

# How unfair is private learning ?

Amartya Sanyal<sup>\*1,3</sup>, Yaxi Hu<sup>\*2</sup>, and Fanny Yang<sup>3</sup>

<sup>1</sup>ETH AI Center, ETH Zürich, Zürich, Switzerland.

<sup>2</sup>Department of Mathematics, ETH Zürich, Zürich, Switzerland.

<sup>3</sup>Department of Computer Science, ETH Zürich, Zürich, Switzerland.

## Abstract

As machine learning algorithms are deployed on sensitive data in critical decision making processes, it is becoming increasingly important that they are also private and fair. In this paper, we show that, when the data has a long-tailed structure, it is not possible to build accurate learning algorithms that are both private and results in higher accuracy on minority subpopulations. We further show that relaxing overall accuracy can lead to good fairness even with strict privacy requirements. To corroborate our theoretical results in practice, we provide an extensive set of experimental results using a variety of synthetic, vision (CIFAR-10 and CelebA), and tabular (Law School) datasets and learning algorithms.

## 1 Introduction

In recent years, reliability of machine learning algorithms have become ever more important due to their widespread use in daily life. Fairness and privacy are two instances of such reliability traits that are desirable but often absent in modern machine learning algorithms [13, 17, 47]. As a result, there has been a flurry of recent works that aim to improve these properties in commonly used learning algorithms. However, most of these works discuss these two properties individually with relatively less attention paid to how they affect each other.

There is a multitude of definitions for privacy and fairness in their respective literatures. Perhaps the most widespread statistical notion of privacy is that of *Differential Privacy* [20] and its slightly relaxed variant, *Approximate Differential Privacy* [21]. Despite its marginally weaker privacy guarantees, *Approximate Differential Privacy* enjoys better theoretical guarantees in terms of statistical complexity for learning [10, 26]. It is also more widely used in practice [1, 57]. Thus, we always use approximate differential privacy in this paper and for the sake of brevity, refer to it as differential privacy (DP). Intuitively, DP limits the amount of influence any single data point has on the output of the DP algorithm. This ensures that DP algorithms do not leak information about whether any particular data point was given as input to the algorithm. While DP was initially popular as a theoretical construct, it has recently been put to practical use by large companies [24, 49] and governments [40] alike. Its popularity is largely due to its strong privacy guarantees, ease of implementation, and the quantitative nature of differential privacy.

There are many notions of fairness in machine learning [19, 22, 31, 33]. Minority or worst group accuracy [18, 36] and its difference from the overall accuracy is a common notion of fairness used in recent works. We define this difference as *accuracy discrepancy* and use it to measure the degree of unfairness in this paper. However, we expect our results to translate to other related fairness metrics as well. Sagawa<sup>\*</sup> et al. [45] observed that robust optimisation methods (under appropriate regularisations) obtain higher minority group accuracy but at the cost of a lower overall accuracy compared to vanilla training. Several other works [?] have also observed this behavior in practice thereby suggesting a possible trade-off between overall accuracy and fairness metrics. Subsequent works including Goel et al. [30] and Menon et al. [42] have tried to minimise this trade-off.

While there have been a large body of works that aim to minimising the trade-off between accuracy and fairness [18, 30, 45] and between accuracy and privacy [21, 29], there is relatively few works that investigate the intersection of privacy, fairness, and accuracy. In this paper, we provide theoretical and experimental results to show that private and accurate algorithms are necessarily unfair. We further show that achieving privacy and fairness simultaneously leads to inaccurate algorithms.

**Contributions** Our main contributions can be stated as —

- In Theorem 1 and 2, we provide asymptotic lower bounds for unfairness (accuracy discrepancy) of DP algorithms, that are accurate, showing that privacy and accuracy comes at the cost of fairness.
- In Theorem 3, we show that in a very strict privacy regime, fairness can be achieved at the cost of accuracy.
- In Section 3, we provide experimental results using multiple architectures on synthetic as well as real world datasets (CelebA [39], CIFAR-10 [38], and Law School [55]) to validate our theory.

**Related works** It is now well understood that by imposing these additional conditions of DP more data are required to achieve high accuracy. A string of theoretical works [3, 8, 12, 26] have shown that the sample complexity of learning certain concept classes privately and accurately can be arbitrarily larger than learning the same classes non-privately (i.e. with high accuracy but without privacy). On the other hand, it is easy to guarantee any arbitrary level of differential privacy if high accuracy is not desired. This can be achieved by simply composing the output of an accurate classifier with a properly calibrated *randomised response* mechanism [53]. This allows for a tradeoff between differential privacy and accuracy.

One of the most popular notions of fairness is *group fairness* that compares the performance of the algorithm on a minority group with other groups in the data. A popular instantiation of this, especially for deep learning algorithms, is comparing the accuracy on the minority group against the entire population [18, 30, 36, 45]. In the fairness literature, Buolamwini and Gebru [13], Raz et al. [44] shows extensively that this discrepancy is large between different groups of people for popular facial recognition systems.

DP-SGD [1] is widely used algorithm for implementing differentially private deep learning models. Bagdasaryan et al. [7] provides some experimental evidence that DP-SGD can have disparate impact on accuracy. Conversely, Chang and Shokri [14] shows, experimentally, that fairness aware machine learning algorithms suffer from less privacy. However, unlike these works, we provide theoretical results that are model agnostic and that discuss the dependance of the trade-off on the subpopulation sizes and frequencies.

Cummings et al. [16] and Agarwal [2] were one of the first to consider the impact of privacy on fairness theoretically. They construct a distribution where any algorithm that is always fair and private will necessarily output a trivial constant classifier, thereby suggesting a tradeoff between fairness and privacy. However, there are multiple drawbacks with their work. First, their work only discusses pure differential privacy which is not only theoretically more restrictive than approximate differential privacy [9, 10, 26] but also rarely used in practice. Second, their proof heavily relies on it being *pure* differential privacy and the algorithm being *always* fair; and their proofs are not amenable to relaxations of these assumptions. Further, they do not provide experiments to corroborate their theory perhaps due to the unrealistic requirements of pure DP. On the other hand, we look at approximate DP (which is a stronger result than pure DP), construct bounds for both fairness and error, and provide experimental results to support our theory. Perhaps, most closely related to our work is that of Feldman [25], who studies, mainly, the impact of memorisation on test accuracy for long-tailed distributions. However, neither does their work foray into differential privacy nor into fairness.

## 2 Theoretical results

The main contribution of our work is to provide a simple explanation for why and when differentially private algorithms cannot be simultaneously accurate and fair. Real world data distributions often contain a large number of subpopulations with very few examples in each of them and a few subpopulations with a large

number of examples. Mathematically, the large number of subpopulations with few examples in each of them constitute the *long tail* of the distribution. We use this structure of data distributions to illustrate the tension between accuracy and fairness of private algorithms.

This structure is common in many datasets commonly used in machine learning. For example, in Figure 1 (left), we show this for CelebA. Using the 40 attributes of the CelebA dataset [39], we partition the trainset (of size  $m = 160k$ ) into  $2^{40}$  subpopulations bin. The blue shaded area shows the group of the subpopulations with large number of examples in them (probability mass greater than  $\frac{1}{m}$ ) and the red shaded area corresponds to subpopulations which contain just one example in them. We refer to the subpopulations in the red area as *minority subpopulations* and the subpopulations in the blue area as *majority subpopulations*. We also refer to the group of minority subpopulations as the *minority group* and the group of majority subpopulations as the *majority group*. Zhu et al. [59] observes a similar pattern in other vision datasets like the SUN [56] and PASCAL [23] datasets. Babbar and Schölkopf [5, 6] observes this in extreme multilabel classification datasets like Amazon-670K [41] and Wikipedia-31k [11] datasets. Various other works [15, 37, 50, 52] have shown this in a range of datasets including eBird [48], Visual Genome [37], Pasadena trees [54], and iNaturalist [51].

## 2.1 Problem setup

For our theoretical results, we model this by viewing each subpopulation as an element of a discrete set  $X$  without any intrinsic structure such as distance. Next, we define a distribution over the subpopulations to enforce the *long-tailed* structure.

**Definition 1** ( $(p, N, k)$ -long-tailed distribution on  $X$ ). *Given  $p \in (0, 1)$ ,  $N \in \mathbb{N}$ , and  $1 < k \ll N$ , define two groups (i) the group of majority subpopulations  $X_1 \subset X$  where  $|X_1| = (1 - p)k$  and (ii) the group of minority subpopulations  $X_2 \subset X \setminus X_1$ , where  $|X_2| = N$ .<sup>1</sup> Now, define the distribution  $\Pi_{p,N,k}$  as*

$$\Pi_{p,N,k}(x) = \begin{cases} \frac{1}{k} & x \in X_1 \\ \frac{p}{N} & x \in X_2. \end{cases} \quad (1)$$

We provide an illustration of the distribution in Figure 1 (left). In the rest of the text, we will use the terms *group of majority subpopulations* and *majority group* interchangeably to denote  $X_1$  and the terms *group of minority subpopulations* and *minority group* to refer to  $X_2$  respectively. Intuitively,  $p$  denotes the total probability mass of the group of minority subpopulations under  $\Pi_{p,N,k}$  and  $N$  denotes the number of minority subpopulations. We let  $N$  go to  $\infty$  and treat  $k$  as a constant. Thus, for the sake of simplicity, we remove  $k$  from the notation of the distribution. Note that each minority subpopulation i.e. the element in  $X_2$  has a probability mass of the order of  $O(\frac{1}{N})$  which is much smaller than  $\frac{1}{k} = \Omega(1)$  i.e. the probability mass of the element in  $X_1$ .

Note that the distribution for CelebA (Figure 1 (left)) does not exactly look like the distribution in Definition 1 due to the probability masses of the majority subpopulations not being exactly equal to  $\frac{1}{k}$ . However, all our results hold true even if different majority subpopulations have different probability masses as long as they satisfy  $\Pi_{p,N}(x) = \Omega(\frac{1}{k})$  for some  $k = O(1)$  and all  $x \in X_1$ . We set them to  $\frac{1}{k}$  for simplicity of the theoretical results. This distribution is inspired by the use of a similar distribution in Feldman [25].

As we deal with a multiclass classification setup, we also define a label space  $\mathcal{Y}$  and a function space  $F$  of labelling functions. We also use  $\mathcal{F}$  to represent a distribution on the function space  $F$  and refer to this distribution as the label prior. Our results do not restrict the size of  $\mathcal{Y}$  and can hence, explain both binary and multi-class classification settings.

Finally, to generate a dataset of size  $m$  from a  $(p, N)$ -long-tailed distribution  $\Pi_{p,N}$  on  $X$ , first sample an unlabelled dataset  $S = \{x_1, \dots, x_m\}$  of size  $m$  from  $\Pi_{p,N}$ . Then, generate the labelled dataset  $S_f = \{(x_1, f(x_1)), \dots, (x_m, f(x_m))\}$  using a labelling function  $f \sim \mathcal{F}$ . In all our theoretical results, we consider an asymptotic regime where  $\frac{N}{m} \rightarrow c$  as  $N, m \rightarrow \infty$ . This is common in high-dimensional statistics where the

<sup>1</sup>WLOG we will assume that  $k$  is such that  $(1 - p)k$  is an integer and if not, replace  $k$  with the closest number such that  $(1 - p)k$  is an integer.

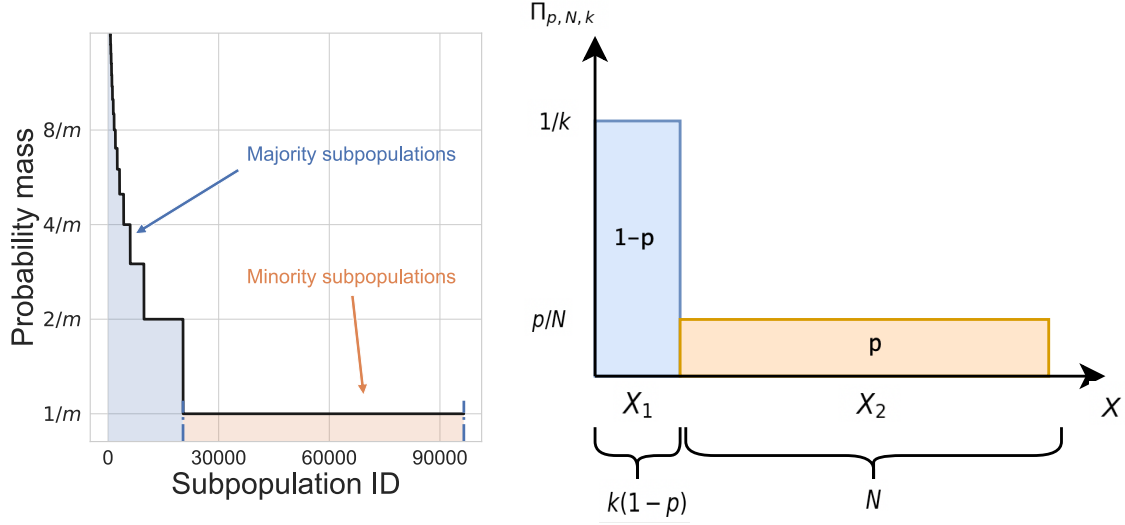


Figure 1: **(Left)** Illustration of the distribution of majority and minority subpopulations of CelebA. Here  $m = 160k$  is the total size of the training set of CelebA. **(Right)** Illustration of  $\Pi_{p,N,k}$ .

number of dimensions often grows to  $\infty$  along with the sample size. Intuitively,  $c$  quantifies the hardness of the learning problem as it is proportional to the number of data points observed per minority subpopulation.

Next, we define the error and fairness measure of an algorithm on the distribution defined above. Consider a domain  $X$ , a label space  $\mathcal{Y}$ , a set of labelling functions  $F$ , a label prior  $\mathcal{F}$ , and a distribution  $\Pi_{p,N}$  on  $X$  as defined above.

