

# Noisy Population Recovery from Unknown Noise

Shachar Lovett\*

SLOVETT@UCSD.EDU

Jiapeng Zhang†

JPENG.ZHANG@GMAIL.COM

*Computer Science and Engineering department  
University of California, San Diego*

## Abstract

The noisy population recovery problem is a statistical inference problem, which is a special case of the problem of learning mixtures of product distributions. Given an unknown distribution on  $n$ -bit strings with support of size  $k$ , and given access only to noisy samples from it, where each bit is flipped independently with some unknown noise probability, estimate from a few samples the underlying parameters of the model. Previous work [De et al., FOCS 2016] designed polynomial time algorithms which work under the assumption that the noise parameters are known exactly. In this work, we remove this assumption, and show how to recover the underlying parameters, even when the noise is unknown, in quasi-polynomial time.

**Keywords:** Noisy recovery, Partial information, learning mixtures of product distributions

## 1. Introduction

Consider a database of patients in a hospital, where for each patient the database lists a large number of traits. Researchers are interested in obtaining this database to perform various statistical studies, but due to privacy concerns the database cannot be released. A possible solution (other than deleting identifying parameters of patients, such as their name) is to delete information at random from the database, or even better, add randomness to the information, with the goal that this will maintain the privacy of the original database, but would still provide researchers with useful information. The question is: does this process ensure privacy, or can the original database be recovered (up to its row order) from a lossy or noisy version of it?

The problem of recovery of data from lossy or noisy samples was studied extensively in statistics in the context of continuous distributions, and was introduced to computer science by Kearns et al. [Kearns et al. \(1994\)](#) who focused on discrete distributions. In this paper, we focus on the binary setting. Setting parameters, the goal is as follows: given a mixture of  $k$  product distributions supported on  $\{0, 1\}^n$ , recover the parameters within accuracy  $\varepsilon$  as efficiently as possible, by an algorithm which has access to samples from the overall mixed distribution. As far as we know, this may be possible to achieve in time  $\text{poly}(n, k, 1/\varepsilon)$  (although, there is some evidence that  $n^{\log k}$  might be a lower bound, see [Feldman et al. \(2008\)](#)). However, the best known algorithm to date is by [Feldman et al. \(2008\)](#), whose

---

\* Supported by NSF CAREER award 1350481 and NSF CCF award 1614023.

† Supported by NSF CAREER award 1350481 and NSF CCF award 1614023.

algorithm requires running time of  $\text{poly}(n, k, 1/\varepsilon)^{O(k^3)}$ . Thus, it can only recover mixtures of essentially a fixed number of product distributions.

A special case of this problem, termed “population recovery”, regained attention recently in a work by Dvir et al. [Dvir et al. \(2012\)](#), who related it to the problem of learning DNFs from partial information. The main reason is that in this restricted settings, much better dependence on the sparsity  $k$  is possible. It was further studied in [Moitra and Saks \(2013\)](#); [Wigderson and Yehudayoff \(2016\)](#); [Batman et al. \(2013\)](#); [Lovett and Zhang \(2015\)](#); [De et al. \(2016\)](#). This problem is better described as that of recovering an unknown sparse distribution given noisy samples.

Formally, suppose there is an unknown distribution  $\pi$  over  $k$  unknown elements in  $\{0, 1\}^n$  (which we think of as “centers”), and a vector of noise parameters  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n) \in [0, 1]^n$ . Samples from the noisy distribution are obtained as follows:

- Sample a string  $v \in \{0, 1\}^n$  according to  $\pi$ .
- Output  $x = v \oplus e \in \{0, 1\}^n$ , obtained by sampling each  $e_i \in \{0, 1\}$  independently with probability  $\Pr[e_i = 0] = (1 + \mu_i)/2$ .

Note that  $\mu_i = 1$  corresponds to no error in coordinate  $i$ , while  $\mu_i = 0$  corresponds to fully randomizing coordinate  $i$ . We denote this noisy distribution by  $\sigma(\pi, \boldsymbol{\mu})$

The goal is to recover all of the underlying parameters of the model, as efficiently as possible, given samples from the noisy distribution. Most of the work so far [Dvir et al. \(2012\)](#); [Wigderson and Yehudayoff \(2016\)](#); [Lovett and Zhang \(2015\)](#); [De et al. \(2016\)](#) focused on the case where the noise distribution is fully known, and moreover  $\mu_1 = \dots = \mu_n = \mu$  (we note that all these works easily extend to the case where  $\mu_1, \dots, \mu_n$  are exactly known, even if they are not equal). While an unrealistic assumption in practice, the main benefit is that the dependence on the sparsity of the mixture  $k$  can be greatly improved. The best results to date are by [De et al. \(2016\)](#) who show that in this case, all underlying parameters can be learned in time  $\text{poly}(n, 1/\varepsilon, k)$ , where we are suppressing the exact dependence on  $\mu$ . One exception is the work of Batman et al. [Batman et al. \(2013\)](#), who are able to recover in the same time complexity a superset of points in the support of  $\pi$ , even under much weaker assumptions on the noise. However, they cannot recover the remaining parameters (that is, the distribution  $\pi$  and the noise parameters  $\boldsymbol{\mu}$ ).

### 1.1. Our results

The goal of the current work is to relax the conditions under which we can learn mixtures of product distributions. Our main focus is on removing the assumption that the noise parameters have to be known. This question was posed by [Wigderson and Yehudayoff \(2016\)](#) as an open problem; we resolve it (at least partially) in this work. To do that, we need to assume that the noise is somewhat bounded. For  $\mu > 0$ , we say that the noise parameters  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)$  are  $\mu$ -bounded if  $\mu_i \geq \mu$  for all  $i \in [n]$ .

It is easy to see that, when the noise is a parameter of the model, it is impossible in general to recover the “true” parameters  $\pi, \boldsymbol{\mu}$  from the noisy distribution. For example, let  $n = 1$ ,  $\pi(0) = p, \pi(1) = 1 - p$  and  $\mu_1 = \mu$ . Then the noisy distribution is

$$\sigma(0) = p \frac{1 + \mu}{2} + (1 - p) \frac{1 - \mu}{2} = \frac{1 + (2p - 1)\mu}{2}, \quad \sigma(1) = 1 - \sigma(0).$$

So, even given perfect knowledge of  $\sigma$ , we can only recover  $(2p - 1)\mu$ . So for example, we cannot distinguish  $p = 2/3, \mu = 1$  from  $p = 3/4, \mu = 2/3$ . Thus, the best that we can do is to recover some  $\hat{\pi}, \hat{\boldsymbol{\mu}}$  such that the noisy distribution  $\sigma(\hat{\pi}, \hat{\boldsymbol{\mu}})$  is close in statistical distance to the observable noisy distribution. Thus, it gives a succinct representation to the observable data (and moreover, it is proper, as it comes from the same class of distributions).

