

Collaborative Learning of Distributions under Heterogeneity and Communication Constraints

Xinmeng Huang,^{*†} Donghwan Lee,^{*‡} Edgar Dobriban,[§] and Hamed Hassani[¶]

June 8, 2022

Abstract

In modern machine learning, users often have to collaborate to learn the distribution of the data. Communication can be a significant bottleneck. Prior work has studied homogeneous users—i.e., whose data follow the same discrete distribution—and has provided optimal communication-efficient methods for estimating that distribution. However, these methods rely heavily on homogeneity, and are less applicable in the common case when users’ discrete distributions are heterogeneous. Here we consider a natural and tractable model of heterogeneity, where users’ discrete distributions only vary sparsely, on a small number of entries. We propose a novel two-stage method named SHIFT: First, the users collaborate by communicating with the server to learn a central distribution; relying on methods from robust statistics. Then, the learned central distribution is fine-tuned to estimate their respective individual distribution. We show that SHIFT is minimax optimal in our model of heterogeneity and under communication constraints. Further, we provide experimental results using both synthetic data and n -gram frequency estimation in the text domain, which corroborate its efficiency.

1 Introduction

Research on learning from data distributed over multiple computational units (machines, users, devices) has grown in recent years, as data is commonly generated by multiple users, such as smart devices and wireless sensors. Communication costs and bandwidth are often bottlenecks on the performance of learning algorithms in these settings [22, 8, 19]. The bottlenecks become even more severe in federated learning and analytics [31], where many users coordinate with a server to learn a central model, while communication is expensive and operates at low rates.

This paper considers learning high-dimensional discrete distributions from user data in the distributed setting. In this setting, several communication-efficient methods have been proposed, and their optimality under communication constraints has been established under various models [26, 9, 27, 2, 3, 16, 17]. However, the key challenge of *heterogeneity*, i.e., that users’ distributions can differ, is rarely considered. Heterogeneity is common, as users inevitably have unique characteristics [15]. Meanwhile, heterogeneity can cause a significant performance drop for learning algorithms designed only for i.i.d data [38, 35, 20]. To use all the data, one needs to learn some central structure, transferable to all individual users. Then one may locally learn—or, finetune—some unique components for each user [21, 50, 18, 54].

To study this paradigm, we first need to introduce a suitable model of heterogeneity. We consider, as an example, the heterogeneous frequencies of words across different texts, e.g., news articles, books, plays (tragedies and comedies), viewed as users. Most words appear with nearly the same probabilities in different texts, however, a few can have very different probabilities, such as “sorrowful” being common in

^{*}Equal Contribution.

[†]Graduate Group in Applied Mathematics and Computational Science, Univ. of Pennsylvania. xinmengh@sas.upenn.edu.

[‡]Graduate Group in Applied Mathematics and Computational Science, Univ. of Pennsylvania. dh7401@sas.upenn.edu.

[§]Department of Statistics and Data Science, Univ. of Pennsylvania. dobriban@wharton.upenn.edu.

[¶]Department of Electrical and Systems Engineering, Univ. of Pennsylvania. hassani@seas.upenn.edu.

Table 1: Estimation error $\mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_2^2]$ of various methods when n is sufficiently large: \mathbf{p}^t is the test distribution, $\hat{\mathbf{p}}^t$ is the estimator, $\boldsymbol{\delta}^t = T^{-1} \sum_{t' \in [T]} \mathbf{p}^{t'} - \mathbf{p}^t$ is a non-vanishing measure of heterogeneity. See Section 1.1 for other notations. Constants and logarithmic factors are omitted for clarity. The “data usage” column indicates whether the estimate is obtained for each cluster separately or by pooling data.

Method	Estimation Error	Data Usage	Bound Type
Unif. Group./Hash. [26]	$O\left(\frac{d}{2^b n}\right)$	Separate	Upper
Unif. Group./Hash. [26]	$O\left(\ \boldsymbol{\delta}^t\ _2^2 + \frac{d}{2^b T n}\right)$	Pool	Upper
Localize-then-Refine* [17]	$O\left(\frac{\ \mathbf{p}^t\ _{1/2}}{2^b n}\right)$	Separate	Upper
Localize-then-Refine* [17]	$O\left(\ \boldsymbol{\delta}^t\ _2^2 + \frac{\frac{1}{T} \sum_{t' \in [T]} \ \mathbf{p}^{t'}\ _{1/2}}{2^b T n}\right)$	Pool	Upper
SHIFT (Theorem 3.1)	$\tilde{O}\left(\frac{\max\{2^b, s\}}{2^b n} + \frac{d}{2^b T n}\right)$	—	Upper
SHIFT (Theorem 4.1)	$\Omega\left(\frac{\max\{2^b, s\}}{2^b n} + \frac{d}{2^b T n}\right)$	—	Lower

* This method [17] requires interactive communication protocols, while other methods are non-interactive.

tragedies and “convivial” being common in comedies. Motivated by this, we formulate a model of sparse heterogeneity. Specifically, suppose that the discrete distributions of all users differ from an underlying central distribution in at most $s \geq 0$ entries, where s is much smaller than the dimension $d \geq 0$. Sparse heterogeneity is relevant to applications such as recommendation systems [28, 42, 36, 10] and medical risk scoring [47, 43, 41].

However, given data generated by multiple distributions with sparse heterogeneity, previous works [26, 2, 3, 17] either do not use all the data, or suffer from bias due to heterogeneity that does not vanish as the sample size increases. Here we propose a novel sparse heterogeneity-inspired collaboration and fine-tuning method (SHIFT) where we first collaboratively learn the central distribution, and then fine-tune the central estimate to individual distributions. Our method makes full use of heterogeneous data, leading to a significant improvement in error rates compared to prior methods. See Table 1 for an overview, explained in detail later.

1.1 Contributions

We consider the problem of learning d -dimensional distributions with s -sparse heterogeneity. We assume there are T clusters of user datapoints, and allow each datapoint to be transmitted in a message with at most b bits of information to the server. Our setting embraces heterogeneous data and thus is a significant generalization of the models from [26, 9, 27, 2]. Our technical contributions are as follows:

- We propose the SHIFT method to learn heterogeneous distributions with collaboration and tuning, in a sample-efficient manner. Our method can, in principle, be used with an arbitrary robust estimate of the probability of each entry/coordinate. When entry-wise median and trimmed mean are used, we provide upper bounds on the estimation error of individual distributions in the ℓ_2 and ℓ_1 norms. We show a factor of (the order of) $\min\{T, d/\max\{s, 2^b\}\}$ improvement in sample complexity compared to previous works; showing the benefit of collaboration (large T) and sparsity (small s), despite communication constraints and heterogeneity.
- To justify the optimality of our method, we prove minimax lower bounds on the estimation errors of individual distributions in the ℓ_2 and ℓ_1 norms, holding for all, possibly interactive, learning methods. These lower bounds, combined with our upper bounds, imply that our median-based method is minimax optimal.
- We support our method with experiments on both synthetic and empirical datasets, showing a significant improvement over previous methods.

1.2 Related Works

Learning with Heterogeneity. Learning with heterogeneity is commonly found in the broader context of multi-task learning [14, 49, 11] and federated learning [5, 44, 24], where a central model or representation is learned from multiple heterogeneous datasets. These central representations can be useful for few-shot learning, i.e., for new problems with a small sample size [48, 23] due to their ability to adapt to new tasks efficiently. In heterogeneous linear regression, [50, 21] show improved sample complexities by assuming a low dimensional central representation, compared to the i.i.d. setting [24, 37]. Related results are proved in [18] for personalized federated learning. [54] study a bandit problem where the unknown parameter in each dataset equals a global parameter plus a sparse instance-specific term. We study a different setting, learning distributions with sparse heterogeneity under communication constraints.

Estimating Distributions under Communication Constraints. Estimating discrete distributions has a rich literature [32, 6]. Under communication constraints, [26, 9, 27, 2] consider the non-interactive scenario, and establish the minimax optimal rates, in terms of data dimension and communication budget, via potentially shared randomness, when all users' data is homogeneous. The optimality for the general interactive (or, blackboard) methods is developed by [1]. A few works study the estimation of sparse distributions. In particular, [3] consider s -sparse distributions and establish minimax optimal rates under communication and privacy constraints, which are further improved by localization strategies in [16]. Complementary to minimax rates, [17] provides pointwise rates, governed by the half-norm of the distribution to be learned, instead of its dimension. Our setting embraces heterogeneous data, and thus is a generalization of the one studied in above works.

Robust Estimation & Learning. Robust statistics and learning study algorithms resilient to unknown data corruption [30, 25, 7, 51, 29]. The median-of-means method [33, 39, 40] partitions the data into subsets, computes an estimate (such as the mean) from each, and takes their median. Similarly, some works study robustness from the optimization perspective, proposing to robustly aggregate gradients of the loss functions [45, 46, 12, 55]. We adapt some analysis techniques from [40, 55] to our significantly different setting of estimation with heterogeneity and communication constraints.

1.3 Notations

Throughout the paper, for an integer $d \geq 1$, we write $[d]$ for both $\{1, \dots, d\}$ and $\{e_1, \dots, e_d\} \subseteq \mathbb{R}^d$, where e_k is the k -th canonical basis vector of \mathbb{R}^d . For a vector $\mathbf{v} \in \mathbb{R}^d$, we refer to the entries of \mathbf{v} by both $[\mathbf{v}]_1, \dots, [\mathbf{v}]_d$ and v_1, \dots, v_d . We denote $\|\mathbf{v}\|_p = (\sum_{k \in [d]} |v_k|^p)^{\frac{1}{p}}$ for all $p > 0$ with $\|\mathbf{v}\|_0$ defined additionally as the number of non-zero entries. We let $\mathcal{P}_d := \{\mathbf{p} = (p_1, \dots, p_d) \in [0, 1]^d : p_1 + \dots + p_d = 1\}$ be the simplex of all d -dimensional discrete probability distributions. For $\mathbf{p} \in \mathcal{P}_d$, we denote by $\mathbb{B}_s(\mathbf{p})$ the s -distinct neighborhood $\{\mathbf{p}' \in \mathcal{P}_d : \|\mathbf{p}' - \mathbf{p}\|_0 \leq s\}$. For a random variable X , we denote n i.i.d. copies of X by $X^{[n]}$. Given any index set \mathcal{I} , we write $|\mathcal{I}|$ for its cardinality and denote by $[\mathbf{v}]_{\mathcal{I}}$ the sub-vector $([\mathbf{v}]_k)_{k \in \mathcal{I}}$ indexed by \mathcal{I} . We use the Bachmann-Landau asymptotic notations $\Omega(\cdot)$, $\Theta(\cdot)$, $O(\cdot)$ to hide constant factors, and use $\tilde{\Omega}(\cdot)$, $\tilde{O}(\cdot)$ to also hide logarithmic factors. We denote the categorical distribution with class probability vector $\mathbf{p} \in \mathcal{P}_d$ by $\text{Cat}(\mathbf{p})$. In our algorithm to be introduced shortly, we use $\tilde{\cdot}$ and $\hat{\cdot}$ to indicate the intermediate estimate and the final estimate, respectively.

2 Problem Setup

We consider the problem of collaboratively learning distributions defined according to the following model of heterogeneity (see Figure 1 for an illustration). There are $T \geq 1$ clusters $\{\mathcal{C}^t \triangleq (X^{t,j})_{j \in [n]} : t \in [T]\}$ of user datapoints, each of which contains n i.i.d. local datapoints. Each datapoint $X^{t,j}$ is in a one-hot format, i.e., $X^{t,j} \in \{e_1, e_2, \dots, e_d\}$, and follows the categorical distribution $\text{Cat}(\mathbf{p}^t)$ where $\mathbf{p}^t \in \mathcal{P}_d$ is unknown. Thus, user datapoints in the same cluster \mathcal{C}^t have an identical distribution \mathbf{p}^t , while the distribution \mathbf{p}^t can vary, i.e., be heterogeneous, across clusters $t \in [T]$. The datapoint $X^{t,j}$ is encoded by

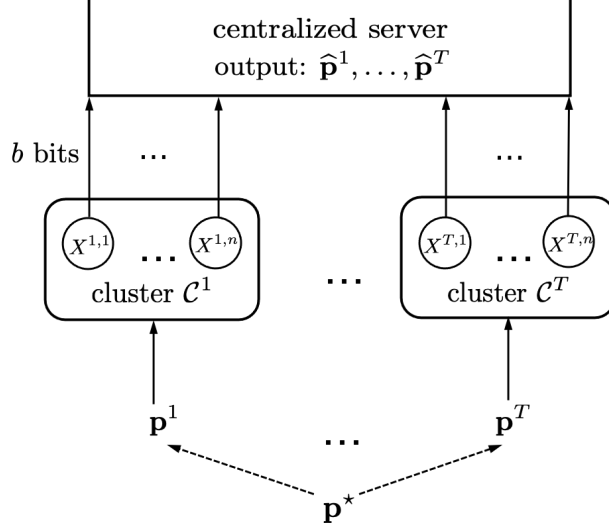


Figure 1: Learning distributions with heterogeneity and communication constraints.

its user into a message $Y^{t,j}$, and then transmitted to a central server. We assume that the message sent by each datapoint is encoded into no more than b bits. We also assume b is significantly smaller than $\log_2 d$ so that the communication is efficient.

The goal here is to collaboratively learn the distributions \mathbf{p}^t from the collection of messages $\{Y^{t',[n]} \triangleq (Y^{t',j})_{j \in [n]} : t' \in [T]\}$ despite heterogeneity. More precisely, we aim to design per-cluster estimators $\hat{\mathbf{p}}^t$, $t \in [T]$, with $\hat{\mathbf{p}}^t : \{Y^{t',[n]} : t' \in [T]\} \rightarrow \mathcal{P}_d$ to minimize the ℓ_2 errors

$$\mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_2^2], \quad \text{for all } t \in [T].$$

We also study the widely-used ℓ_1 error metric (in addition to the ℓ_2 metric). When $T = 1$, i.e., all user datapoints are homogeneous and there is a single distribution to learn, the problem reduces to the one studied by [26, 9, 27, 2, 1].

Model of Heterogeneity. In heterogeneous settings, collaboration among the users is most beneficial if the local distributions are related. We model this by assuming that the local distributions are sparse perturbations of an unknown *central distribution* $\mathbf{p}^* \in \mathcal{P}_d$. The distribution \mathbf{p}^t of each cluster t differs from \mathbf{p}^* in at most $s \geq 0$ entries:

$$\|\mathbf{p}^t - \mathbf{p}^*\|_0 \leq s, \quad \forall t \in [T]. \quad (1)$$

The central distribution \mathbf{p}^* can be viewed as the central structure across heterogeneous clusters of datapoints. The level of heterogeneity is controlled by the parameter s . When s is much smaller than d , the local distributions differ from the center in a small number of entries.

While motivated by word frequencies of different texts, our model of sparse heterogeneity is also relevant for recommendation systems, where the high-dimensional item-preference vectors of users can vary sparsely [28, 42, 36, 10]; and for medical risk scoring, where hospitals can have similar characteristics, with a few systematic differences in diagnosis behavior, healthcare utilization, etc. [47, 43, 41].

3 Algorithm

We now introduce our method for leveraging heterogeneous data to improve per-cluster estimation accuracy. We first discuss a hashing-based method to handle the communication constraint. Since the communication between each user datapoint and the server is restricted to at most $b > 0$ bits (where we may have $2^b \ll d$), the datapoint $X^{t,j}$ needs to be encoded by an encoding function $W^{t,j} : \mathcal{X} \triangleq [d] \rightarrow \mathcal{Y}$.

Algorithm 1 SHIFT: Sparse Heterogeneity Inspired collaboration and Fine-Tuning

input: individual hashed estimators $\check{\mathbf{b}}^1, \dots, \check{\mathbf{b}}^T$, threshold parameter α
▷ Stage I: Collaborative Learning
Estimate \mathbf{b}^* via robust statistical methods: $\check{\mathbf{b}}^* \leftarrow \text{robust_estimate}(\{\check{\mathbf{b}}^t : t \in [T]\})$
▷ Stage II: Fine-Tuning
for $k = 1, \dots, d$ **do**
 for $t = 1, \dots, T$ **do**
 $[\hat{\mathbf{b}}^t]_k \leftarrow [\check{\mathbf{b}}^*]_k$ **if** $|[\check{\mathbf{b}}^*]_k - [\check{\mathbf{b}}^t]_k| \leq \sqrt{\frac{\alpha[\check{\mathbf{b}}^t]_k}{n}}$ **else** $[\check{\mathbf{b}}^t]_k$
 $[\hat{\mathbf{p}}^t]_k \leftarrow \text{Proj}_{[0,1]}(\frac{2^b[\hat{\mathbf{b}}^t]_k - 1}{2^b - 1})$
 end for
end for
output: estimates $\hat{\mathbf{p}}^1, \dots, \hat{\mathbf{p}}^T$

The b -bit constraint enforces $|\mathcal{Y}| \leq 2^b$. Then the encoded message $Y^{t,j} := W^{t,j}(X^{t,j})$ is sent to the server, where it is decoded and used.