## 2.2 Privacy, error, and fairness

In the context of this paper, a differentially private (randomised) learning algorithm generates similar distributions over classifiers when trained on *neighbouring datasets*. Two datasets are neighbouring when they differ in one entry<sup>2</sup>. Formally,

**Definition 2** (Approximate Differential Privacy [20, 21]). *Given any two neighbouring datasets  $S$  and  $S'$  and  $\epsilon > 0, \delta \in (0, 1)$  a randomised algorithm  $\mathcal{A}$  is called  $(\epsilon, \delta)$ -differentially private if for all sets of outputs  $\mathcal{Z}$ , the following holds*

$$\mathbb{P}[\mathcal{A}(S) \in \mathcal{Z}] \leq e^\epsilon \mathbb{P}[\mathcal{A}(S') \in \mathcal{Z}] + \delta.$$

Next, we define the error of an algorithm in our problem setup. For a randomised learning algorithm  $\mathcal{A}$ , a distribution  $\Pi_{p,N}$ , a label prior  $\mathcal{F}$ , we can define the error of the algorithm as follows

**Definition 3** (Error measure on  $\Pi_{p,N}$ ). *The error of the algorithm  $\mathcal{A}$  trained on a dataset of size  $m$  from the distribution  $\Pi_{p,N}$  with respect to a label prior  $\mathcal{F}$  is*

$$\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = \mathbb{E}[\mathbb{I}\{h(x) \neq f(x)\}] \quad (2)$$

where  $\mathbb{I}\{\cdot\}$  is the indicator function and the expectation is over  $S \sim \Pi_{p,N}^m, f \sim \mathcal{F}, h \sim \mathcal{A}(S_f)$ , and  $x \sim \Pi_{p,N}$ .

While it is slightly unconventional to take an expectation over  $\mathcal{F}$ , this was previously used in Feldman [25]. In fact, for the purpose of lower bounds on fairness, this is a stronger notion than the worst case  $f \in F$  as, here, the lower bound is on the expectation which is stronger than a lower bound on the worst case. Next,

<sup>2</sup>However, the definition can be extended to difference in more entries using *group DP* which entails a gradual decline in the privacy guarantees.

Notation	Description
$m$	Size of dataset
$N$	Number of minority subpopulations
$p$	Probability of minority group
$k$	Reciprocal of the probability of individual subpopulations in $X_2$
$c$	Ratio of $N$ and $m$
$\epsilon, \delta$	Privacy parameters of Approximate Differential Privacy
$p_1, p_2$	Parameters used in Assumption A1 and A2 respectively
$\mathcal{A}$	Randomised learning algorithm
$X, X_1, X_2$	Entire data domain, majority, and minority group respectively
$\Pi_{p,N}, \Pi_{p,N}^2$	Data distribution on $X$ and marginal distributions on $X_2$ respectively
$S, S_f$	$m$ -sized unlabelled and labelled (with $f \in F$ ) dataset respectively
$S^\ell$	All points that appear $\ell$ times in $S$
$F, \mathcal{F}$	Set of labelling functions and distribution over the set respectively

Table 1: A table of notations frequently used in the text

we define the *accuracy discrepancy* of an algorithm, represented by  $\Gamma$  over the distribution  $\Pi_{p,N}$ . For any  $\Pi_{p,N}$ , define the marginal distribution on the group of *minority* subpopulations  $X_2$  as

$$\Pi_{p,N}^2(x) = \begin{cases} \frac{\Pi_{p,N}(x)}{\sum_{x \in X_1} \Pi_{p,N}(x)} = \frac{\Pi_{p,N}(x)}{p} & x \in X_2 \\ 0 & x \notin X_2 \end{cases} \quad (3)$$

**Definition 4** (Accuracy discrepancy on  $\Pi_{p,N}$ ). *For  $X, \mathcal{F}, \Pi_{p,N}$ , and  $\Pi_{p,N}^2$  as defined above, the accuracy discrepancy of the algorithm  $\mathcal{A}$  trained on a dataset of size  $m$  on the distribution  $\Pi_{p,N}$ , with respect to the label prior  $\mathcal{F}$ , is*

$$\Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = \text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) - \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}). \quad (4)$$

where  $\text{err}_m(\cdot)$  is as defined in Definition 3.

This notion of group fairness is similar to the notion of subgroup performance gap used in Goel et al. [30] and has also been implicitly used in multiple works Du et al. [18], Koh et al. [36], Sagawa\* et al. [45] as discussed before. It has also been used in works related to the privacy [7, 14] and fairness literature [13, 44].

## 2.3 Privacy and accuracy at the cost of fairness

Next, we state two assumptions that we use in our results. For any  $\ell \in \mathbb{N}$ , define  $S^\ell$  to denote the set of examples that appear exactly  $\ell$  times in  $S$ . Given  $s_0 \in \mathbb{N}$  and  $p_1, p_2 \in (0, 1)$ , we state that  $\mathcal{A}$  satisfies the assumptions A1 and A2 if the following conditions are satisfied by the algorithm  $\mathcal{A}$  for all datasets  $S$ .

### Assumption on algorithm

- (Accuracy) For all  $\ell > s_0$  and  $x \in S^\ell$ ,

$$\mathbb{P}_{f \sim \mathcal{F}, h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \leq p_1 \quad (A1)$$

- (Privacy) For all  $\ell \leq s_0$  and  $x \in S^\ell$ ,

$$\mathbb{P}_{f \sim \mathcal{F}, h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] > 1 - p_2 \quad (A2)$$

Assumption A1 essentially requires the algorithm to have reasonable accuracy. Note that since our domain is discrete, high training accuracy translates to high test accuracy in particular as the sample size approaches infinity. When  $p_1$  is small, the algorithm  $\mathcal{A}$  obtains low training (and hence test) error on frequently occurring or *typical* data (i.e.  $\ell \geq s_0$ ).

On the other hand, Assumption A2 implies that the algorithm is likely to be incorrect on subpopulations that are rare in the training set  $\ell \leq s_0$ . The algorithm errs on them with a probability of at least  $(\geq 1 - p_2)$ . We refer to this assumption as the privacy assumption because for  $(\epsilon, \delta)$ -DP algorithms, it holds true for certain  $s_0, p_2$  that depend on  $\epsilon, \delta$ .

Intuitively, the parameter  $s_0$  represents the smallest frequency (in the observed data) of a subpopulation that is, with high probability, correctly classified by the algorithm. In this paper we distinguish between a small (Theorem 2) and large regime (Theorem 3) for  $s_0$ . Choosing  $s_0 = o(m)$  imposes a high overall accuracy by ensuring that the algorithm is accurate on majority subpopulations by virtue of Assumption A1. On the other hand, in this regime of  $s_0$ , Assumption A2 implies that the algorithm is still likely inaccurate on minority subpopulations. In contrast, large  $s_0 = \Omega(m)$  implies large overall error of the algorithm: Assumption A2 implies that the algorithm even errs on majority subpopulations.

The theoretical results below use the definitions and notations described above and summarised in Table 1. First, in Theorem 1, we show that there are distributions (within the family of distributions defined in Definition 1) where any accurate and approximately differentially private algorithm (with additional assumptions) is necessarily unfair. In Theorem 2, we relax some of the stronger assumptions and present a more general result.

**Theorem 1.** *For  $\epsilon \in (0, 1)$  and  $\delta \in (0, 0.01)$ , consider any  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$  that does not make mistakes on subpopulations occurring more than once in the dataset. Then, there exists a family of label priors  $\mathcal{F}$  where for any  $\alpha \in (0, 1)$ , there exists  $p \in (0, 1/2), c > 0$  such that,*

$$\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq \alpha \quad \text{and} \quad \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq 0.5.$$

where  $\frac{N}{m} \rightarrow c$  as  $N, m \rightarrow \infty$ .

A detailed version along with its full proof is presented in Appendix A.1. First, we note that the assumption in Theorem 1, that the algorithm does not make mistakes on subpopulations that appear more than once is an instantiation of assumption A1 and A2 with the parameters  $s_0 = 1, p_1 = 0$ , and  $p_2 = 1$ . Further, the consequence of unfairness being greater than 0.5, coupled with the very small error  $\alpha$ , is that the algorithm essentially behaves worse than random chance on the minority subpopulations thereby rendering the algorithm useless for these subpopulations.

**Proof sketch** Here we present a proof sketch and discuss the results. By Definition 1, the probability mass of each majority subpopulation is  $\Omega(1)$  whereas the probability mass of each minority subpopulation is  $O(\frac{1}{m})$ . Thus, for a large enough dataset (i.e. large  $m$ ), we show that almost all majority subpopulations appear more than once and consequently, the algorithm in Theorem 1 makes very less mistakes on the majority subpopulations. As a result, the error and the accuracy discrepancy are both caused by mistakes, majorly, on the minority subpopulations. We, then, count the number of subpopulations that do not appear or appear just once among the minority subpopulations and use that to provide the upper bounds for error and lower bound for unfairness (accuracy discrepancy). As these bounds are expressed in terms of  $p$  and  $c$ , the proof then follows by showing the existence of  $p, c$  that satisfy the inequalities in the theorem.

While Theorem 1 shows the existence of distributions under which private and accurate algorithms are necessarily unfair, in Theorem 2, we provide a quantitative lower bound for unfairness of private algorithms. In particular, we present a detailed result showing how unfairness of a DP algorithm varies with respect to the parameters in assumptions A1 and A2, privacy parameters, and distributional parameters. For easier interpretation, we show a simplified version in Theorem 2 and highlight the key takeaways, and provide a detailed version in Appendix A.2. Consider  $X, \mathcal{F}$  as defined in Section 2.1. As discussed above, the dataset size  $m$  and the number of minority subpopulations  $N$  both simultaneously go to  $\infty$  in the ratio  $c = \frac{N}{m}$ .

**Theorem 2.** For any  $p \in (0, 1/2)$ ,  $c > 0$  such that  $p/c \leq 1$ , consider the distribution  $\Pi_{p,N}$  where  $\frac{N}{m} \rightarrow \infty$  as  $N, m$  goes to  $\infty$ . Also, for any  $\epsilon, \delta > 0$ , consider an  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$  that satisfies assumptions A1 and A2 with  $s_0 = o(m)$  and  $p_1, p_2 \in (0, 1)$ . Then,

$$\Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq (1 - p)\gamma_0$$

where  $\gamma_0$  is some constant depending on  $c, s_0, p, \epsilon, \delta$ , and  $\mathcal{F}$ . Further, the error of the algorithm is upper-bounded as  $\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq (1 - p_1)p\alpha_0 + p_1$  and  $\alpha_0$  depends on  $c, s_0$  and  $p$ . As  $c, s_0 \rightarrow \infty$ ,  $\gamma_0$  increases as  $1 - O\left(\frac{e^{-cs_0}}{s_0\sqrt{c}}\right)$ . As  $p \rightarrow 0$ ,  $\alpha_0$  increases as  $1 - O\left(\sqrt{pe}^{-1/p}\right)$ .

The detailed expressions of  $c_0, \alpha_0$ , and  $\gamma_0$  can be found in Theorem 5 in Appendix A.2. We now briefly discuss how the theorem characterizes the effect of privacy and accuracy (via  $s_0$ ) of the algorithm and the distribution of subpopulations (via  $c, p$ ) on the accuracy discrepancy. First note that when the ratio  $c$  of the number of minority subpopulations with respect to the sample size is large, more minority subpopulations are less likely to appear, or appear infrequently, in the observed dataset.

Overall, when  $c, s_0$  are both large, minority subpopulations occur infrequently and most subpopulations are misclassified. Indeed, Theorem 2 indicates that as  $c, s_0$  increase,  $\gamma_0$  and hence unfairness increases, while the average error decreases. Intuitively, this is because minority subpopulations appear infrequently and the algorithm is less likely to memorise infrequent subpopulations. Therefore, the lower bound on unfairness increases with  $c$  and  $s_0$ , and in the extreme case of  $c, s_0 \rightarrow \infty$  approaches  $(1 - p)$ . While the discussion here treats  $c, s_0$ , and  $p$  asymptotically, our results in Appendix A.2 shows non-asymptotic dependence on these terms.

Further, recall that  $p$  quantifies the total probability mass of the group of minority subpopulations (see Figure 1). Hence, for a small  $p$ , error on minority subpopulations do not contribute significantly to the overall error despite causing a disproportionate increase to the marginal error of the minority group. As a result, as  $p$  decreases, Theorem 3 states that the lower bound on unfairness increases as  $1 - p - O\left(\sqrt{pe}^{-1/p}\right)$  while the upper bound for error decreases as  $O\left(p - p^{3/2}e^{-1/p}\right) = O(p)$ .

**Discussion of the assumptions** We now discuss how the privacy parameters  $\epsilon, \delta$  of the DP algorithm lead to feasible parameters  $s_0, p_2$  that appear in Assumptions A1 and A2 used in Theorem 2. We also provide an example of an algorithm that satisfies these assumptions. First of all, Lemma 1 shows that for all values of  $\epsilon, \delta$ , and  $s_0$ , there exists a value of  $p_2$  that satisfies Assumption A2.

**Lemma 1.** Let  $S$  be any dataset. For any  $(\epsilon, \delta)$ -differentially private algorithm  $\mathcal{A}$ , and for all  $s_0 \in \mathbb{N}$ , we have that with high probability over  $\tilde{f} \sim \mathcal{F}$ , for all subpopulations  $x \in X$  that appear fewer than  $s_0$  times in the dataset  $S$ ,

$$\mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})} \left[ h(x) \neq \tilde{f}(x) \right] > 1 - p_2$$

for  $p_2 = \frac{1 + s_0 e^{-\epsilon} \delta}{1 + e^{-s_0 \epsilon}}$ .

Please see Lemma 3 in the Appendix for a detailed version of Lemma 1. First, note that for fixed privacy parameters  $\epsilon$  and  $\delta$ , a small  $s_0$  leads to a small  $p_2$ . Thus, for a fixed privacy level, the algorithm is less likely to be correct on less frequent subpopulations. Similarly, for a fixed  $s_0$ , more privacy, i.e. small  $\epsilon$  and  $\delta$ , also leads to a small  $p_2$ . This means that, for subpopulations of a given frequency, the algorithm is less likely to be correct on those subpopulations with increasing privacy levels. Finally, for a fixed  $p_2$ , there is an inverse relationship between  $s_0$  and  $\epsilon$ , due to which, we view  $s_0$  as a “proxy” for the privacy of the algorithm.