**Theorem 1 (Main theorem)** *Fix  $\mu > 0$ . Let  $\pi$  be an unknown distribution over  $\{0, 1\}^n$  of support size  $k$ . Let  $\boldsymbol{\mu} \in [\mu, 1]^n$  be unknown noise parameters. Then, for any  $\varepsilon > 0$  there exists an algorithm which, given samples from  $\sigma(\pi, \boldsymbol{\mu})$ , returns*

- A distribution  $\hat{\pi}$  over  $\{0, 1\}^n$  of support size at most  $k$ .
- A vector  $\hat{\boldsymbol{\mu}} \in [\mu, 1]^n$  of noise parameters.

*Such that the noisy distributions  $\sigma(\pi, \boldsymbol{\mu}), \sigma(\hat{\pi}, \hat{\boldsymbol{\mu}})$  are  $\varepsilon$ -close in statistical distance. The algorithm requires time  $\text{poly}(n^{\log k}, (1/\varepsilon)^{(\log k)^2}, k^{(\log k)^3})$ , where we suppress the exact dependence on  $\mu$ .*

## 1.2. Proof overview

Our proof has three steps, the first step is to recover a list of candidates of noise parameters, such that one of them is very close to the real noise parameters. Our main contribution is the following theorem.

**Theorem 2** *Fix  $\mu > 0$ . Let  $V \subset \{0, 1\}^n$  be an unknown set of size  $|V| = k$ ,  $\pi$  an unknown probability distribution on  $V$ ,  $\boldsymbol{\mu} \in [\mu, 1]^n$  unknown noise parameters. Then, for any  $\varepsilon > 0$  there exists an algorithm which, given samples from the noisy distribution  $\sigma = \sigma(\pi, \boldsymbol{\mu})$ , returns a list  $L \subset [\mu, 1]^n$  of potential noise parameters such that*

- There exist  $\hat{\boldsymbol{\mu}} \in L$  for which  $\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_\infty \leq \varepsilon$ .
- $|L| \leq (k/\varepsilon\mu)^{O(\log^2 k)} \cdot n^{\log k}$ .

*Moreover, the algorithm runs in time  $\text{poly}(|L|)$ .*

Once we recovered almost perfect noise parameters, we will run the algorithm of [Wigderson and Yehudayoff \(2016\)](#), which can recover the underlying parameters assuming that the noise parameters are perfectly known. As part of the analysis, we show that the algorithm is robust, in the sense that if we know the noise parameters up to a small error, it still succeeds in recovering the remaining parameters  $V, \pi$ .

**Theorem 3 (Recovering mixture from almost perfect noise parameters)** *Fix  $\mu > 0$ . Let  $\boldsymbol{\mu} \in [\mu, 1]^n$  be a vector of noise parameters, and let  $\boldsymbol{\mu}' \in [\mu, 1]^n$  be assumed noise parameters, where  $\|\boldsymbol{\mu} - \boldsymbol{\mu}'\|_\infty \leq \delta$  for  $\delta = \varepsilon \cdot k^{-O(\log k)}$ . Let  $\pi$  be an unknown distribution over  $\{0, 1\}^n$  of support size  $k$ . Given access to samples from the noisy distribution corresponding to  $\sigma(\pi, \boldsymbol{\mu})$ , the algorithm given in [Wigderson and Yehudayoff \(2016\)](#), which assumes (erroneously) that  $\boldsymbol{\mu}'$  are the noise parameters, still recovers  $\pi$  within accuracy  $\varepsilon$  in the same running time.*

We have so far built a short list of potential candidates for the distribution  $\pi$  and the noise parameters  $\boldsymbol{\mu}$ . In the last stage, we prune any distribution which does not match the actual observed noisy distribution. This part is generic and does not rely on the specific properties of our parametric model, except that we can both sample and compute probability of individual elements. See Lemma 12 for the details.

### 1.3. Open problems

An obvious open problem is to improve the parameters in Theorem 2. Concretely, if the methods from Lovett and Zhang (2015); De et al. (2016) can be extended to the framework of this paper, then it will improve our algorithm runtime to  $\text{poly}((n \cdot k/\varepsilon)^{\log k})$ .

Another open problem is to extend our current algorithm to the more general model of mixtures of product distributions. We note that in this case, the best algorithms have exponential dependency on  $k$ , in contrast with the current work where the dependence is quasi-polynomial (e.g. the exponent is poly-logarithmic in  $k$ ).

## 2. Preliminaries

### 2.1. Noisy distributions

We recall and slightly adapt some definitions from the introduction. Let  $V \subseteq \{0, 1\}^n$  be a collection of  $|V| = k$  binary vectors, and let  $\pi$  be a distribution on  $V$ , called the *mixing distribution*. Both  $V$  and  $\pi$  are unknown, and our goal is to recover  $V$  and  $\pi$  from noisy samples. Let  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n) \in [0, 1]^n$  be a vector of (also unknown) noise parameters. We sample from the noisy distribution as follows:

- Sample a string  $v \in V$  according to  $\pi$ .
- Output  $x = v \oplus e \in \{0, 1\}^n$ , obtained by sampling each  $e_i \in \{0, 1\}$  independently with probability  $\Pr[e_i = 0] = (1 + \mu_i)/2$ .

We denote the resulting noisy distribution by  $\sigma(\pi, \boldsymbol{\mu})$ . For  $0 < \mu < 1$ , we say that the noise is  $\mu$ -bounded if  $\mu_i \geq \mu$  for all  $i \in [n]$ .

It will be useful to view this as applying a noise operator to the original distribution. Slightly abusing notations, identify a distribution  $\pi$  on  $\{0, 1\}^n$  with the function  $\pi : \{0, 1\}^n \rightarrow \mathbb{R}$  where  $\pi(x)$  is the probability that  $\pi$  assigns to  $x$ . Let  $T_{\boldsymbol{\mu}}$  denote the noise operator operating on functions  $f : \{0, 1\}^n \rightarrow \mathbb{R}$ , defined as

$$(T_{\boldsymbol{\mu}} \circ f)(x) = \mathbb{E}_e[f(x + e)],$$

where  $e \in \{0, 1\}^n$  is sampled as  $\Pr[e_i = 0] = (1 + \mu_i)/2$  independently for each  $i \in [n]$ . Then  $\sigma(\pi, \boldsymbol{\mu}) = T_{\boldsymbol{\mu}} \circ \pi$ .