There are relatively sophisticated communication protocols under which the design of encoding functions can be *interactive* [17, 16, 1], i.e., depend on previously sent messages. Here we adopt a *non-interactive* encoding-decoding scheme, based on uniform hashing [4, 16] where $W^{t,j}$ depends only on $X^{t,j}$ and is independent of other messages. Specifically, each datapoint $X^{t,j}$ is encoded via an independent random hash function $h^{t,j} : [d] \rightarrow [2^b]$. Upon receiving all messages, the server counts the empirical frequencies of all symbols, leading to hashed estimates $\check{\mathbf{b}}^t$. The communication scheme based on uniform hashing is summarized below.

(Encoding) : Send the message $Y^{t,j} = h^{t,j}(X^{t,j})$ encoded by a hash function $h^{t,j} : [d] \rightarrow [2^b]$;

(Decoding) : Count $N_k^t(Y^{t,[n]}) = |\{j \in [n] : h^{t,j}(k) = Y^{t,j}\}|$ and return $[\check{\mathbf{b}}^t]_k = N_k^t/n$.

One can readily verify that $\mathbb{E}[\check{\mathbf{b}}^t] = [(2^b - 1)\mathbf{p}^t + 1]/2^b \triangleq \mathbf{b}^t$; and thus the hashed estimate $\check{\mathbf{b}}^t$ is biased for \mathbf{p}^t . We also write $\mathbf{b}^* = [(2^b - 1)\mathbf{p}^* + 1]/2^b$ for the mean of a hashed datapoint sampled from the central distribution. More details on the hashed estimator $\check{\mathbf{b}}^t$ are given in Appendix 8.

3.1 The SHIFT Method

We now introduce the SHIFT method, which consists of two stages: *collaborative learning* and *fine-tuning*. The first stage estimates the central hashed distribution \mathbf{b}^* using all hashed estimates $\{\check{\mathbf{b}}^t : t \in [T]\}$. This is achieved via methods from robust statistics such as the median or trimmed mean. The key insight here is that, since the heterogeneity is sparse, for each entry where the individual distributions mostly match with the central one, most datapoints (used to estimate that entry) are sampled from the probability of the central distribution. Hence, to estimate those entries of the central distribution, we can treat the datapoints generated by heterogeneous users as outliers, and leverage robust statistical methods to mitigate their influence.

In the second stage—fine-tuning—we detect mismatched entries between individual hashed estimates $\check{\mathbf{b}}^t$ and the central estimate $\check{\mathbf{b}}^*$. Recall that the central and individual distributions differ in only a few entries. For entries $k \in [d]$ such that $|[\check{\mathbf{b}}^*]_k - [\check{\mathbf{b}}^t]_k|$ is below $(\alpha[\check{\mathbf{b}}^*]_k/n)^{1/2}$ for some threshold parameter α , we may expect that $p_k^t = p_k^*$. As a result, we expect the estimate $[\check{\mathbf{b}}^*]_k$ of b_k^* to be more accurate, as it is learned collaboratively using a larger sample size. Thus, we assign $[\check{\mathbf{b}}^*]_k$ as the final estimate $[\hat{\mathbf{b}}^t]_k$ of b_k^* .

On the other hand, for the entries where the central and individual distributions differ, i.e., $p_k^t \neq p_k^*$, the threshold is more likely to be exceeded. In this case, we keep the individual estimate $[\check{\mathbf{b}}^t]_k$ as $[\hat{\mathbf{b}}^t]_k$. Finally, since the hashed distributions \mathbf{b}^t are biased, we debias them in the final estimates of \mathbf{p}^t where

$\text{Proj}_{[0,1]}(\cdot)$ in Algorithm (1) indicates truncating the input to the $[0, 1]$ interval. Our method does not require sample splitting, despite using two stages, leading to increased sample-efficiency.

Knowledge Transfer to New Clusters. The collaboratively learned central distribution from Algorithm 1 is adaptable to new clusters, which possibly only have a few datapoints. This can be particularly beneficial for sample efficiency, because most entries of the target distribution are well-estimated through collaborative learning. One can transfer those entries, and it suffices to estimate the few remaining entries, instead of the whole distribution. See Theorem 3.2 for the details. The knowledge transfer utility further motivates the importance of collaborative learning.

3.2 Median-Based SHIFT

In this section, we provide upper bounds on the error for the median-based SHIFT method, where $\text{robust_estimate}(\{\tilde{\mathbf{b}}^t : t \in [T]\})$ in Algorithm 1 is the entry-wise median. Specifically, we let

$$[\tilde{\mathbf{b}}^\star]_k = \text{median}(\{[\tilde{\mathbf{b}}^t]_k : t \in [T]\}), \quad \text{for each } k \in [d].$$

When there is no ambiguity, we write $\tilde{\mathbf{b}}^\star = \text{median}(\{\tilde{\mathbf{b}}^t\}_{t \in [T]})$. We also provide results for the trimmed-mean-based SHIFT method, see Appendix 11.

By setting the threshold parameter α in Algorithm 1 as $\alpha = \Theta(\ln(n))$, we prove the following upper bounds on the final individual ℓ_2 estimation errors. The results for the ℓ_1 error are in Appendix 10.

Theorem 3.1. *Suppose $n = \Omega(2^b \ln(n))$ and $\alpha = \Theta(\ln(n))^1$. Then, for the median-based SHIFT method, for any $t \in [T]$,*

$$\mathbb{E} [\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_2^2] = \tilde{O} \left(\frac{\max\{2^b, s\}}{2^b n} + \frac{d}{2^b T n} + \frac{d}{n^2} \right).$$

When $n = \Omega(2^b \max\{T, \ln(n)\})$, the rate further becomes

$$\mathbb{E} [\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_2^2] = \tilde{O} \left(\frac{\max\{2^b, s\}}{2^b n} + \frac{d}{2^b T n} \right). \quad (2)$$

The upper bound in (2) consists of two terms. The first term— $\max\{2^b, s\}/(2^b n)$ —is independent of the dimension d , and is smaller than the rate $d/(2^b n)$ obtained by the minimax optimal method using only homogeneous datapoints [26] by a factor $d/\max\{2^b, s\}$. Thus, it brings a significant benefit under sparse heterogeneity and reasonable communication restrictions, i.e., when $\max\{2^b, s\} \ll d$. Meanwhile, the second, dimension-dependent, term $d/(2^b T n)$ is T times smaller than $d/(2^b n)$, since it depends on the *total sample-size* Tn used collaboratively, despite heterogeneity. Therefore, our method shows a factor of $\min\{T, d/\max\{2^b, s\}\}$ improvement in sample efficiency, compared to previous work designed for homogeneous datapoints.

For completeness, we also consider a heuristic application of estimators from prior works [26, 17], in which datapoints from all clusters are pooled to learn a global distribution $T^{-1} \sum_{t \in [T]} \mathbf{p}^t$, which is then used by each cluster. While this uses all datapoints, it inevitably introduces a non-vanishing bias $\delta^t = \mathbf{p}^t - T^{-1} \sum_{t' \in [T]} \mathbf{p}^{t'}$ in estimating individual distributions, and can behave poorly when the bias is large. See Table 1 for more details.

Finally, we discuss our results on knowledge transfer. The central estimator $\tilde{\mathbf{b}}^\star$ is adaptable to a new cluster \mathcal{C}^{T+1} of size \tilde{n} in the following way. We adjust the fine-tuning procedure in Algorithm 1 to $[\hat{\mathbf{b}}^{T+1}]_k \leftarrow [\tilde{\mathbf{b}}^\star]_k$ if $|[\tilde{\mathbf{b}}^\star]_k - [\tilde{\mathbf{b}}^{T+1}]_k| \leq \sqrt{\alpha [\tilde{\mathbf{b}}^{T+1}]_k / \tilde{n}}$, and $[\hat{\mathbf{b}}^{T+1}]_k \leftarrow [\tilde{\mathbf{b}}^{T+1}]_k$ otherwise. We then show the following result.

Theorem 3.2. *Let $\tilde{\mathbf{b}}^{T+1}$ be the hashed estimate of any new cluster \mathcal{C}^{T+1} with \tilde{n} datapoints such that $n \geq \tilde{n} = \Omega(2^b \max\{T, \ln(\tilde{n})\})$. Let the threshold parameter be $\alpha = \Theta(\ln(\tilde{n}))$. Then, the median-based*

¹To be precise, we require $\alpha = O(\ln(n))$ and $\alpha \geq c \ln(n)$ for some absolute constant c . The analogous statement applies in Theorem 3.2.

SHIFT method has error bounded by

$$\mathbb{E} [\|\hat{\mathbf{p}}^{T+1} - \mathbf{p}^{T+1}\|_2^2] = \tilde{O} \left(\frac{\max\{2^b, s\}}{2^b \tilde{n}} + \frac{d}{2^b T n} \right).$$

Similarly, one can see that the adaptation to new clusters with the median-based SHIFT method achieves a factor of $\min\{Tn/\tilde{n}, d/\max\{2^b, s\}\}$ improvement in sample-efficiency compared to estimating the distribution of the new cluster without knowledge transfer.

3.2.1 Highlights of Theoretical Analysis

In this section, we introduce the key ideas behind the proof of Theorems 3.1 and 3.2. Our analysis is novel compared to previous analyses for methods with homogeneous datapoints. The final individual estimation errors relate to the error of estimating the central hashed distribution \mathbf{b}^* . However, we only expect high accuracy at the center for entries with few individual misalignments. To quantify the influence of heterogeneity, for any $0 < \eta \leq 1$, we define the set of η -well-aligned entries as

$$\mathcal{I}_\eta := \{k \in [d] : |\mathcal{B}_k| < \eta T\}, \quad \text{where} \quad \mathcal{B}_k \triangleq \{t \in [T] : b_k^t \neq b_k^*, \text{ i.e., } p_k^t \neq p_k^*\}$$

is the set of clusters whose distribution differs from \mathbf{p}^* in the k -th entry. We aim to estimate the η -well-aligned entries accurately by using robust statistical methods.

Further, we argue that there are few poorly-aligned entries, and they affect the final per-cluster error only mildly. By the pigeonhole principle, the number of entries that are not η -well-aligned is upper bounded by $|\mathcal{I}_\eta^c| \triangleq |[d] \setminus \mathcal{I}_\eta| \leq \frac{sT}{\eta T} = \frac{s}{\eta}$. Therefore, given an estimator $\tilde{\mathbf{b}}^*$ that is accurate for the η -well-aligned entries, the entries of \mathbf{b}^* can be estimated accurately except for at most s/η entries. The following technical lemma bounds the error for each entry $k \in \mathcal{I}_\eta$.

Lemma 3.3. *Suppose $\tilde{\mathbf{b}}^* = \text{median}(\{\tilde{\mathbf{b}}^t\}_{t \in [T]})$. Then for any $0 < \eta \leq 1/5$ and $k \in \mathcal{I}_\eta$, it holds that*

$$\mathbb{E}[(\tilde{\mathbf{b}}^*)_k - (\mathbf{b}^*)_k]^2 = \tilde{O} \left(\frac{|\mathcal{B}_k|^2 b_k^*(1 - b_k^*)}{T^2 n} + \frac{b_k^*(1 - b_k^*)}{T n} + \frac{1}{n^2} \right).$$

Lemma 3.3 provides an upper bound, in terms of the frequency $|\mathcal{B}_k|/T$ of misalignment (which is smaller than η), and a variance term $b_k^*(1 - b_k^*)$. This result cannot be obtained by directly applying the standard Chernoff or Hoeffding bounds to random variables distributed in $[0, 1]$ as in previous works [17] for two reasons: 1) the datapoints are heterogeneous, 2) the variance $b_k^*(1 - b_k^*)$ here can be small, compared to what it can be for general random variables in $[0, 1]$, implying more concentration than what follows from Hoeffding's inequality. To address these issues, we analyze the concentration of the empirical $(1/2 \pm |\mathcal{B}_k|/T)$ -quantiles to mitigate the influence of heterogeneity, and we also use Bernstein's inequality, which is variance-dependent [52], to obtain bounds relying on both the sample size Tn and the variance $b_k^*(1 - b_k^*)$.

Also, the constant $1/5$, upper bounding the heterogeneity, is not essential and is chosen for clarity. It can be replaced with any number less than half; in which case estimating the central probability distribution becomes possible, as the information conveyed by homogeneous datapoints dominates.

Lemma 3.3 reveals that well-aligned entries of the central distribution are accurately estimated. Thus one can use the central estimate for the entries where the central distribution \mathbf{p}^* aligns with the target distribution \mathbf{p}^t . The remaining entries, that are neither well-aligned nor satisfy $p_k^* = p_k^t$, can be estimated by the local estimator. We argue that a properly chosen threshold parameter α filters out the only desired entries to be estimated individually, with high probability, leading to Theorems 3.1 and 3.2.

While estimating \mathbf{p}^* is not our main goal, one can readily obtain from Lemma 3.3 the following bound for estimating \mathbf{p}^* by summing up the errors for all entries $k \in [d] = \mathcal{I}_\eta$ with $\eta = \max_{k \in [d]} |\mathcal{B}_k|/T$. Corollary 3.4 reveals that the central distribution can be accurately estimated if the mismatch of distributions happens uniformly across all entries, i.e., each entry differs in $O(sT/d)$ clusters.

Corollary 3.4. *Let $\hat{\mathbf{p}}^* = \text{Proj}_{[0,1]}(\frac{2^b \tilde{\mathbf{b}}^* - 1}{2^b - 1})$ be obtained by the debiasing operation from Algorithm 1. Suppose $|\mathcal{B}_k| = O(sT/d)$ for any $k \in [d]$, with $\eta = \max_{k \in [d]} |\mathcal{B}_k|/T$. Then the median-based SHIFT*

method enjoys

$$\mathbb{E}[\|\hat{\mathbf{p}}^* - \mathbf{p}^*\|_2^2] = \tilde{O}\left(\frac{s^2}{d2^{bn}} + \frac{d}{2^b T n} + \frac{d}{n^2}\right).$$

4 Lower Bounds

To complement our upper bounds, we now provide minimax lower bounds for estimating distributions under heterogeneity. Since our setting contains T heterogeneous clusters of datapoints, our minimax error metric is slightly different from the one studied in [26, 9, 27, 2]. Using the ℓ_2 error as the loss, the lower bound metric is defined as

$$\inf_{\substack{(W^{t',[n]})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E} \left[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_2^2 \right], \quad (3)$$

where the supremum is taken over all possible central distributions $\mathbf{p}^* \in \mathcal{P}_d$ and individual distributions $\{\mathbf{p}^t : t \in [T]\}$ in $\mathbb{B}_s(\mathbf{p}^*) \triangleq \{\mathbf{p} \in \mathcal{P}_d : \|\mathbf{p} - \mathbf{p}^*\|_0 \leq s\}$, and the infimum is taken over all estimation methods $\hat{\mathbf{p}}^t$ that use all heterogeneous messages $\{Y^{t',j} \triangleq W^{t',j}(X^{t',j}) : j \in [n], t' \in [T]\}$ encoded (possibly interactively) by any encoding functions $\{W^{t',j} : j \in [n], t' \in [T]\}$ with output in $[2^b]$, e.g., the random hashing maps. The measure (3) characterizes the best possible worst-case performance of estimating distributions under our model of heterogeneity.

Since the supremum is taken over all distributions $\mathbf{p}^*, \mathbf{p}^1, \dots, \mathbf{p}^T$ in \mathcal{P}_d such that $\|\mathbf{p}^t - \mathbf{p}^*\|_0 \leq s$ for all $t \in [T]$, we consider two representative cases therein: First, in the *homogeneous* case where $\mathbf{p}^1 = \dots = \mathbf{p}^T = \mathbf{p}^* \in \mathcal{P}_d$, the setting reduces to the single-cluster problem but with nT datapoints, and with the goal of estimating \mathbf{p}^* , leading to the lower bound $\Omega(d/(2^b T n))$. Second, in the *s/2-sparse* case where $\|\mathbf{p}^*\|_0 \leq s/2$ and $\|\mathbf{p}^t\|_0 \leq s/2$ for all $t \in [T]$, we have $\{\mathbf{p}^t : t \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)$. By constructing independent priors for $\{\mathbf{p}^t : t \in [T]\}$ and \mathbf{p}^* , one can show that only datapoints generated by \mathbf{p}^t itself are informative for estimating \mathbf{p}^t . In this case, we show the lower bound $\Omega(\max\{2^b, s\}/(2^b n))$. Combining the two cases, we find the following lower bound. The formal argument is provided in Appendix 12.

Theorem 4.1. *For any—possibly interactive—estimation method, and for any $t \in [T]$, we have*

$$\inf_{\substack{(W^{t',[n]})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_2^2] = \Omega\left(\frac{\max\{2^b, s\}}{2^b n} + \frac{d}{2^b T n}\right). \quad (4)$$

By a similar argument but with an additional $(T+1)$ -st cluster of \tilde{n} users, we obtain a lower bound for adapting to a new cluster.

