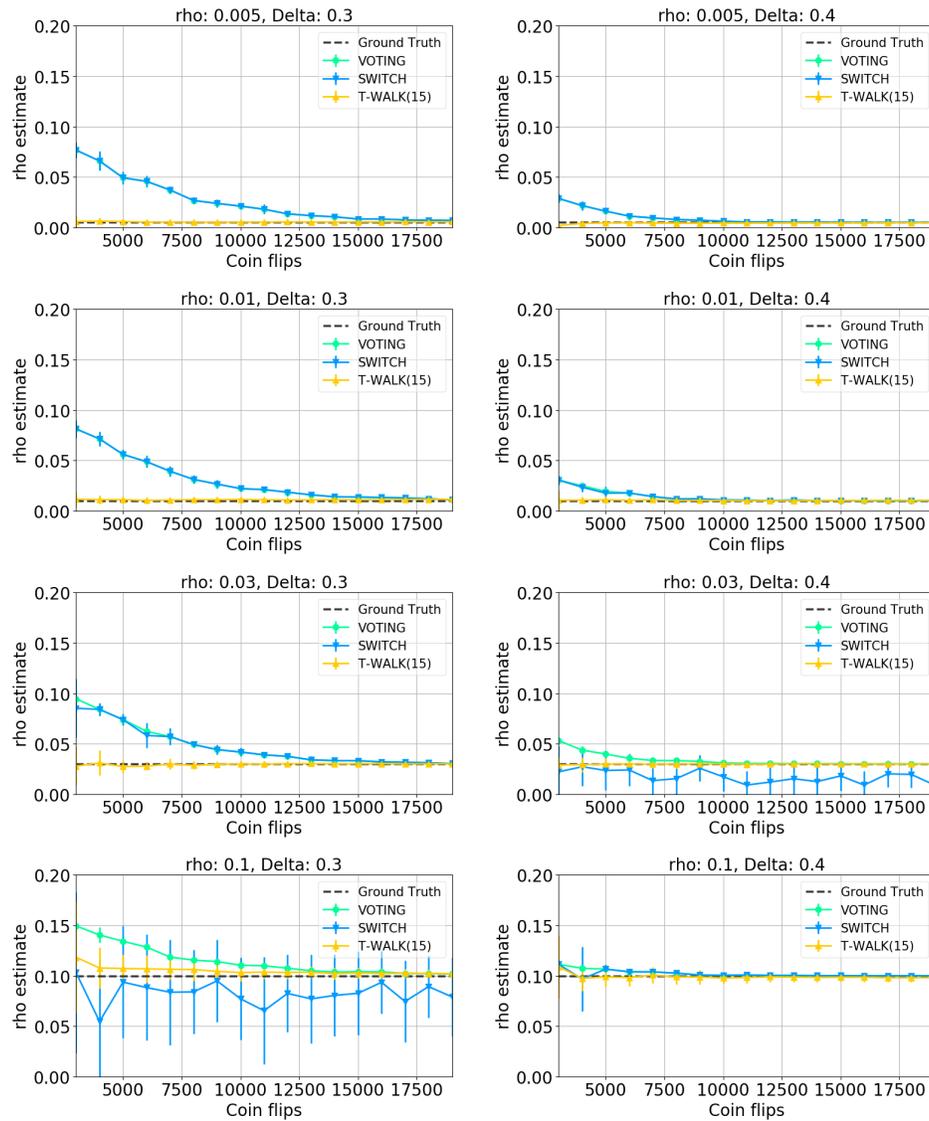


Figure 4.1: Experimental Results

coin flips as given by the x -coordinate, are computed according to Section 4.2.1.

In all cases, our algorithm (plotted in yellow) performs close to the ground truth (horizontal black line), while the alternative algorithms take longer to converge, or have high variance, as depicted by the error bars. In particular, as discussed in the introduction, our adaptive methods have the most potential for improvement in the more challenging and more practical regime where ρ is small (top few plots), and where Δ is smaller (left column).

4.8 Algorithm for Known Conditional Distributions of Coins

We now present our algorithms for the scenario where we know 1) the conditional distribution h^+ of the biases of positive coins, 2) analogously the distribution h^- for the negative coins, as well as, for now circularly 3) the mixture parameter ρ itself. In practice, of course, we would only have an estimate ρ for the mixture parameter itself, with the goal being to refine the estimate. Assuming for the sake of analysis that our knowledge of the two conditional distributions as well as the mixture parameter are perfect (even if in practice they are only guesses), we derive a simple method based on linear and quadratic programming tools for computing the triangular walk linear estimator (an instantiation of Algorithm 7 in Section 4.4) with the *minimum variance* subject to the constraints that 1) the estimator has expected output exactly 0 when given a randomly chosen negative coin, and 2) expected output exactly 1 for a randomly chosen positive coin. That is, we enforce that the estimator is unbiased no matter what the true mixture parameter is, but we optimize its variance given our (assumed to be perfect) knowledge of the mixture parameter.

This method is practically relevant as a bootstrapping approach. If our estimates of the conditional distributions and mixture parameter are indeed close to the ground truth, then it is easy to show bounds on the decrease in the estimator's performance as our estimates deviate from the truth. As such, we focus on the analysis of the method when our knowledge of the parameters are assumed to be perfect. The sample complexity of our algorithm is given in Theorem 4.5.

To complement the above upper bound result, we show that the linear estimator constructed from perfect knowledge of the relevant parameters is essentially an optimal estimator (Theorem 4.6) up to constant factors in sample complexity, under those exact same parameters. This gives strong evidence for the unique algorithmic challenges presented by the "uncertainty about uncertainty" regime of our problem, as discussed at the beginning of the paper.

4.8.1 A Quadratic+Linear Programming Approach

In this section we shall use extensively the notations $\alpha_{n,k}$, $\beta_{n,k}$ and $\gamma_{n,k}$ defined in Definition 4.14.