For a distribution  $\pi$  on  $\{0, 1\}^n$  and  $S \subseteq [n]$ , define  $\pi_S$  to be the marginal distribution of  $\pi$  on  $\{0, 1\}^S$ . For  $\boldsymbol{\mu}$  let  $\boldsymbol{\mu}_S = (\mu_i : i \in S)$ . Observe that the noisy distribution on  $S$  can be obtained by applying noise to the original distribution on  $S$ ,

$$(T_{\boldsymbol{\mu}} \circ \pi)_S = T_{\boldsymbol{\mu}_S} \circ \pi_S.$$

We denote by  $T_{\boldsymbol{\mu}}^{-1}, T_{\boldsymbol{\mu}_S}^{-1}$  the inverse operators, which exist whenever the noise is  $\mu$ -bounded for any  $\mu > 0$ .

## 2.2. PID graphs

We follow definitions from [Wigderson and Yehudayoff \(2016\)](#). Let  $V \subset \{0, 1\}^n$ . A *PID graph* defined over  $V$  is a labeled directed graph  $G = (V, E, \{S_v : v \in V\})$  defined as follows. The vertices of  $G$  are identified with  $V$ . For every  $v \in V$ , we associate a subset  $S_v \subseteq [n]$ , called the *partial ID* (abbreviated PID) of  $v$ . For any  $S \subseteq [n]$ , denote by  $v[S] \in \{0, 1\}^S$  the restriction of  $v$  to the coordinates in  $S$ . The *imposters* of  $v$  are the set of vertices which agree with  $v$  on  $S_v$ ,

$$I(v) = \{u \in V \setminus \{v\} : u[S_v] = v[S_v]\}.$$

The edges of  $G$  are  $E = \{(v, u) : v \in V, u \in I(v)\}$ . A partial ID  $v_S$  is said to be a *unique ID* (abbreviated UID) for  $v \in V$  if  $I(v) = \emptyset$ . That is,  $v$  is uniquely identified by its restriction to  $S_v$ .

Wigderson and Yehudayoff proved that for any choice of  $V \subset \{0, 1\}^n$  there exists a choice of PIDs such that the resulting PID graph is efficient in two ways: the size of the sets  $S_v$ , as well as the depth of the tree, are both at most logarithmic in  $V$ . In this paper, all logarithms are in base two.

**Theorem 4 ([Wigderson and Yehudayoff \(2016\)](#))** *For any subset  $V \subset \{0, 1\}^n$  of size  $|V| = k$ , there exists a choice of  $S_v \subseteq [n]$  for each  $v \in V$ , such that*

- (i)  $|S_v| \leq \log k$  for all  $v \in V$ .
- (ii) The PID graph  $G = (V, E, \{S_v : v \in V\})$  has depth  $\leq \log k$ .

Moreover, the choices of  $\{S_v : v \in V\}$  and the PID graph can be constructed in time  $\text{poly}(n, k)$ .

## 2.3. Previous work: learning when the noise parameters are known

Several clever ideas introduced in previous works allow to learn the underlying parameters of the distribution, under the assumption that the noise parameters are known exactly. To recall, the underlying parameters are as follows: an unknown distribution  $\pi$  over an unknown set in  $\{0, 1\}^n$  of size  $k$ . We denote by  $V \subset \{0, 1\}^n$  the support of  $\pi$ .

First, [Dvir et al. \(2012\)](#) showed that one can reduce to the easier problem, where the support  $V$  of  $\pi$  is known, but the probabilities assigned by  $\pi$  are still unknown. That is, assuming the existence of an algorithm which can recover  $\pi$  given  $V$  from noisy samples, they designed an algorithm which recovers both  $V$  and  $\pi$ . The high level idea is to recover them one coordinate at a time.

As a base case, if  $n = 1$  then we can clearly set  $V = \{0, 1\}$  and recover  $\pi$ . Note that as we allow  $\pi(v) = 0$  for some  $v \in V$ , so its sufficient that  $V$  is a superset of the support of  $\pi$ . Given the restriction of  $V, \pi$  to the first  $t < n$  coordinates, we extend them to  $t + 1$  coordinates as follows. First, extend  $V$  to all possible values of the  $t + 1$  bit; this at most doubles the size of  $V$ . Then, run our assumed algorithm, which learns  $\pi$  given  $V$ . Then, delete from  $V$  any element whose probability under  $\pi$  is negligible. Thus, we may assume from now on that  $V$  is known while  $\pi$  is unknown; this will only introduce an overhead of  $n$  in the overall running time. See the paper [Dvir et al. \(2012\)](#) for further details.

Next, we introduce the learning algorithm of [Wigderson and Yehudayoff \(2016\)](#), under the assumption that both  $V$  and the noise parameters are known. It starts by constructing a PID graph for  $V$  of depth  $\log k$ , where each PID has size  $\leq \log k$ . We assume below that  $\mu_1, \dots, \mu_n \geq \mu$ . They do as follows:

- Estimate the restriction of the distribution  $\pi$  to each  $S_v$ . As the noise is known exactly, this can be achieved as follows. First, estimate the noisy distribution  $\sigma_{S_v}$ . Then compute

$$\pi_{S_v} = T_{\mu_{S_v}}^{-1} \circ \sigma_{S_v}.$$

It can be shown that in order to estimate  $\pi_{S_v}$  within accuracy  $\varepsilon$ , we need to estimate  $\sigma_{S_v}$  within accuracy  $\varepsilon \mu^{|S_v|}$ . Recalling that  $|S_v| = \log k$  and surpassing the dependence on  $\mu$ , this takes  $\text{poly}(k, 1/\varepsilon)$  samples and time.