Theorem 4.2. *For any—possibly interactive—estimation method, and a new cluster \mathcal{C}^{T+1} , we have*

$$\inf_{\substack{(W^{t',[n]})_{t' \in [T]} \\ W^{T+1,[\tilde{n}]}, \hat{\mathbf{p}}^{T+1}}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T+1]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^{T+1} - \mathbf{p}^{T+1}\|_2^2] = \Omega\left(\frac{\max\{2^b, s\}}{2^b \tilde{n}} + \frac{d}{2^b T n}\right). \quad (5)$$

Theorem 4.1 and 4.2, combined with the upper bounds in Section 3, imply that our method is minimax optimal up to logarithmic terms. We provide similar lower bounds for the ℓ_1 error in Appendix 12.

5 Experiments

We test SHIFT on synthetic data as well as on the Shakespeare dataset [13]. As a baseline method, we use the estimator in [26] that is minimax optimal under homogeneity. We apply the baseline method both locally and globally. In the local case, the estimator $\hat{\mathbf{p}}^t$ for each cluster is computed without datapoints from other clusters. In the global case, we pool data from all clusters, and compute estimators $\hat{\mathbf{p}} = \hat{\mathbf{p}}^1 = \dots = \hat{\mathbf{p}}^T$. The performance measure for estimating \mathbf{p}^t , $t \in [T]$ is taken as $T^{-1} \sum_{t=1}^T \|\mathbf{p}^t - \hat{\mathbf{p}}^t\|_2^2$.

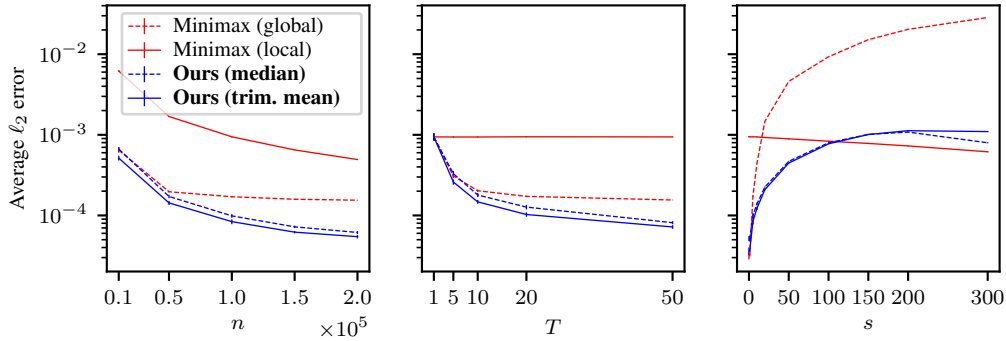


Figure 2: Average ℓ_2 estimation error in synthetic experiment. (Left): Fixing $s = 5$, $T = 30$ and varying n . (Middle): Fixing $s = 5$, $n = 100,000$ and varying T . (Right): Fixing $T = 30$, $n = 100,000$ and varying s . The standard error bars are obtained from 10 independent runs.

5.1 Synthetic Data

We set the uniform distribution, $\mathbf{p}^* = (1/d, \dots, 1/d)$ as the central distribution. In Appendix 13, we also experiment on the truncated geometric distribution, and compare our method with the localization-refinement method [17]. Among the d entries of \mathbf{p}^* , we draw s entries uniformly at random and assign new values for them uniformly at random over $[0, 1]$, with re-normalization to preserve their sum. We repeat this procedure T times to obtain sparsely perturbed distributions $\mathbf{p}^1, \dots, \mathbf{p}^T \in \mathcal{P}_d$. Then, n i.i.d. datapoints $X^{t,1}, \dots, X^{t,n} \sim \text{Cat}(\mathbf{p}^t)$ are generated for each cluster $t \in [T]$. We set the dimension to $d = 300$ and run the simulation by varying n, T, s . As we see from (2), the error of our method depends on s only when $2^b < s$. For this reason, we let $b = 2$ in our experiments.

We run SHIFT with the entry-wise median and entry-wise trimmed mean as the robust estimate. We set the threshold parameter $\alpha = \ln(n)$ and the trimming proportion $\omega = 0.1$. In Appendix 13, we provide results for other choices of the hyperparameters α, ω and discuss a heuristic for choosing α . Figure 2 illustrates that our method outperforms the baseline method for most choices of n, T, s . Specifically, as Theorem 3.1 reveals, the ℓ_2 error of our method decreases as the number of clusters T increases. On the other hand, when the baseline methods are applied globally, without considering heterogeneity, they show a bias that does not disappear as the sample size n or the number of clusters T increases. This shows that the fine-tuning step in SHIFT is crucial for the estimation of heterogeneous distributions. Finally, the right panel of Figure 2 shows that our method is effective only when s is small compared to the dimension d ; which highlights the crucial role of sparsity. When s is close to d , the distributions $\mathbf{p}^1, \dots, \mathbf{p}^T$ could be considerably different without any meaningful central structure, making collaboration less useful than local estimation.

5.2 Shakespeare Dataset

The Shakespeare dataset was proposed as a benchmark for federated learning in [13]. The dataset consists of dialogues of 1,129 speaking roles in Shakespeare’s 35 different plays. In our experiment, we study the distribution of k -grams (k -tuples of consecutive letters from the 26-letter English alphabet, see Chapter 3 of [34]) appearing in the dialogues. We consider each play as a cluster \mathcal{C}^t and estimate the distribution $\mathbf{p}^t \in \mathcal{P}_d$, $d = 26^k$ of k -grams. Since the ground-truth distribution \mathbf{p}^t is unknown, we regard the empirical frequency as \mathbf{p}^t .

To verify the heterogeneity, we run the chi-squared goodness-of-fit test for each pair of distributions from distinct clusters \mathbf{p}^u and \mathbf{p}^v . Resulting p-values are essentially zero within machine precision, which suggests that the distributions of k -grams are strongly heterogeneous. We also perform entry-wise tests comparing $[\mathbf{p}^u]_i$ and $[\mathbf{p}^v]_i$ for all $u \neq v \in [T]$, $i \in [d]$. In total, 25.8% of the tests are rejected at the 5% level of significance. This is again consistent with heterogeneity.

$k = 2$	$b = 2$	$b = 4$	$b = 6$	$b = 8$
Minimax (local)	640 ± 6.0	142 ± 1.2	40 ± 0.40	14 ± 0.13
Minimax (global)	33 ± 1.8	17 ± 0.37	14 ± 0.081	13 ± 0.037
SHIFT (median)	47 ± 2.4	21 ± 0.66	14 ± 0.17	11 ± 0.10
SHIFT (trimmed mean)	36 ± 2.2	19 ± 0.51	13 ± 0.24	10 ± 0.062
$k = 3$	$b = 2$	$b = 4$	$b = 6$	$b = 8$
Minimax (local)	15000 ± 21	3000 ± 5.9	720 ± 2.1	180 ± 0.39
Minimax (global)	4400 ± 5.7	100 ± 1.4	38 ± 0.35	23 ± 0.090
SHIFT (median)	7300 ± 9.6	180 ± 2.1	53 ± 1.0	20 ± 0.18
SHIFT (trimmed mean)	5100 ± 6.3	140 ± 2.3	43 ± 0.66	18 ± 0.18

Table 2: Average ℓ_2 error for estimating distributions of k -grams in the Shakespeare dataset. Numbers are scaled by 10^{-5} .

We draw $n = 20,000$ datapoints with replacement from each cluster and test SHIFT with communication budgets $b \in \{2, 4, 6, 8\}$. We set the fine-tuning threshold $\alpha = \ln(n)$ and the trimming proportion $\omega = 0.1$, which we choose following the heuristic discussed in Appendix 13. We repeat the experiment ten times by randomly drawing different datapoints, and report the average ℓ_2 error of estimation in Table 2. The standard deviations are small even over ten repetitions. SHIFT shows competitive performance on the empirical dataset, even though we do not rigorously know if the sparse heterogeneity model (1) applies.

6 Conclusion and Future Directions

We formulate the problem of learning distributions under sparse heterogeneity and communication constraints. We propose the SHIFT method, which first learns a central distribution, and then fine-tunes the estimate to adapt to individual distributions. We provide both theoretical and experimental results to show its sample-efficiency improvement compared to classical methods that target only homogeneous data. Many interesting directions remain to be explored, including investigating if there is a point-wise optimal method with rate depending on $\{\mathbf{p}^t : t \in [T]\}$ and \mathbf{p}^* ; and designing methods for other information constraints, such as local differential privacy constraints.

Acknowledgements

During this work, Xinmeng Huang was supported in part by the NSF TRIPODS 1934960, NSF DMS 2046874 (CAREER), NSF CAREER award CIF-1943064; Donghwan Lee was supported in part by ARO W911NF-20-1-0080, DCIST, Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) #FA9550-20-1-0111 award.

References

- [1] J. Acharya, C. L. Canonne, and H. Tyagi. General lower bounds for interactive high-dimensional estimation under information constraints. *arXiv: Data Structures and Algorithms*, 2020.
- [2] J. Acharya, C. L. Canonne, and H. Tyagi. Inference under information constraints ii: Communication constraints and shared randomness. *IEEE Transactions on Information Theory*, 66:7856–7877, 2020.
- [3] J. Acharya, P. Kairouz, Y. Liu, and Z. Sun. Estimating sparse discrete distributions under privacy and communication constraints. In V. Feldman, K. Ligett, and S. Sabato, editors, *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*, volume 132 of *Proceedings of Machine Learning Research*, pages 79–98. PMLR, 16–19 Mar 2021.

- [4] J. Acharya, Y. Liu, and Z. Sun. Estimating sparse discrete distributions under local privacy and communication constraints. *Algorithmic Learning Theory*, 2021.
- [5] I. Achituve, A. Shamsian, A. Navon, G. Chechik, and E. Fetaya. Personalized federated learning with gaussian processes. *ArXiv*, abs/2106.15482, 2021.
- [6] A. Agresti. *Categorical data analysis*. John Wiley & Sons, 2003.
- [7] F. J. Anscombe. Rejection of outliers. *Technometrics*, 2:123–146, 1960.
- [8] M. F. Balcan, A. Blum, S. Fine, and Y. Mansour. Distributed learning, communication complexity and privacy. In *Proceedings of the 25th Annual Conference on Learning Theory*, volume 23, pages 26.1–26.22, 2012.
- [9] L. P. Barnes, Y. Han, and A. Ozgur. Lower bounds for learning distributions under communication constraints via fisher information. *arXiv: Information Theory*, 2019.
- [10] H. Bastani. Predicting with proxies: Transfer learning in high dimension. *Management Science*, 67:2964–2984, 2021.
- [11] J. Baxter. A model of inductive bias learning. *Journal of Artificial Intelligence Research*, 12:149–198, 2000.
- [12] P. Blanchard, E. M. E. Mhamdi, R. Guerraoui, and J. Stainer. Byzantine-tolerant machine learning. *ArXiv*, abs/1703.02757, 2017.
- [13] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.
- [14] R. Caruana. Multitask learning. *Machine Learning*, 28:41–75, 1997.
- [15] W.-N. Chen, P. Kairouz, and A. Ozgur. Breaking the communication-privacy-accuracy trilemma. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3312–3324. Curran Associates, Inc., 2020.
- [16] W.-N. Chen, P. Kairouz, and A. Özgür. Breaking the dimension dependence in sparse distribution estimation under communication constraints. *arXiv preprint arXiv:2106.08597*, 2021.
- [17] W.-N. Chen, P. Kairouz, and A. Özgür. Pointwise bounds for distribution estimation under communication constraints. *ArXiv*, abs/2110.03189, 2021.
- [18] L. Collins, H. Hassani, A. Mokhtari, and S. Shakkottai. Exploiting shared representations for personalized federated learning. *ArXiv*, abs/2102.07078, 2021.
- [19] H. Daumé, J. M. Phillips, A. Saha, and S. Venkatasubramanian. Efficient protocols for distributed classification and optimization. *ArXiv*, abs/1204.3523, 2012.
- [20] E. Dobriban and Y. Sheng. Distributed linear regression by averaging. *The Annals of Statistics*, 2018.
- [21] S. S. Du, W. Hu, S. M. Kakade, J. Lee, and Q. Lei. Few-shot learning via learning the representation, provably. *ArXiv*, abs/2002.09434, 2021.
- [22] A. Garg, T. Ma, and H. L. Nguyen. On communication cost of distributed statistical estimation and dimensionality. In *Advances in Neural Information Systems*, 2014.
- [23] P. Goyal, D. K. Mahajan, A. K. Gupta, and I. Misra. Scaling and benchmarking self-supervised visual representation learning. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6390–6399, 2019.

- [24] F. Grimberg, M.-A. Hartley, S. P. Karimireddy, and M. Jaggi. Optimal model averaging: Towards personalized collaborative learning. *ArXiv*, abs/2110.12946, 2021.
- [25] F. R. Hampel, E. M. Ronchetti, P. J. Rousseeuw, and W. A. Stahel. *Robust statistics: the approach based on influence functions*, volume 196. John Wiley & Sons, 2011.
- [26] Y. Han, P. Mukherjee, A. Özgür, and T. Weissman. Distributed statistical estimation of high-dimensional and nonparametric distributions. *2018 IEEE International Symposium on Information Theory (ISIT)*, pages 506–510, 2018.
- [27] Y. Han, A. Özgür, and T. Weissman. Geometric lower bounds for distributed parameter estimation under communication constraints. *IEEE Transactions on Information Theory*, 67:8248–8263, 2021.
- [28] L. Hu, S. Jian, L. Cao, Z. Gu, Q. Chen, and A. Amirbekyan. Hers: Modeling influential contexts with heterogeneous relations for sparse and cold-start recommendation. In *The Thirty-Third AAAI Conference on Artificial Intelligence*, 2019.
- [29] P. J. Huber. Robust estimation of a location parameter. *Annals of Mathematical Statistics*, 35:492–518, 1964.
- [30] P. J. Huber. Robust statistics. *Wiley Series in Probability and Mathematical Statistics*, 1981.
- [31] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. B. Charles, G. Cormode, R. Cummings, R. G. L. D’Oliveira, S. Y. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, O. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. X. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao. Advances and open problems in federated learning. *Found. Trends Mach. Learn.*, 14:1–210, 2021.
- [32] E. L. Lehmann and G. Casella. *Theory of point estimation*. Springer Science & Business Media, 2006.
- [33] M. Lerasle and R. I. Oliveira. Robust empirical mean estimators. *arXiv: Statistics Theory*, 2011.
- [34] J. Leskovec, A. Rajaraman, and J. D. Ullman. *Mining of massive data sets*. Cambridge university press, 2020.
- [35] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang. On the convergence of fedavg on non-iid data. *ArXiv*, abs/1907.02189, 2020.
- [36] T. Liu, Z. Wang, J. Tang, S. Yang, G. Y. Huang, and Z. Liu. Recommender systems with heterogeneous side information. *The World Wide Web Conference*, 2019.
- [37] A. Maurer, M. Pontil, and B. Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17:81:1–81:32, 2016.
- [38] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, 2017.
- [39] S. Minsker. Geometric median and robust estimation in banach spaces. *Bernoulli*, 21:2308–2335, 2015.
- [40] S. Minsker and N. Strawn. Distributed statistical estimation and rates of convergence in normal approximation. *ArXiv*, abs/1704.02658, 2019.
- [41] S. Mullainathan and Z. Obermeyer. Does machine learning automate moral hazard and error? *The American Economic Review*, 107(5):476–480, 2017.

- [42] M. Qian, L. Hong, Y. Shi, and S. Rajan. Structured sparse regression for recommender systems. *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, 2015.
- [43] J. Quionero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence. *Dataset Shift in Machine Learning*. MIT Press, 2009.
- [44] A. Shamsian, A. Navon, E. Fetaya, and G. Chechik. Personalized federated learning using hypernetworks. In *International Conference on Machine Learning*, 2021.
- [45] L. Su and N. H. Vaidya. Fault-tolerant multi-agent optimization: Optimal iterative distributed algorithms. In *Proceedings of the 2016 ACM Symposium on Principles of Distributed Computing*, 2016.
- [46] L. Su and N. H. Vaidya. Non-bayesian learning in the presence of byzantine agents. In *International Symposium on Distributed Computing*, pages 414–427, 2016.
- [47] A. Subbaswamy and S. Saria. From development to deployment: dataset shift, causality, and shift-stable models in health ai. *Biostatistics*, 2019.
- [48] C. Sun, A. Shrivastava, S. Singh, and A. K. Gupta. Revisiting unreasonable effectiveness of data in deep learning era. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 843–852, 2017.
- [49] S. Thrun and L. Y. Pratt. Learning to learn: Introduction and overview. In *Learning to Learn*. Springer, 1998.
- [50] N. Tripuraneni, C. Jin, and M. I. Jordan. Provable meta-learning of linear representations. In *International Conference on Machine Learning*, 2021.
- [51] J. W. Tukey. A survey of sampling from contaminated distributions. *Contributions to probability and statistics*, pages 448–485, 1960.
- [52] J. V. Uspensky. *Introduction to mathematical probability*. McGraw-Hill Book Company, 1937.
- [53] R. Vershynin. *High-Dimensional Probability*. 2018.
- [54] K. Xu and H. Bastani. Learning across bandits in high dimension via robust statistics. *ArXiv*, abs/2112.14233, 2021.
- [55] D. Yin, Y. Chen, K. Ramchandran, and P. L. Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. *ArXiv*, abs/1803.01498, 2018.