We now provide a differentially private algorithm that satisfies Assumptions A1 and A2. In particular, consider an algorithm  $\mathcal{A}_\eta$  that accepts an  $m$ -sized dataset  $S_f \in (X \times \mathcal{Y})^m$  and a noise rate  $\eta \in (0, \frac{1}{2})$  as input and outputs a dictionary matching every subpopulation in  $X$  to a label in  $\mathcal{Y}$ . The algorithm first creates a dictionary where the set  $X$  is the set of keys. In order to assign values to every key, it first randomly flips the label of every element in  $S_f$  with probability  $\eta$ , then for every unique key in  $S_f$ , the algorithm computes the majority label of that key in the flipped dataset and assigns that majority label to the corresponding key.

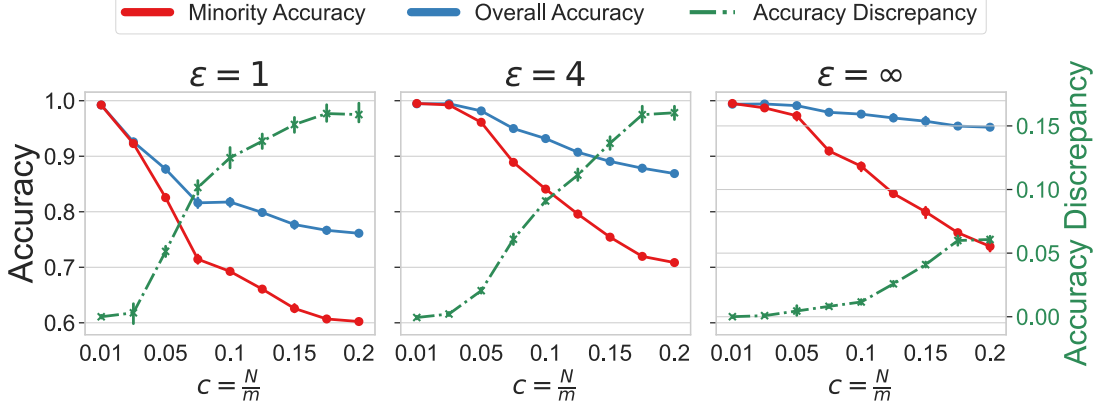


Figure 2: Each figure plots the accuracy discrepancy ( $\Gamma$ ; higher is less fair) in green dashed line, the accuracy of the minority group with red, and the overall accuracy with blue on the y-axis and the parameter  $c$  in the X-axis. The left most ( $\epsilon = 1$ ) achieves the strictest level of privacy and the right most ( $\epsilon = \infty$ ) is vanilla training without any privacy constraints. The two figures in between achieve intermediate levels of privacy. Here  $p = 0.2$ . Experiment for  $p = 0.5$  is in Appendix B.1

For elements in  $X$  not present in  $S_f$ , it assigns a random element from  $\mathcal{Y}$ . Lemma 2 provides privacy and accuracy guarantees for this algorithm.

**Lemma 2.** *The algorithm  $A_\eta$  is  $\left(O\left(\log\left(\frac{1}{\eta}\right)\right), 0\right)$  differentially private as  $\eta \rightarrow 0$ . Further for any dataset  $S_f$  and  $s_0 \in \mathbb{N}$ ,*

- *if a subpopulation  $x$  appears more than  $s_0$  times in  $S$ , then  $\mathbb{P}_{h \sim A_\eta(S_f)}[h(x) \neq f(x)] \leq e^{-s_0(1-2\eta)^2/8(1-\eta)}$  and*
- *if a subpopulation  $x$  appears less than  $s_0$  times in  $S$ , then  $\mathbb{P}_{h \sim A_\eta(S_f)}[h(x) \neq f(x)] \geq (4\eta(1-\eta))^{s_0/2} e^{-s_0}$ .*

*Equivalently, algorithm  $A_\eta$  satisfies Assumption A1 with  $p_1 = e^{-s_0(1-2\eta)^2/8(1-\eta)}$  and Assumption A2 with  $p_2 = 1 - (4\eta(1-\eta))^{s_0/2} e^{-s_0}$ .*

Lemma 1 and 2 are proved in Appendix A.2. Lemma 2 shows that for all  $\epsilon > 0$ , we can find an  $\eta = O(e^{-\epsilon})$  such that  $A_\eta$  is  $(\epsilon, 0)$ -differentially private. Further, this algorithm is more accurate on frequently occurring subpopulations and inaccurate on rare subpopulations, which aligns with Lemma 1. Hence, for any  $\epsilon > 0$ , there is an  $\eta = O(e^{-\epsilon})$  such that the algorithm  $A_\eta$  is  $(\epsilon, 0)$ -differentially private and is accurate on points appearing more than  $\ell$  times with probability  $1 - O(e^{-(1-2\eta)^\ell})$ .

## 2.4 Privacy and fairness at the cost of accuracy

So far we have shown that under strict privacy and high average accuracy requirements on the algorithm, fairness necessarily suffers. A natural question to ask is whether it is possible to sacrifice accuracy for fairness. As discussed in Section 2.3, increasing  $s_0$  leads to higher error – in particular, we consider  $s_0 = \Omega(m)$ .

We present a simplified theorem statement here for easier interpretation and prove a more precise version in Appendix A.3 along with a discussion. In words, the theorem states that for very strict privacy parameters, fairness can be achieved at the cost of accuracy.



**Theorem 3.** For any  $p \in (0, 1/2)$ ,  $c > 0$  such that  $p/c \leq 1$ , consider the distribution  $\Pi_{p,N}$  where  $N$  is the number of minority subpopulations. For any  $\epsilon, \delta, \alpha > 0$ , consider an  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$  that satisfies assumptions (A1) and (A2) with  $s_0 = \left(\frac{2-p}{2k(1-p)}\right)m + \alpha\sqrt{m}$  and some  $p_2 \in (0, 1)$  where  $\frac{N}{m} \rightarrow c$  as both  $m, N \rightarrow \infty$ . Then, ,

$$\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq c_1 p + (1 - p_2)(1 - p) \left(1 - e^{-\frac{4(1-p)\alpha^2}{(2-p)^2}}\right)$$

$$\Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq (1 - p)(1 - p_2) e^{-\frac{4(1-p)\alpha^2}{(2-p)^2}} + p_2$$

where  $c_1$  is a constant depending on  $\epsilon, \delta$ , and  $\mathcal{F}$ .

The detailed expression of  $c_1$  and the full proof can be found in Theorem 6 in Appendix A.3. We now briefly discuss how the theorem characterises the effect of privacy and accuracy (via  $\alpha$ ) on fairness. When  $\alpha$  is small, the algorithm is incorrect only on minority subpopulations. Thus,  $\alpha$ , essentially, characterises what fraction of majority subpopulations the algorithm is incorrect on.

Theorem 3 shows that when  $\alpha$  is large, the error increases and the unfairness decreases. Intuitively, this is because, with increasing  $\alpha$ , the algorithm is incorrect, not only on minority subpopulations, but also on *majority subpopulations* (due to Assumption A2). Thus, as subpopulations of larger frequency gets misclassified due to increasing  $\alpha$ , the overall accuracy as well as the unfairness decreases.

### 3 Experimental results

In this section, we look at experiments to support our theoretical arguments from Section 2. We also note that our theoretical results are model-agnostic and to demonstrate the universality of our result, we conduct a broad set of experiments on both synthetic (in Section 3.1), and real world datasets (in Section 3.2 and 3.3), using multiple machine learning models including deep neural networks and random forests.

#### 3.1 Synthetic experiments

First, we look at a synthetic data distribution that closely emulates the data distribution we use in our theoretical results in Theorem 1 and 2. Given  $N, k \in \mathbb{N}, c \in \mathbb{R}_+$ , and  $p \in (0, 0.5)$ , we construct a continuous version of the long-tailed distribution  $\Pi_{p,N,k}$  (Definition 1) on a domain  $X$ . First of all, since the domain  $X$  is discrete, we can place each element on a vertex of a  $O(\log(N))$ -dimensional hypercube. The continuous distribution we use in our experiments is a mixture of Gaussians where each Gaussian is centered around the vertices of the hypercube. In the experiments, we choose  $k = 64, m = 10^4$ , vary the ratio  $c$  from 0.01 to 0.2, set the number of minority subpopulations to  $N = mc$ , and choose  $p \in \{0.2, 0.5\}$ . We train a five-layer fully connected neural network with ReLU activations using DP-SGD [1] for varying levels of  $\epsilon$  while setting  $\delta = 10^{-3}$ . We refer to Appendix B for a more detailed description of the data distribution and the training algorithm.

**Unfairness aggravates with increasing number of minority subpopulations** As discussed in Section 2.1, increasing the number of subpopulations compared to the number of samples via  $c$  decreases accuracy on the minority subpopulations while the majority subpopulations remain unaffected. Figure 2 shows how increasing  $c$  hurts fairness since the accuracy discrepancy (green dashed line) increases, most pronounced for small values  $\epsilon$  (i.e. more private algorithms). This corroborates our theoretical results from Theorem 2 regarding the dependence of accuracy discrepancy on  $c$ . We further observe that the increase in unfairness is almost entirely due to the drop in the minority accuracy (red solid) whereas the overall accuracy (blue) stays relatively constant. This highlights our claim that, in the presence of strong privacy, fairness can be poor even when overall error is low.

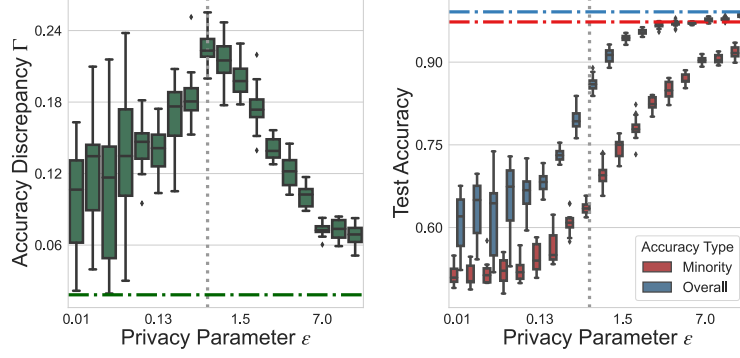


Figure 3: **Left:** Accuracy discrepancy (green) where the box plots reflect the variance when run several times. **Right:** Overall (blue boxes) and minority (red boxes) accuracies for varying  $\epsilon$ . The horizontal dashed line of different colors show the respective metrics for vanilla training without privacy constraints. The gray vertical dashed line marks the privacy parameter for which significant ( $\geq 80\%$ ) overall test accuracy is achieved.

**Privacy constraints hurt fairness for accurate models** In this section, we analyse the dependence of fairness on the privacy parameter  $\epsilon$  for a fixed  $c$ . In Figure 3 (left), we plot the disparate accuracies  $\Gamma$  for varying privacy parameter  $\epsilon$  and Figure 3 (right) depicts the minority and overall accuracy as a function of  $\epsilon$ .

There are two distinct phases in the development of the accuracy discrepancy with increasing  $\epsilon$  separated by the gray dashed line: For a very small  $\epsilon$ , the learned classifier is essentially a trivial classifier as evidenced by the very low overall accuracy ( $\approx 60\%$ ). This is a trivial way of achieving fairness without learning an accurate classifier and is explained by Theorem 3 in our theoretical section. As the privacy restrictions are relaxed, the classifier becomes more accurate and less fair in the first phase.

The interesting regime is when classifier obtains decent overall accuracy ( $\approx 80\%$ ) and is marked by the vertical gray dashed line. In the region to the right of the vertical dashed line, assumption  $m = o(m)$  is fulfilled and Figure 3 (left) reflects the behavior as predicted in Theorem 2: *loosening privacy increases fairness or smaller  $\epsilon$  implies larger accuracy discrepancy*.

### 3.2 Experiments on vision datasets

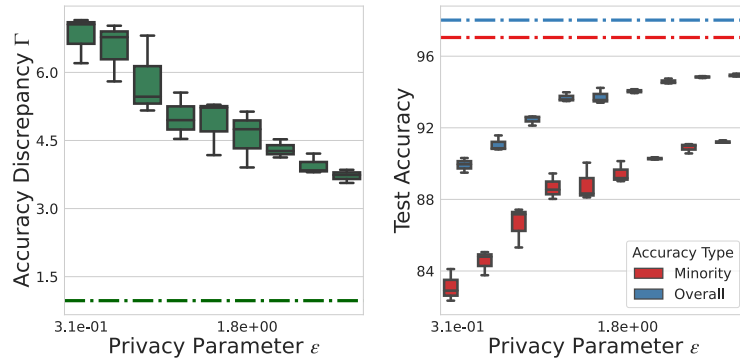


Figure 4: CelebA: **Left:** Accuracy discrepancy (green) where the box plots reflect the variance when run several times. **Right:** Overall (blue boxes) and minority (red boxes) accuracies for varying  $\epsilon$ . The horizontal dashed line of different colors show the respective metrics for vanilla training without privacy constraints.

In this section, we show that our claims resulting from Theorem 2 and 3 do not only hold in synthetic settings

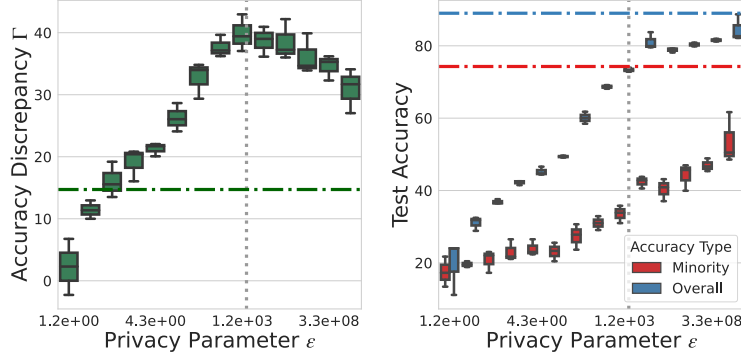


Figure 5: CIFAR-10 **Left:** Accuracy discrepancy (green) where the box plots reflect the variance when run several times. **Right:** Overall (blue boxes) and minority (red boxes) accuracies for varying  $\epsilon$ . The horizontal dashed line of different colors show the respective metrics for vanilla training without privacy constraints. The vertical dashed line marks the  $\epsilon$  for which significant ( $\geq 75\%$ ) overall test accuracy is achieved.

but can also be observed in real-world computer vision datasets. In particular, we conduct experiments on two popular computer vision datasets — CelebA [39] and CIFAR-10 [38]. CelebA is a dataset of approximately 160k training images of dimension  $178 \times 218$  and another 20k of the same dimension for testing. CIFAR-10 is a 10-class classification dataset where there are 50k training images and 10k test images of dimension  $3 \times 32 \times 32$ . For CIFAR-10, we use a ResNet-18 [32] and for CelebA, we use a ResNet-50 architecture.