- Solve equations for  $\pi(v)$  for each  $v \in V$ . The main idea is to traverse the PID graph from bottom to top. Given a node  $v \in V$ , assume that we already recovered  $\pi(u)$  for all his children, namely for all  $u \in I(v)$ . We also know  $\Pr[\pi_{S_v} = v[S_v]]$  by the first step. Then we can solve

$$\pi(v) = \Pr[\pi_{S_v} = v[S_v]] - \sum_{u \in I(v)} \pi(u).$$

Due to approximation errors, if we are aiming for an overall error of  $\varepsilon$  in approximating  $\pi$ , it is needed to approximate  $\pi_{S_v}$  to within error of  $\varepsilon \cdot k^{-O(\log k)}$ , where we suppress the exact dependence on  $\mu$ . The overall running time of the algorithm is thus  $\text{poly}(n, 1/\varepsilon, k^{\log k})$ . See the paper [Wigderson and Yehudayoff \(2016\)](#) for further details.

We summarize these in the following theorem.

**Theorem 5 (Recovering mixture from exact noise parameters)** *Fix  $\mu > 0$ . Let  $\boldsymbol{\mu} \in [\mu, 1]^n$  be a known vector of noise parameters. Let  $\pi$  be an unknown distribution over  $\{0, 1\}^n$  of support size  $k$ . Given access to samples from the noisy distribution corresponding to  $(\pi, \boldsymbol{\mu})$ , it is possible to recover  $\pi$  within accuracy  $\varepsilon$  in time*

$$\text{poly}(n, 1/\varepsilon, k^{\log k})$$

where we suppress the exact dependence on  $\mu$ .

We would need a variant of Theorem 5 where we assume that the noise parameters are known not exactly, but up to a very small error, so that the existing algorithms would work as is. That is, if we are given noisy samples according to  $(\pi, \boldsymbol{\mu})$ , but instead we believe that the noise parameters are  $\boldsymbol{\mu}'$  where  $|\mu_i - \mu'_i| \leq \delta$  for all  $i \in [n]$  where  $\delta$  is small enough, then the algorithm given in Theorem 5 would still succeed in producing an  $\varepsilon$  approximation of  $\pi$ .

To calculate how small  $\delta$  needs to be, observe that the noisy distribution is used in the algorithm of [Wigderson and Yehudayoff \(2016\)](#) to estimate  $\pi_{S_v}$  for all PIDs  $S_v, v \in V$ . So, the algorithm would estimate  $\sigma_{S_v}$  and calculate  $\pi_{S_v} = T_{\mu'_{S_v}}^{-1} \circ \sigma_{S_v}$  instead of  $\pi_{S_v} = T_{\mu_{S_v}}^{-1} \circ \sigma_{S_v}$ . We would like these two distributions to be close enough, so that the remainder of the algorithm would still succeed.

**Claim 6** Let  $\boldsymbol{\mu}, \boldsymbol{\mu}' \in [\mu, 1]^n$  be two noise distribution such that  $\|\boldsymbol{\mu} - \boldsymbol{\mu}'\|_\infty \leq \delta$ . Then for any  $S \subset [n]$  and any distribution  $\sigma_S$  on  $\{0, 1\}^S$ ,

$$\|T_{\boldsymbol{\mu}_S}^{-1} \circ \sigma_S - T_{\boldsymbol{\mu}'_S}^{-1} \circ \sigma_S\|_1 \leq \delta |S| (2/\mu)^{2|S|}.$$

We defer the proof to the appendix (Claim 13). Note that as  $|S_v| \leq \log k$  for all  $v \in V$ , the statistical distance introduced is  $\delta \cdot \text{poly}(k)$  where we suppress the dependence on  $\mu$ . To recall, the algorithm of [Wigderson and Yehudayoff \(2016\)](#) requires an approximation of each  $\pi_{S_v}$  to within  $\varepsilon \cdot k^{-O(\log k)}$  accuracy; this requires us to also choose  $\delta = \varepsilon \cdot k^{-O(\log k)}$ . We summarize this in the following theorem.

**Theorem 3 (restated)** Fix  $\mu > 0$ . Let  $\boldsymbol{\mu} \in [\mu, 1]^n$  be a vector of noise parameters, and let  $\boldsymbol{\mu}' \in [\mu, 1]^n$  be assumed noise parameters, where  $\|\boldsymbol{\mu} - \boldsymbol{\mu}'\|_\infty \leq \delta$  for  $\delta = \varepsilon \cdot k^{-O(\log k)}$ . Let  $\pi$  be an unknown distribution over  $\{0, 1\}^n$  of support size  $k$ . Given access to samples from the noisy distribution corresponding to  $\sigma(\pi, \boldsymbol{\mu})$ , the algorithm given in [Theorem 5](#), which assumes (erroneously) that  $\boldsymbol{\mu}'$  are the noise parameters, still recovers  $\pi$  within accuracy  $\varepsilon$  in the same time.

### 3. Recovering mixtures with unknown noise parameters

Let  $V \subset \{0, 1\}^n$  be an unknown set of  $|V| = k$  vectors,  $\pi$  an unknown distribution on  $V$ , and  $\boldsymbol{\mu} \in [0, 1]^n$  an unknown set of noise parameters which are assumed to be  $\mu$ -bounded for some  $0 < \mu < 1$ . Our goal is to recover the noise parameters, and then apply [Theorem 3](#) to recover the remaining parameters  $V, \pi$ .

Let  $\sigma = \sigma(\pi, \boldsymbol{\mu}) = T_{\boldsymbol{\mu}} \circ \pi$  denote the observed distribution. Recall that  $\sigma_S$  for  $S \subseteq [n]$  is the marginal of  $\sigma$  on  $S$ .

#### 3.1. Recovering the noise parameters

Our main contribution is an algorithm which recovers a small list of noise parameters, one of which is guaranteed to be close to the true noise parameters.

**Theorem 2 (restated)** Fix  $\mu > 0$ . Let  $V \subset \{0, 1\}^n$  be an unknown set of size  $|V| = k$ ,  $\pi$  an unknown probability distribution on  $V$ ,  $\boldsymbol{\mu} \in [\mu, 1]^n$  unknown noise parameters. Then, for any  $\varepsilon > 0$  there exists an algorithm which, given samples from the noisy distribution  $\sigma = \sigma(\pi, \boldsymbol{\mu})$ , returns a list  $L \subset [\mu, 1]^n$  of potential noise parameters such that

- There exist  $\hat{\boldsymbol{\mu}} \in L$  for which  $\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_\infty \leq \varepsilon$ .
- $|L| \leq (k/\varepsilon\mu)^{O(\log^2 k)} \cdot n^{\log k}$ .