We now give an overview on the steps required to derive the minimum variance unbiased estimator (in the form of Algorithm 7), as described at the beginning of the section. First, we assume that we are given a fixed stopping rule, and derive output coefficients for the corresponding linear estimator that has minimum variance. We formulate a quadratic program (Figure 4.8.1) with the output coefficients $\{v_{n,k}\}$ as the variables, fixing $\alpha_{n,k}$ as constants. The quadratic program can be solved analytically, which allows us to derive for $\{v_{n,k}\}$ closed form expressions that makes an unbiased estimator with minimum variance assuming the given stopping rule, as well as perfect knowledge of the conditional distribution of biases and the mixture parameter. Furthermore, the objective value (a function in $\alpha_{n,k}$) of the quadratic program turns out (Lemma 4.34) to be the reciprocal of a *linear* function in terms of $\alpha_{n,k}$. With this representation of the objective, then, we can use the structural observations in Section 4.4 to formulate a *linear* program that solves for the optimum stopping rule given the conditional distributions of biases (conditional on a positive coin, or a negative coin) and mixture parameter. In practice, the linear program is first solved to give the stopping rule, then the output coefficients can be calculated from the first step in the analysis.

Having the above overview in mind, we describe the details of the derivation. To simplify notation, let $h_{n,k}^-$ be shorthand for $\mathbb{E}_{p \leftarrow h^-} (p^k(1-p)^{n-k})$ (a generalization of the notation from Section 4.5), and similarly for $h_{n,k}^+$. Thus $\alpha_{n,k}h_{n,k}^-$ is the probability that if we randomly choose a negative coin, executing the triangular walk with that coin will stop at state (n, k) . Similarly, $\alpha_{n,k}h_{n,k}^+$ is the analogous probability using a randomly chosen positive coin instead.

The quadratic program mentioned above is given in Figure 4.8.1. We use variables $\{\tilde{v}_{n,k}\}$, constraining them such that the expected output over a randomly chosen positive coin (from distribution h^+) has value 1 greater than that over a randomly chosen negative coin (from distribution h^-). Under this constraint, we minimize the second moment of the output when items are drawn from the mixture $\rho h^+ + (1 - \rho)h^-$. Any optimal solution to

$\begin{aligned} &\text{minimize} && \sum_{n,k} \alpha_{n,k} (\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-) \tilde{v}_{n,k}^2 \\ &\text{subject to} && \sum_{n,k} \alpha_{n,k} h_{n,k}^+ \tilde{v}_{n,k} = 1 + \sum_{n,k} \alpha_{n,k} h_{n,k}^- \tilde{v}_{n,k} \end{aligned}$

Figure 4.2: A QP formulation for computing the output coefficients in terms of the stopping rule

this optimization will choose the variables $\{\tilde{v}_{n,k}\}$ such that the expected output of an item drawn from the universe is 0, implying that $\sum_{n,k} \alpha_{n,k} h_{n,k}^+ \tilde{v}_{n,k} = 1 - \rho$ and $\sum_{n,k} \alpha_{n,k} h_{n,k}^- \tilde{v}_{n,k} = -\rho$. Therefore, we can compute $\{v_{n,k}\}$ using $\{\tilde{v}_{n,k}\}$ by setting $v_{n,k} = \tilde{v}_{n,k} + \rho$. As a consequence, $\sum_{n,k} \alpha_{n,k} h_{n,k}^+ v_{n,k} = 1$ and $\sum_{n,k} \alpha_{n,k} h_{n,k}^- v_{n,k} = 0$, satisfying the unbiasedness requirement as desired.

The quadratic program in Figure 4.8.1 can be solved analytically using Lagrange multipliers. We give the results as Lemma 4.34, and defer the calculations to Section 4.10.

Lemma 4.34. *For the quadratic program in Figure 4.8.1, the optimal assignments to $\{\tilde{v}_{n,k}\}$ are*

$$\tilde{v}_{n,k} = \frac{\frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-}}{\sum_{m,j} \alpha_{m,j} \frac{(h_{m,j}^+ - h_{m,j}^-)^2}{\rho h_{m,j}^+ + (1-\rho)h_{m,j}^-}}$$

(and we choose $v_{n,k} = \tilde{v}_{n,k} + \rho$), giving an objective value of

$$\frac{1}{\sum_{n,k} \alpha_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-}}$$

As mentioned at the beginning of the section, the optimal objective value of the quadratic program, namely the minimum variance achievable given a stopping rule, is the reciprocal of a *linear* function in $\{\alpha_{n,k}\}$. Note that the total sample complexity of the linear estimator, if we use the median-of-means method to estimate its expectation, is proportional to product of the variance of the linear estimator and the expected sample complexity of one run of the random walk. Therefore, if we fix the expected sample complexity of one run to be n_0 , we can in fact optimize the total sample complexity by minimizing the variance over all possible stopping rules with the expected sample complexity of n_0 . Observe that the reciprocal of the variance, divided by n_0 , is simply the reciprocal of the total sample complexity of the stopping rule, that we would therefore like to *maximize*. Moreover, such function is a linear function in $\{\alpha_{n,k}\}$. Thus, we can write the optimization

maximize	$\frac{1}{n_0} \sum_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-} \alpha_{n,k}$
subject to	$\beta_{0,0} = 1$ $\beta_{n+1,k+1} = \beta_{n,k+1} - \alpha_{n,k+1} + \beta_{n,k} - \alpha_{n,k}$ $\alpha_{n,k} \leq \beta_{n,k}$ $\alpha_{n_{\max},k} = 1$ for all k (Max depth constraint) $\sum_{n,k} n \cdot \alpha_{n,k} (\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-) \leq n_0$ (Bounding expected sample complexity)
where	$\alpha_{n,k}, \beta_{n,k} \geq 0$

Figure 4.3: An LP formulation for finding the best stopping rule given an expected sample complexity

problem as the linear program in Figure 4.3, by taking the objective to maximize the reciprocal of the quadratic program solution, divided by n_0 . The program includes (slightly adapted versions of) the recurrence relations introduced in Equation 4.1 as constraints. Moreover, in order to control the sample complexity of the algorithm, the program also contains constraints enforcing that 1) the expected number of responses solicited for a random item is bounded by n_0 and 2) the maximum depth of the triangle is bounded by some parameter n_{\max} . In addition to the interpretation as the maximum amount of resources we would ever invest on a single coin/item, the maximum depth constraint can also be interpreted as a computational constraint on how much time we can spend on computing the description of the linear estimator.