### 3.2.1 Minority and majority subpopulations

In practice, datasets like CelebA and CIFAR-10 often do not come with a label of what constitutes a subpopulation. In this section, we describe how we define the minority and majority subpopulations for CIFAR-10 and CelebA. Note that the creation of the subpopulation is done before any private training in order to avoid problems like fairness gerrymandering [35].

**CelebA** The CelebA dataset provides 40 attributes for each image including characteristics like gender, hair color, facial hair etc. We create a binary classification problem by using the gender attribute as the target label. In addition, we use 11 of the remaining 39 binary attributes to create  $2^{11}$  subpopulations and categorise each example into one of these  $2^{11}$  subpopulations. Then, we create various groups of minority subpopulations by aggregating the samples of all the unique subpopulations that appear less than  $s \in \{5, 10, 20, 40, 60, 80, 100\}$  times in the test set. The remaining examples constitute the majority group. In this section, we run experiments using  $s = 40$ . We report results for the other values of  $s$  in Appendix B.2.

**CIFAR-10** Unlike the synthetic distribution and CelebA as described above, CIFAR-10 cannot be readily grouped into subpopulations using explicit attributes. However, recent works [46, 58] have shown the presence of subpopulations in CIFAR-10 in the context of influence functions and adversarial training respectively. We use the influence score estimates from Zhang and Feldman [58] to create the minority and majority subpopulations. Intuitively, we treat examples that are atypical i.e. unlike any other examples in the dataset as minority examples belonging to minority subpopulations; and examples that are *typical* i.e. similar to a significant number of other examples in the dataset as examples belonging to majority subpopulations.

To define these subpopulations, first, we sort the examples in the training set according to their self-influence [58]. We define all of those that surpass a threshold  $\rho$  as minority populations. In order to find the samples belonging to each subpopulation  $x$  in the test set, we search for images that are heavily influenced (influence score is greater than the threshold) by at least one of the samples in  $x$  in the training

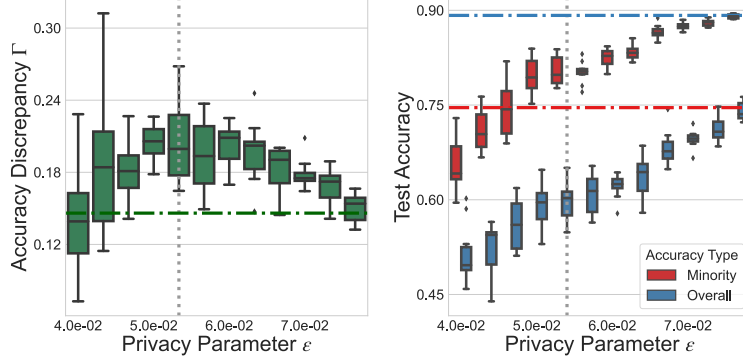


Figure 6: Law School—**Left:** Accuracy discrepancy (green) where the box plots reflect the variance when run several times. **Right:** Overall (blue boxes) and minority (red boxes) accuracies for varying  $\epsilon$ . The horizontal dashed line of different colors show the respective metrics for vanilla training without privacy constraints. The vertical dashed line marks the  $\epsilon$  for which significant ( $\geq 80\%$ ) overall test accuracy is achieved.

set. In this section, we report results with  $\rho = 0.1$ . Other values of  $\rho$  show a similar trend and we plot results using  $\rho = 0.01$  in Appendix B.3.

### 3.2.2 Privacy leads to worse fairness for accurate models

In this section, we use the above definitions of minority and majority groups to measure the impact of privacy on fairness (using Definition 4). Like Section 3.1, we measure both the accuracy discrepancy and the individual minority and overall accuracies.

**CelebA** Figure 4 plots the change in accuracy discrepancy, with respect to the  $\epsilon$  parameter of differential privacy (smaller  $\epsilon$  indicates stricter privacy). Figure 4 (left) shows that smaller  $\epsilon$ , with high accuracy (see Figure 4 (right)) implies a larger accuracy discrepancy. This aligns with our theoretical results from Section 2. Note that unlike Figure 3 (left), the accuracy discrepancy here monotonically decreases with increasing  $\epsilon$  without exhibiting a two-phase behavior. Figure 4 (right) shows that, the reason why we do not observe the two-phase behavior is that throughout the range of observed  $\epsilon$ , we are in the regime of high accuracy.

**CIFAR-10** Figure 5 (left) plots the change of accuracy discrepancy  $\Gamma$  with respect to the privacy parameter  $\epsilon$ . Interestingly, the results here exactly mimic those from the synthetic experiments in Figure 3, which are based on our theoretical setting. This indicates that our theoretical setting is indeed relevant for real world observations. Similar to the synthetic experiments, we observe two distinct phases in how the accuracy discrepancy changes with  $\epsilon$ .

For small values of  $\epsilon$ , Figure 5 (right) shows that the learned classifier is highly inaccurate. As discussed in Theorem 3, this is a trivial way to achieve fairness and this is reflected in Figure 5 (left). However, if we restrict ourselves to classifiers with high average accuracy, marked by the area to the right of the vertical gray dashed line, Figure 5 (left) shows that accuracy discrepancy increases with decreasing  $\epsilon$ . This corresponds to the  $s_0 = o(m)$  assumption, explored in Theorem 2.

## 3.3 Experiments on tabular data

To show that our observations hold across a wider range of publicly used datasets, we next conduct similar experiments using tabular data. We run our experiments on the the Law school dataset [55] that has previously been used in fairness-awareness studies like Quy et al. [43]. It is a binary classification dataset with 21k data points and 12 dimensional features. Out of the 12 attributes, two binary attributes are used to obtain the minority group as defined in Quy et al. [43].

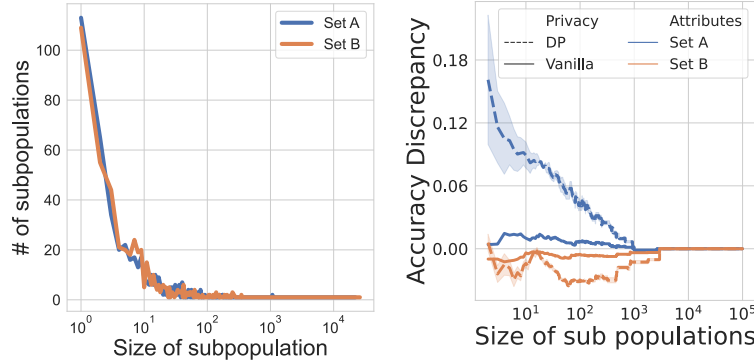


Figure 7: **Left** Both sets of attributes induces a similar distribution over sizes of subpopulations on CelebA. **Right** Set A has high accuracy discrepancy for small sized subpopulations whereas Set B does not.

In contrast to previous experiments, we use random forest model from Fletcher and Islam [28] instead of neural networks as in Section 3.1 and 3.2. For our implementation, we use the publicly available code in Holohan et al. [34] with 10 trees and of a maximum depth 50. The results are plotted in Figure 6 and they show a similar trend as in our previous experiments. The gray vertical dotted line indicates the smallest  $\epsilon$  for which an overall accuracy of 80% is achieved. The results in Figure 6 (left) shows a similar two-phase behavior as previously observed in CIFAR-10 with the highest accuracy discrepancy at the 80% overall accuracy mark.

Thus, all our experiments provide empirical evidence in support of the theoretical arguments in Section 2. The behavior is consistent across multiple kinds of datasets, machine learning models, and learning algorithms.

## 4 Future work

The experimental results on CelebA in Section 3.2 shows that when the minority group is composed of small sized subpopulations, differential privacy requirements hurt the fairness of the algorithm. Here, we highlight that not all small-sized subpopulations are hurt equally in this process. In Figure 7, we show that a different partition of CelebA composed of similar sized populations do not show similar behaviours in terms of how accuracy discrepancy changes with sizes of subpopulations. We will refer to the 11 attributes we chose to partition the testset for our experiments so far as *Set A* and here we choose another set of 11 attributes and refer to them as *Set B*. Figure 7 (left) shows that both Set A and Set B induces a very similar distribution over sizes of subpopulations on the test set. However, Figure 7 (right) shows that while the group of minority subpopulations induced by Set A suffers very high accuracy discrepancy from private training compared to vanilla training, Set B shows very minor difference in accuracy discrepancy between private and vanilla training (see Appendix B.2 for more details on Set A and Set B). This indicates that, irrespective of sizes, private training hurts fairness disproportionately more for certain subpopulations compared to others. In particular, an interesting direction of further research is to investigate where these minority subpopulations that are worse-affected by private training intersects with the subpopulations that are relevant for the specific domain. While most past works [7, 14] have also used sizes of subpopulations to differentiate between disparately impacted subpopulations, this suggests that that is not always the case. In future we would like to understand the assumptions, for which, certain subpopulations are worse impacted than others.

In this paper, we have shown theoretically that when the minority group in the data is composed of multiple subpopulations, a DP algorithm can achieve very low error but necessarily incurs worse fairness. Further, we corroborated our theoretical results with experimental evidence on synthetic and real world computer vision datasets. However, our model-agnostic results, that shed a rather pessimistic light on algorithmic fairness and differential privacy, only apply under certain distributional assumptions. It is possible that in some real-world datasets there are fair and private algorithms that achieve a more optimistic trade-off. This begs further research to develop fair and private algorithms that are closer to the pareto optimal frontier.

## 5 Acknowledgements

AS is partially supported by the ETH AI Center postdoctoral fellowship. AS also acknowledges the support of Hasler Stiftung.

## References

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security*, 2016.
- [2] S. Agarwal. Trade-offs between fairness and privacy in machine learning. 2021.
- [3] N. Alon, R. Livni, M. Malliaris, and S. Moran. Private PAC learning implies finite Littlestone dimension. In *ACM Symposium on Theory of Computing*, 2019.
- [4] R. B. Ash. *Information theory*. 1967.
- [5] R. Babbar and B. Schölkopf. Dismec: Distributed sparse machines for extreme multi-label classification. In *Proceedings of the tenth ACM international conference on web search and data mining*, 2017.
- [6] R. Babbar and B. Schölkopf. Data scarcity, robustness and extreme multi-label classification. *Machine Learning*, 2019.
- [7] E. Bagdasaryan, O. Poursaeed, and V. Shmatikov. Differential privacy has disparate impact on model accuracy. In *Conference on Neural Information Processing Systems*, 2019.
- [8] A. Beimel, S. P. Kasiviswanathan, and K. Nissim. Bounds on the sample complexity for private learning and private data release. In *Theory of Cryptography*. 2010.
- [9] A. Beimel, K. Nissim, and U. Stemmer. Characterizing the sample complexity of private learners. In *Innovations in Theoretical Computer Science*, 2013.
- [10] A. Beimel, K. Nissim, and U. Stemmer. Private learning and sanitization: Pure vs. approximate differential privacy. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*. 2013.
- [11] K. Bhatia, K. Dahiya, H. Jain, P. Kar, A. Mittal, Y. Prabhu, and M. Varma. The extreme classification repository: Multi-label datasets and code, 2016.
- [12] M. Bun, R. Livni, and S. Moran. An equivalence between private classification and online prediction. In *Annual Symposium on Foundations of Computer Science*, 2020.
- [13] J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, 2018.
- [14] H. Chang and R. Shokri. On the privacy risks of algorithmic fairness. In *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2021.
- [15] Y. Cui, Y. Song, C. Sun, A. Howard, and S. Belongie. Large scale fine-grained categorization and domain-specific transfer learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [16] R. Cummings, V. Gupta, D. Kimpara, and J. Morgenstern. On the compatibility of privacy and fairness. In *27th Conference on User Modeling, Adaptation and Personalization - (UMAP)*, 2019.

- [17] T. De Vries, I. Misra, C. Wang, and L. Van der Maaten. Does object recognition work for everyone? In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [18] M. Du, S. Mukherjee, G. Wang, R. Tang, A. Awadallah, and X. Hu. Fairness via representation neutralization. In *Conference on Neural Information Processing Systems*, 2021.
- [19] J. C. Duchi and H. Namkoong. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 2021.
- [20] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor. Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology - EUROCRYPT*. 2006.
- [21] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. *Theory of Cryptography*, 2006.
- [22] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel. Fairness through awareness. In *Innovations in Theoretical Computer Science*, 2012.
- [23] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes (VOC) challenge. *International Journal of Computer Vision*, 2010.
- [24] G. Fanti, V. Pihur, and Ú. Erlingsson. Building a RAPPOR with the unknown: Privacy-preserving learning of associations and data dictionaries. *Proceedings on Privacy Enhancing Technologies (PoPETs)*, 2016.
- [25] V. Feldman. Does learning require memorization? A short tale about a long tail. In *ACM Symposium on Theory of Computing*, 2019.
- [26] V. Feldman and D. Xiao. Sample complexity bounds on differentially private learning via communication complexity. In *Conference on Learning Theory*, 2014.
- [27] V. Feller. *An introduction to probability theory and its applications. Vol. 1*. J. Wiley, 1970.
- [28] S. Fletcher and M. Z. Islam. Differentially private random decision forests using smooth sensitivity. *Expert systems with applications*, 2017.
- [29] B. Ghazi, N. Golowich, R. Kumar, P. Manurangsi, and C. Zhang. Deep learning with label differential privacy. In *Conference on Neural Information Processing Systems*, 2022.
- [30] K. Goel, A. Gu, Y. Li, and C. Ré. Model patching: Closing the subgroup performance gap with data augmentation. In *International Conference on Learning Representations*, 2021.
- [31] M. Hardt, E. Price, and N. Srebro. Equality of opportunity in supervised learning. In *Conference on Neural Information Processing Systems*, 2016.
- [32] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [33] U. Hébert-Johnson, M. Kim, O. Reingold, and G. Rothblum. Multicalibration: Calibration for the (computationally-identifiable) masses. In *International Conference on Machine Learning*, 2018.
- [34] N. Holohan, S. Braghin, P. Mac Aonghusa, and K. Levacher. Diffprivlib: the IBM differential privacy library. 2019.
- [35] M. Kearns, S. Neel, A. Roth, and Z. S. Wu. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In *International Conference on Machine Learning*, 2018.