Moreover, the algorithm runs in time  $\text{poly}(|L|, n)$ .

Let  $V$  be the unknown population set, and let  $\{S_v \subseteq [n] : v \in V\}$  be the (unknown) corresponding PIDs guaranteed by [Theorem 4](#), so that  $|S_v| \leq \log k$  and the resulting PID graph  $G = (V, E, \{S_v : v \in V\})$  has depth  $\leq \log k$ .

For  $v \in V$ , let  $d(v) \in \{1, \dots, \log k\}$  denote the *depth* of  $v$  in  $G$ , where the root has depth 1 and leaves have depth  $\leq \log k$ . Intuitively, the probabilities of nodes with larger depth are easier to recover, as they have less imposters. To formalize this, we introduce the notion of dominant elements.

**Definition 7 (Dominant elements)** Let  $\gamma := \varepsilon/(16k)$ . An element  $v \in V$  is said to be dominant if  $\pi(v) \geq \gamma^{d(v)}$ , and for each  $u \in I(v)$ ,  $\pi(u) < \gamma^{d(u)}$ .

**Claim 8** There exist at least one dominant element  $v \in V$ .

**Proof** Let  $A = \{v \in V : \pi(v) \geq \gamma^{d(v)}\}$ . If  $A$  is nonempty, then any node  $v \in A$  of maximum depth is dominant. The set  $A$  cannot be empty since this would imply that

$$1 = \sum_{v \in V} \pi(v) < \sum_{v \in V} \gamma^{d(v)} \leq k\gamma = \varepsilon/16 < 1.$$

■

Assume that  $v \in V$  is dominant, and that furthermore we know  $(\mu_i : i \in S_v)$ . We next show how to use this to learn all the other noise parameters.

**Claim 9** Let  $v \in V$  be a dominant element and denote  $S := S_v$ . Let  $\rho := T_{\mu_S}^{-1} \circ \sigma$ . For each  $i \notin S$  define

$$\hat{\mu}_i := 2 \frac{\rho_{S \cup \{i\}}[v_{S \cup \{i\}}]}{\rho_S(v[S])} - 1.$$

Then  $\|\hat{\mu} - \mu\|_\infty \leq \varepsilon/2$ , where to recall  $\mu$  are the real noise parameters.

**Proof** Let  $d = d(v)$ . Initially, note that  $\rho_S = \pi_S$  and hence

$$\rho_S(v[S]) = \pi_S(v[S]) = \pi(v) + \sum_{u \in I(v)} \pi(u).$$

By our assumption that  $v$  is dominant, we have  $\pi(v) \geq \gamma^d$  and

$$\sum_{u \in I(v)} \pi(u) \leq k\gamma^{d+1} \leq (k\gamma)\pi(v).$$

Thus,

$$1 \leq \frac{\rho_S(v[S])}{\pi(v)} \leq 1 + k\gamma.$$

Next, as  $\rho_{S \cup \{i\}} = T_{\mu_i} \circ \pi_{S \cup \{i\}}$  then we have

$$\begin{aligned} \rho_{S \cup \{i\}}(v[S \cup \{i\}]) &= (T_{\mu_i} \circ \pi_{S \cup \{i\}})(v[S \cup \{i\}]) \\ &= \sum_{u \in \{v\} \cup I_v, u_i = v_i} \pi(u) \left( \frac{1 + \mu_i}{2} \right) + \sum_{u \in I_v, u_i \neq v_i} \pi(u) \left( \frac{1 - \mu_i}{2} \right). \end{aligned}$$

We thus have

$$\left| \rho_{S \cup \{i\}}(v[S \cup \{i\}]) - \pi(v) \left( \frac{1 + \mu_i}{2} \right) \right| \leq \sum_{u \in I(v)} \pi(u) \leq (k\gamma)\pi(v).$$



Combining the estimates we obtain that

$$|\rho_{S \cup \{i\}}(v[S \cup \{i\}]) - \rho_S(v[S]) \cdot \left(\frac{1 + \mu_i}{2}\right)| \leq 2k\gamma\pi(v) \leq 4k\gamma\rho_S(v[S]).$$

Hence

$$\left| \frac{\rho_{S \cup \{i\}}(v[S \cup \{i\}])}{\rho_S(v[S])} - \frac{1 + \mu_i}{2} \right| \leq 4k\gamma.$$

This implies that  $|\hat{\mu}_i - \mu_i| \leq 8k\gamma \leq \varepsilon/2$ . ■

In order to apply Claim 9, we would need to first guess  $v \in V$  to be dominant, and then enumerate over  $\mu_i, i \in S$ . The next claim shows that it suffices to approximate  $\mu_i$  within a fine enough accuracy.

**Claim 10** *Let  $v \in V$  be a dominant element and denote  $S := S_v$ . Let  $\mu'_i \in [\mu, 1]$  for  $i \in S$  be such that  $|\mu'_i - \mu_i| \leq \delta$  for  $\delta = (k/\varepsilon\mu)^{-O(\log k)}$ . Define  $\rho' := T_{\mu'_S}^{-1} \circ \sigma$ . For each  $i \notin S$  define*

$$\hat{\mu}_i := 2 \frac{\rho'_{S \cup \{i\}}[v_{S \cup \{i\}}]}{\rho'_S(v[S])} - 1.$$

Then  $|\hat{\mu}_i - \mu_i| \leq \varepsilon$  for all  $i \in [n]$ .

**Proof** Let  $\rho = T_{\mu_S}^{-1} \circ \sigma$  be the distribution obtain after removing the correct noises in the coordinates of  $S$ . We will show that as expected,  $\rho'_{S \cup \{i\}}[v_{S \cup \{i\}}] \approx \rho_{S \cup \{i\}}[v_{S \cup \{i\}}]$  and  $\rho'_S[v_S] \approx \rho_S[v_S]$ , and then apply Claim 9. Formally, we will show for a small enough  $\eta > 0$  that

$$\|\rho' - \rho\|_1 \leq \eta.$$

In the proof of Claim 9 we showed that

$$\rho_S(v[S]) \geq \pi(v) \geq \gamma^{|S|}.$$

Thus, if we choose  $\eta = O(\varepsilon \cdot \gamma^{|S|})$  then we would obtain that both the numerator and denominator in the definition of  $\hat{\mu}_i$  approximate up to a multiplicative factor of  $1 \pm O(\varepsilon)$  the corresponding quantities with  $\rho$  instead of  $\rho'$ ; the claim then follows from Claim 9.