Since, ultimately, we wish to optimize over all possible values in n_0 , such a linear program formulation (in Figure 4.3) cannot be used directly. However, consider the following rewriting of the program. We can always divide the $\{\alpha_{n,k}, \beta_{n,k}\}$ variables by n_0 and not change the meaning of the program, if we rescale the constraints and objective correspondingly. This modification has the following effects: it 1) changes the n_0 in the objective and the fifth constraint into 1, 2) preserves the second and third constraints as well as the non-negativity constraints and 3) changes the first and fourth constraints into “variable = $1/n_0$ ”. The first and fourth constraints are now the only components in the new program that depend on n_0 , and since we ultimately wish to optimize over all possible n_0 , we can replace these constraints with the weaker constraint that they all equal to each other without specifying what they are equal to. This results in the linear program in Figure 4.4, which by the above reasoning is equivalent to optimizing the total sample complexity for

maximize	$\sum_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-} \alpha_{n,k}$
subject to	$\beta_{n+1,k+1} = \beta_{n,k+1} - \alpha_{n,k+1} + \beta_{n,k} - \alpha_{n,k}$ $\alpha_{n,k} \leq \beta_{n,k}$ $\alpha_{n_{\max},k} = \beta_{0,0}$ for all k (Max depth constraint) $\sum_{n,k} n \cdot \alpha_{n,k} (\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-) \leq 1$
where	$\alpha_{n,k}, \beta_{n,k} \geq 0$

Figure 4.4: An LP formulation for finding the best stopping rule independent of the expected sample complexity for a single coin

a stopping rule.

To obtain the optimal stopping rule $\{\gamma_{n,k}\}$, we solve the linear program in Figure 4.4, rescale every variable such that $\beta_{0,0} = 1$, and calculate $\gamma_{n,k} = \alpha_{n,k}/\beta_{n,k}$. If the solution to the linear program (in Figure 4.4) is $1/S$, then the expected sample complexity is $O(\frac{S}{\epsilon^2} \log \frac{1}{\delta})$ to estimate ρ to within an additive ϵ with probability at least $1 - \delta$. This can be achieved by taking the median-of-means of $O(\log \frac{1}{\delta})$ groups of samples of size $O(S/\epsilon^2)$, each of which has a constant probability concentration to within additive ϵ by Chebyshev's inequality. Summarizing the above gives the following theorem.

Theorem 4.5. *Suppose we are given 1) the distribution of coin biases conditioned on being a positive coin, 2) the analogous distribution for negative coins and 3) the mixture parameter ρ (which, again, is a circular assumption but useful for a bootstrapping approach). Suppose further that we are given 4) the parameter n_{\max} , which controls the maximum depth of the triangular walk.*

Then, following the method described earlier in this section, we can find the linear estimator for ρ that minimizes variance, subject to a) the expected output of the estimator on input a random positive coin is 1 and b) the analogous expected output for a random negative coin is 0.

Moreover, if the objective of the linear program in Figure 4.4 is $1/S$, then the expected sample complexity of the constructed linear estimator is $O(\frac{S}{\epsilon^2} \log \frac{1}{\delta})$, which will estimate ρ to within an additive error of ϵ with probability at least $1 - \delta$.

4.8.2 Optimality of such linear estimators

In this section, we show that in fact, the linear estimators produced by the linear program in Figure 4.3 are optimal compared with any single-coin adaptive but possibly non-linear

estimators, subject to the same maximum depth constraints.

Our approach for lower bounding the sample complexity is to fix the distributions h^+ and h^- of positive and negative coin biases respectively, and show that with a small number of samples, it is impossible to distinguish the case between A) a ρ and $(1 - \rho)$ mixture of positive and negative coins and B) a $(\rho + \epsilon)$ and $(1 - \rho - \epsilon)$ mixture. To show indistinguishability, we again use the notion of Hellinger distance. Since each stopping rule induces different distribution on the Pascal triangle, under randomly chosen coins from each of the A and B scenarios, we will upper bound the (squared) Hellinger distance between the scenarios.

Lemma 4.35 shows that the squared Hellinger distance is in fact a linear function in $\{\alpha_{n,k}\}$ and furthermore, in the regime where $\epsilon \ll \rho$, is within a constant factor of the objective in the linear program in Figure 4.4. The coincidence will allow us to show matching lower bounds.

Lemma 4.35. *Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ giving coefficients $\{\alpha_{n,k}\}$. If ϵ/ρ is smaller than some universal constant, then the squared Hellinger distance between 1) a coin randomly chosen as in case A (described in the paragraphs above) inducing a distribution on the Pascal triangle given the stopping rule and 2) a coin randomly chosen as in case B instead, is*

$$\Theta(\epsilon^2) \sum_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} \alpha_{n,k}$$

We defer the proof and calculations to Section 4.10, but it is completely analogous to that of Lemma 4.23.

With Lemma 4.35, we now prove Theorem 4.6.

Theorem 4.6. *As in Theorem 4.5, suppose we are given 1) the distribution of coin biases conditioned on being a positive coin, 2) the analogous distribution for negative coins, 3) the mixture parameter ρ , as well as the parameter n_{\max} , which controls the maximum depth of the triangular walk.*

The linear estimator produced from solving the linear program in Figure 4.4, as described in Theorem 4.5, has total expected sample complexity that is within a constant factor of any optimal single-coin adaptive algorithm with $\geq \frac{2}{3}$ probability of success, subject to the same maximum depth constraint.

Combining with a corollary of Lemma 4.19, restricted to fully-adaptive algorithms that invests at most n_{\max} flips on any single coin, this shows that our linear estimator in fact has sample complexity within a constant factor of any fully-adaptive algorithm satisfying the maximum depth constraint for every single coin.