- [36] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga, R. L. Phillips, I. Gao, et al. WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, 2021.
- [37] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 2017.
- [38] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [39] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [40] G. Long. Formal privacy methods for the 2020 census. 2020.
- [41] J. McAuley and J. Leskovec. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, 2013.
- [42] A. K. Menon, A. S. Rawat, and S. Kumar. Overparameterisation and worst-case generalisation: friend or foe? In *International Conference on Learning Representations*, 2021.
- [43] T. L. Qu, A. Roy, V. Iosifidis, W. Zhang, and E. Ntoutsis. A survey on datasets for fairness-aware machine learning. *WIREs Data Mining and Knowledge Discovery*, 2022.
- [44] D. Raz, C. Bintz, V. Guetler, A. Tam, M. Katell, D. Dailey, B. Herman, P. M. Krafft, and M. Young. Face Mis-ID: An interactive pedagogical tool demonstrating disparate accuracy rates in facial recognition. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, 2021.
- [45] S. Sagawa\*, P. W. Koh\*, T. B. Hashimoto, and P. Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020.
- [46] A. Sanyal, P. K. Dokania, V. Kanade, and P. H. S. Torr. How benign is benign overfitting? In *International Conference on Learning Representations*, 2021.
- [47] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, 2017.
- [48] B. L. Sullivan, C. L. Wood, M. J. Iliff, R. E. Bonney, D. Fink, and S. Kelling. ebird: A citizen-based bird observation network in the biological sciences. *Biological conservation*, 2009.
- [49] J. Tang, A. Korolova, X. Bai, X. Wang, and X. Wang. Privacy loss in Apple’s implementation of differential privacy on MacOS 10.12. 2017.
- [50] G. Van Horn and P. Perona. The devil is in the tails: Fine-grained classification in the wild. *arXiv:1709.01450*, 2017.
- [51] G. Van Horn, O. Mac Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and S. Belongie. The inaturalist species classification and detection dataset. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [52] Y.-X. Wang, D. Ramanan, and M. Hebert. Learning to model the tail. In *Conference on Neural Information Processing Systems*, 2017.
- [53] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 1965.



- [54] J. D. Wegner, S. Branson, D. Hall, K. Schindler, and P. Perona. Cataloging public objects using aerial and street-level images-urban trees. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [55] L. F. Wightman. LSAC national longitudinal bar passage study. lsac research report series. 1998.
- [56] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2010.
- [57] A. Yousefpour, I. Shilov, A. Sablayrolles, D. Testuggine, K. Prasad, M. Malek, J. Nguyen, S. Ghosh, A. Bharadwaj, J. Zhao, G. Cormode, and I. Mironov. Opacus: User-friendly differential privacy library in PyTorch. *arXiv preprint arXiv:2109.12298*, 2021.
- [58] C. Zhang and V. Feldman. What neural networks memorize and why: Discovering the long tail via influence estimation. In *Conference on Neural Information Processing Systems*, 2020.
- [59] X. Zhu, D. Anguelov, and D. Ramanan. Capturing long-tail distributions of object subcategories. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.

## A Proofs

### A.1 Proof for Theorem 1

**Theorem 4** (Detailed version of Theorem 1). *For any  $\alpha \in (0, 1)$ ,  $\epsilon \in (0, 1)$ ,  $\delta \in (0, 0.01)$ , consider any  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$  that does not make mistakes on points occurring more than once in the dataset. Then, there exists  $p \in (0, 1/2)$ ,  $c > 0$  such that*

$$\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq \alpha$$

and

$$\Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq 0.5$$

where  $N/m \rightarrow c$  as  $N, m \rightarrow \infty$  and  $\mathcal{F}$  belongs to a family of label priors such that  $\max_{x \in X_2} \|\mathcal{F}(x)\|_\infty \leq 0.1$  and  $\min_{x \in X_2} \min_{y \in \mathcal{Y}} \mathbb{P}_{f \sim \mathcal{F}}[f(x) = y] \geq 0.05$ .

*Proof.* Recall the definition of the majority and minority subpopulations  $X_1$  and  $X_2$  from Definition 4. Given a dataset  $S$ , define  $S_1$  to be the partition of  $S$  that belongs to  $X_1$  and  $m_1 = |S_1|$  to be the number of examples in  $S$  that belong to the majority subpopulation. Similarly, define  $S_2$  and  $m_2 = m - m_1$  as the set of minority examples and the size of the set of minority examples respectively. We also use  $S_i^\ell$  to denote the set of  $x \in X$  that appears  $\ell$  times in  $S_i$ .

First, we expand the expression for error, defined in Definition 3 as follows.

$$\begin{aligned} \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \mathbb{E}_{h \sim \mathcal{A}(S_f), f \sim \mathcal{F}} \left[ \sum_{x \in X} \Pi_{p,N}(x) \mathbb{1}\{h(x) \neq f(x)\} \right] \\ &= \mathbb{E}_{f \sim \mathcal{F}} \left[ \sum_{\ell=0, i=0}^{\ell=m, i=1} \sum_{x \in S_{i+1}^\ell} \Pi_{p,N}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] \\ &= \mathbb{E}_{f \sim \mathcal{F}} \left[ \sum_{\ell, i=0}^1 \sum_{x \in S_{i+1}^\ell} \Pi_{p,N}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] \\ &\leq \frac{1}{k} \mathbb{E}[|S_1^0| + |S_1^1|] + \frac{c_1 p}{N} \mathbb{E}[|S_2^0| + |S_2^1|] \end{aligned}$$

where  $\mathbb{P}_{f \sim \mathcal{F}, h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \leq \max_{x \in X_2} \max_{y \in \mathcal{Y}} \mathbb{P}_{f \sim \mathcal{F}}[f(x) \neq y] := c_1 \leq 0.1$ . To see why, note that this upper bound can be achieved by an algorithm that ignores the labels in the dataset  $S_f$  and returns a deterministic classifier that predicts a fixed label for an example, possibly different for different examples. Similarly, we decompose the expression of accuracy discrepancy, defined in Definition 4, as follows

$$\begin{aligned} \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) - \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\ &= (1-p) [\text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) - \text{err}_m(\mathcal{A}, \Pi_{p,N}^1, \mathcal{F})] \\ &= (1-p) \sum_{\ell, i=0}^1 \mathbb{E}_{f \sim \mathcal{F}} \left[ (-1)^{i+1} \sum_{x \in S_{i+1}^\ell} \Pi_{p,N}^{i+1}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] \\ &\geq \frac{c_2(1-p)}{N} \mathbb{E}[|S_2^1| + |S_2^0|] - \frac{(1-p)}{k} \mathbb{E}[|S_1^1| + |S_1^0|] \end{aligned}$$

where we define  $c_2$  as follows. Let  $\tilde{f}$  be sampled from a class of functions that only differ from  $f$  at  $x$ , then  $S_{\tilde{f}}$  and  $S_f$  are neighboring datasets differing only at  $x$ . As  $\mathcal{A}$  is  $(\epsilon, \delta)$ -differentially private, for all  $S \sim (\Pi_{p,N})^m$

and for all  $x \in S_2$ , we have the following inequality by the definition of differential privacy

$$\begin{aligned}
\mathbb{P}_{h \sim \mathcal{A}(S, f)}[h(x) \neq f(x)] &= 1 - \mathbb{P}_{h \sim \mathcal{A}(S, f)}[h(x) = f(x)] \\
&\geq 1 - e^\epsilon \mathbb{P}_{h \sim \mathcal{A}(S, \bar{f})}[h(x) = f(x)] - \delta \\
&\geq \min_{x \in X_2} 1 - e^\epsilon \max_{y \in \mathcal{Y}} \mathbb{P}_{f \sim \mathcal{F}}[f(x) = y] - \delta \\
&= \min_{x \in X_2} 1 - e^\epsilon \|\mathcal{F}(x)\|_\infty - \delta := c_2
\end{aligned} \tag{5}$$

where  $\|\mathcal{F}(x)\|_\infty = \max_{y \in \mathcal{Y}} \mathbb{P}_{f \sim \mathcal{F}}[f(x) = y]$ . Now, we will bound the terms  $|S_1^0|, |S_1^1|, |S_2^0|, |S_2^1|$  individually to obtain the relevant upper and lower bounds. First, we define the following random event over the sampling of the  $m$ -sized dataset from  $\Pi_{p, N}$ .

$$\mathcal{E} = \left\{ \frac{p}{2} \leq \frac{m_2}{m} \leq \frac{1+p}{2} \right\}$$

As  $m_2$  is a binomial distribution with parameter  $(m, p)$ , as  $m \rightarrow \infty$  it could be estimated with a Gaussian distribution with mean  $mp$  and variance  $mp(1-p)$  by Central Limit Theorem (CLT). Using Chebyshev's inequality and the assumption  $p < \frac{1}{2}$ , we lower bound the probability of  $\mathcal{E}$  as follows:

$$\begin{aligned}
\mathbb{P}[\mathcal{E}] &= 1 - \mathbb{P}\left[m_2 \leq \frac{mp}{2}\right] - \mathbb{P}\left[m_2 \geq \left(\frac{1+p}{2}\right)m\right] \\
&\geq 1 - \mathbb{P}\left[|m_2 - \mathbb{E}(m_2)| \geq \sqrt{\frac{mp}{4(1-p)} \text{Var}(m_2)}\right] \\
&= 1 - o\left(\frac{1}{m^2}\right) \rightarrow 1 \quad \text{as } m \rightarrow \infty.
\end{aligned} \tag{6}$$

Then, by law of total expectation, we have that

$$\begin{aligned}
\lim_{m \rightarrow \infty} \mathbb{E}[|S_1^0|] &= \lim_{m \rightarrow \infty} \mathbb{E}[|S_1^0| | \mathcal{E}] \mathbb{P}[\mathcal{E}] + \lim_{m \rightarrow \infty} \underbrace{\mathbb{E}[|S_1^0| | \mathcal{E}^c]}_{O(m)} \underbrace{\mathbb{P}[\mathcal{E}^c]}_{o(\frac{1}{m^2})} \\
&= \lim_{m \rightarrow \infty} \sum_{x \in X_1} \mathbb{P}_{S_1 \sim (\Pi_{p, N}^1)^{m_1}} [x \text{ occurs 0 times in } S_1 | \mathcal{E}] \mathbb{P}[\mathcal{E}] \\
&\leq \lim_{m \rightarrow \infty} k(1-p) \left(1 - \frac{1}{k(1-p)}\right)^{\frac{1-p}{2}m} \mathbb{P}[\mathcal{E}] = 0
\end{aligned} \tag{7}$$

where the last step follows because the event  $\mathcal{E}$  requires  $m_2 \geq \frac{(1+p)m}{2}$  and  $\frac{1}{k(1-p)} \in (0, 1)$ . It follows from similar arguments that  $\lim_{m \rightarrow \infty} \mathbb{E}[|S_1^1|] = 0$ .

Next, we compute  $\mathbb{E}[S_2^0]$  and  $\mathbb{E}[S_2^1]$ . By simple counting argument, we have that

$$\mathbb{E}[S_2^\ell] = \sum_{x \in X_2} \mathbb{P}_{S_2 \sim (\Pi_{p, N}^2)^{m_2}} [x \text{ occurs } \ell \text{ times in } S_2] = N \binom{m_2}{\ell} \left(\frac{1}{N}\right)^\ell \left(1 - \frac{1}{N}\right)^{m_2 - \ell} \tag{8}$$

Plugging in  $\ell = 1$  and  $\ell = 0$ , we get

$$\mathbb{E}[|S_2^1|] = m_2 \left(1 - \frac{1}{N}\right)^{m_2 - 1} \quad \text{and} \quad \mathbb{E}[|S_2^0|] = N \left(1 - \frac{1}{N}\right)^{m_2} \tag{9}$$

respectively. Plugging these expressions back into the expression of error in Equation (5), we obtain the

following upper bound

$$\begin{aligned}
\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \lim_{m, N \rightarrow \infty} \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\
&\leq \underbrace{\frac{1}{k} \lim_{m \rightarrow \infty} \mathbb{E}[|S_1^0| + |S_1^1|]}_{\rightarrow 0} + \lim_{m, N \rightarrow \infty} \frac{c_1 p}{N} \mathbb{E}[|S_2^0| + |S_2^1|] \\
&= \lim_{m, N \rightarrow \infty} c_1 p \left( \frac{(1+p)m}{2N} \left(1 - \frac{1}{N}\right)^{pm/2-1} + \left(1 - \frac{1}{N}\right)^{pm/2} \right) \\
&= c_1 p \left( \frac{(1+p)}{2c} e^{p/2c} + e^{-p/2c} \right)
\end{aligned} \tag{10}$$

where the last step follows from limit rules because  $\frac{N}{m} \rightarrow c$  as  $m, N \rightarrow \infty$ .

Similarly, we simplify the expression of accuracy discrepancy from Equation (5).

$$\begin{aligned}
\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \lim_{m, N \rightarrow \infty} \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\
&\geq \lim_{m, N \rightarrow \infty} \frac{c_2(1-p)}{N} \mathbb{E}[|S_2^1| + |S_2^0|] - \underbrace{\lim_{m, N \rightarrow \infty} \frac{(1-p)}{k} \mathbb{E}[|S_1^1| + |S_1^0|]}_{=0} \\
&\geq c_2(1-p) \left( \frac{p}{2c} e^{-(1+p)/2c} + e^{-(1+p)/2c} \right)
\end{aligned} \tag{11}$$

Next, we show that for any  $\alpha \in [0, 1]$ , there is some  $p, c$  such that  $\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq \alpha$  and  $\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq 0.5$ . Setting  $c = 10, p = 0.26$ , the minority group accuracy discrepancy evaluates to  $\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq 0.5$ . For  $p \in (0, 0.26)$ , the lower bound of accuracy discrepancy increases with decreasing  $p$  as its derivative is negative. Therefore,  $\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq 0.5$  is always satisfied for all  $p \leq 0.26$  and  $c = 10$ .