7 Appendix

Additional Notations. In the appendix, we use the following additional notations. For an integer $d \geq 1$, and a vector $\mathbf{v} \in \mathbb{R}^d$, the support $\text{supp}(\mathbf{v}) = \{j \in [d] : \mathbf{v}_j \neq 0\}$ denotes the indices of non-zero entries. For an event A on a probability space (Ω, B, P) (which is usually self-understood from the context), we denote by $I(A)$, $\mathbb{1}\{A\}$, or $\mathbb{1}(A)$ its indicator function, such that $I(A)(\omega) = 1$ if $\omega \in A$, and zero otherwise. We denote by Φ the cumulative distribution function of the standard normal random variable. For two scalars $a, b \in \mathbb{R}$, we write $a \wedge b = \min(a, b)$.

Algorithm 2 Encoding-Decoding via Uniform Hashing

input: cluster \mathcal{C}^t with $n \geq 1$ users having data $X^{t,j}$, $j = 1, \dots, n$
for $j = 1, \dots, n$ **do**
 Generate a uniformly random hash function $h^{t,j} : [d] \rightarrow [2^b]$ using shared randomness
 Encode message $Y^{t,j} = h^{t,j}(X^{t,j})$ and send it to the server ▷ Encoding
end for
for $k = 1, \dots, d$ **do**
 Count $N_k^t(Y^{t,[n]}) \leftarrow |\{j \in [n] : h^{t,j}(k) = Y^{t,j}\}|$ ▷ Decoding
 Estimate $\tilde{b}_k^t \leftarrow N_k^t/n$
end for
output: $\hat{\mathbf{b}}^t$

8 Properties of Uniform Hashing

Recall that for all $t \in [T]$ and $k \in [d]$, $b_k^t = \frac{p_k^t(2^b-1)+1}{2^b} \in [\frac{1}{2^b}, 1]$.

Proposition 8.1 (Properties of Hashed Estimates). *For each $t \in [T]$, suppose $\check{\mathbf{b}}^t$ is computed in cluster \mathcal{C}^t as in Algorithm 2 with i.i.d datapoints $X^{t,j} \sim \text{Cat}(\mathbf{p}^t)$, $\forall j \in [n]$. Then, it holds that*

1. $\check{\mathbf{b}}^1, \dots, \check{\mathbf{b}}^T \in [0, 1]$ are independent;
2. for any $t \in [T]$ and $k \in [d]$, $N_k^t \sim \text{Binom}(n, b_k^t)$;
3. $\text{supp}(\mathbf{p}^t - \mathbf{p}^*) = \text{supp}(\mathbf{b}^t - \mathbf{b}^*)$ and $p_k^* = 1$ (or 0) is equivalent to $b_k^* = 1$ (or $\frac{1}{2^b}$, respectively).

Proof. Property 1 holds because $\hat{\mathbf{b}}^1, \dots, \hat{\mathbf{b}}^T$ are obtained by cluster-wise encoding-decoding of independent datapoints. To see property 2, we have for any $j \in [n]$ and $k \in [d]$ that

$$\begin{aligned} \mathbb{P}(h^{t,j}(k) = Y^{t,j}) &= \mathbb{P}(k = X^{t,j}) + \mathbb{P}(k \neq X^{t,j} \text{ and } h^{t,j}(k) = h^{t,j}(X^{t,j})) \\ &= p_k^t + (1 - p_k^t) \cdot \frac{1}{2^b} = b_k^t \in \left[\frac{1}{2^b}, 1\right]. \end{aligned}$$

Thus, $I(h^{t,j}(k) = Y^{t,j})$ is a Bernoulli variable with success probability b_k^t . Since each datapoint is encoded with an independent hash function, N_k^t has a binomial distribution with n trials and parameter b_k^t . Property 3 directly follows from $\mathbf{b}^t - \mathbf{b}^* = (\mathbf{p}^t - \mathbf{p}^*)(2^b - 1)/2^b$ and as $b > 0$. \square

Proposition 8.2 (Property of Debiasing). *For any $\mathbf{y}, \mathbf{y}^* \in \mathbb{R}^d$, let $\mathbf{x} = \text{Proj}_{[0,1]}(\frac{2^b \mathbf{y} - 1}{2^b - 1})$ and $\mathbf{x}^* = \text{Proj}_{[0,1]}(\frac{2^b \mathbf{y}^* - 1}{2^b - 1})$. Then it holds that for $q = 1, 2$, $\mathbb{E}[\|\mathbf{x} - \mathbf{x}^*\|_q^q] = O(\mathbb{E}[\|\mathbf{y} - \mathbf{y}^*\|_q^q])$. In particular, we have for $q = 1, 2$ and any $t \in [T]$, $\mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] = O(\mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_q^q])$, where $\hat{\mathbf{p}}^t = \text{Proj}_{[0,1]}(\frac{2^b \hat{\mathbf{b}}^t - 1}{2^b - 1})$ is the final per-cluster estimate obtained in Algorithm 1.*

Proof. Using the inequality that $|\text{Proj}_{[0,1]}(x) - \text{Proj}_{[0,1]}(y)| \leq |x - y|$ for any $x, y \in \mathbb{R}$, we have

$$\begin{aligned} \mathbb{E}[\|\mathbf{x} - \mathbf{x}^*\|_q^q] &= \sum_{k \in [d]} \mathbb{E} \left[\left| \text{Proj}_{[0,1]} \left(\frac{2^b y_k - 1}{2^b - 1} \right) - \text{Proj}_{[0,1]} \left(\frac{2^b y_k^* - 1}{2^b - 1} \right) \right|^q \right] \\ &\leq \sum_{k \in [d]} \mathbb{E} \left[\left| \frac{2^b (y_k - y_k^*)}{2^b - 1} \right|^q \right] = \left(\frac{2^b}{2^b - 1} \right)^q \mathbb{E}[\|\mathbf{y} - \mathbf{y}^*\|_q^q] = O(\mathbb{E}[\|\mathbf{y} - \mathbf{y}^*\|_q^q]). \end{aligned}$$

In the last step, we used that $2^b/(2^b - 1) \leq 2$ for all $b \geq 1$, and thus the $O(\cdot)$ only depends on universal constants. \square

9 General Lemmas

In this section, we state some general lemmas that will be used in the analysis.

Lemma 9.1 (Berry-Esseen Theorem; [53]). *Assume that Z_1, \dots, Z_n are i.i.d. copies of a random variable Z with mean μ , variance $\sigma^2 > 0$, and such that $\mathbb{E}[|Z - \mu|^3] < \infty$. Then,*

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P} \left\{ \sqrt{n} \frac{\bar{Z} - \mu}{\sigma} \leq x \right\} - \Phi(x) \right| \leq 0.4748 \frac{\gamma(Z)}{\sqrt{n}}.$$

where $\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i$ and $\gamma(Z) = \mathbb{E}[|Z - \mu|^3]/\sigma^3$ is the absolute skewness of Z .

Lemma 9.2 (Hoeffding's Inequality; [53]). *Let $Z_1, \dots, Z_n \in [l, r]$, $l < r$, be independent random variables and let $\bar{Z} = \frac{1}{n} \sum_{j=1}^n Z_j$. Then for any $\delta \geq 0$,*

$$\max\{\mathbb{P}(\bar{Z} - \mathbb{E}[\bar{Z}] > \delta), \mathbb{P}(\bar{Z} - \mathbb{E}[\bar{Z}] < -\delta)\} \leq \exp\left(-\frac{2n\delta^2}{(r-l)^2}\right).$$

Lemma 9.3 (Bernstein's Inequality; [52]). *Let Z_1, \dots, Z_n be i.i.d. copies of a random variable Z with $|Z - \mathbb{E}[Z]| \leq M$, $M > 0$ and $\text{Var}(Z_1) = \sigma^2 > 0$, and let $\bar{Z} = \frac{1}{n} \sum_{j=1}^n Z_j$. Then for any $\delta \geq 0$,*

$$\mathbb{P}(|\bar{Z} - \mathbb{E}[\bar{Z}]| > \delta) \leq 2 \exp\left(-\frac{n\delta^2}{2(\sigma^2 + M\delta)}\right) \leq 2 \exp\left(-\frac{n}{4} \min\left\{\frac{\delta^2}{\sigma^2}, \frac{\delta}{M}\right\}\right). \quad (6)$$

The second inequality above directly follows from $\frac{1}{a+b} \geq \frac{1}{2} \min\{\frac{1}{a}, \frac{1}{b}\}$ for any $a, b > 0$. Note that (6) also allows $\sigma = 0$ because $\mathbb{P}(|\bar{Z} - \mathbb{E}[\bar{Z}]| > \delta) = 0$ and $\min\left\{\delta^2/\sigma^2 \triangleq +\infty, \delta/M\right\} = \frac{\delta}{M}$. Therefore, we use this lemma for all $\sigma \geq 0$ below.

9.1 Analysis Framework

For each $t \in [T]$, we denote by

$$\mathcal{K}_\alpha^t = \{k \in [d] : (\check{b}_k^\star - \check{b}_k^t)^2 \leq \alpha \check{b}_k^t/n\} \quad (7)$$

the set of entries in which the central estimate $[\hat{\mathbf{b}}^\star]_k$ is adapted to cluster \mathcal{C}^t . In this language, the final estimates can be expressed as $\hat{b}_k^t = \check{b}_k^\star \mathbb{1}\{k \in \mathcal{K}_\alpha^t\} + \check{b}_k^t \mathbb{1}\{k \notin \mathcal{K}_\alpha^t\}$ for $t \in [T]$. Therefore, it holds that, for $q = 1, 2$,

$$\mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_q^q] = \sum_{k \in [d]} \mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\} |\check{b}_k^\star - \check{b}_k^t|^q] + \sum_{k \in [d]} \mathbb{E}[\mathbb{1}\{k \notin \mathcal{K}_\alpha^t\} |\check{b}_k^t - b_k^t|^q]. \quad (8)$$

Let $\mathcal{I}^t \triangleq \{k \in [d] : b_k^t = b_k^\star, \text{ i.e., } p_k^t = p_k^\star\}$ be set of entries at which the t -th cluster's distribution \mathbf{p}^t aligns with the central distribution \mathbf{p}^\star . We next bound the two terms from (8) in Lemmas 9.4 and 9.5. These do not need the independence of $\check{\mathbf{b}}^\star$ and $\check{\mathbf{b}}^t$, and hence do not require sample splitting despite the division between stages.

Lemma 9.4. *For any $t \in [T]$, $\alpha \geq 1$ and $\eta \in (0, 1]$, with \mathcal{K}_α^t from (7), we have, for $q = 1, 2$*

$$\sum_{k \in [d]} \mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\} |\check{b}_k^\star - \check{b}_k^t|^q] = O\left(\mathbb{E}[\|\check{b}_{\mathcal{I}_\eta \cap \mathcal{I}^t}^\star - \check{b}_{\mathcal{I}_\eta \cap \mathcal{I}^t}^t\|_q^q] + \sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \left(\frac{\alpha b_k^t}{n}\right)^{q/2}\right).$$

Proof. We first take $q = 1$. For any $k \in [d]$, clearly

$$\mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\} |\check{b}_k^\star - \check{b}_k^t|] \leq \mathbb{E}[|\check{b}_k^\star - b_k^t|]. \quad (9)$$

We use this bound for $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$. For $k \notin \mathcal{I}_\eta \cap \mathcal{I}^t$, we instead bound

$$\mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}|\check{b}_k^* - b_k^t|] \leq \mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}|\check{b}_k^* - \check{b}_k^t|] + \mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}|\check{b}_k^t - b_k^t|].$$

If $k \in \mathcal{K}_\alpha^t$, it holds by definition that $|\check{b}_k^* - \check{b}_k^t| \leq \sqrt{\alpha \check{b}_k^t/n}$, thus we further have

$$\begin{aligned} \mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}|\check{b}_k^* - b_k^t|] &\leq \mathbb{E}\left[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}\sqrt{\alpha \check{b}_k^t/n}\right] + \mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}|\check{b}_k^t - b_k^t|] \\ &\leq \mathbb{E}\left[\sqrt{\alpha \check{b}_k^t/n}\right] + \mathbb{E}[|\check{b}_k^t - b_k^t|]. \end{aligned} \quad (10)$$

By Jensen's inequality and since $n\check{b}_k^t \sim \text{Binom}(n, b_k^t)$, we have

$$\mathbb{E}\left[\sqrt{\check{b}_k^t}\right] \leq \sqrt{\mathbb{E}[\check{b}_k^t]} = \sqrt{b_k^t} \quad (11)$$

and

$$\mathbb{E}[|\check{b}_k^t - b_k^t|] \leq \sqrt{\mathbb{E}[(\check{b}_k^t - b_k^t)^2]} = \sqrt{\frac{b_k^t(1-b_k^t)}{n}} \leq \sqrt{\frac{b_k^t}{n}}. \quad (12)$$

Plugging (11) and (12) into (10), we find

$$\mathbb{E}[\mathbb{1}\{k \in \mathcal{K}_\alpha^t\}|\check{b}_k^* - b_k^t|] \leq (\sqrt{\alpha} + 1)\sqrt{\frac{b_k^t}{n}} = O\left(\sqrt{\frac{\alpha b_k^t}{n}}\right). \quad (13)$$

Summing up (9) over all entries in $\mathcal{I}_\eta \cap \mathcal{I}^t$ and summing up (13) over all entries not in $\mathcal{I}_\eta \cap \mathcal{I}^t$ leads to the claim for $q = 1$ in Lemma 9.4. The case $q = 2$ follows by a similar argument. \square

Lemma 9.5. *For any $t \in [T]$, $\alpha \geq 1$ and $\eta \in (0, 1]$, with \mathcal{K}_α^t from (7), we have, for $q = 1, 2$*

$$\begin{aligned} &\sum_{k \in [d]} \mathbb{E}[\mathbb{1}\{k \notin \mathcal{K}_\alpha^t\}|\check{b}_k^t - b_k^t|^q] \\ &= O\left(\sum_{k \in \mathcal{I}_\eta \cap \mathcal{I}^t} \mathbb{P}(k \notin \mathcal{K}_\alpha^t) \wedge \left(\frac{b_k^t(1-b_k^t)}{n}\right)^{q/2} + \sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \left(\frac{b_k^t}{n}\right)^{q/2}\right). \end{aligned}$$

Proof. For $q = 1$, note that

$$\mathbb{E}[\mathbb{1}\{k \notin \mathcal{K}_\alpha^t\}|\check{b}_k^t - b_k^t|] \leq \mathbb{P}(k \notin \mathcal{K}_\alpha^t) \quad (14)$$

and

$$\mathbb{E}[\mathbb{1}\{k \notin \mathcal{K}_\alpha^t\}|\check{b}_k^t - b_k^t|] \leq \mathbb{E}[|\check{b}_k^t - b_k^t|]. \quad (15)$$

Combining (14), (15) with the first inequality in (12) for $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$, and using the last inequality in (12) for $k \notin \mathcal{I}_\eta \cap \mathcal{I}^t$ leads to the claim with $q = 1$. We can similarly obtain the bound with $q = 2$. \square

Combing Lemma 9.4 and 9.5 with (8), we find the following proposition:

Proposition 9.6. *For any $\alpha \geq 1$, and $q = 1, 2$, it holds that*

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_q^q] &= O\left(\sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \left(\frac{\alpha b_k^t}{n}\right)^{q/2} \right. \\ &\quad \left. + \sum_{k \in \mathcal{I}_\eta \cap \mathcal{I}^t} \mathbb{P}(k \notin \mathcal{K}_\alpha^t) \wedge \left(\frac{b_k^t(1-b_k^t)}{n}\right)^{q/2} + \mathbb{E}[\|\check{\mathbf{b}}_{\mathcal{I}_\eta \cap \mathcal{I}^t}^* - \mathbf{b}_{\mathcal{I}_\eta \cap \mathcal{I}^t}^*\|_q^q]\right). \end{aligned}$$

Proposition 9.6 does not rely on how $\check{\mathbf{b}}^*$ is obtained. The next part is devoted to proving that when $\check{\mathbf{b}}^*$ is obtained via a certain robust estimate, the bounds in Proposition 9.6 are small for certain values of α and η .

10 Median-Based Method

In this section, we provide the proofs for the median-based SHIFT method. We first re-state the detailed version of some key results that apply to both the ℓ_2 and ℓ_1 errors.

Below, we use $\sigma_k = \sqrt{b_k^*(1-b_k^*)}$ to denote the standard deviation of the Bernoulli variable with success probability $b_k^* = p_k^* + (1-p_k^*)/2^b$. We also recall that \mathcal{B}_k is defined as the set of clusters with distributions mismatched with the central distribution at the k -th entry, i.e., $\mathcal{B}_k = \{t \in [T] : p_k^t \neq p_k^*\}$, and \mathcal{I}_η is defined as the η -well-aligned entries, i.e., $\mathcal{I}_\eta = \{k \in [d] : |\mathcal{B}_k| < \eta T\}$.