Proof. Given an arbitrary stopping rule, if it induces a squared Hellinger distance of H^2 between the two cases with a single random walk, then we can lower bound the number of random walks needed in the single-coin adaptive algorithm in order to solve the distinguishing task with constant probability of success, by $\Theta(1/H^2)$, using the subadditivity of squared Hellinger distance, and that the total Hellinger distance needs to be at least constant to solve the distinguishing task. Thus, if n_0 is the expected number of coin flips for a random walk, the overall expected sample complexity is lower bounded by $\Omega(n_0/H^2)$. Since we need to find a lower bound that applies to *all* single-coin adaptive algorithms, we need to find the smallest n_0/H^2 over all the possible stopping rules (subject to the same max-depth constraint), or equivalently, maximize H^2/n_0 (which can alternatively be interpreted as the squared Hellinger distance per expected sample). Lemma 4.35 tells us that we can replace H^2 with the expression in the lemma and lose no more than multiplicative constants. Thus, if we fix n_0 , finding the best lower bound up to multiplicative constants is equivalent to solving the optimization problem that is exactly the one in Figure 4.3, except for an extra factor of $\Theta(\epsilon^2)$ in the objective. We again wish to maximize the H^2/n_0 over all possible choices of n_0 as well, and therefore, following the same reasoning as before, we arrive at the linear program that is essentially the one in Figure 4.4, again except for the factor of $\Theta(\epsilon^2)$ in the objective. This linear program has no n_0 dependency, and has objective that is $\Theta(\epsilon^2)$ times that of the one in Figure 4.4, which is the reciprocal of the (expected) total sample complexity of the optimal linear estimator produced as described in Section 4.8.1. Summarizing, if the solution to the linear program in Figure 4.4 is $1/S$, then the maximum H^2/n_0 over all possible stopping rules would be within a constant factor of ϵ^2/S , giving a lower bound of $\Omega(S/\epsilon^2)$ on the expected sample complexity (under case A) for a constant probability of success in the task of distinguishing between case A and case B. This lower bound matches the upper bound of $O(S/\epsilon^2)$ on the total sample complexity of the linear estimator we produce according to Theorem 4.5.

As given in the theorem statement, combining this result with a corollary of Lemma 4.19 shows that our linear estimator is in fact competitive to within a constant factor in sample complexity with *fully-adaptive* algorithms that invest at most n_{\max} flips on any single coin. \square

4.9 Non-Adaptive Lower Bound

Here we give the remaining calculations for the non-adaptive lower bound of $O(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\rho})$.

Recall that, to show a non-adaptive lower-bound, consider a random variable S that uniformly chooses between scenarios " ρ " and " $\rho + \epsilon$ " respectively, where coins will have bias $\frac{1}{2} + \Delta$ with probability ρ or $\rho + \epsilon$ respectively depending on the outcome of S , and bias $\frac{1}{2} - \Delta$ otherwise. We will show that the mutual information between n flips of a single coin and the scenario variable S is at most $O(n \frac{\epsilon^2 \Delta^2}{\rho \log(1/\rho)})$, and thus that, even when combining information from several coins, at least $\Omega(\frac{\rho}{\epsilon^2 \Delta^2} \log \frac{1}{\rho})$ samples are needed to distinguish the two scenarios with constant probability.

Let $\text{Bin}(n, p, k)$ denote the probability that a Binomial distribution with n trials and bias p has value k .

The mutual information is exactly represented as

$$\begin{aligned} & \frac{1}{2} \sum_{k=0}^n (\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho) \text{Bin}(n, \frac{1}{2} - \Delta, k)) \log \frac{(\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho) \text{Bin}(n, \frac{1}{2} - \Delta, k))}{((\rho + \frac{\epsilon}{2}) \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho - \frac{\epsilon}{2}) \text{Bin}(n, \frac{1}{2} - \Delta, k))} \\ & + ((\rho + \epsilon) \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho - \epsilon) \text{Bin}(n, \frac{1}{2} - \Delta, k)) \log \frac{((\rho + \epsilon) \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho - \epsilon) \text{Bin}(n, \frac{1}{2} - \Delta, k))}{((\rho + \frac{\epsilon}{2}) \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho - \frac{\epsilon}{2}) \text{Bin}(n, \frac{1}{2} - \Delta, k))} \end{aligned}$$

Claim is that, for $x, y \geq 0$, we have $x \log \frac{x}{(x+y)/2} + y \log \frac{y}{(x+y)/2} \leq \frac{(x-y)^2}{x+y}$.

Letting x be $\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho) \text{Bin}(n, \frac{1}{2} - \Delta, k)$ and y be the $\rho + \epsilon$ mixture analogue, the mutual information is less than or equal to:

$$\begin{aligned} & \frac{\epsilon^2}{4} \sum_{k=0}^n \frac{(\text{Bin}(n, \frac{1}{2} + \Delta, k) - \text{Bin}(n, \frac{1}{2} - \Delta, k))^2}{((\rho + \frac{\epsilon}{2}) \text{Bin}(n, \frac{1}{2} + \Delta, k) + (1 - \rho - \frac{\epsilon}{2}) \text{Bin}(n, \frac{1}{2} - \Delta, k))} \\ & \leq \epsilon^2 \sum_{k=0}^n \frac{(\text{Bin}(n, \frac{1}{2} + \Delta, k) - \text{Bin}(n, \frac{1}{2} - \Delta, k))^2}{(\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + \text{Bin}(n, \frac{1}{2} - \Delta, k))} \end{aligned}$$

Since $(x - y)^2 \leq 2(x^2 + y^2)$, and $\frac{1}{x+y} \leq \min\{\frac{1}{x}, \frac{1}{y}\}$, we also have

$$\sum_{k=0}^n \frac{(\text{Bin}(n, \frac{1}{2} + \Delta, k) - \text{Bin}(n, \frac{1}{2} - \Delta, k))^2}{(\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + \text{Bin}(n, \frac{1}{2} - \Delta, k))}$$

$$\begin{aligned}
&\leq 2 \sum_{k=0}^n \frac{\text{Bin}(n, \frac{1}{2} + \Delta, k)^2 + \text{Bin}(n, \frac{1}{2} - \Delta, k)^2}{(\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + \text{Bin}(n, \frac{1}{2} - \Delta, k))} \\
&\leq 2 \min \left\{ \sum_{k=0}^n \frac{\text{Bin}(n, \frac{1}{2} + \Delta, k)^2}{\rho \text{Bin}(n, \frac{1}{2} + \Delta, k)}, \sum_{k=0}^n \frac{\text{Bin}(n, \frac{1}{2} + \Delta, k)^2}{\text{Bin}(n, \frac{1}{2} - \Delta, k)} \right\} + 2 \sum_{k=0}^n \frac{\text{Bin}(n, \frac{1}{2} - \Delta, k)^2}{\text{Bin}(n, \frac{1}{2} - \Delta, k)} \\
&= 2 \min \left\{ \frac{1}{\rho}, \left(\frac{1 + 12\Delta^2}{1 - 4\Delta^2} \right)^n \right\} + 2
\end{aligned}$$