So, our goal is to show that for a small enough  $\delta > 0$  we obtain  $\|\rho' - \rho\|_1 \leq \eta$ . We have

$$\rho = T_{\mu_S}^{-1} \circ \sigma = \left( \prod_{i \in S} T_{\mu_i}^{-1} \right) \circ \sigma.$$

Explicitly computing this, the inverse noise operator  $T_{\mu}^{-1}$  corresponds to the following  $2 \times 2$  matrix

$$\frac{1}{2} \begin{pmatrix} 1 + \mu & 1 - \mu \\ 1 - \mu & 1 + \mu \end{pmatrix}^{-1} = \frac{1}{2\mu} \begin{pmatrix} 1 + \mu & -1 + \mu \\ -1 + \mu & 1 + \mu \end{pmatrix}$$

Thus we have

$$\rho(x) = \sum_{y \in \{0,1\}^n: y_{S^c} = x_{S^c}} \prod_{i \in S} \left( \frac{(-1)^{x_i + y_i} + \mu_i}{2\mu_i} \right) \sigma(y)$$

and similarly

$$\rho'(x) = \sum_{y \in \{0,1\}^n: y_{S^c} = x_{S^c}} \prod_{i \in S} \left( \frac{(-1)^{x_i + y_i} + \mu'_i}{2\mu'_i} \right) \sigma(y).$$

We can bound

$$\begin{aligned} & \left| \prod_{i \in S} \left( \frac{(-1)^{x_i + y_i} + \mu_i}{2\mu_i} \right) - \prod_{i \in S} \left( \frac{(-1)^{x_i + y_i} + \mu'_i}{2\mu'_i} \right) \right| \\ & \leq \frac{1}{\mu^{|S|-1}} \sum_{i \in S} \frac{|\mu'_i - \mu_i|}{2\mu'_i \mu_i} \\ & \leq \frac{|S|\delta}{\mu^{|S|+1}}. \end{aligned}$$

where we assume  $\mu'_i, \mu_i \geq \mu$  and  $|\mu'_i - \mu_i| \leq \delta$ . Thus

$$\|\rho' - \rho\|_1 \leq \frac{|S|\delta 2^{|S|}}{\mu^{|S|+1}} \leq \delta(k/\mu)k^{\log(1/\mu)} \log k.$$

we to obtain  $\|\rho' - \rho\|_1 \leq O(\varepsilon\gamma^{|S|})$  we need to choose  $\delta = (k/\varepsilon\mu)^{-O(\log k)}$ .  $\blacksquare$

In final step, we need to argue that given some noise parameters  $\mu'$  we can estimate  $\rho'$ .

**Claim 11** Fix  $S \subseteq [n]$ , noise parameters  $\mu'_S \in [\mu, 1]^S$ . Let  $\rho'(x) = T_{\mu'_S}^{-1}\sigma$ . Then we can estimate  $\rho'_S$  and  $\rho'_{S \cup \{i\}}$  for any  $i$  within error  $\eta$ , with success probability  $1 - \delta$ , using  $(2/\mu)^{O(|S|)} \text{poly}(1/\eta) \log(1/\delta)$  samples.

**Proof** Let  $S' = S$  or  $S' = S \cup \{i\}$ . Estimate  $\sigma_{S'}$  within error  $\mu^{|S'|}\eta$  with probability  $1 - \delta$  using standard estimation techniques. As each value of  $\rho'_{S'}$  is the linear combination of the elements of  $\sigma_{S'}$  with coefficients bounded by  $(1/\mu)^{|S'|}$ , the claim follows.  $\blacksquare$

We can now describe the algorithm, that would generate a list of potential suggestions for the noise parameters. The main observation is that in order to compute  $\hat{\mu}$  given in Claim 10, we don't really need to know  $v \in V$ . We only need to know:

- A PID  $S_v \subset [n]$  of size  $|S_v| \leq \log k$ ,
- The value  $v[S_v]$ ,
- A good enough approximation for  $\mu_{S_v}$ .

Note that Claim 10 requires us to know  $v[S_v \cup \{i\}]$  for each  $i \in [n]$ . However, as we assume that  $v$  is dominant, this can be easily found from samples. Indeed, let  $S = S_v$ , fix  $i \notin S$  and let  $S' = S \cup \{i\}$  and  $v' = v \oplus e_i$ . Then

$$\pi_{S'}(v'[S']) \leq (k\gamma) \cdot \pi_{S'}(v[S']) \leq \varepsilon \cdot \pi_{S'}(v[S']).$$

That is, if we sample  $x \sim \pi$  and condition that  $x_S = v_S$ , then  $\Pr[x_i = v_i | x_S = v_S] \geq 1 - \varepsilon$ . Now, we do not have access to  $\pi$ , but we do have access to  $\rho = T_{\mu_S}^{-1} \circ \sigma$ . Note that  $\rho_{S'} = T_{\mu_i} \circ \pi_{S'}$ , and hence

$$\begin{aligned} \Pr_{x \sim \rho}[x_i = v_i | x_S = v_S] &= \Pr_{x \sim \pi}[x_i = v_i | x_S = v_S] \cdot \left(\frac{1 + \mu_i}{2}\right) + \Pr_{x \sim \pi}[x_i \neq v_i | x_S = v_S] \cdot \left(\frac{1 - \mu_i}{2}\right) \\ &\geq \frac{1 + \mu_i}{2} - \varepsilon \\ &\geq \frac{1 + \mu}{2} - \varepsilon. \end{aligned}$$

Thus, we can learn all the bits of  $v$  by samples. We thus obtain the following algorithm.

**Algorithm: Recover-Noise-Parameters**

- **Input:** Samples from  $\sigma = \sigma(\pi, \mu)$ .
  - **Output:** A list  $L$  of candidates for the noise parameters.
1. Initialize empty list  $L = \emptyset$ .
  2. Enumerate  $S \subseteq [n]$  with  $|S| \leq \log k$  and  $v_S \in \{0, 1\}^S$ . For each choice:
    - 2.1 Enumerate all possible values of  $\mu'_S \in [\mu, 1]^S$  within accuracy  $\delta = (k/\varepsilon\mu)^{-O(\log k)}$  in each coordinate. For each choice:
      - 2.1.1 For each  $i \in [n]$ , estimate  $T_{\mu'_S}^{-1} \circ \sigma_{S \cup \{i\}}$  and compute  $\mu'_i$  for all  $i \notin S$  as in Claim 10.
      - 2.1.2 Add  $\mu' \in [\mu, 1]^n$  to the list  $L$ .
  3. Output  $L$ .