Lemma 10.1 (Detailed statement of Lemma 3.3). *Suppose $\check{\mathbf{b}}^* = \text{median}(\{\check{\mathbf{b}}^t\}_{t \in [T]})$. Then for any $0 < \eta \leq \frac{1}{5}$, $k \in \mathcal{I}_\eta$, and $q = 1, 2$, it holds that*

$$\mathbb{E}[|\check{b}_k^* - b_k^*|^q] = \tilde{O}\left(\left(\frac{|\mathcal{B}_k|\sigma_k}{T\sqrt{n}}\right)^q + \left(\frac{\sigma_k}{\sqrt{Tn}}\right)^q + \left(\frac{1}{n}\right)^q\right).$$

Let us define, for $q = 1, 2$,

$$E(q) \triangleq E(q; n, d, b, T) := \frac{d}{(2^b T n)^{q/2}} + \frac{d}{n^q}.$$

Proposition 10.2. *Suppose $\check{\mathbf{b}}^* = \text{median}(\{\check{\mathbf{b}}^t\}_{t \in [T]})$. Then for any $0 < \eta \leq \frac{1}{5}$ and $q = 1, 2$, it holds that*

$$\mathbb{E}[\|\check{\mathbf{b}}_{\mathcal{I}_\eta}^* - \mathbf{b}_{\mathcal{I}_\eta}^*\|_q^q] = \tilde{O}\left(\sum_{k \in \mathcal{I}_\eta} \left(\frac{|\mathcal{B}_k|\sigma_k}{T\sqrt{n}}\right)^q + E(q)\right).$$

We omit the proofs of Proposition 10.2 and Theorem 10.4 (below), as Proposition 10.2 is a direct corollary of Lemma 10.1 by using $\sum_{k \in [d]} \sigma_k^q = O(d/2^{bq/2})$ for $q = 1, 2$, and Theorem 10.4 follows from the same analysis as Theorem 10.3.

Theorem 10.3 (Detailed statement of Theorem 3.1). *Suppose $n \geq 2^{b+6} \ln(n)$ and $\alpha \geq 2(8 + \sqrt{8 \ln(n)})^2$ with $\alpha = O(\ln(n))$. Then for the median-based SHIFT method, for any $0 < \eta \leq \frac{1}{5}$, $q = 1, 2$, and $t \in [T]$,*

$$\mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] = \tilde{O}\left(\sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \left(\frac{b_k^t}{n}\right)^{q/2} + \sum_{k \in \mathcal{I}_\eta \cap \mathcal{I}^t} \left(\frac{|\mathcal{B}_k|^2 b_k^*}{T^2 n}\right)^{q/2} + E(q)\right).$$

Furthermore, by setting $\eta = \Theta(1)$ with $\eta \leq \frac{1}{5}$, we have

$$\mathbb{E}[\|\check{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] = \tilde{O}\left(s^{1-q/2} \left(\frac{\max\{2^b, s\}}{2^b n}\right)^{q/2} + E(q)\right).$$

Theorem 10.4 (Detailed statement of Theorem 3.2). *Suppose $n \geq \tilde{n} \geq 2^{b+6} \ln(\tilde{n})$ and $\alpha \geq 2(8 + \sqrt{8 \ln(\tilde{n})})^2$ with $\alpha = O(\ln(\tilde{n}))$. Then the median-based SHIFT method for predicting the distribution of the new cluster with \tilde{n} users achieves, for $q = 1, 2$,*

$$\mathbb{E}[\|\check{\mathbf{p}}^{T+1} - \mathbf{p}^{T+1}\|_q^q] = \tilde{O}\left(s^{1-q/2} \left(\frac{\max\{2^b, s\}}{2^b \tilde{n}}\right)^{q/2} + E(q)\right).$$

10.1 Proof of Lemma 10.1

We first assume $T = O(\ln(n))$. In this case, by Bernstein's inequality (Lemma 9.3) with $M = 1$, we have for any $t \in [T] \setminus \mathcal{B}_k$ that for any $\delta \geq 0$,

$$\mathbb{P}\left(|\check{b}_k^t - b_k^*| > \delta\right) \leq 2e^{-\frac{\eta}{4} \min\{\delta^2/\sigma_k^2, \delta\}}. \quad (16)$$

Taking $\delta = \max\{\sigma_k \sqrt{8 \ln(n)/n}, 8 \ln(n)/n\}$ in (16), we find

$$\mathbb{P} \left(|\check{b}_k^t - b_k^*| > \max \left\{ \sigma_k \sqrt{\frac{8 \ln(n)}{n}}, \frac{8 \ln(n)}{n} \right\} \right) \leq \frac{2}{n^2}. \quad (17)$$

Since $|[T] \setminus \mathcal{B}_k| > \frac{T}{2}$ for any $k \in \mathcal{I}_\eta$ with $\eta \leq \frac{1}{5}$, we have, since $\check{b}_k^* = \text{median}(\{\check{b}_k^t\}_{t \in [T]})$, that there are $t_-, t_+ \in [T] \setminus \mathcal{B}_k$ with $\check{b}_k^{t_-} \leq \check{b}_k^* \leq \check{b}_k^{t_+}$. Hence, $|\check{b}_k^* - b_k^*| \leq \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*|$.

Recall that for any random variable $0 \leq X \leq 1$ and any $\delta \geq 0$, $\mathbb{E}[X] \leq \delta + \mathbb{P}(X \geq \delta)$. Therefore, by taking the union bound of (17) over $k \in [T] \setminus \mathcal{B}_k$, and by the assumption that $T = O(\ln(n))$, we have

$$\begin{aligned} \mathbb{E}[|\check{b}_k^* - b_k^*|] &\leq \mathbb{E}[\max_{k \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*|] \leq \sigma_k \sqrt{\frac{8 \ln(n)}{n}} + \frac{8 \ln(n)}{n} + \frac{2T}{n^2} \\ &= O \left(\sigma_k \sqrt{\frac{\ln(n)}{n}} + \frac{\ln(n)}{n} \right) = O \left(\sigma_k \frac{\ln(n)}{\sqrt{Tn}} + \frac{\ln(n)}{n} \right) = \tilde{O} \left(\frac{\sigma_k}{\sqrt{Tn}} + \frac{1}{n} \right). \end{aligned} \quad (18)$$

Similarly, we have

$$\mathbb{E}[(\check{b}_k^* - b_k^*)^2] \leq \sigma_k^2 \frac{8 \ln(n)}{n} + \frac{64 \ln(n)^2}{n^2} + \frac{2T}{n^2} = \tilde{O} \left(\frac{\sigma_k^2}{Tn} + \frac{1}{n^2} \right). \quad (19)$$

For each $k \in [d]$ with $b_k^* \neq 1$ (recall that $b_k^* \geq 1/2^b$ by definition), let $\gamma_k = (1 - 2b_k^*(1 - b_k^*)) / \sqrt{b_k^*(1 - b_k^*)}$, and let $\tilde{F}_k(x) := \frac{1}{T - |\mathcal{B}_k|} \sum_{t \in [T] \setminus \mathcal{B}_k} \mathbf{1}(\check{b}_k^t \leq x)$ be the empirical distribution function of $\{\check{b}_k^t : b_k^t = b_k^*\}$. Let $\varepsilon \in (0, 1/2)$ and $C_\varepsilon = \sqrt{2\pi} \exp((\Phi^{-1}(1 - \varepsilon))^2/2)$. For $\delta \geq 0$, define, recalling $\eta T > |\mathcal{B}_k|$ for all $k \in \mathcal{I}_\eta$,

$$G_{k,T,\delta} = \frac{|\mathcal{B}_k|}{T} + \frac{10^{-8}}{Tn} + \sqrt{\frac{\delta}{T - |\mathcal{B}_k|}}$$

where the term $\frac{10^{-8}}{Tn}$ is used to overcome some challenges due to the discreteness of empirical distributions, and can be replaced with other suitably small terms (see the proof of Lemma 10.6). Further, define

$$G'_{k,T,\delta} = G_{k,T,\delta} + 0.4748 \frac{\gamma_k}{\sqrt{n}}.$$

To prove Lemma 3.3 for $T = \Omega(\ln(n))$, we need the following additional lemmas:

Lemma 10.5. *For any $\delta \geq 0$ such that*

$$G'_{k,T,\delta} \leq \frac{1}{2} - \varepsilon, \quad (20)$$

it holds with probability at least $1 - 4e^{-2\delta}$ that

$$\tilde{F}_k \left(b_k^* + C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \geq \frac{1}{2} + \frac{|\mathcal{B}_k|}{T} + \frac{10^{-8}}{Tn}$$

and

$$\tilde{F}_k \left(b_k^* - C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \leq \frac{1}{2} - \frac{|\mathcal{B}_k|}{T} - \frac{10^{-8}}{Tn}.$$

Proof. The proof essentially follows Lemma 1 of [55]. We provide the proof for the sake of being self-contained.

Let $Z_k^t = (\check{b}_k^t - b_k^t) / \sqrt{\text{Var}(\check{b}_k^t)}$ be a standardized version of \check{b}_k^t for each $t \in [T]$ and $k \in [d]$, with $b_k^* \neq 1$. Let $\tilde{\Phi}_k(z) = \frac{1}{T - |\mathcal{B}_k|} \sum_{t \in [T] \setminus \mathcal{B}_k} \mathbf{1}(Z_k^t \leq z)$ be the empirical distribution of $\{Z_k^t : t \in [T] \setminus \mathcal{B}_k\}$. The distribution of Z_k^t is identical $t \in [T] \setminus \mathcal{B}_k$, and we denote by Φ_k their common cdf.

By definition, $\mathbb{E}[\tilde{\Phi}_k(z)] = \Phi_k(z)$ for any $z \in \mathbb{R}$. Let $z_1 > 0 > z_2$ be such that $\Phi(z_1) = \frac{1}{2} + G'_{k,T,\delta}$ and $\Phi(z_2) = \frac{1}{2} - G'_{k,T,\delta}$, which exist due to (20). Then, by Lemma 9.1, we have

$$\Phi_k(z_1) \geq \frac{1}{2} + G_{k,T,\delta} \quad \text{and} \quad \Phi_k(z_2) \leq \frac{1}{2} - G_{k,T,\delta}. \quad (21)$$

Further, by the Hoeffding's inequality, we have for any $\delta \geq 0$ and $z \in \mathbb{R}$,

$$\left| \tilde{\Phi}_k(z) - \Phi_k(z) \right| \leq 0.4748 \sqrt{\frac{\delta}{T - |\mathcal{B}_k|}} \quad (22)$$

with probability at least $1 - 2e^{-2\delta}$. Then, by a union bound of (22) for $z = z_1, z_2$, and by (21), it holds with probability at least $1 - 4e^{-2\delta}$ that

$$\tilde{\Phi}_k(z_1) \geq \frac{1}{2} + \frac{|\mathcal{B}_k|}{T} + \frac{10^{-8}}{Tn} \quad \text{and} \quad \tilde{\Phi}_k(z_2) \leq \frac{1}{2} - \frac{|\mathcal{B}_k|}{T} - \frac{10^{-8}}{Tn}. \quad (23)$$

Finally, we bound the values of z_1 and z_2 . By the mean value theorem, there exists $\xi \in [0, z_1]$ such that

$$G'_{k,T,\delta} = z_1 \Phi'(\xi) = \frac{z_1}{\sqrt{2\pi}} e^{-\frac{\xi^2}{2}} \geq \frac{z_1}{\sqrt{2\pi}} e^{-\frac{z_1^2}{2}}. \quad (24)$$

By (20) and the definition of z_1 , we have $z_1 \leq \Phi^{-1}(1 - \varepsilon)$, and thus, by (24), we have

$$z_1 \leq \sqrt{2\pi} G'_{k,T,\delta} \exp\left(\frac{1}{2}(\Phi^{-1}(1 - \varepsilon))^2\right). \quad (25)$$

Similarly, we have

$$z_1 \geq -\sqrt{2\pi} G'_{k,T,\delta} \exp\left(\frac{1}{2}(\Phi^{-1}(1 - \varepsilon))^2\right). \quad (26)$$

Since for all z , $\tilde{\Phi}_k(z) = \tilde{F}_k(\sigma_k z / \sqrt{n} + b_k^*)$, plugging (25) and (26) into (23), we find the conclusion of this lemma. \square

This leads to our next result.

Lemma 10.6. *For any $k \in [d]$ such that condition (20) holds, we have with probability at least $1 - 4e^{-2\delta}$ that*

$$\left| \check{b}_k^t - b_k^t \right| \leq C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G_{k,T,\delta} + \frac{0.4748 C_\varepsilon}{n}. \quad (27)$$

Proof. Let \hat{F}_k be the empirical distribution function of $\{\check{b}_k^t : t \in [T]\}$, such that for all $x \in \mathbb{R}$, $\hat{F}_k(x) := \frac{1}{T} \sum_{t \in [T]} \mathbb{1}(\check{b}_k^t \leq x)$. We have

$$\begin{aligned} |\hat{F}_k(x) - \tilde{F}_k(x)| &= \left| \frac{1}{T} \sum_{t \in [T]} \mathbb{1}(\check{b}_k^t \leq x) - \frac{1}{T - |\mathcal{B}_k|} \sum_{t \in [T] \setminus \mathcal{B}_k} \mathbb{1}(\check{b}_k^t \leq x) \right| \\ &= \left| \frac{1}{T} \sum_{t \in \mathcal{B}_k} \mathbb{1}(\check{b}_k^t \leq x) - \frac{|\mathcal{B}_k|}{T(T - |\mathcal{B}_k|)} \sum_{t \in [T] \setminus \mathcal{B}_k} \mathbb{1}(\check{b}_k^t \leq x) \right| \\ &\leq \max \left\{ \frac{1}{T} \cdot |\mathcal{B}_k|, \frac{|\mathcal{B}_k|}{T(T - |\mathcal{B}_k|)} \cdot (T - |\mathcal{B}_k|) \right\} = \frac{|\mathcal{B}_k|}{T}. \end{aligned} \quad (28)$$

Define $\tilde{F}_k^-(x) := \frac{1}{T-|\mathcal{B}_k|} \sum_{t \in [T] \setminus \mathcal{B}_k} \mathbb{1}(\tilde{b}_k^t < x) \leq \tilde{F}_k(x)$. Then by (28) and Lemma 10.5, we have, with probability at least $1 - 4e^{-2\delta}$ that

$$\hat{F}_k \left(b_k^* + C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \geq \frac{1}{2} + \frac{10^{-8}}{Tn} \quad \text{and} \quad \hat{F}_k^- \left(b_k^* - C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \leq \frac{1}{2} - \frac{10^{-8}}{Tn}. \quad (29)$$

Let $\tilde{b}_k^{(j)}$, $\forall j \in [T]$ be the j -th smallest element in $\{\tilde{b}_k^t : t \in [T]\}$. Recalling the definition of the median, if T is odd, then $\tilde{b}_k^* = \tilde{b}_k^{((T+1)/2)}$. Therefore, $b_k^* + C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} < \tilde{b}_k^*$ implies $\hat{F}_k \left(b_k^* + C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \leq \frac{1}{2} - \frac{1}{2T}$ and $b_k^* - C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} > \tilde{b}_k^*$ implies $\hat{F}_k^- \left(b_k^* - C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \geq \frac{1}{2} + \frac{1}{2T}$, leading to a contradiction with (29).

On the other hand, if T is even, $\tilde{b}_k^* = (\tilde{b}_k^{(T/2)} + \tilde{b}_k^{(T/2+1)})/2$. Therefore, $b_k^* + C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} < \tilde{b}_k^*$ implies $\hat{F}_k \left(b_k^* + C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \leq \frac{1}{2}$ and $b_k^* - C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} > \tilde{b}_k^*$ implies $\hat{F}_k^- \left(b_k^* - C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta} \right) \geq \frac{1}{2}$, which is also contradictory to (29).