For Δ bounded below by any universal positive constant, this last expression is $O(\min\{\frac{1}{\rho}, e^{O(\Delta^2 n)}\})$. Since the components of the minimum are 1) equal for $n = O(\frac{1}{\Delta^2} \log \frac{1}{\rho})$ and 2) constant and convex in n respectively, we can bound the minimum by a linear function that goes through this intersection point: for $n \geq \frac{1}{\Delta^2}$ the minimum is bounded by $O(n \frac{\Delta^2}{\rho \log \frac{1}{\rho}})$. Multiplying by ϵ^2 gets a bound on the mutual information, and dividing by n gets a bound on mutual information per sample of $O(\frac{\epsilon^2 \Delta^2}{\rho \log \frac{1}{\rho}})$.

For the remaining regime of $n \leq \frac{1}{\Delta^2}$:

$$\sum_{k=0}^n \frac{(\text{Bin}(n, \frac{1}{2} + \Delta, k) - \text{Bin}(n, \frac{1}{2} - \Delta, k))^2}{\rho \text{Bin}(n, \frac{1}{2} + \Delta, k) + \text{Bin}(n, \frac{1}{2} - \Delta, k)} \leq \sum_{k=0}^n \frac{(\text{Bin}(n, \frac{1}{2} + \Delta, k) - \text{Bin}(n, \frac{1}{2} - \Delta, k))^2}{\text{Bin}(n, \frac{1}{2} - \Delta, k)} = \left(\frac{1 + 12\Delta^2}{1 - 4\Delta^2} \right)^n - 1$$

This last expression is $O(\Delta^2 n)$ for $n \leq \frac{1}{\Delta^2}$, and thus we can bound the mutual information per sample by $O(\epsilon^2 \Delta^2)$ here.

Combining the two bounds, we conclude the mutual information per sample is at most $O(\frac{\epsilon^2 \Delta^2}{\rho \log \frac{1}{\rho}})$ for all n , and thus its inverse, $O(\frac{\rho \log \frac{1}{\rho}}{\epsilon^2 \Delta^2})$ lower-bounds the number of non-adaptive samples needed for our task.

4.10 Remaining Proofs/Calculations of Results

Lemma 4.23. *Consider the two probability distributions in Proposition 4.22 over locations (n, k) in the Pascal triangle of depth $10^{-8}/\Delta^2$ and bias $p \in \{\frac{1}{2} \pm \Delta\}$, generated by the given stopping rule $\{\gamma_{n,k}\}$ in the two cases of 1) a coin with bias $\frac{1}{2} + \Delta$ is used with probability ρ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise versus 2) a coin with bias $\frac{1}{2} + \Delta$ is used with probability $\rho + \epsilon$ and a coin with bias $\frac{1}{2} - \Delta$ is used otherwise. If ϵ/ρ is smaller than some universal constant, then the squared*

Hellinger distance between these two distributions can be written as

$$\Theta(\epsilon^2) \left[\sum_{n < \frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)^2} \right. \\ \left. \sum_{n = \frac{10^{-8}}{\Delta^2}, k \in [0..n]} \alpha_{n,k} \left(\left(\rho + \frac{\epsilon}{2} \right) h_{n,k}^+ + \left(1 - \rho - \frac{\epsilon}{2} \right) h_{n,k}^- \right) \frac{\frac{h_{n,k}^+}{\rho} + h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} \right]$$

Proof. For $n < 10^{-8}/\Delta^2$, the probability of that the location (n, k) is revealed, for the two distributions we consider in Proposition 4.22, are $\alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)$ and $\alpha_{n,k}((\rho + \epsilon) h_{n,k}^+ + (1 - \rho - \epsilon) h_{n,k}^-)$ respectively. Thus, the contribution by these locations to the squared Hellinger distance is proportional to:

$$\sum_{n,k} \left(\sqrt{\alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)} - \sqrt{\alpha_{n,k}((\rho + \epsilon) h_{n,k}^+ + (1 - \rho - \epsilon) h_{n,k}^-)} \right)^2 \\ = \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \left(1 - \sqrt{\frac{(\rho + \epsilon) h_{n,k}^+ + (1 - \rho - \epsilon) h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-}} \right)^2 \\ = \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \left(1 - \sqrt{1 + \epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-}} \right)^2 \\ = \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \left(1 - \left(1 + \frac{1}{2} \epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} + \Theta \left(\left(\epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} \right)^2 \right) \right) \right)^2 \\ = \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \left(\frac{1}{2} \epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} + \Theta \left(\left(\epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} \right)^2 \right) \right)^2$$

Note that the multiplier to ϵ is upper bounded by $1/\rho$, and therefore if ϵ/ρ is sufficiently small, we have the last line being equal to

$$\sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \left(\Theta \left(\epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-} \right) \right)^2 \\ = \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \Theta \left(\epsilon^2 \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)^2} \right) \\ = \Theta(\epsilon^2) \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-) \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho) h_{n,k}^-)^2}$$

Finally, note that ϵ/ρ is a small constant, then ρ and $1 - \rho$ are respectively within a small constant factor of $\rho + \frac{\epsilon}{2}$ and $1 - \rho - \frac{\epsilon}{2}$, meaning that $(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)$ is within a constant factor of $((\rho + \frac{\epsilon}{2})h_{n,k}^+ + (1 - \rho - \frac{\epsilon}{2})h_{n,k}^-)$.