Now, we show that for  $\alpha \in [0, 1]$ , there exists a distribution  $\Pi_{p,N}$  with  $p \leq 0.26, c = 10$  such that  $\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq \alpha$ . Specifically, we first argue that the upper bound of error monotonically increases with  $p$  for  $p \in [0, 0.26]$  for  $c = 10$ , then we show that the upper bound achieves its maximum and minimum values at 0 and 0.26 respectively and invoke the intermediate value theorem to complete the proof.

Taking the derivative of the upper bound of the error, for  $p \in [0, 0.26]$  and  $c = 10$ , we observe that the upper bound on error monotonically increases from  $p = 0$  to  $p = 0.26$  for  $c = 10$ . Plugging in  $p = 0.26$  and  $p = 0$ , we obtain the maximum and minimum values to be 0.259 and 0 respectively. Invoking the intermediate value theorem on this completes the proof.  $\square$

## A.2 Proof for Theorem 2

In this section, we will look at the detailed version and the proof of Theorem 1.

**Theorem 5** (Detailed version of Theorem 2). *For any  $p \in (0, \frac{1}{2})$ ,  $c > 0$  such that  $\frac{p}{c} \leq 1$ , consider the distribution  $\Pi_{p,k,N}$  where  $N/m \rightarrow \infty$  as the number of minority subpopulations  $N$  and the sample size  $m$  goes to  $\infty$  and  $k$  is a constant satisfying  $k = O(1)$  and  $k(1-p) \geq 2$ . For any  $\epsilon, \delta > 0$ , consider an  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$  that satisfies assumptions A1 and (A2) with  $s_0 = o(m)$  and some  $p_1, p_2 \in (0, 1)$ . Then,*

$$\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq (1-p_1)p \left( 1 - \frac{(s_0 - \frac{p}{2c})^2 \frac{2c}{p} - 1}{\sqrt{2\pi}(s_0 - \frac{p}{2c})^3 \left(\sqrt{\frac{2c}{p}}\right)^3} e^{-(s_0 - \frac{p}{2c})^2 \frac{c}{p}} \right) + p_1$$

and

$$\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq c_3(1-p) \left( 1 - \frac{\sqrt{\frac{1+p}{2c}}}{\sqrt{2\pi}(s_0 - \frac{1+p}{2c})} e^{-(s_0 - \frac{1+p}{2c})^2 \frac{c}{1+p}} \right) - (1-p)p_1$$

where  $c_3 = \min_{x \in X} 1 - e^{s_0 \epsilon} \|\mathcal{F}(x)\|_\infty - s_0 e^{(s_0-1)\epsilon} \delta$  and  $\|\mathcal{F}(x)\|_\infty = \max_{y \in \mathcal{Y}} [f(x) = y]$ . Further,  $\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = \lim_{m, N \rightarrow \infty} \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F})$  and  $\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = \lim_{m, N \rightarrow \infty} \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F})$ .

*Proof.* We use the same extra notation for  $\{m_i, S_i, S_i^j\}$  for  $i \in \{1, 2\}$  and  $j \in \mathbb{N}$ ,  $j \leq m$  as in the proof of Theorem 3. First, we decompose the terms for error. Note that all expectations here is with respect to  $S \sim \Pi_{p,N}^m$  and  $f \sim \mathcal{F}$ .

$$\begin{aligned}
\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \mathbb{E} \left[ \sum_{\ell=0}^{s_0} \sum_{x \in S^\ell} \Pi_{p,N}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] + \sum_{\ell=s_0+1}^m \sum_{x \in S^\ell} \Pi_{p,N}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] \right] \\
&\leq \sum_{\ell=0}^{s_0} \mathbb{E} \left[ \sum_{x \in S^\ell} \Pi_{p,N}(x) \right] + p_1 \mathbb{E} \left[ \sum_{\ell=0}^m \sum_{x \in S^\ell} \Pi_{p,N}(x) \right] - p_1 \sum_{\ell=0}^{s_0} \mathbb{E} \left[ \sum_{x \in S^\ell} \Pi_{p,N}(x) \right] \\
&\stackrel{(a)}{=} (1 - p_1) \sum_{\ell=0}^{s_0} \mathbb{E} \left[ \sum_{x \in S^\ell} \Pi_{p,N}(x) \right] + p_1 \\
&= (1 - p_1) \sum_{\ell=0}^{\ell=s_0} (\mathbb{E}_{S \sim (\Pi_{p,N})^m} [|S_1^\ell| + |S_2^\ell|]) + p_1
\end{aligned} \tag{12}$$

where in step (a), we get rid of the second expectation using the fact  $\sum_{\ell=0}^m \sum_{x \in S^\ell} \Pi_{p,N}(x) = 1$ . Then, recalling the decomposition of accuracy discrepancy into the error of the minority and majority groups,

$$\Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = (1 - p) (\text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) - \text{err}_m(\mathcal{A}, \Pi_{p,N}^1, \mathcal{F})), \tag{13}$$

we expand the two terms as follows.

$$\begin{aligned}
\text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) &= \mathbb{E} \left[ \sum_{\ell=0}^{m_2} \sum_{x \in S_2^\ell} \Pi_{p,N}^2(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] \right] \\
&\geq \sum_{\ell=1}^{s_0} \mathbb{E} \left[ \sum_{x \in S_2^\ell} \Pi_{p,N}^2(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] \right] \\
&\geq \frac{c_3}{N} \sum_{\ell=0}^{s_0} \mathbb{E} [|S_2^\ell|]
\end{aligned} \tag{14}$$

where we define  $c_3$  as follows using group differential privacy. Let  $\tilde{f}$  be sampled from a class of functions that only differ from  $f$  at  $x \in S^\ell$ , then  $S_f$  and  $S_{\tilde{f}}$  are neighboring datasets differing in  $\ell$  entries. As  $\mathcal{A}$  is  $(\epsilon, \delta)$ -differentially private, we have the following inequality using group differential privacy for all  $x \in \bigcup_{\ell \leq s_0} S^\ell$ .

$$\begin{aligned}
\mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] &= 1 - \mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) = f(x)] \\
&\geq 1 - e^{s_0 \epsilon} \mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})} [h(x) = f(x)] - \ell e^{(\ell-1)\epsilon} \delta \\
&\geq 1 - e^{s_0 \epsilon} \max_{y \in \mathcal{Y}} \mathbb{P}_{f \sim \mathcal{F}} [f(x) = y] - s_0 e^{(s_0-1)\epsilon} \delta \\
&\geq \min_{x \in X} 1 - e^{s_0 \epsilon} \|\mathcal{F}(x)\|_\infty - s_0 e^{(s_0-1)\epsilon} \delta := c_3
\end{aligned} \tag{15}$$

where  $\|\mathcal{F}(x)\|_\infty = \max_{y \in \mathcal{Y}} \mathbb{P}_{f \sim \mathcal{F}}[f(x) = y]$ . Next, we upper bound the error of the majority group.

$$\begin{aligned}
\text{err}_m(\mathcal{A}, \Pi_{p,N}^1, \mathcal{F}) &= \mathbb{E} \left[ \sum_{\ell=0}^{s_0} \sum_{x \in S_1^\ell} \Pi_{p,N}^1(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] + \mathbb{E} \left[ \sum_{\ell=s_0+1}^{m_1} \sum_{x \in S_1^\ell} \Pi_{p,N}^1(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] \\
&\stackrel{(a)}{\leq} \mathbb{E} \left[ \sum_{\ell=0}^{s_0} \sum_{x \in S_1^\ell} \Pi_{p,N}^1(x) \right] + p_1 \mathbb{E} \left[ \sum_{\ell=s_0+1}^{m_1} \sum_{x \in S_1^\ell} \Pi_{p,N}^1(x) \right] \\
&\stackrel{(b)}{\leq} \frac{1}{k(1-p)} \mathbb{E} \left[ \sum_{\ell=s_0}^{m_1} |S_1^\ell| \right] + p_1
\end{aligned} \tag{16}$$

where step (a) follows from Assumption A1 and step (b) follows because  $\sum_{\ell=s_0}^{m_1} \sum_{x \in S_1^\ell} \frac{1}{k(1-p)} \Pi_{p,N}^1(x) \leq 1$  and  $\Pi_{p,N}^1(x) = \frac{1}{k(1-p)}$ . Now, we will evaluate the terms  $S_1^\ell$  and  $S_2^\ell$  respectively. We recall the definition of the random event  $\mathcal{E}$  from the proof of Theorem 4

$$\mathcal{E} = \left\{ \frac{p}{2} \leq \frac{m_2}{m} \leq \frac{1+p}{2} \right\} \quad \text{and} \quad \mathbb{P}[\mathcal{E}] = 1 - o\left(\frac{1}{m^2}\right) \quad \text{as } m \rightarrow \infty$$

Then, by law of total expectation, we have that

$$\begin{aligned}
\lim_{m \rightarrow \infty} \mathbb{E}[|S_1^\ell|] &= \lim_{m \rightarrow \infty} \mathbb{E}[|S_1^\ell| | \mathcal{E}] \mathbb{P}[\mathcal{E}] + \lim_{m \rightarrow \infty} \underbrace{\mathbb{E}[|S_1^\ell| | \mathcal{E}^c]}_{O(m)} \underbrace{\mathbb{P}[\mathcal{E}^c]}_{o(\frac{1}{m^2})} \\
&= \lim_{m \rightarrow \infty} \sum_{x \in X_1} \mathbb{P}_{S_1 \sim (\Pi_{p,N}^1)^{m_1}}[x \text{ occurs } \ell \text{ times in } S_1 | \mathcal{E}] \mathbb{P}[\mathcal{E}] \\
&\stackrel{(a)}{=} \lim_{m \rightarrow \infty} k(1-p) \binom{m_1}{\ell} \left( \frac{1}{k(1-p)} \right)^\ell \left( 1 - \frac{1}{k(1-p)} \right)^{m_1-\ell} \mathbb{P}[\mathcal{E}] \\
&\stackrel{(b)}{=} \lim_{m \rightarrow \infty} k(1-p) \left( \frac{m_1}{\ell} \right)^\ell \left( \frac{1}{k(1-p)} \right)^\ell \left( 1 - \frac{1}{k(1-p)} \right)^{m_1-\ell} \mathbb{P}[\mathcal{E}] \\
&\stackrel{(c)}{=} 0
\end{aligned} \tag{17}$$

where step (a) follows from a simple counting argument, step (b) follows from the well known inequality  $\binom{m_1}{\ell} \leq \left( \frac{m_1}{\ell} \right)^\ell$  and step (c) follows from limit rules and the assumption that  $\ell \leq s_0 = o(m)$ ,  $k = o(m)$  and conditioned on the event  $\mathcal{E}$  which requires  $m_1 \geq \frac{(1-p)m}{2}$ .

Using a similar technique, we next lower bound  $\lim_{m \rightarrow \infty} \mathbb{E}[|S_2^\ell|]$

$$\begin{aligned}
\lim_{m, N \rightarrow \infty} \mathbb{E}[|S_2^\ell|] &= \lim_{m, N \rightarrow \infty} N \binom{m_2}{\ell} \left( \frac{1}{N} \right)^\ell \left( 1 - \frac{1}{N} \right)^{m_2-\ell} \mathbb{P}[\mathcal{E}] \\
&= \lim_{m, N \rightarrow \infty} \mathbb{P}_{\ell \sim \text{binom}(m_2, \frac{1}{N})}[l \leq s_0] \mathbb{P}[\mathcal{E}] \\
&\stackrel{(a)}{\rightarrow} \lim_{m, N \rightarrow \infty} \mathbb{P}_{\ell \sim \mathcal{N}(\frac{m_2}{N}, \frac{m_2}{N}(1-\frac{1}{N}))}[l \leq s_0] \mathbb{P}[\mathcal{E}] \\
&\stackrel{(b)}{\geq} \Phi \left( \frac{s_0 - \frac{1+p}{2c}}{\sqrt{\frac{1+p}{2c}}} \right)
\end{aligned} \tag{18}$$

where step (a) follows due to CLT,  $\Phi$  is the Gaussian CDF function, and step (b) follows from the condition on  $\mathcal{E}$  which enforces  $m_2 \leq \frac{(1+p)m}{2}$  and from the limit  $\frac{N}{m} \rightarrow c$  as  $m, N \rightarrow \infty$ . Using a similar argument, but

using the fact that conditioning on  $\mathcal{E}$  also enforces  $m_2 \geq \frac{pm}{2}$ , we get an upper bound on  $\lim_{m \rightarrow \infty} \mathbb{E} [|S_2^\ell|]$  as follows

$$\lim_{m, N \rightarrow \infty} \mathbb{E} [|S_2^\ell|] \leq \Phi \left( \left( s_0 - \frac{p}{2c} \right) \sqrt{\frac{2c}{p}} \right) \quad (19)$$

We use the following well known tail bounds for Gaussian distributions to bound the CDF function in the expression for the upper (Equation (19)) and lower (Equation (18)) bounds on  $\lim_{m \rightarrow \infty} \mathbb{E} [|S_2^\ell|]$ .