To summarize, it holds that

$$|\tilde{b}_k^* - b_k^*| \leq C_\varepsilon \frac{\sigma_k}{\sqrt{n}} G'_{k,T,\delta}$$

with probability at least $1 - 4e^{-2\delta}$. \square

If $T = O(\ln(n))$, Lemma 10.1 follows directly from (18) and (19). Now, given Lemma 10.5 and Lemma 10.6, we turn to prove Lemma 10.1 with $T \geq 20 \ln(n)$. We first check condition (20). Since $|\mathcal{B}_k| \leq \eta T$ for any $k \in \mathcal{I}_\eta$, $\eta \leq \frac{1}{5}$, and $\gamma_k \sigma_k \leq 1$, we have for each $k \in [d]$ that

$$G'_{k,T,\delta} \leq \eta + \frac{10^{-8}}{Tn} + \sqrt{\frac{5\delta}{4T}} + \frac{0.4748}{\sqrt{n}\sigma_k}.$$

When $T \geq 20 \ln(n)$, for any $k \in [d]$ such that $\sigma_k \geq \frac{20}{\sqrt{n}(1-2\eta)}$, taking $\delta = \ln(n)$ above, we have

$$G'_{k,T,\delta} \leq \eta + 10^{-8} + \frac{1}{4} + 0.4748 \frac{1-2\eta}{20} \leq \frac{1}{2} - 0.035755.$$

Therefore, condition (20) in Lemma 10.6 is satisfied with $\varepsilon = 0.035755$, for which we can check that and $C_\varepsilon \leq 13$. Thus, for any $\delta \leq \ln(n)$,

$$\mathbb{P} \left(|\tilde{b}_k^* - b_k^*| \geq 13 \frac{\sigma_k}{\sqrt{n}} G_{k,T,\delta} + \frac{13}{n} \right) \leq 4e^{-2\delta}. \quad (30)$$

Therefore, by (30), we have, using that for any random variable $X \geq 0$ and any $0 \leq \delta \leq 1$, $\mathbb{E}[X] \leq \delta + \mathbb{P}(X \geq \delta)$, and that for $\delta = (\ln n)/2$, one has $4e^{-2\delta} = 4/n$, we find

$$\mathbb{E}[|\tilde{b}_k^* - b_k^*|] \leq 13 \frac{\sigma_k}{\sqrt{n}} G_{k,T,(\ln n)/2} + \frac{17}{n} = \tilde{O} \left(\frac{\sigma_k}{\sqrt{n}} \frac{|\mathcal{B}_k|}{T} + \frac{\sigma_k}{\sqrt{nT}} + \frac{1}{n} \right). \quad (31)$$

Similarly, by the Cauchy-Schwarz inequality, we also have

$$\begin{aligned} \mathbb{E}[(\tilde{b}_k^* - b_k^*)^2] &= O \left(\frac{\sigma_k^2}{n} \left(\frac{|\mathcal{B}_k|^2}{T^2} + \frac{\ln(n)}{T-|\mathcal{B}_k|} \right) + \frac{1}{n^2} + e^{-2\ln(n)} \right) \\ &= \tilde{O} \left(\frac{\sigma_k^2}{n} \frac{|\mathcal{B}_k|^2}{T^2} + \frac{\sigma_k^2}{nT} + \frac{1}{n^2} \right). \end{aligned} \quad (32)$$

On the other hand, for any $k \in [d] \setminus \mathcal{B}_k$ such that $\sigma_k < \frac{20}{\sqrt{n}(1-2\eta)}$, by Bernstein's inequality and a union bound, we have

$$\mathbb{P} \left(\max_{k \in [T] \setminus \mathcal{B}_k} |\tilde{b}_k^t - b_k^*| > \delta \right) \leq 2(T-|\mathcal{B}_k|) e^{-\frac{n}{4} \min\{\delta^2/\sigma_k^2, \delta\}} \leq 2Te^{-\frac{n}{4} \min\{\frac{n(1-2\eta)^2\delta^2}{400}, \delta\}}. \quad (33)$$

Since $|[T] \setminus \mathcal{B}_k| > \frac{T}{2}$, we have as before that $|\check{b}_k^\star - b_k^\star| \leq \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^\star|$. Taking

$$\delta = 4 \max\{\ln(Tn^2), 10\sqrt{\ln(Tn^2)}\}/n$$

in (33), with the same steps as above, we find

$$\begin{aligned} \mathbb{E}[|\check{b}_k^\star - b_k^\star|] &\leq \mathbb{E}[\max_{k \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^\star|] \leq \delta + 2Te^{-\frac{n}{4} \min\{\frac{(1-2\eta)^2 n \delta^2}{400}, \delta\}} \\ &\leq \frac{4 \max\{\ln(Tn^2), 10\sqrt{\ln(Tn^2)}\} + 2}{n} = \tilde{O}\left(\frac{1}{n}\right) \end{aligned} \quad (34)$$

and

$$\mathbb{E}[(\check{b}_k^\star - b_k^\star)^2] \leq \delta^2 + 2Te^{-\frac{n}{4} \min\{\frac{(1-2\eta)^2 n \delta^2}{400}, \delta\}} = \tilde{O}\left(\frac{1}{n^2}\right). \quad (35)$$

To summarize, combining (31), (32) with (34), (35), we complete the proof when $T = \Omega(\ln(n))$.

Furthermore, by using $\sum_{k \in [d]} \sigma_k^q = O(d/2^{bq/2})$ for $q = 1, 2$, we directly reach Proposition 10.2.

10.2 Proof of Theorem 10.3

We first consider the case where $T = O(\ln(n))$. By definition, \hat{b}_k^t is either equal to \check{b}_k^t or \check{b}_k^\star , and the latter happens only when $k \in \mathcal{K}_\alpha^t$, i.e., $|\check{b}_k^\star - \check{b}_k^t| \leq \sqrt{\alpha \check{b}_k^t/n}$. In this case, we have

$$|\hat{b}_k^t - b_k^t| = |\check{b}_k^\star - b_k^t| \leq |\check{b}_k^t - b_k^t| + |\check{b}_k^\star - \check{b}_k^t| \leq |\check{b}_k^t - b_k^t| + \sqrt{\frac{\alpha \check{b}_k^t}{n}}.$$

Therefore, we have $|\hat{b}_k^t - b_k^t| \leq |\check{b}_k^t - b_k^t| + \sqrt{\alpha \check{b}_k^t/n}$ for all $k \in [d]$. This leads to

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_1] &\leq \mathbb{E}[\|\check{\mathbf{b}}^t - \mathbf{b}^t\|_1] + \sqrt{\frac{\alpha}{n}} \sum_{k \in [d]} \mathbb{E}\left[\sqrt{\check{b}_k^t}\right] \\ &\leq \mathbb{E}[\|\check{\mathbf{b}}^t - \mathbf{b}^t\|_1] + \sqrt{\frac{\alpha}{n}} \sum_{k \in [d]} \sqrt{\mathbb{E}[\check{b}_k^t]}, \end{aligned} \quad (36)$$

where (36) holds by Jensen's inequality. By further using the Cauchy-Schwarz inequality, we have

$$\mathbb{E}[\|\check{\mathbf{b}}^t - \mathbf{b}^t\|_1] \leq \sqrt{d \mathbb{E}[\|\check{\mathbf{b}}^t - \mathbf{b}^t\|_2^2]} = O\left(\frac{d}{\sqrt{2^b n}}\right) \quad (37)$$

and

$$\sum_{k \in [d]} \sqrt{\mathbb{E}[\check{b}_k^t]} = \sum_{k \in [d]} \sqrt{b_k^t} \leq \sqrt{d \sum_{k \in [d]} b_k^t} = O\left(\frac{d}{\sqrt{2^b}}\right). \quad (38)$$

Plugging (37) and (38) into (36), we find

$$\mathbb{E}[\|\check{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O}\left(\frac{d}{\sqrt{2^b n}}\right) = \tilde{O}\left(\frac{d}{\sqrt{2^b T n}}\right).$$

We can similarly prove

$$\mathbb{E}[\|\check{\mathbf{b}}^t - \mathbf{b}^t\|_2^2] = \tilde{O}\left(\frac{d}{2^b n}\right) = \tilde{O}\left(\frac{d}{2^b T n}\right).$$

Next we prove the case where $T \geq 20 \ln(n) = \Omega(\ln(n))$. We first consider the estimation errors over $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$ such that $\sigma_k \geq \frac{20}{\sqrt{n(1-2\eta)}}$. Let $\mathcal{E}_k^t := \{\check{b}_k^t \geq \frac{1}{2}b_k^t \text{ and } |\check{b}_k^* - b_k^*| \leq 8\sqrt{b_k^*/n}\}$. If $n \geq 2^{b+6} \ln(n)$ and $0 < \eta \leq 1/5$, then since $b_k^* \geq \frac{1}{2^b}$ for any $k \in [d]$, we have

$$\begin{aligned} 13 \frac{\sigma_k}{\sqrt{n}} G_{k,T,\ln n} + \frac{13}{n} &= 13 \frac{\sigma_k}{\sqrt{n}} \left(\frac{|\mathcal{B}_k|}{T} + \frac{10^{-8}}{Tn} + \sqrt{\frac{\ln(n)}{T - |\mathcal{B}_k|}} \right) + \frac{13}{n} \\ &\leq 13 \frac{\sigma_k}{\sqrt{n}} \left(\frac{|\mathcal{B}_k|}{T} + \frac{10^{-8}}{Tn} + \sqrt{\frac{5 \ln(n)}{4T}} \right) + \frac{13}{n} \leq 13 \frac{\sigma_k}{\sqrt{n}} \left(\frac{1}{5} + 10^{-8} + \frac{1}{4} \right) + \frac{13}{\sqrt{n 2^{b+6} \ln(n)}} \\ &\leq 13 \frac{\sigma_k}{\sqrt{n}} \left(\frac{1}{5} + 10^{-8} + \frac{1}{4} \right) + \frac{13 \sqrt{b_k^*}}{\sqrt{n 64 \ln(n)}} \leq 8 \sqrt{\frac{b_k^*}{n}}. \end{aligned}$$

Hence, by (30), it holds that

$$\mathbb{P} \left(|\check{b}_k^* - b_k^*| \geq 8 \sqrt{\frac{b_k^*}{n}} \right) \leq \frac{4}{n^2}. \quad (39)$$

By Bernstein's inequality and as $b_k^* \geq \frac{1}{2^b}$, we have

$$\mathbb{P} \left(|\check{b}_k^t - b_k^t| > \frac{b_k^t}{2} \right) \leq 2e^{-\frac{n}{4} \min\{\frac{b_k^t}{4(1-b_k^t)}, \frac{b_k^t}{2}\}} \leq 2e^{-\frac{nb_k^t}{16}} \leq 2e^{-\frac{n}{16 \cdot 2^b}} \leq \frac{2}{n^2}, \quad (40)$$

where the last inequality holds because $n \geq 2^{b+6} \ln(n)$. Combining (40) with (39), we find $\mathbb{P}((\mathcal{E}_k^t)^c) \leq \frac{6}{n^2}$. By definition, $k \notin \mathcal{K}_\alpha^t$ implies $|\check{b}_k^* - \check{b}_k^t| > \sqrt{\alpha \check{b}_k^t/n}$. On the event \mathcal{E}_k^t , this further implies $|\check{b}_k^* - \check{b}_k^t| > \sqrt{\alpha b_k^t/2n}$. Combined with (39) and that $b_k^* = b_k^t$ for any $k \in \mathcal{I}^t$, we have on the event \mathcal{E}_k^t

$$|\check{b}_k^t - b_k^t| = |\check{b}_k^t - b_k^*| \geq |\check{b}_k^t - \check{b}_k^*| - |\check{b}_k^* - b_k^*| > \sqrt{\frac{b_k^t}{n}} \left(\sqrt{\frac{\alpha}{2}} - 8 \right). \quad (41)$$

Let $\zeta \triangleq \sqrt{\alpha/2} - 8 \geq \sqrt{8 \ln(n)}$ and $\mathcal{F}_k^t := \left\{ |\check{b}_k^t - b_k^t| \geq \zeta \sqrt{b_k^t/n} \right\}$. By Bernstein's inequality, and using $n \geq 2^{b+6} \ln(n)$, we have

$$\mathbb{P}(\mathcal{F}_k^t) \leq 2e^{-\frac{n}{4} \min\{\frac{\zeta^2}{n(1-b_k^t)}, \zeta \sqrt{\frac{b_k^t}{n}}\}} \leq 2e^{-\min\{\frac{\zeta^2}{4}, \zeta \sqrt{\frac{n}{2^b}}\}} \leq \frac{2}{n^2}. \quad (42)$$

Combining (41) with (42), we find that for any $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$ with $\sigma_k \geq \frac{20}{\sqrt{n(1-2\eta)}}$, it holds that

$$\begin{aligned} \mathbb{P}(k \notin \mathcal{K}_\alpha^t) &\leq \mathbb{P}((\mathcal{E}_k^t)^c) + \mathbb{P}(\mathcal{E}_k^t \cap \{k \notin \mathcal{K}_\alpha^t\}) \leq \mathbb{P}((\mathcal{E}_k^t)^c) + \mathbb{P}(\mathcal{E}_k \cap \mathcal{F}_k^t) \\ &\leq \mathbb{P}((\mathcal{E}_k^t)^c) + \mathbb{P}(\mathcal{F}_k^t) \leq \frac{8}{n^2}. \end{aligned}$$

On the other hand for any $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$ with $\sigma_k < \frac{20}{\sqrt{n(1-2\eta)}}$, we have

$$\sqrt{\frac{b_k^t(1-b_k^t)}{n}} = \sqrt{\frac{b_k^*(1-b_k^*)}{n}} = \frac{\sigma_k}{\sqrt{n}} = O\left(\frac{1}{n}\right).$$

Therefore, we have for all $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$, and $q = 1, 2$

$$\min \left\{ \mathbb{P}(k \notin \mathcal{K}_\alpha^t), \left(\frac{b_k^t(1-b_k^t)}{n} \right)^{q/2} \right\} = O\left(\frac{1}{n^q}\right). \quad (43)$$

Since $\alpha = O(\ln(n))$, by (43) and Proposition 9.6, we obtain

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O} \left(\sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \sqrt{\frac{b_k^t}{n}} + \mathbb{E}[\|\check{\mathbf{b}}_{\mathcal{I}_\eta \cap \mathcal{I}^t}^* - \mathbf{b}_{\mathcal{I}_\eta \cap \mathcal{I}^t}^*\|_1] + \frac{d}{n} \right). \quad (44)$$

Combining (44) with Proposition 10.2 and using that $\sigma_k \leq \sqrt{b_k^*} = \sqrt{b_k^t}$ for any $k \in \mathcal{I}^t$, we have

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O} \left(\sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \sqrt{\frac{b_k^t}{n}} + \sum_{k \in \mathcal{I}_\eta \cap \mathcal{I}^t} \frac{|\mathcal{B}_k|}{T} \sqrt{\frac{b_k^t}{n}} + E(1) \right). \quad (45)$$

Since $|(\mathcal{I}^t)^c| = \|\mathbf{p}^t - \mathbf{p}^*\|_0 \leq s$, by the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \sum_{k \notin \mathcal{I}_\eta \cap \mathcal{I}^t} \sqrt{\frac{b_k^t}{n}} &\leq \sum_{k \notin \mathcal{I}_\eta} \sqrt{\frac{b_k^t}{n}} + \sum_{k \notin \mathcal{I}^t} \sqrt{\frac{b_k^t}{n}} \leq \sum_{k \notin \mathcal{I}_\eta} \sqrt{\frac{b_k^t}{n}} + \sqrt{\frac{s \sum_{k \notin \mathcal{I}^t} b_k^t}{n}} \\ &\leq \sum_{k \notin \mathcal{I}_\eta} \sqrt{\frac{b_k^t}{n}} + \sqrt{\frac{s \sum_{k \notin \mathcal{I}^t} ((2^b - 1)p_k^t + 1)}{2^b n}} \leq \sum_{k \notin \mathcal{I}_\eta} \sqrt{\frac{b_k^t}{n}} + \sqrt{\frac{s(2^b - 1 + s)}{2^b n}}. \end{aligned} \quad (46)$$

Plugging (46) into (44), we further obtain

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O} \left(\sum_{k \notin \mathcal{I}_\eta} \sqrt{\frac{b_k^t}{n}} + \sum_{k \in \mathcal{I}_\eta} \frac{|\mathcal{B}_k|}{T} \sqrt{\frac{b_k^t}{n}} + \sqrt{\frac{s \max\{2^b, s\}}{2^b n}} + E(1) \right). \quad (47)$$

Similarly, we can reach the following ℓ_2 counterpart:

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_2^2] = \tilde{O} \left(\sum_{k \notin \mathcal{I}_\eta} \frac{b_k^t}{n} + \sum_{k \in \mathcal{I}_\eta} \frac{|\mathcal{B}_k|^2}{T^2} \frac{b_k^t}{n} + \frac{\max\{2^b, s\}}{2^b n} + E(2) \right). \quad (48)$$

Note that $\sum_{k \in [d]} |\mathcal{B}_k|/T \leq s$ and for any set \mathcal{I} with $|\mathcal{I}| = \lceil \frac{s}{\eta} \rceil$,

$$\sum_{k \in \mathcal{I}} \sqrt{\frac{b_k^t}{n}} \leq \sqrt{\frac{|\mathcal{I}| \sum_{k \in \mathcal{I}} ((2^b - 1)p_k^t + 1)}{2^b n}} = O \left(\sqrt{\frac{s/\eta \max\{2^b, s/\eta\}}{2^b n}} \right).$$

Now, recalling the definition of \mathcal{I}_η , we apply Lemma 10.7 in (47) with $(r_k, x_k) = (\sqrt{b_k^t/n}, |\mathcal{B}_k|/T)$ for all $k \in [d]$, to find

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O} \left(\sqrt{\frac{s/\eta \max\{2^b, s/\eta\}}{2^b n}} + E(1) \right).$$

Therefore, for any $\eta = \Theta(1)$ with $\eta \leq \frac{1}{5}$, we finally have

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O} \left(\sqrt{\frac{s \max\{2^b, s\}}{2^b n}} + E(1) \right).$$

Similarly, by combining (48) with Lemma 10.7, we have for any $\eta = \Theta(1)$ with $\eta \leq \frac{1}{5}$,

$$\mathbb{E}[\|\widehat{\mathbf{b}}^t - \mathbf{b}^t\|_2^2] = \tilde{O} \left(\frac{\max\{2^b, s\}}{2^b n} + E(2) \right).$$

The result directly follows Proposition 8.2.