This concludes the proof of Theorem 2: some choice of  $v \in V$  is dominant by Claim 8. For this  $v$ , some choice of  $\mu'_S$  is  $\delta$ -close to the true noise parameters on  $S$ . By Claim 10, this suffices to learn the remaining noise parameters. The running time is dominated by the enumeration of  $\mu'$ , which takes  $(1/\delta)^{\log k} = (k/\varepsilon\mu)^{-O(\log^2 k)}$  time.

### 3.2. The full algorithm

Given the algorithms given in Theorem 3 and Theorem 2, the full algorithm follows by simply composing the two. To recall, the algorithm given in Theorem 3 guarantees to approximate  $\pi$  within an error of  $\varepsilon$ , assuming knowledge of all the noise parameters to within accuracy of  $\delta = \varepsilon \cdot k^{-O(\log k)}$ . Its run time is  $\text{poly}(n^{\log k}, 1/\varepsilon, k^{\log k})$ . The algorithm given in Theorem 2 outputs a list  $L$  of potential noise parameters. If we want one of them to be a  $\delta$ -approximation of the true noise parameters, we have  $|L| = (k/\delta)^{O(\log^2 k)}$ . Thus, we obtain a list of

$$T := |L| = \text{poly}(n^{\log k}, (1/\varepsilon)^{\log^2 k}, k^{\log^3 k})$$

potential choices of  $\{(\hat{\pi}^{(i)}, \hat{\mu}^{(i)}) : i \in [T]\}$ , where it is guaranteed that there exists some  $i \in [T]$  for which

$$\|\hat{\pi}^{(i)} - \pi\|_1 \leq \varepsilon, \quad \|\hat{\mu}^{(i)} - \mu\|_\infty \leq \delta$$

and hence  $\|\sigma(\pi, \boldsymbol{\mu}) - \sigma(\hat{\pi}^{(i)}, \hat{\boldsymbol{\mu}}^{(i)})\|_1 \leq O(\varepsilon)$ .

To conclude, we have a list of potential noisy distributions described by  $\hat{\pi}^{(i)}, \hat{\boldsymbol{\mu}}^{(i)}$ . Let  $\hat{\sigma}^{(i)} = \sigma(\hat{\pi}^{(i)}, \hat{\boldsymbol{\mu}}^{(i)})$ . Observe that we can both sample from  $\hat{\sigma}^{(i)}$ , as well as calculate the probability of each specific element. The next lemma shows that this is sufficient to prune any distribution  $\hat{\sigma}^{(i)}$  which is far from the true distribution  $\sigma$ .

**Lemma 12 (Pruning lemma)** *Let  $\mathcal{S} = \{\sigma^{(i)} : i \in [T]\}$  be a set of distributions over a universe  $X$ , given by some succinct representation, such that we can:*

- *Sample efficiently from each  $\sigma^{(i)}$ .*
- *Calculate  $\sigma^{(i)}(x)$  for each  $x \in X$ .*

*Let  $\sigma$  be an (unknown) observable distribution over  $X$  that we can get samples from. Then, for any  $\varepsilon > 0$  we can find a subset  $\mathcal{S}' \subset \mathcal{S}$  such that:*

- *If the statistical distance between  $\sigma$  and  $\sigma^{(i)}$  is  $\geq 4\varepsilon$ , then  $\sigma^{(i)} \notin \mathcal{S}'$ .*
- *If the statistical distance between  $\sigma$  and  $\sigma^{(i)}$  is  $\leq \varepsilon$ , then  $\sigma^{(i)} \in \mathcal{S}'$ .*

*Moreover, we can find  $\mathcal{S}'$  with high probability in time  $O((T/\varepsilon)^2 \cdot \log T)$ .*

**Proof** The main idea is as follows. Assume that  $\sigma^{(i)}, \sigma^{(j)} \in \mathcal{S}$  are two distributions whose statistical distance is  $\geq 4\varepsilon$ . Then it cannot be the case that both are  $\varepsilon$ -close to  $\sigma$ . The following procedure will reject at least one of them (one which is not  $\varepsilon$ -close to  $\sigma$ ). Define

$$A = \{x : \sigma^{(i)}(x) > \sigma^{(j)}(x)\}.$$

By definition of statistical distance,

$$\Pr_{x \sim \sigma^{(i)}} [x \in A] - \Pr_{x \sim \sigma^{(j)}} [x \in A] \geq 4\varepsilon.$$

We will estimate  $\Pr[\sigma \in A]$ ,  $\Pr[\sigma^{(i)} \in A]$ ,  $\Pr[\sigma^{(j)} \in A]$ , and reject either  $\sigma^{(i)}$  or  $\sigma^{(j)}$  (or both) for which the probability is  $2\varepsilon$  far from that of  $\sigma$ . To estimate  $\Pr[\sigma \in A]$ , sample  $x \sim \sigma$ , calculate  $\sigma^{(i)}(x), \sigma^{(j)}(x)$  and decide whether  $x \in A$  or not. Similarly estimate  $\Pr[\sigma^{(i)} \in A]$  and  $\Pr[\sigma^{(j)} \in A]$ . Reject  $i$  if the estimated  $\Pr[\sigma \in A] - \Pr[\sigma^{(i)} \in A]$  exceed  $2\varepsilon$ . To make sure that we make the correct decision with probability  $1 - \delta$  we need  $O((1/\varepsilon)^2 \log(1/\delta))$  samples.

So, set  $\delta = 1/|T|^2$ , enumerate all pairs  $i, j \in [T]$  and apply the above procedure. Let  $\mathcal{S}'$  be the non-rejected distributions. By the above analysis, with high probability we will keep all distributions which are  $\varepsilon$ -close to  $\sigma$ , and reject all distributions which are  $4\varepsilon$ -far from  $\sigma$ . ■

The full algorithm is as follows.