For $n = 10^{-8}/\Delta^2$, the probability that $((n, k), \frac{1}{2} + \Delta)$ is revealed is $\alpha_{n,k} \rho h_{n,k}^+$ and $\alpha_{n,k}(\rho + \epsilon)h_{n,k}^+$ for the two scenarios respectively. A similar calculation above gives a squared Hellinger distance contribution of

$$\Theta(\epsilon^2) \alpha_{n,k} \left(\rho + \frac{\epsilon}{2} \right) \frac{h_{n,k}^+}{\rho}$$

As for the contribution from the revealing of $((n, k), \frac{1}{2} - \Delta)$, the respective probabilities are $\alpha_{n,k}(1 - \rho)h_{n,k}^-$ and $\alpha_{n,k}(1 - \rho - \epsilon)h_{n,k}^-$, and similar calculations give a squared Hellinger distance contribution of

$$\Theta(\epsilon^2) \alpha_{n,k} h_{n,k}^-$$

which with algebraic manipulation and approximations as in the $n < 10^{-8}/\Delta^2$ case completes the proof of the lemma. \square

Lemma 4.34. *For the quadratic program in Figure 4.8.1, the optimal assignments to $\{\tilde{v}_{n,k}\}$ are*

$$\tilde{v}_{n,k} = \frac{\frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-}}{\sum_{m,j} \alpha_{m,j} \frac{(h_{m,j}^+ - h_{m,j}^-)^2}{\rho h_{m,j}^+ + (1-\rho)h_{m,j}^-}}$$

(and we choose $v_{n,k} = \tilde{v}_{n,k} + \rho$), giving an objective value of

$$\frac{1}{\sum_{n,k} \alpha_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-}}$$

Proof. We use the method of Lagrangian multiplier to find the optimal assignment to $\{\tilde{v}_{n,k}\}$.

The Lagrangian of the program is

$$L = \sum_{m,j} \alpha_{m,j} (\rho h_{m,j}^+ + (1 - \rho)h_{m,j}^-) \tilde{v}_{m,j}^2 + \lambda \left(\left(\sum_{m,j} \alpha_{m,j} (h_{m,j}^+ - h_{m,j}^-) \tilde{v}_{m,j} \right) - 1 \right)$$

where λ is the Lagrange multiplier.

We need to find assignments to $\{\tilde{v}_{n,k}\}$ and λ such that $\nabla_{\{\tilde{v}_{n,k}\},\lambda} L = 0$. Computing the partial derivatives gives the following system of equations:

$$2\alpha_{n,k}(\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-)\tilde{v}_{n,k} + \lambda\alpha_{n,k}(h_{n,k}^+ - h_{n,k}^-) = 0 \text{ for all } n, k \quad (4.9)$$

$$\sum_{m,j} \alpha_{m,j}(h_{m,j}^+ - h_{m,j}^-)\tilde{v}_{m,j} = 1 \quad (4.10)$$

Rearranging Equation 4.9 gives

$$\tilde{v}_{n,k} = \frac{-\lambda(h_{n,k}^+ - h_{n,k}^-)}{2(\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-)} \quad (4.11)$$

and substituting this into Equation 4.10 gives

$$-\lambda \sum_{m,j} \frac{\alpha_{m,j}(h_{m,j}^+ - h_{m,j}^-)^2}{2(\rho h_{m,j}^+ + (1-\rho)h_{m,j}^-)} = 1$$

which lets us solve for λ

$$\lambda = -1 / \sum_{m,j} \frac{\alpha_{m,j}(h_{m,j}^+ - h_{m,j}^-)^2}{2(\rho h_{m,j}^+ + (1-\rho)h_{m,j}^-)}$$

which when substituted back into Equation 4.11 gives

$$\tilde{v}_{n,k} = \frac{\frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-}}{\sum_{m,j} \frac{\alpha_{m,j}(h_{m,j}^+ - h_{m,j}^-)^2}{\rho h_{m,j}^+ + (1-\rho)h_{m,j}^-}}$$

as desired.

The optimal value of the program can be calculated by substituting the assignment to the objective function. \square

Lemma 4.35. Consider an arbitrary stopping rule $\{\gamma_{n,k}\}$ giving coefficients $\{\alpha_{n,k}\}$. The squared Hellinger distance between 1) a random coin drawn in case A inducing a distribution on the Pascal triangle given the stopping rule and 2) a random coin drawn in case B instead, is

$$\Theta(\epsilon^2) \sum_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1-\rho)h_{n,k}^-} \alpha_{n,k}$$

assuming that ϵ/ρ is smaller than some universal constant.

Proof. In scenario 1, the distribution induced by a random coin on the Pascal triangle is

$$\alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)$$

and similarly for scenario 2,

$$\alpha_{n,k}((\rho + \epsilon)h_{n,k}^+ + (1 - \rho - \epsilon)h_{n,k}^-)$$

The squared Hellinger distance is therefore proportional to

$$\begin{aligned} & \sum_{n,k} \left(\sqrt{\alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)} - \sqrt{\alpha_{n,k}((\rho + \epsilon)h_{n,k}^+ + (1 - \rho - \epsilon)h_{n,k}^-)} \right)^2 \\ &= \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-) \left(1 - \sqrt{\frac{(\rho + \epsilon)h_{n,k}^+ + (1 - \rho - \epsilon)h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-}} \right)^2 \\ &= \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-) \left(1 - \sqrt{1 + \epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-}} \right)^2 \\ &= \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-) \left(1 - \left(1 + \frac{1}{2} \epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} + \Theta \left(\left(\epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \right)^2 \right) \right) \right)^2 \\ &= \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-) \left(\frac{1}{2} \epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} + \Theta \left(\left(\epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \right)^2 \right) \right)^2 \end{aligned}$$

Note that the multiplier to ϵ is upper bounded by $1/\rho$, and therefore if ϵ/ρ is sufficiently small, we have the last line being equal to

$$\begin{aligned} & \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-) \left(\Theta \left(\epsilon \frac{h_{n,k}^+ - h_{n,k}^-}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \right) \right)^2 \\ &= \sum_{n,k} \alpha_{n,k}(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-) \Theta \left(\epsilon^2 \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{(\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-)^2} \right) \\ &= \Theta(\epsilon^2) \sum_{n,k} \frac{(h_{n,k}^+ - h_{n,k}^-)^2}{\rho h_{n,k}^+ + (1 - \rho)h_{n,k}^-} \alpha_{n,k} \end{aligned}$$