Lemma 10.7. Given $\eta \in (0, 1]$, $r_k \geq 0$ for all $k \in [d]$, and for $q = 1, 2$, consider the functions $f_q : \{x \in \mathbb{R}^d : 0 \leq x_k \leq 1, \forall k \in [d] \text{ and } \sum_{k \in [d]} x_k \leq s\} \rightarrow \mathbb{R}$, $f_q(x_1, \dots, x_d) := \sum_{k \in [d]} r_k^q (\mathbb{1}\{x_k \geq \eta\} + x_k^q \mathbb{1}\{x_k < \eta\})$. Then it holds that

$$\max_{x_1, \dots, x_d} f_q(x_1, \dots, x_d) \leq \sum_{k=1}^{\lceil s/\eta \rceil} r_{(k)}^q, \quad (49)$$

where $r_{(1)} \geq \dots \geq r_{(d)}$ is the non-decreasing rearrangement of $\{r_1, \dots, r_d\}$.

Proof. We only prove the result for f_1 , and the result for function f_2 follows similarly. Note that $\mathbb{1}\{r_k \geq \eta\} + x_k \mathbb{1}\{r_k < \eta\}$ is increasing with respect to r_k and x_k . To consider the maximum of the sum in f , by the rearrangement inequality, without loss of generality, we can assume $r_1 \geq r_2 \geq \dots \geq r_d \geq 0$ and $1 \geq x_1 \geq x_2 \geq \dots \geq x_d \geq 0$. In this case, we claim that the maximum is attained at $x_1 = \dots = x_{\lfloor s/\eta \rfloor} = \eta$, $x_{\lfloor s/\eta \rfloor + 1} = s - \eta \lfloor s/\eta \rfloor$, and $x_k = 0$ for all $k > \lfloor s/\eta \rfloor + 1$. Further, the maximum is $\sum_{k=1}^{\lfloor s/\eta \rfloor} r_k + r_{\lfloor s/\eta \rfloor + 1} (s - \eta \lfloor s/\eta \rfloor)$, which is upper bounded by the right-hand side of (49). We now use the exchange argument to prove the claim.

- S. 1 If there is some k such that $x_k > \eta \geq x_{k+1}$, then defining x' by letting $(x'_k, x'_{k+1}) = (\eta, x_k + x_{k+1} - \eta)$ while for other j , $x'_j = x_j$, increases the value of f . Therefore, the maximum is attained by x such that for some j , $x_1 = \dots = x_j = \eta > x_{j+1} \geq \dots \geq x_d$.
- S. 2 If there is some k such that $\eta > x_k \geq x_{k+1} > 0$, then defining x' by letting $(x'_k, x'_{k+1}) = (\min\{\eta, x_k + x_{k+1} - \eta\}, \max\{0, x_k + x_{k+1} - \eta\})$ while for other j , $x'_j = x_j$, increases the value of f . Therefore, combined with Step 1, the maximum is attained by x such that for some j , $x_1 = \dots = x_j = \eta > x_{j+1} \geq 0$ and $x_k = 0$ for all $k > j + 1$. Thus most one element lies in $(0, \eta)$.

Combining S. 1 and S. 2 above, we complete the proof of the claim, which further leads to (49). \square

11 Trimmed-Mean-Based Method

In this section, we study the trimmed-mean-based estimator. Fix $\omega \in (0, 1/2)$. Specifically, for each $k \in [d]$, let \mathcal{U}_k be the subset of $\{\tilde{\mathbf{p}}^t\}_{t \in [T]}$ obtained by removing the largest and smallest ωT elements². Then, the trimmed-mean-based method can be expressed as

$$\check{b}_k^* = \frac{1}{|\mathcal{U}_k|} \sum_{t \in \mathcal{U}_k} \check{b}_k^t. \quad (50)$$

We also write $\check{\mathbf{b}}^* = \text{trmean}(\{\check{\mathbf{b}}^t\}_{t \in [T]}, \omega)$. For any chosen trimming proportion $0 \leq \eta \leq \omega \leq \frac{1}{5}$, we control the estimation error of each η -well aligned entry. Intuitively, this is small because there are at most a fraction of η elements from heterogeneous distributions. These are trimmed if they behave as outliers, and otherwise kept in \mathcal{U}_k . The error control for a single entry $k \in \mathcal{I}_\eta$ is in Lemma 11.1.

Lemma 11.1. Suppose $\check{\mathbf{b}}^* = \text{trmean}(\{\check{\mathbf{b}}^t\}_{t \in [T]}, \omega)$ such that $0 \leq \omega \leq \frac{1}{5}$. Then for each $k \in \mathcal{I}_\eta$ with $0 < \eta \leq \omega$ and any $q = 1, 2$, it holds that

$$\mathbb{E}[|\check{b}_k^* - b_k^*|^q] = \tilde{O} \left(\left(\omega^2 \frac{b_k^*}{n} \right)^{q/2} + \left(\frac{b_k^*}{Tn} \right)^{q/2} + \frac{1}{(Tn)^q} + \left(\frac{\omega}{n} \right)^q \right). \quad (51)$$

Proof. To prove Lemma 11.1, we need the following lemma.

Lemma 11.2. For each $k \in \mathcal{I}_\eta$ with $0 < \eta \leq \omega \leq \frac{1}{5}$, and any $\varepsilon_k, \delta_k \geq 0$, it holds with probability at least $1 - 2e^{-\frac{(T - |\mathcal{B}_k|)n}{4} \min\{\frac{\varepsilon_k^2}{\sigma_k^2}, \varepsilon_k\}} - 2(T - |\mathcal{B}_k|)e^{-\frac{n}{4} \min\{\frac{\delta_k^2}{\sigma_k^2}, \delta_k\}}$ that

$$|\check{b}_k^* - b_k^*| \leq \frac{\varepsilon_k + 3\omega\delta_k}{1 - 2\omega}.$$

²To be precise, one can either trim $\lceil \omega T \rceil$ or $\lfloor \omega T \rfloor$ elements. From now on, we write ωT for conciseness without further notice.

Proof of Lemma 11.2. By Bernstein's inequality and the union bound, we have for any $\varepsilon_k, \delta_k > 0$ that

$$\mathbb{P} \left(\left| \frac{1}{T - |\mathcal{B}_k|} \sum_{t \in [T] \setminus \mathcal{B}_k} \check{b}_k^t - b_k^* \right| > \varepsilon_k \right) \leq 2e^{-\frac{(T - |\mathcal{B}_k|)n}{4} \min\{\frac{\varepsilon_k^2}{\sigma_k^2}, \varepsilon_k\}}$$

and

$$\mathbb{P} \left(\max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*| > \delta_k \right) \leq 2(T - |\mathcal{B}_k|)e^{-\frac{n}{4} \min\{\frac{\delta_k^2}{\sigma_k^2}, \delta_k\}}.$$

By the definition of \check{b}_k^* , we have

$$\begin{aligned} |\check{b}_k^* - b_k^*| &= \frac{1}{T(1 - 2\omega)} \left| \sum_{t \in \mathcal{U}_k} \check{b}_k^t - b_k^* \right| \\ &= \frac{1}{T(1 - 2\omega)} \left| \sum_{t \in [T] \setminus \mathcal{B}_k} (\check{b}_k^t - b_k^*) - \sum_{t \in [T] \setminus (\mathcal{B}_k \cup \mathcal{U}_k)} (\check{b}_k^t - b_k^*) + \sum_{t \in \mathcal{B}_k \cap \mathcal{U}_k} (\check{b}_k^t - b_k^*) \right| \\ &\leq \frac{1}{T(1 - 2\omega)} \left(\left| \sum_{t \in [T] \setminus \mathcal{B}_k} \check{b}_k^t - b_k^* \right| + \left| \sum_{i \notin \mathcal{B}_k \cup \mathcal{U}_k} \check{b}_k^t - b_k^* \right| + \left| \sum_{t \in \mathcal{B}_k \cap \mathcal{U}_k} \check{b}_k^t - b_k^* \right| \right). \end{aligned}$$

It is clear that

$$\left| \sum_{t \in [T] \setminus (\mathcal{B}_k \cup \mathcal{U}_k)} (\check{b}_k^t - b_k^*) \right| \leq |[T] \setminus \mathcal{U}_k| \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*| = 2\omega T \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*|.$$

Then we claim that $\left| \sum_{t \in \mathcal{B}_k \cap \mathcal{U}_k} \check{b}_k^t - b_k^* \right| \leq |\mathcal{B}_k| \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*|$. Let $\mathcal{Q}_{k,l}$ and $\mathcal{Q}_{k,r}$ be the indices of the trimmed elements on the left side and right side, respectively, i.e., the smallest and largest ωT elements among $\{\check{b}_k^t\}_{t \in [T]}$. If $\mathcal{B}_k \cap \mathcal{U}_k \neq \emptyset$, then $|\mathcal{U}_k \setminus \mathcal{B}_k| < T(1 - 2\omega)$. Furthermore, we have $|\mathcal{Q}_{k,l} \cup (\mathcal{U}_k \setminus \mathcal{B}_k)| = |\mathcal{Q}_{k,r} \cup (\mathcal{U}_k \setminus \mathcal{B}_k)| = \omega T + |\mathcal{U}_k \setminus \mathcal{B}_k| < T(1 - \omega) \leq |[T] \setminus \mathcal{B}_k|$, which leads to $([T] \setminus \mathcal{B}_k) \cap \mathcal{Q}_{k,l} \neq \emptyset$ and $(T \setminus \mathcal{B}_k) \cap \mathcal{Q}_{k,r} \neq \emptyset$. In conclusion, we have $\max_{t \in \mathcal{U}_k} |\check{b}_k^t - b_k^*| \leq \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*|$, which completes the proof of the claim. Therefore, we have

$$|\check{b}_k^* - b_k^*| \leq \frac{1}{T(1 - 2\omega)} \left(\left| \sum_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*| \right| + (2\omega T + |\mathcal{B}_k|) \max_{t \in [T] \setminus \mathcal{B}_k} |\check{b}_k^t - b_k^*| \right) \leq \frac{\varepsilon_k + 3\omega\delta_k}{1 - 2\omega}$$

with probability at least $1 - 2e^{-\frac{(T - |\mathcal{B}_k|)n}{4} \min\{\frac{\varepsilon_k^2}{\sigma_k^2}, \varepsilon_k\}} - 2(T - |\mathcal{B}_k|)e^{-\frac{n}{4} \min\{\frac{\delta_k^2}{\sigma_k^2}, \delta_k\}}$. \square

Given Lemma 11.2, by setting

$$\varepsilon_k = \max \left\{ \frac{4\sigma_k \sqrt{\ln(T^2 n^2)}}{\sqrt{(T - |\mathcal{B}_k|)n}}, \frac{8 \ln(T^2 n^2)}{(T - |\mathcal{B}_k|)n} \right\} = \tilde{O} \left(\frac{\sigma_k}{\sqrt{Tn}} + \frac{1}{Tn} \right)$$

and

$$\delta_k = \max \left\{ \frac{4\sigma_k \sqrt{\ln(T^2 (T - |\mathcal{B}_k|)n^2)}}{\sqrt{n}}, \frac{4 \ln(T^2 (T - |\mathcal{B}_k|)n^2)}{n} \right\} = \tilde{O} \left(\frac{\sigma_k}{\sqrt{n}} + \frac{1}{n} \right),$$

using that $1/(1-2\omega) \leq \frac{5}{3}$, and recalling $\sigma_k \leq \sqrt{b_k^*}$, we have with probability at least $1 - \frac{4}{T^2 n^2}$ that

$$\begin{aligned} |\check{b}_k^* - b_k^*| &\leq \frac{\varepsilon_k + 3\omega\delta_k}{1-2\omega} \\ &\leq \frac{5\omega}{3} \max \left\{ \frac{4\sqrt{b_k^* \ln(T^3 n^2)}}{\sqrt{n}}, \frac{4\ln(T^3 n^2)}{n} \right\} + \frac{5}{3} \max \left\{ \frac{4\sqrt{b_k^* \ln(T^2 n^2)}}{\sqrt{(T-|\mathcal{B}_k|)n}}, \frac{4\ln(T^2 n^2)}{(T-|\mathcal{B}_k|)n} \right\} \\ &= \tilde{O} \left(\omega \sqrt{\frac{b_k^*}{n}} + \frac{\omega}{n} + \frac{\sigma_k}{\sqrt{Tn}} + \frac{1}{Tn} \right), \end{aligned} \quad (52)$$

which implies

$$\begin{aligned} \mathbb{E}[|\check{b}_k^* - b_k^*|] &= \tilde{O} \left(\omega \sqrt{\frac{b_k^*}{n}} + \frac{\omega}{n} + \frac{\sigma_k}{\sqrt{Tn}} + \frac{1}{Tn} + \frac{1}{T^2 n^2} \right) \\ &= \tilde{O} \left(\omega \sqrt{\frac{b_k^*}{n}} + \sqrt{\frac{b_k^*}{Tn}} + \frac{1}{Tn} + \frac{\omega}{n} \right). \end{aligned}$$

Similarly, we can obtain

$$\begin{aligned} \mathbb{E}[(\check{b}_k^* - b_k^*)^2] &= \tilde{O} \left(\frac{\omega^2 b_k^*}{n} + \frac{\omega^2}{n^2} + \frac{\sigma_k^2}{Tn} + \frac{1}{T^2 n^2} + \frac{1}{T^2 n^2} \right) \\ &= \tilde{O} \left(\omega^2 \frac{b_k^*}{n} + \frac{b_k^*}{Tn} + \frac{1}{T^2 n^2} + \frac{\omega^2}{n^2} \right). \end{aligned}$$

□

Given these results, we readily establish the following bound on the total error over all η -well-aligned entries.

Proposition 11.3. *Suppose $\check{\mathbf{b}}^* = \text{trmean}(\{\check{\mathbf{b}}^t\}_{t \in [T]}, \omega)$ such that $0 \leq \omega \leq 1/5$. Then for each $k \in \mathcal{I}_\eta$ with $0 < \eta \leq \omega$ and any $q = 1, 2$, it holds that*

$$\mathbb{E}[\|\check{b}_{\mathcal{I}_\eta}^* - b_{\mathcal{I}_\eta}^*\|_q^q] = \tilde{O} \left(d \left(\frac{\omega^2}{2^b n} \right)^{q/2} + \frac{d}{(2^b T n)^{q/2}} + \frac{d}{(T n)^q} + d \left(\frac{\omega}{n} \right)^q \right).$$

By setting $\alpha = \Theta(\ln(Tn))$, we find the following result.

Theorem 11.4. *Suppose $n \geq 2^{b+5} \ln(Tn)$ and $\alpha \geq 2(8 + \sqrt{8 \ln(Tn)})^2$ with $\alpha = O(\ln(Tn))$. Then for the trimmed-mean-based SHIFT method, for any $0 < \omega \leq \frac{1}{5}$, $t \in [T]$ and $q = 1, 2$,*

$$\mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] = \tilde{O} \left(\left(\frac{s}{\omega} \right)^{1-q/2} \left(\frac{\max\{2^b, s/\omega\}}{2^b n} \right)^{q/2} + d \left(\frac{\omega^2}{2^b n} \right)^{q/2} + \frac{d}{(2^b T n)^{q/2}} \right).$$

Proof. To apply Proposition 9.6, we need to bound $\sum_{k \in \mathcal{I}_\eta \cap \mathcal{I}^t} \min\{\mathbb{P}(k \notin \mathcal{K}_\alpha^t), \sqrt{b_k^t(1-b_k^t)/n}\}$ and $\sum_{k \in \mathcal{I}_\eta \cap \mathcal{I}^t} \min\{\mathbb{P}(k \notin \mathcal{K}_\alpha^t), b_k^t(1-b_k^t)/n\}$.