**Algorithm: Recover-Mixture**

- **Input:** Samples from  $\sigma = \sigma(\pi, \boldsymbol{\mu})$ .
  - **Output:**  $\hat{\pi}, \hat{\boldsymbol{\mu}}$  such that  $\sigma(\hat{\pi}, \hat{\boldsymbol{\mu}})$  is  $\varepsilon$ -close to  $\sigma$  in statistical distance.
1. Run algorithm Recover-Noise-Parameters, given in Theorem 2, to obtain a list  $L$  of potential noise parameters.
  2. For each noise parameter  $\hat{\boldsymbol{\mu}}^{(i)} \in L$ , recover  $\hat{\pi}^{(i)}$  by applying the Wigderson-Yehudayoff algorithm given in Theorem 3.
  3. Apply the pruning procedures given in Lemma 12 to prune any pair  $\hat{\pi}^{(i)}, \hat{\boldsymbol{\mu}}^{(i)}$  whose noisy distribution  $\sigma(\hat{\pi}^{(i)}, \hat{\boldsymbol{\mu}}^{(i)})$  is  $\varepsilon$ -far from the observable noisy distribution.
  4. Output an arbitrary  $\hat{\pi}^{(i)}, \hat{\boldsymbol{\mu}}^{(i)}$  which was not pruned.

## References

- Lucia Batman, Russell Impagliazzo, Cody Murray, and Ramamohan Paturi. Finding heavy hitters from lossy or noisy data. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 347–362. Springer, 2013.
- Anindya De, Michael Saks, and Sijian Tang. Noisy population recovery in polynomial time. In *Foundations of Computer Science (FOCS), 2016 IEEE 57th Annual Symposium on*, pages 675–684. IEEE, 2016.
- Zeev Dvir, Anup Rao, Avi Wigderson, and Amir Yehudayoff. Restriction access. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 19–33. ACM, 2012.
- Jon Feldman, Ryan O’Donnell, and Rocco A Servedio. Learning mixtures of product distributions over discrete domains. *SIAM Journal on Computing*, 37(5):1536–1564, 2008.
- Michael Kearns, Yishay Mansour, Dana Ron, Ronitt Rubinfeld, Robert E Schapire, and Linda Sellie. On the learnability of discrete distributions. In *Proceedings of the twenty-sixth annual ACM symposium on Theory of computing*, pages 273–282. ACM, 1994.
- Shachar Lovett and Jiapeng Zhang. Improved noisy population recovery, and reverse bonami-beckner inequality for sparse functions. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing*, pages 137–142. ACM, 2015.
- Ankur Moitra and Michael Saks. A polynomial time algorithm for lossy population recovery. In *Foundations of Computer Science (FOCS), 2013 IEEE 54th Annual Symposium on*, pages 110–116. IEEE, 2013.
- Avi Wigderson and Amir Yehudayoff. Population recovery and partial identification. *Machine Learning*, 102(1):29–56, 2016.

## Appendix A. Technical claims

**Claim 13** Let  $\boldsymbol{\mu}, \boldsymbol{\mu}' \in [\mu, 1]^n$  be two vectors with  $\|\boldsymbol{\mu} - \boldsymbol{\mu}'\|_\infty \leq \epsilon$ . Then for any  $f : \{0, 1\}^n \rightarrow \mathbb{R}$  and  $S \subseteq [n]$ ,

$$\|T_{\boldsymbol{\mu}_S}^{-1} \circ f - T_{\boldsymbol{\mu}'_S}^{-1} \circ f\|_1 \leq (1/\mu)^{2|S|} \cdot \epsilon \cdot 2^{|S|+1} \cdot |S| \cdot \|f\|_1.$$

**Proof** To prove the claim, we first need the following claim.

**Claim 14** Let  $\boldsymbol{\mu}, \boldsymbol{\mu}' \in [0, \mu]^n$  be two vectors with  $\|\boldsymbol{\mu} - \boldsymbol{\mu}'\|_\infty \leq \epsilon$ . Then for any  $f : \{0, 1\}^n \rightarrow \mathbb{R}$  and  $S \subseteq [n]$ ,

$$\|(T_{\boldsymbol{\mu}_S} - T_{\boldsymbol{\mu}'_S}) \circ f\|_1 \leq \epsilon \cdot 2^{|S|+1} \cdot |S| \cdot \|f\|_1.$$

**Proof** By the definitions, we have that

$$\begin{aligned}
\|(T_{\mu_S} - T_{\mu'_S}) \circ f\|_1 &= \sum_x |(T_{\mu_S} \circ f)(x) - (T_{\mu'_S} \circ f)(x)| \\
&= \sum_x |\mathbb{E}_{e \sim \mu_S}[f(x+e)] - \mathbb{E}_{e \sim \mu'_S}[f(x+e)]| \\
&\leq \sum_x \max_{e \in \{0,1\}^S} \{|\Pr_{e' \sim \mu_S}[e' = e] - \Pr_{e' \sim \mu'_S}[e' = e]|\} \cdot \left( \sum_{e \in \{0,1\}^S} |f(x+e)| \right) \\
&= 2^{|S|} \|f\|_1 \cdot \max_{e \in \{0,1\}^S} \{|\Pr_{e' \sim \mu_S}[e' = e] - \Pr_{e' \sim \mu'_S}[e' = e]|\} \\
&\leq \epsilon \cdot 2^{|S|+1} \cdot |S| \cdot \|f\|_1.
\end{aligned}$$

■

By Claim 14 we have

$$\begin{aligned}
\|T_{\mu_S}^{-1} \circ f - T_{\mu'_S}^{-1} \circ f\|_1 &= \|T_{\mu_S}^{-1} \circ (f - T_{\mu_S} T_{\mu'_S}^{-1} \circ f)\|_1 \\
&\leq \|T_{\mu_S}^{-1}\|_{1 \rightarrow 1} \cdot \|f - T_{\mu_S} T_{\mu'_S}^{-1} \circ f\|_1 \\
&\leq (1/\mu)^{|S|} \cdot \|f - (T_{\mu'_S} - T_{\mu'_S} + T_{\mu_S}) T_{\mu'_S}^{-1} \circ f\|_1 \\
&= (1/\mu)^{|S|} \cdot \|(T_{\mu'_S} - T_{\mu_S})(T_{\mu'_S}^{-1} \circ f)\|_1 \\
&\leq (1/\mu)^{|S|} \cdot \|(T_{\mu'_S} - T_{\mu_S})\|_{1 \rightarrow 1} \cdot \|T_{\mu'_S}^{-1}\|_{1 \rightarrow 1} \cdot \|f\|_1 \\
&\leq (1/\mu)^{2|S|} \cdot \epsilon \cdot 2^{|S|+1} \cdot |S| \cdot \|f\|_1
\end{aligned}$$

This completes the proof.

■