**Inequality 1** (From Feller [27]). *Let  $W$  be a random variable following the normal distribution  $N(\mu, \sigma^2)$ . Then, we have the following lower and upper tail bounds on  $W$  for all  $x > 0$  :*

$$\left( \frac{1}{x} - \frac{1}{x^3} \right) \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} \leq \mathbb{P}[W \geq \mu + \sigma x] \leq \frac{e^{-\frac{x^2}{2}}}{x\sqrt{2\pi}}$$

Plugging Equations (17) and (19) in Equation (12), we get the upper bound on error.

$$\begin{aligned} \text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \lim_{m, N \rightarrow \infty} \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\ &\leq (1 - p_1)p \left( 1 - \frac{(s_0 - \frac{p}{2c})^2 \frac{2c}{p} - 1}{\sqrt{2\pi} (s_0 - \frac{p}{2c})^3 \left( \sqrt{\frac{2c}{p}} \right)^3} e^{-(s_0 - p/2c)^2 c/p} \right) + p_1 \end{aligned} \quad (20)$$

Similarly, plugging Equation (18) into Equation (14) and Equation (17) into Equation (16), we get the required expressions for accuracy discrepancy in Equation (13), we get the desired result.

$$\begin{aligned} \Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \lim_{m, N \rightarrow \infty} \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\ &\geq (1 - p)c_3 \left( 1 - \frac{\sqrt{\frac{1+p}{2c}}}{\sqrt{2\pi} (s_0 - \frac{1+p}{2c})} e^{-(s_0 - \frac{1+p}{2c})^2 \frac{c}{1+p}} \right) - (1 - p)p_1 \end{aligned} \quad (21)$$

□

**Remark 1.** In the constant  $c_3 = \min_{x \in X} 1 - e^{s_0 \epsilon} \|\mathcal{F}(x)\|_\infty - s_0 e^{(s_0-1)\epsilon} \delta$ , the term  $s_0 e^{(s_0-1)\epsilon} \delta$  appears prohibitively large when  $s_0$  is large. However, note that for the privacy guarantees to be reasonable,  $\delta$  is usually of the order  $O\left(\frac{1}{\sqrt{m}}\right)$  because if it is not, then the privacy guarantees are somewhat vacuous. Therefore,  $\delta$  decreases faster than  $s_0$  grows with respect to  $m$  and hence the term is actually decreases with increasing  $m$ .

### A.3 Proof for Theorem 3

**Theorem 6** (Detailed version of Theorem 3). *For any  $p \in (0, 1/2)$ ,  $c > 0$  such that  $p/c \leq 1$ , consider the distribution  $\Pi_{p,N}$  where  $N$  is the number of minority subpopulations. The number of majority subpopulations is  $k(1 - p)$  which satisfies  $k(1 - p) \geq 2$  and  $k = o(m)$ . For any  $\epsilon, \delta, \alpha > 0$ , consider any  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$  that satisfies assumptions A1 and A2 with  $s_0 = \left( \frac{2-p}{2k(1-p)} \right) m + \alpha\sqrt{m}$  and some  $p_1, p_2 \in (0, 1)$  where  $N/m \rightarrow c$  as  $m, N \rightarrow \infty$ . Then,*

$$\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \geq c_4 p + (1 - p_2)(1 - p) \left( 1 - e^{-4(1-p)\alpha^2/(2-p)^2} \right)$$

$$\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \leq (1 - p) \left[ 1 - (1 - e^{-4(1-p)\alpha^2/(2-p)^2})(1 - p_2) \right]$$

where  $c_4 := \min_{x \in X} 1 - e^{s_0 \epsilon} \|\mathcal{F}(x)\|_\infty - s_0 e^{(s_0-1)\epsilon} \delta$ ,  $\text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = \lim_{m, N \rightarrow \infty} \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F})$ , and  $\Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = \lim_{m, N \rightarrow \infty} \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F})$ .

*Proof.* We use the same extra notations as in the proof of Theorem 4 and 5. As in the previous proofs, we first decompose the expression for error as follows.

$$\begin{aligned}
\text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \mathbb{E} \left[ \sum_{\ell=0}^{s_0} \sum_{x \in S^\ell} \Pi_{p,N}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] \\
&\quad + \mathbb{E} \left[ \sum_{\ell=s_0+1}^m \sum_{x \in S^\ell} \Pi_{p,N}(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] \right] \\
&\geq c_4 p + \sum_{\ell=s_0+1}^{m_2} \mathbb{E}[|S_2^\ell|] + \frac{(1-p_2)}{k} \sum_{\ell=0}^{s_0} \mathbb{E}[|S_1^\ell|]
\end{aligned} \tag{22}$$

where  $c_4 := \min_{x \in X} 1 - e^{s_0 \epsilon} \|\mathcal{F}(x)\|_\infty - s_0 e^{(s_0-1)\epsilon} \delta$  follows from similar arguments as Equation (15) and the third term exploits the assumption A2. We again recall the decomposition of accuracy discrepancy.

$$\Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) = (1-p) (\text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) - \text{err}_m(\mathcal{A}, \Pi_{p,N}^1, \mathcal{F})) \tag{23}$$

We use the simple upper bound  $\text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) \leq 1$  for error in the minority subpopulations. As we discuss later, this is a legitimate approximation of this term. We next lower bound the error of the majority subpopulations.

$$\begin{aligned}
\text{err}_m(\mathcal{A}, \Pi_{p,N}^1, \mathcal{F}) &= \sum_{x \in X_1} [\Pi_{p,N}^1(x) \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)]] \\
&\stackrel{(a)}{\geq} \mathbb{E} \left[ \sum_{\ell=1}^{s_0} \sum_{x \in S_1^\ell} \frac{1}{k(1-p)} (1-p_2) \right] \\
&= \frac{(1-p_2)}{k(1-p)} \sum_{\ell=1}^{s_0} \mathbb{E}[|S_1^\ell|]
\end{aligned} \tag{24}$$

We again recall the definition of the random event  $\mathcal{E}$  from the proof of Theorem 4

$$\mathcal{E} = \left\{ \frac{p}{2} \leq \frac{m_2}{m} \leq \frac{1+p}{2} \right\} \quad \text{and} \quad \mathbb{P}[\mathcal{E}] = 1 - o\left(\frac{1}{m^2}\right) \quad \text{as } m \rightarrow \infty$$

We first show that all points in the minority group occur less than  $s_0$  times in the dataset, i.e.  $\sum_{\ell=s_0+1}^{m_2} \mathbb{E}[|S_2^\ell|] \rightarrow 0$  as  $m \rightarrow \infty$ .

$$\begin{aligned}
\lim_{m, N \rightarrow \infty} \sum_{\ell=s_0+1}^{m_2} \mathbb{E}[|S_2^\ell|] &= \lim_{m, N \rightarrow \infty} \sum_{\ell=s_0+1}^{m_2} \mathbb{E}[|S_2^\ell| | \mathcal{E}] \mathbb{P}[\mathcal{E}] + \mathbb{E}[|S_2^\ell| | \mathcal{E}^c] \mathbb{P}[\mathcal{E}^c] \\
&= \lim_{m, N \rightarrow \infty} \sum_{\ell=s_0+1}^{m_2} N \binom{m_2}{\ell} \left(\frac{1}{N}\right)^\ell \left(1 - \frac{1}{N}\right)^{m_2-\ell} \mathbb{P}[\mathcal{E}] \\
&\stackrel{(a)}{\leq} \lim_{m, N \rightarrow \infty} N e^{\left(-2m_2 \left(1 - \frac{1}{N} - \frac{m_2 - s_0}{m_2}\right)\right)^2} \\
&\stackrel{(b)}{\leq} c m e^{\left(\frac{1+p}{c} - \frac{pm}{k(1-p)} - \alpha\sqrt{m}\right)^2} = 0
\end{aligned} \tag{25}$$

where step (a) is an instantiation of the Hoeffding's inequality on the upper tail of a bernoulli random variable, and step (b) is due to the limit  $\frac{N}{m} \rightarrow c$  as  $m, N \rightarrow \infty$ , the definition  $s_0 = \frac{2-p}{2k(1-p)}m + \alpha\sqrt{m}$ , and  $m_2 \geq \frac{(1+p)m}{2}$  due to the conditioning on the event  $\mathcal{E}$ . Note that this justifies the upper bound  $\text{err}_m(\mathcal{A}, \Pi_{p,N}^2, \mathcal{F}) \leq 1$ .



Similarly, using the law of total expectation and a simple counting argument, we can lower bound  $\sum_{\ell=1}^{s_0} \mathbb{E} [|S_1^\ell|]$

$$\begin{aligned} \sum_{\ell=1}^{s_0} \mathbb{E} [|S_1^\ell|] &\geq \sum_{\ell=1}^{s_0} k(1-p) \binom{m_1}{\ell} \left( \frac{1}{k(1-p)} \right)^\ell \left( 1 - \frac{1}{k(1-p)} \right)^{m_1-\ell} \\ &= k(1-p) \mathbb{P}_{\ell \sim \text{binom}(m_1, \frac{1}{k(1-p)})} [\ell \leq s_0] \end{aligned} \quad (26)$$

To calculate  $\mathbb{P}_{\ell \sim \text{binom}(m_1, \frac{1}{k(1-p)})} [\ell \leq s_0]$ , we apply Hoeffding's inequality on Bernoulli distribution to get the following,

$$\begin{aligned} \mathbb{P}_{\ell \sim \text{binom}(m_1, \frac{1}{k(1-p)})} [\ell \leq s_0] &= 1 - \mathbb{P}_{\ell \sim \text{binom}(m_1, \frac{1}{k(1-p)})} [\ell \geq s_0] \\ &\stackrel{(a)}{\geq} 1 - e^{-2m_1 \left( 1 - \frac{1}{k(1-p)} - \frac{m_1 - s_0}{m_1} \right)^2} \\ &\stackrel{(b)}{=} 1 - e^{-(1-p)m \left( \frac{2\alpha}{(2-p)\sqrt{m}} \right)^2} = 1 - e^{-\frac{4(1-p)\alpha^2}{(2-p)^2}} \end{aligned} \quad (27)$$

where, step (a) is an instantiation of Hoeffding's inequality and step (b) follows from conditioning on  $\mathcal{E}$  and due to  $s_0 = \frac{2-p}{2k(1-p)}m + \alpha\sqrt{m}$ . Plugging, Equation (27) into Equation (26), we get

$$\lim_{m, N \rightarrow \infty} \sum_{\ell=1}^{s_0} \mathbb{E} [|S_1^\ell|] \geq k(1-p) \left( 1 - e^{-\frac{4(1-p)\alpha^2}{(2-p)^2}} \right). \quad (28)$$

Plugging Equations (25) and (28) into Equation (22), we get the lower bound for error

$$\begin{aligned} \text{err}(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \lim_{m, N \rightarrow \infty} \text{err}_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\ &\geq c_4 p + (1-p_2)(1-p) \left( 1 - e^{-\frac{4(1-p)\alpha^2}{(2-p)^2}} \right) \end{aligned} \quad (29)$$

Plugging Equations (24) and (25), in Equation (23) we get the upper bound for fairness

$$\begin{aligned} \Gamma(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) &= \lim_{m, N \rightarrow \infty} \Gamma_m(\mathcal{A}, \Pi_{p,N}, \mathcal{F}) \\ &\leq (1-p) \left[ 1 - (1-p_2) \left( 1 - e^{-4(1-p)\alpha^2/(2-p)^2} \right) \right] \end{aligned} \quad (30)$$

which completes the proof of Theorem 6. □

## A.4 Proof of Lemma 1

Next, we restate and prove Lemma 1. We will define some extra notation for this. First, we define the concept of vertex cover and independent set.

**Definition 5** (Vertex Cover and Independent Set). *Consider any undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  where  $\mathcal{V}$  is the set of vertices and  $\mathcal{E}$  is the set of edges.*

- A set  $\mathcal{C} \subseteq \mathcal{V}$  is called the vertex cover of  $\mathcal{G}$  if for all edges  $(u, v) \in \mathcal{E}$  at least one of  $u$  and  $v$  belongs to  $\mathcal{C}$ . The smallest such set is known as the minimum vertex cover. If a vertex cover does not remain a vertex cover upon the removal of any of its vertices, then it is known as a minimal vertex cover.

- A set  $\mathcal{I} \subseteq \mathcal{V}$  is called an independent set of  $\mathcal{G}$  if for any two vertices  $u, v \in \mathcal{I}$ , the edge  $(u, v) \notin \mathcal{E}$ . The largest such set is known as the maximum independent set. If an independent set does not remain an independent set upon including a vertex from outside the set, then it is known as a maximal independent set.

- All minimum vertex covers and maximum independent sets are minimal vertex covers and maximal independent sets respectively but not vice versa.

- The complement of any vertex cover  $\mathcal{C}$  is an independent set  $\mathcal{I} = \mathcal{V} \setminus \mathcal{C}$ . Therefore, the complement of a minimum vertex cover is a maximum independent set.

Let  $Q \in \mathcal{Y}^{m-1}$  represent an arbitrary labeling of the points in  $S^\ell \setminus \{x\}$  and use  $\mathcal{F}_{|S^\ell \setminus \{x\}, Q}$  to denote the marginal distribution of  $\mathcal{F}$  over the set of functions that satisfies the labeling  $Q$  on  $S^\ell \setminus \{x\}$ .

**Lemma 3.** Consider a label prior  $\mathcal{F}$ , an  $m$ -sized unlabeled  $S \sim \Pi_{p,N}^m$ , and an  $(\epsilon, \delta)$ -DP algorithm  $\mathcal{A}$ . For any  $p \in (0, 1)$ , let  $s_0$  be the largest integer that satisfies  $p \geq \frac{1+s_0 e^{-\epsilon} \delta}{1+e^{-s_0 \epsilon}}$ . Then for any  $x \in S^\ell$ , where  $\ell \in \mathbb{N}$ ,  $\ell \leq s_0$ , there exists a set of labeling functions  $\hat{F} \subset F$  such that for all functions  $f \in \hat{F}$ ,

$$\mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] > 1 - p \quad (31)$$

and the size of  $\hat{F}$  is large

$$|\hat{F}| \geq |F| - \sum_{Q \in \mathcal{Y}^{m-1}} \max_i |\{f \in \mathcal{F}_{|S^\ell \setminus \{x\}, Q} : f(x) = \mathcal{Y}_i\}| \quad (32)$$

In particular, if  $\mathcal{F}$  defines a uniform distribution over  $\mathcal{Y}^\mathcal{X}$ , then with probability at least  $1 - \frac{1}{|\mathcal{Y}|}$ , any function  $f \in \mathcal{F}$  satisfies Equation (31).