□

Bibliography

- [1] Jayadev Acharya, Clément L Canonne, and Gautam Kamath. A chasm between identity and equivalence testing with conditional queries. *Theory Comput.*, 14:1–46, 2018.
- [2] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci*, 58(1):137–147, 1999.
- [3] Tugkan Batu, Sanjoy Dasgupta, Ravi Kumar, and Ronitt Rubinfeld. The complexity of approximating the entropy. *SIAM J. Comput*, 35(1):132–150, 2005.
- [4] Graham Bell, Martin J Lechowicz, and Marcia J Waterway. Environmental heterogeneity and species diversity of forest sedges. *Journal of Ecology*, 88(1):67–87, 2000.
- [5] Aleksandrs Belovs and Eric Blais. A polynomial lower bound for testing monotonicity. In *Proc. STOC '16*, pages 1021–1032, 2016.
- [6] Mark Braverman, Ankit Garg, Tengyu Ma, Huy L. Nguyen, and David P. Woodruff. Communication lower bounds for statistical estimation problems via a distributed data processing inequality. In *Proc. STOC'16*, page 1011–1020, 2016.
- [7] Jennifer Brennan, Ramya Korlakai Vinayak, and Kevin Jamieson. Estimating the number and effect sizes of non-null hypotheses. In *Proc. ICML'20*, 2020.
- [8] Christian Brownlees, Emilien Joly, and Gábor Lugosi. Empirical risk minimization for heavy-tailed losses. *Ann. Stat.*, 43(6):2507–2536, 2015.
- [9] Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *IEEE Trans. Inf. Theory*, 59(11):7711–7717, 2013.

- [10] Clément Canonne. A Survey on Distribution Testing. Your Data is Big. But is it Blue? *Theory of Computing Graduate Surveys*, 9:1–100, 2017.
- [11] Clément Canonne, Dana Ron, and Rocco A Servedio. Testing equivalence between distributions using conditional samples. In *Proc. SODA '14*, pages 1174–1192, 2014.
- [12] Clément Canonne and Ronitt Rubinfeld. Testing probability distributions underlying aggregated data. In *Proc. ICALP '14*, pages 283–295, 2014.
- [13] Clément L Canonne. Big data on the rise? In *Proc. ICALP '15*, pages 294–305, 2015.
- [14] Clément L Canonne, Dana Ron, and Rocco A Servedio. Testing probability distributions using conditional samples. *SIAM J. Comput*, 44(3):540–616, 2015.
- [15] Olivier Catoni. Challenging the empirical mean and empirical variance: a deviation study. *Ann. I. H. Poincaré-PR*, 48(4):1148–1185, 2012.
- [16] Olivier Catoni and Ilaria Giulini. Dimension-free PAC-Bayesian bounds for matrices, vectors, and linear least squares regression. *arXiv:1712.02747*, 2017.
- [17] Sourav Chakraborty, Eldar Fischer, Yonatan Goldhirsh, and Arie Matsliah. On the power of conditional samples in distribution testing. *SIAM J. Comput*, 45(4):1261–1296, 2016.
- [18] Karthekeyan Chandrasekaran and Richard Karp. Finding a most biased coin with fewest flips. In *Proc. COLT'14*, pages 394–407, 2014.
- [19] Xi Chen, Erik Waingarten, and Jinyu Xie. Beyond Talagrand Functions: New Lower Bounds for Testing Monotonicity and Unateness. In *Proc. STOC '17*, pages 523–536, 2017.
- [20] Yeshwanth Cherapanamjeri, Nicolas Flammarion, and Peter L Bartlett. Fast mean estimation with sub-gaussian rates. In *Proc. COLT '20*, pages 786–806, 2019.
- [21] Yeounoh Chung, Sanjay Krishnan, and Tim Kraska. A Data Quality Metric (DQM): How to Estimate the Number of Undetected Errors in Data Sets. *Proc. VLDB '17*, 10(10):1094–1105, 2017.

- [22] Robert K Colwell and Jonathan A Coddington. Estimating terrestrial biodiversity through extrapolation. *Phil. Trans. R. Soc. Lond. B*, 345(1311):101–118, 1994.
- [23] Luc Devroye, Matthieu Lerasle, Gabor Lugosi, and Roberto I. Oliveira. Sub-gaussian mean estimators. *Ann. Stat.*, 44(6):2695–2725, 2016.
- [24] Ilias Diakonikolas, Daniel M Kane, and Ankit Pensia. Outlier robust mean estimation with subgaussian rates via stability. In *Proc. NeurIPS '20*, 2020.
- [25] John C Duchi, Michael I Jordan, and Martin J Wainwright. Local privacy and statistical minimax rates. In *Proc. FOCS '13*, pages 429–438, 2013.
- [26] John C Duchi, Michael I Jordan, and Martin J Wainwright. Minimax optimal procedures for locally private estimation. *J. Am. Stat. Assoc.*, 113(521):182–201, 2018.
- [27] Moein Falahatgar, Yi Hao, Alon Orlitsky, Venkatadheeraj Pichapati, and Vaishakh Ravindrakumar. Maxing and ranking with few assumptions. In *Proc. NeurIPS '17*, pages 7060–7070, 2017.
- [28] Moein Falahatgar, Ashkan Jafarpour, Alon Orlitsky, Venkatadheeraj Pichapati, and Ananda Theertha Suresh. Faster algorithms for testing under conditional sampling. In *Proc. COLT '15*, pages 607–636, 2015.
- [29] Moein Falahatgar, Alon Orlitsky, Venkatadheeraj Pichapati, and Ananda Theertha Suresh. Maximum selection and ranking under noisy comparisons. In *Proc. ICML '17*, pages 1088–1096, 2017.
- [30] Sudipto Guha, Andrew McGregor, and Suresh Venkatasubramanian. Streaming and sublinear approximation of entropy and information distances. In *Proc. SODA '06*, pages 733–742, 2006.
- [31] Samuel B Hopkins. Mean estimation with sub-gaussian rates in polynomial time. *Ann. Stat.*, 48(2):1193–1213, 2020.
- [32] Daniel Hsu and Sivan Sabato. Loss minimization and parameter estimation with heavy tails. *J. Mach. Learn. Res.*, 17(1):543–582, 2016.