Let $\mathcal{E}_k^t := \{\check{b}_k^t \geq \frac{1}{2}b_k^t \text{ and } |\check{b}_k^t - b_k^t| \leq 8\sqrt{b_k^* \ln(T^3 n^2)/n}\}$. For each entry $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$, since $n \geq 2^b \ln(T^3 n^2)$ and $b_k^* \leq \frac{1}{2^b}$, we have $\frac{1}{n} \leq \sqrt{\frac{b_k^*}{n \ln(T^3 n^2)}}$. By (52), we have with probability at least $1 - \frac{4}{T^2 n^2}$ that

$$\begin{aligned} |\check{b}_k^t - b_k^t| &\leq \frac{5\omega}{3} \max \left\{ \frac{4\sqrt{b_k^* \ln(T^3 n^2)}}{\sqrt{n}}, \frac{4\ln(T^3 n^2)}{n} \right\} + \frac{5}{3} \max \left\{ \frac{4\sqrt{b_k^* \ln(T^2 n^2)}}{\sqrt{(T-|\mathcal{B}_k|)n}}, \frac{4\ln(T^2 n^2)}{(T-|\mathcal{B}_k|)n} \right\} \\ &\leq \frac{4}{3} \sqrt{\frac{b_k^* \ln(T^3 n^2)}{n}} + \frac{20}{3} \sqrt{\frac{b_k^* \ln(T^2 n^2)}{(T-|\mathcal{B}_k|)n}} \leq 8 \sqrt{\frac{b_k^* \ln(T^3 n^2)}{n}}. \end{aligned} \quad (53)$$

By Bernstein's inequality and as $b_k^* \geq \frac{1}{2^b}$, we have

$$\mathbb{P}\left(\left|\check{b}_k^t - b_k^t\right| > \frac{b_k^t}{2}\right) \leq 2e^{-\frac{n}{4} \min\{\frac{b_k^t}{4(1-b_k^t)}, \frac{b_k^t}{2}\}} \leq 2e^{-\frac{nb_k^t}{16}} \leq 2e^{-\frac{n}{16 \cdot 2^b}} \leq \frac{2}{T^2 n^2}, \quad (54)$$

where the last inequality is because $n \geq 2^{b+5} \ln(Tn)$. Combining (53) with (54), we find $\mathbb{P}((\mathcal{E}_k^t)^c) \leq \frac{6}{T^2 n^2}$. Now following the argument from (41)-(43), we can obtain that for all $k \in \mathcal{I}_\eta \cap \mathcal{I}^t$,

$$\mathbb{P}(k \notin \mathcal{K}_\alpha^t) = O\left(\frac{1}{T^2 n^2}\right).$$

Since $\alpha = O(\ln(Tn))$, by applying (43) to Proposition 9.6 with $\eta = \omega$ and using Proposition 11.3 with $n = \Omega(2^b)$, we find

$$\mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_1] = \tilde{O}\left(\sum_{k \notin \mathcal{I}_\omega \cap \mathcal{I}^t} \sqrt{\frac{b_k^t}{n}} + \frac{d\omega}{\sqrt{2^b n}} + \frac{d}{\sqrt{2^b T n}}\right) \quad (55)$$

and

$$\mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_2^2] = \tilde{O}\left(\sum_{k \notin \mathcal{I}_\omega \cap \mathcal{I}^t} \frac{b_k^t}{n} + \frac{d\omega}{2^b n} + \frac{d}{2^b T n}\right).$$

Note that $|(\mathcal{I}_\omega \cap \mathcal{I}^t)^c| \leq |\mathcal{I}_\omega^c| + |(\mathcal{I}^t)^c| \leq s/\omega + s = O(s/\omega)$ and

$$\begin{aligned} \sum_{k \notin \mathcal{I}_\omega \cap \mathcal{I}^t} \sqrt{\frac{b_k^t}{n}} &\leq \sqrt{|(\mathcal{I}_\omega \cap \mathcal{I}^t)^c| \sum_{k \notin \mathcal{I}_\omega \cap \mathcal{I}^t} \frac{b_k^t}{n}} = \sqrt{\frac{|(\mathcal{I}_\omega \cap \mathcal{I}^t)^c| \max\{2^b, |(\mathcal{I}_\omega \cap \mathcal{I}^t)^c|\}}{2^b n}} \\ &= \sqrt{\frac{s/\omega \max\{2^b, s/\omega\}}{2^b n}}. \end{aligned} \quad (56)$$

Plugging (56) into (55) and using $\mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_1] = O(\mathbb{E}[\|\hat{\mathbf{b}}^t - \mathbf{b}^t\|_1])$, we find the conclusion in terms of the ℓ_1 error. The results in terms of the ℓ_2 error can be obtained similarly. \square

12 Lower Bounds

In this section, we provide the proofs for the minimax lower bounds for estimating distributions under our heterogeneity model. We first re-state the detailed version the lower bounds that apply to both the ℓ_2 and ℓ_1 errors.

Theorem 12.1 (Detailed statement of Theorem 4.1). *For any—possibly interactive—estimation method, and for any $t \in [T]$ and $q = 1, 2$, we have*

$$\inf_{\substack{(W^{t'}, [n])_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] = \Omega\left(s^{1-q/2} \left(\frac{\max\{2^b, s\}}{2^b n}\right)^{q/2} + \frac{d}{(2^b T n)^{q/2}}\right). \quad (57)$$

Theorem 12.2 (Detailed statement of Theorem 4.2). *For any—possibly interactive—estimation method, and a new cluster \mathcal{C}^{T+1} , we have*

$$\begin{aligned} &\inf_{\substack{(W^{t'}, [n])_{t' \in [T]} \\ W^{T+1, [\tilde{n}]}, \hat{\mathbf{p}}^{T+1}}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T+1]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^{T+1} - \mathbf{p}^{T+1}\|_q^q] \\ &= \Omega\left(s^{1-q/2} \left(\frac{\max\{2^b, s\}}{2^b \tilde{n}}\right)^{q/2} + \frac{d}{(2^b T n)^{q/2}}\right). \end{aligned}$$

We omit the proof of Theorem 12.2 since it follows from the same analysis as Theorem 12.1.

12.1 Proof of Theorem 12.1

As discussed in Section 4, we will prove (57) by considering two special cases of our sparse heterogeneity model:

1. The *homogeneous* case where $\mathbf{p}^1 = \dots = \mathbf{p}^T = \mathbf{p}^* \in \mathcal{P}_d$.
2. The *s/2-sparse* case where $\|\mathbf{p}^*\|_0 \leq s/2$ and $\|\mathbf{p}^t\|_0 \leq s/2$ for all $t \in [T]$.

Therefore, it naturally holds that

$$\inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] \geq \inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^*}} \sup_{\mathbf{p}^* \in \mathcal{P}_d} \mathbb{E}[\|\hat{\mathbf{p}}^* - \mathbf{p}^*\|_q^q] \quad (58)$$

and

$$\inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] \geq \inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^{t'} \in \mathcal{P}_d \\ \|\mathbf{p}^{t'}\|_0 \leq s/2 \\ \forall t' \in [T]}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q]. \quad (59)$$

For the first case, combining (58) with the existing lower bound result [9, Cor 7] and [26, Thm 2] for the homogeneous setup, where all datapoints are generated by a single distribution, that for any estimation method (possibly based on interactive encoding),

$$\inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^*}} \sup_{\mathbf{p}^* \in \mathcal{P}_d} \mathbb{E}[\|\hat{\mathbf{p}}^* - \mathbf{p}^*\|_q^q] = \Omega \left(\frac{d}{(2^b T n)^{q/2}} \right),$$

we prove that the lower bound is at least of the order of the second term in (57).

For the second case, without loss of generality, we assume s is even. This can be achieved by considering $s - 1$ instead of s , if necessary. Recall that $\text{supp}(\cdot)$ denotes the indices of non-zero entries of a vector. Fixing any $t \in [T]$, we further consider the scenario where

$$\text{supp}(\mathbf{p}^t) \cap \left(\bigcup_{t' \neq t} \text{supp}(\mathbf{p}^{t'}) \right) = \emptyset. \quad (60)$$

One example where (60) holds is when $\text{supp}(\mathbf{p}^t) \subseteq [s/2]$ and $\text{supp}(\mathbf{p}^{t'}) \subseteq \{s/2 + 1, \dots, d\}$ for all $t' \neq t$. If (60) holds, then the support of the datapoints generated by $\{\mathbf{p}^{t'} : t' \neq t\}$ does not overlap with the support of those generated by \mathbf{p}^t , and hence former are not informative for estimating \mathbf{p}^t . Therefore, by further combining (59) with the existing lower bound result [16, Thm 2] for the $s/2$ -sparse homogeneous setup, where all datapoints are generated by a single $s/2$ -sparse distribution, that for any estimation method (possibly based on interactive encoding),

$$\begin{aligned} \inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \|\mathbf{p}^t\|_0 \leq s/2 \\ \mathbf{p}^t \in \mathcal{P}_d}} \sup_{\mathbf{p}^t \in \mathcal{P}_d} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] &= \Omega \left((s/2)^{1-q/2} \left(\frac{\max\{2^b, s/2\}}{2^b n} \right)^{q/2} \right) \\ &= \Omega \left(s^{1-q/2} \left(\frac{\max\{2^b, s\}}{2^b n} \right)^{q/2} \right). \end{aligned}$$

Thus, we have

$$\begin{aligned} &\inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^* \in \mathcal{P}_d \\ \{\mathbf{p}^{t'} : t' \in [T]\} \subseteq \mathbb{B}_s(\mathbf{p}^*)}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] \geq \inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^{t'} \in \mathcal{P}_d \\ \|\mathbf{p}^{t'}\|_0 \leq s/2 \\ \forall t' \in [T]}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] \\ &\geq \inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^{t'} \in \mathcal{P}_d \\ \|\mathbf{p}^{t'}\|_0 \leq s/2, \forall t' \in [T] \\ (60) \text{ holds}}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] = \inf_{\substack{(W^{t'}, [n]_{t'})_{t' \in [T]} \\ \hat{\mathbf{p}}^t}} \sup_{\substack{\mathbf{p}^t \in \mathcal{P}_d \\ \|\mathbf{p}^t\|_0 \leq s/2}} \mathbb{E}[\|\hat{\mathbf{p}}^t - \mathbf{p}^t\|_q^q] \\ &= \Omega \left(s^{1-q/2} \left(\frac{\max\{2^b, s\}}{2^b n} \right)^{q/2} \right). \end{aligned}$$

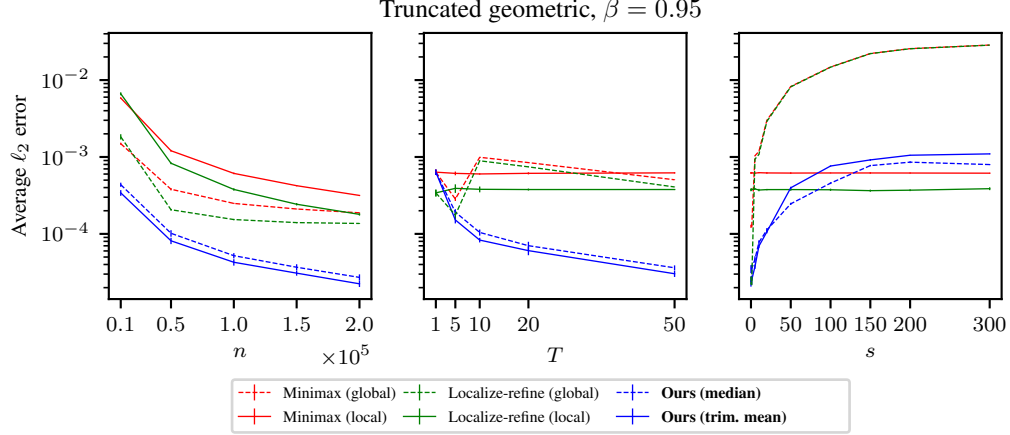


Figure 3: Average ℓ_2 estimation error in synthetic experiment using the truncated geometric distribution. (Left): Fixing $s = 5$, $T = 30$ and varying n . (Middle): Fixing $s = 5$, $n = 100,000$ and varying T . (Right): Fixing $T = 30$, $n = 100,000$ and varying s . The standard error bars are obtained from 10 independent runs.

This proves that the lower bound is at least of the order of the first term in (57). Overall, we conclude the desired result.

13 Supplementary Experiments

Truncated Geometric Distribution. We consider the truncated geometric distribution with parameter $\beta \in (0, 1)$, $\mathbf{p}^* = \frac{1-\beta}{1-\beta^d}(1, \beta, \dots, \beta^{d-1})$, as the central distribution and repeat the experiment in Section 5.1. We use $d = 300$, $\beta = 0.95$, $b = 2$ and vary n, T, s . Figure 3 summarizes the results. As in Section 5.1, we observe that our methods outperform the baseline methods in most cases, especially when s is small. Also, we see the benefit of collaboration, i.e., decreasing trend of the error as T increases, only in our methods.

Hyperparameter Selection. We provide additional experiments using different hyperparameters α and ω from discussed in Section 5.1. All other settings are identical to Section 5.1. We test the hyperparameters $(\alpha, \omega) = (2^r \ln(n), 0.1)$ and $(\alpha, \omega) = (\ln(n), \omega)$, where $r \in \{-5, -4, \dots, 4\}$ and $\omega \in \{0.05, 0.1, \dots, 0.25\}$, respectively. Figure 4 summarizes the results.

We find that setting the threshold α too small leads to replacing almost all coordinates of the central estimate $\hat{\mathbf{p}}^*$ with local ones. In the extreme case of $\alpha \approx 0$, our method is essentially returns the local minimax estimates. On the other hand, we observe that the performance of our method is less sensitive to the trimming proportion ω .

While the choice of α is crucial to the performance of our method, we argue that it is possible to select a reasonably good α by checking the number of fine-tuned entries, i.e.,

$$\frac{1}{T} \sum_{t=1}^T \left| \left\{ k \in [d] : |\hat{\mathbf{b}}^t_k - \hat{\mathbf{b}}^t_k| > \sqrt{\frac{\alpha[\hat{\mathbf{b}}^t]_k}{n}} \right\} \right|.$$

In Figure 5, we observe that more than half ($d/2 = 150$) of the entries are fine-tuned when $r \in \{-5, -4, -3\}$. These correspond to the three curves in the top left of Figure 4 that perform no better than the baseline methods. In conclusion, by selecting α such that the number of fine-tuned entries are small enough compared to d , it is possible to reproduce the results in Section 5.

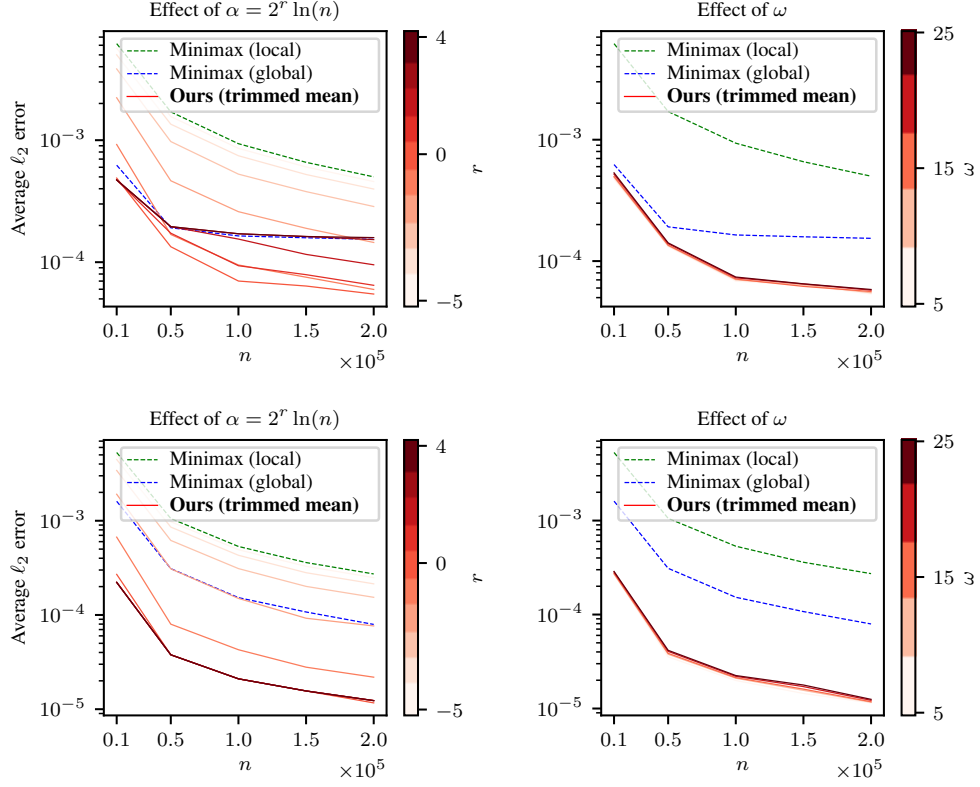


Figure 4: Effect of the hyperparameters α and ω . The top row shows results for the uniform distribution and the bottom row shows the results for the truncated geometric distribution with $\beta = 0.8$.

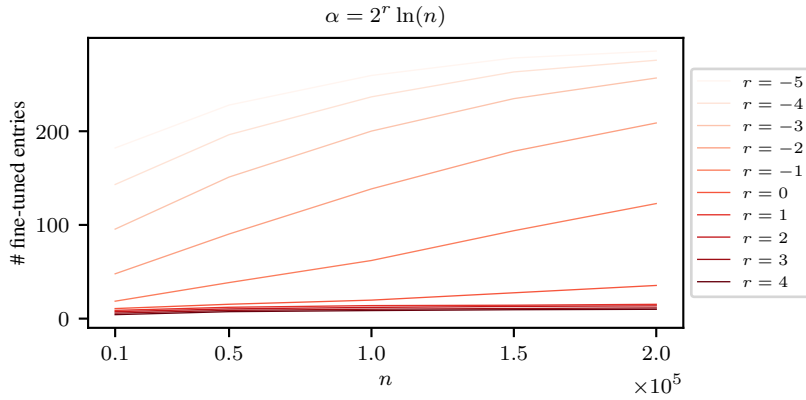


Figure 5: Average number of fine-tuned entries for different values of $\alpha = 2^r \ln(n)$. We use the trimmed mean with $\omega = 0.1$ and the uniform distribution with $d = 300$. This corresponds to the top left of Figure 4.