- [33] Kevin Jamieson, Daniel Haas, and Ben Recht. The power of adaptivity in identifying statistical alternatives. In *Proc. NIPS'16*, pages 775–783, 2016.
- [34] Mark R Jerrum, Leslie G Valiant, and Vijay V Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theor. Comput. Sci*, 43:169–188, 1986.
- [35] Jiantao Jiao, Kartik Venkat, Yanjun Han, and Tsachy Weissman. Minimax estimation of functionals of discrete distributions. *IEEE Trans. Inf. Theory*, 61(5):2835–2885, 2015.
- [36] Gautam Kamath, Vikrant Singhal, and Jonathan Ullman. Private mean estimation of heavy-tailed distributions. In *Proc. COLT '20*, volume 125, pages 2204–2235. PMLR, 09–12 Jul 2020.
- [37] Jasper C.H. Lee and Paul Valiant. Uncertainty about uncertainty: Optimal adaptive algorithms for estimating mixtures of unknown coins. In *Proc. SODA '21*, pages 394–413. SIAM, 2021.
- [38] Zhixian Lei, Kyle Luh, Prayaag Venkat, and Fred Zhang. A fast spectral algorithm for mean estimation with sub-gaussian rates. In *Proc. COLT '20*, pages 2598–2612, 2020.
- [39] Reut Levi, Dana Ron, and Ronitt Rubinfeld. Testing properties of collections of distributions. *Theory Comput.*, 9(1):295–347, 2013.
- [40] Reut Levi, Dana Ron, and Ronitt Rubinfeld. Testing similar means. *SIAM J. Discrete Math*, 28(4):1699–1724, 2014.
- [41] Frederic M Lord. A strong true-score theory, with applications. *Psychometrika*, 30(3):239–270, 1965.
- [42] Frederic M Lord and Noel Cressie. An empirical bayes procedure for finding an interval estimate. *Sankhyā: The Indian Journal of Statistics, Series B*, pages 1–9, 1975.
- [43] Gabor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions—a survey. *Found Comput Math*, 19:1145–1190, 2019.
- [44] Gabor Lugosi and Shahar Mendelson. Robust multivariate mean estimation: the optimality of trimmed mean. *Ann. Stat*, 2019. To appear.

- [45] Gábor Lugosi and Shahar Mendelson. Sub-Gaussian estimators of the mean of a random vector. *Ann. Stat.*, 47(2):783–794, 2019.
- [46] Matthew L Malloy, Gongguo Tang, and Robert D Nowak. Quickest search for a rare distribution. In *Proc. CISS'12*, pages 1–6, 2012.
- [47] Wayne J Millar. Distribution of body weight and height: comparison of estimates based on self-reported and observed measures. *Journal of Epidemiology & Community Health*, 40(4):319–323, 1986.
- [48] Stanislav Minsker. Uniform bounds for robust mean estimators. *arXiv:1812.03523*, 2018.
- [49] Ankur Moitra and Michael Saks. A polynomial time algorithm for lossy population recovery. In *Proc. FOCS'13*, pages 110–116, 2013.
- [50] A.S. Nemirovsky and D.B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. Wiley, 1983.
- [51] Roberto I. Oliveira and Paulo Orenstein. The sub-gaussian property of trimmed means estimators. *Technical Report, IMPA*, 2019.
- [52] Michael W Palmer and Philip M Dixon. Small-scale environmental heterogeneity and the analysis of species distributions along gradients. *Journal of Vegetation Science*, 1(1):57–65, 1990.
- [53] Yury Polyanskiy, Ananda Theertha Suresh, and Yihong Wu. Sample complexity of population recovery. In *Proc. COLT'17*, volume 65, 2017.
- [54] Yury Polyanskiy and Yihong Wu. Dualizing Le Cam’s method, with applications to estimating the unseens. *arXiv:1902.05616*, 2019.
- [55] Ronitt Rubinfeld and Rocco A Servedio. Testing monotone high-dimensional distributions. *Random Struct Algor*, 34(1):24–44, 2009.

- [56] N. B. Shah, S. Balakrishnan, and M. J. Wainwright. Feeling the bern: Adaptive estimators for bernoulli probabilities of pairwise comparisons. In *Proc. ISIT '16*, pages 1153–1157, 2016.
- [57] Nihar Shah, Dengyong Zhou, and Yuval Peres. Approval voting and incentives in crowdsourcing. In *Proc. ICML '15*, pages 10–19, 2015.
- [58] Nihar B Shah and Martin J Wainwright. Simple, robust and optimal ranking from pairwise comparisons. *JMLR*, 18(1):7246–7283, 2017.
- [59] Maurice Sion. On general minimax theorems. *Pac. J. Math*, 8(1):171–176, 1958.
- [60] Leonard A. Stefanski and Dennis D. Boos. The calculus of M-estimation. *The American Statistician*, 56(1):29–38, 2002.
- [61] Kevin Tian, Weihao Kong, and Gregory Valiant. Learning populations of parameters. In *Proc. NeurIPS '17*, pages 5778–5787, 2017.
- [62] Alexandre B Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.
- [63] Gregory Valiant and Paul Valiant. The power of linear estimators. In *Proc. FOCS '11*, pages 403–412, 2011.
- [64] Ramya Korlakai Vinayak, Weihao Kong, Gregory Valiant, and Sham Kakade. Maximum likelihood estimation for learning populations of parameters. In *Proc. ICML '19*, pages 6448–6457, 2019.
- [65] Jingyan Wang and Nihar B Shah. Your 2 is my 1, your 3 is my 9: Handling arbitrary miscalibrations in ratings. In *Proc. AAMAS '19*, 2019. To appear.
- [66] Yihong Wu and Pengkun Yang. Minimax rates of entropy estimation on large alphabets via best polynomial approximation. *IEEE Trans. Inf. Theory*, 62(6):3702–3720, 2016.
- [67] Yihong Wu and Pengkun Yang. Chebyshev polynomials, moment matching, and optimal estimation of the unseen. *Ann. Stat*, 47(2):857–883, 2019.