*Proof of Lemma 3.* Let  $S$  be a dataset of size  $m$  and For any  $x \in S$ , two labelings functions  $f_1, f_2$  satisfy condition

$$C_1 \text{ w.r.t } x \text{ if } \begin{cases} f_1(z) = f_2(z) & z \neq x \\ f_1(z) \neq f_2(z) & z = x \end{cases} \quad (C1)$$

By Lemma 4, if  $p \geq \frac{1+s_0 e^{-\epsilon} \delta}{1+e^{-s_0 \epsilon}}$  and if  $f_1, f_2$  satisfy condition C1 then at least one of them satisfies the following.

$$\mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) \neq f(x)] > 1 - p \quad (C2)$$

Without loss of generality, let  $Q \in \mathcal{Y}^{m-1}$  represent an arbitrary labeling of the points in  $S^\ell \setminus \{x\}$ , and let  $\mathcal{F}_{|S^\ell \setminus \{x\}, Q}$  represent the marginal distribution of  $\mathcal{F}$  over the set of functions that satisfies the labeling  $Q$  on  $S^\ell \setminus \{x\}$ . Let  $\zeta(Q, x)$  denote the total number of functions in  $\text{Supp}(\mathcal{F}_{|S^\ell \setminus \{x\}, Q})$  that satisfies condition C2. We show that  $\zeta(Q, x)$  can be calculated as the size of the minimum vertex cover (see Definition 5) of the following graph.

Construct a graph  $\mathcal{G}_{S, x, Q}$  where every function  $f_i \in \text{Supp}(\mathcal{F}_{|S^\ell \setminus \{x\}, Q})$  represents a vertex. Construct an edge between two vertices  $f_i, f_j$  if  $f_i(x) \neq f_j(x)$ . Note that the two functions constructing any edge in this graph satisfies condition C1. It is easy to see that this is a complete  $|\mathcal{Y}|$ -partite graph where the  $i^{\text{th}}$  partition consists of functions  $\{f \in \mathcal{F}_{|S^\ell \setminus \{x\}, Q} : f(x) = \mathcal{Y}_i\}$  i.e. functions which predict the  $i^{\text{th}}$  label on  $x$ . A vertex cover of a graph, defined in Definition 5, is a set of vertices such that for every edge at least one of its end points are included in the set. As for every edge in  $\mathcal{G}$ , both its endpoints satisfies condition C1, at least one of them should satisfy condition C2. Therefore, the total number of functions that satisfies condition C2 is at least the minimum vertex cover of  $\mathcal{G}$ . Using Lemma 5 to calculate the minimum vertex cover of  $\mathcal{G}$ , we have that

$$\zeta(Q, x) \geq |\mathcal{F}_{|S^\ell \setminus \{x\}, Q}| - \max_i |\{f \in \mathcal{F}_{|S^\ell \setminus \{x\}, Q} : f(x) = \mathcal{Y}_i\}|.$$

Summing this up for all labelings, the total number of labeling functions that satisfy condition C2 is

$$\begin{aligned} \sum_{Q \in \mathcal{Y}^{m-1}} \zeta(Q, x) &\geq \sum_{Q \in \mathcal{Y}^{m-1}} |\mathcal{F}_{|S^\ell \setminus \{x\}, Q}| - \max_i |\{f \in \mathcal{F}_{|S^\ell \setminus \{x\}, Q} : f(x) = \mathcal{Y}_i\}| \\ &= |F| - \sum_{Q \in \mathcal{Y}^{m-1}} \max_i |\{f \in \mathcal{F}_{|S^\ell \setminus \{x\}, Q} : f(x) = \mathcal{Y}_i\}| \end{aligned} \quad (33)$$

This completes the proof of Equation (32). To prove the next part, let  $\mathcal{F}$  be a uniform distribution over all functions in  $\mathcal{Y}^\mathcal{X}$ . Then for all  $Q \in \mathcal{Y}^{m-1}$ ,  $\mathcal{G}_{S, x, Q}$  is a complete multi-partite graph with equal-sized partitions.

Further, in any of the partitions, the total probability density of the functions corresponding to the vertices in that partitions under  $\mathcal{F}_{|S \setminus \{x\}, Q}$  is exactly  $\frac{1}{|\mathcal{Y}|}$ . Hence, the total probability mass, under  $\mathcal{F}_{|S \setminus \{x\}, Q}$  of the minimum vertex cover is  $1 - \frac{1}{|\mathcal{Y}|}$ . The probability of observing any arbitrary  $Q$  under  $\mathcal{F}$  is exactly  $\frac{1}{|\mathcal{Y}|^{m-1}}$  and there are  $|\mathcal{Y}|^{m-1}$  such labelings. Summing this up for all labelings  $Q \in \mathcal{Y}^{m-1}$  and multiplying with the conditional probability of observing  $Q$  under  $\mathcal{F}$ , we get

$$\sum_{Q \in \mathcal{Y}^{m-1}} 1 \cdot \frac{1}{|\mathcal{Y}|^{m-1}} - \frac{1}{|\mathcal{Y}|} \cdot \frac{1}{|\mathcal{Y}|^{m-1}} = |\mathcal{Y}|^{m-1} \cdot \frac{1}{|\mathcal{Y}|^{m-1}} - |\mathcal{Y}|^{m-1} \cdot \frac{1}{|\mathcal{Y}|} \cdot \frac{1}{|\mathcal{Y}|^{m-1}} = 1 - \frac{1}{|\mathcal{Y}|}.$$

This completes the proof.  $\square$

**Lemma 4.** *Let  $S$  be a dataset of size  $m$  and  $p \in (0, 1)$  be a fixed value. Let  $s_0$  be the largest integer that fulfills  $p \geq \frac{1+s_0 e^{-\epsilon} \delta}{1+e^{-s_0 \epsilon}}$ . Then, for an  $(\epsilon, \delta)$ -differentially private algorithm  $\mathcal{A}$ , for any  $x \in S$  where  $\ell \in \mathbb{N}, \ell \leq s_0$ , if two functions  $f_1, f_2$  agree on all points in  $S^\ell \setminus \{x\}$  but disagree on  $x$  then at least one of the two functions, say  $f$  incurs*

$$\mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] > 1 - p$$

*Proof of Lemma 4.* Let  $S$  be a dataset of size  $m$  and For any  $x \in S$ , two labelings functions  $f_1, f_2$  satisfy condition

$$C_1 \text{ w.r.t } x \text{ if } \begin{cases} f_1(z) = f_2(z) & z \neq x \\ f_1(z) \neq f_2(z) & z = x \end{cases} \quad (C1)$$

For a given  $p \in (0, 1)$ , let  $s_0$  be the largest integer that satisfies  $p \geq \frac{1+s_0 e^{-\epsilon} \delta}{1+e^{-s_0 \epsilon}}$ . For any  $x \in S^\ell$  where  $\ell \leq s_0$ , we will show if two labelings functions  $f_1, f_2$  satisfies condition C1, then at least one of them  $f \in \{f_1, f_2\}$  will not satisfy condition C2 as follows

$$\mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) \neq f(x)] \leq 1 - p \quad (C2)$$

We will prove this via contradiction. Let's assume that two labeling functions  $f, \tilde{f}$  satisfying condition C1 also satisfies the condition C2. Then, we lower bound the probability that an  $(\epsilon, \delta)$ -DP algorithm trained on a dataset, labeled with  $\tilde{f}$ , disagrees with  $f$  on  $x \in S^\ell$ .

$$\begin{aligned} \mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})} [h(x) = f(x)] &\geq \left( \mathbb{P}_{h \sim \mathcal{A}(S_f)} [h(x) = f(x)] - \ell e^{(\ell-1)\epsilon} \delta \right) e^{-\ell\epsilon} \\ &\geq \left( p - \ell e^{(\ell-1)\epsilon} \delta \right) e^{-\ell\epsilon} \end{aligned} \quad (34)$$

The first inequality follows from the notion of group privacy and the second inequality follows from our assumption (C2). Next, we provide an upper bound for the same quantity

$$\mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})} [h(x) = f(x)] \leq \mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})} [h(x) \neq \tilde{f}(x)] \leq 1 - p \quad (35)$$

where the first inequality follows using the definition of  $\tilde{f}$  that  $\tilde{f}(x) \neq f(x)$  and the second inequality again uses assumption (C2). Now, we combine the two inequalities (Equations (34) and (35)), to get the following relationship between  $p, \ell, \epsilon$ , and  $\delta$ ,

$$\begin{aligned} \left( p - \ell e^{(\ell-1)\epsilon} \delta \right) e^{-\ell\epsilon} &\leq 1 - p \\ \implies p &\leq \frac{1 + \ell e^{-\epsilon} \delta}{1 + e^{-\ell\epsilon}} \end{aligned} \quad (36)$$

As  $\ell \leq s_0$ , we get

$$\frac{1 + \ell e^{-\epsilon} \delta}{1 + e^{-\ell\epsilon}} \leq \frac{1 + s_0 e^{-\epsilon} \delta}{1 + e^{-s_0 \epsilon}} \leq p$$

which leads to the desired contradiction. Thus, if a pair of functions satisfies condition C1, then at least one of them does not satisfy condition 34. This completes the proof.  $\square$

**Lemma 5.** Given  $k$  sets of vertices  $\{\mathcal{V}_i\}_{i=1}^k$ , consider a complete  $k$ -partite graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  where the  $k^{\text{th}}$  partition is constituted of  $\mathcal{V}_i$ . Then the size of the minimum vertex cover is  $|\mathcal{V}| - \max_i |\mathcal{V}_i|$ .

*Proof of Lemma 5.* A set of vertices is known as an independent set if and only if its complement is a vertex cover. Therefore, we use the identity that the size of the minimum vertex cover and the maximum independent set of a graph equals to the total number of vertices of the graph. Therefore, the problem of finding the minimum vertex cover is equivalent to finding the maximum independent set.

By the definition of independent set, no two vertices in the independent set are adjacent. Without loss of generality, let  $v \in \mathcal{V}_i$  for some  $i$  be in the independent set  $\mathcal{I}$ . Then, by definition of complete  $k$ -partite graph, all vertices belonging to  $\bigcup_{j \neq i} \mathcal{V}_j$  are adjacent to  $v$  and therefore cannot be in  $\mathcal{I}$ . However, all other vertices in  $\mathcal{V}_i$  can be included in  $\mathcal{I}$ . This creates a *maximal independent set*. Thus, there are only  $\mathcal{Y}$  maximal independent sets and the set corresponding to  $i^* = \operatorname{argmax}_i |\mathcal{V}_i|$  is the maximum independent set. Therefore, the size of the minimum vertex cover is  $|\mathcal{V}| - \max_i |\mathcal{V}_i|$ .  $\square$

## A.5 Proof of Lemma 2

We restate a detailed version of Lemma 2 and provide the proof below.

**Lemma 6.** The algorithm  $A_\eta : S \rightarrow \mathcal{Y}^X$  is  $\left(\log\left(\frac{(1-\eta)(|\mathcal{Y}|-1)}{\eta}\right), 0\right)$ -differentially private. Further, for any  $\beta_1, \beta_2 \in (0, 1)$ ,

- if a point appears less than  $\frac{\log(1/\beta_1)}{1 + \frac{1}{2} \log[1/4\eta(1-\eta)]}$  times in the training set, then  $\mathbb{P}_{h \sim \mathcal{A}}[h(x) \neq f(x)] \geq \beta_1$  and
- if a point appears more than  $\frac{8(1-\eta)}{(1-2\eta)^2} \log(1/\beta_2)$  times, then  $\mathbb{P}_{h \sim \mathcal{A}}[h(x) \neq f(x)] \leq \beta_2$ .

*Proof.* We restate the algorithm for the sake of completeness. The algorithm  $A_\eta$  accepts a dataset  $S_f \in (X \times \mathcal{Y})^m$  and a noise rate  $\eta \in (0, \frac{1}{2})$  as input. Then it creates a dictionary where the set  $X$  is the set of keys. In order to assign values to every key, it first randomly flips the label of every element in  $S_f$  with probability  $\eta$ , then for every unique key in  $S_f$ , the algorithm computes the majority label of that key in the flipped dataset and assigns that majority label to the corresponding key. For elements in  $X$  not present in  $S_f$ , it assigns a random element from  $\mathcal{Y}$ .

**Privacy:** We first prove that  $A_\eta$  is  $(\epsilon, 0)$ -differentially private. For any two neighboring datasets  $S_f$  and  $S_{\tilde{f}}$  that differ in one element, the classifiers obtained from the algorithm  $A_\eta$  on the two datasets differs the most if the datasets differ on a singleton point  $x'$  i.e. a point which appears just once in the dataset.

$$\begin{aligned}
& \text{For all } x' \in S, f \in F, \text{ and } \tilde{f} \in F_{|S \setminus \{x'\}, f} \\
& \frac{\mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) = f(x), \forall x \in S]}{\mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})}[h(x) = f(x), \forall x \in S]} = \frac{\prod_{x \in S} \mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x) = f(x)]}{\prod_{x \in S} \mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})}[h(x) = f(x)]} \\
& = \frac{\mathbb{P}_{h \sim \mathcal{A}(S_f)}[h(x') = f(x')]}{\mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})}[h(x') = f(x')]} \\
& = \frac{\mathbb{P}_{h \sim \mathcal{A}(S_f)}[x' \text{ is not flipped}]}{\mathbb{P}_{h \sim \mathcal{A}(S_{\tilde{f}})}[x' \text{ is flipped and flipped to } f(x')]} \\
& = \frac{1-\eta}{\eta \frac{1}{|\mathcal{Y}|-1}} = \frac{1-\eta}{\eta} (|\mathcal{Y}| - 1)
\end{aligned} \tag{37}$$

where  $F_{|S \setminus \{x'\}, f}$  is the subset of  $F$  consisting of all labeling functions that agree with  $f$  at all points in  $S$  except  $x$ . This shows that the algorithm  $A_\eta$  is  $\left(O\left(\log \frac{1}{\eta}\right), 0\right)$ -DP.

**Utility:** Next, we bound the probability that  $A_\eta$  misclassifies a point occurring  $\ell$  times in  $S$  in order to establish the utility of the algorithm. For a point  $x$  occurring  $\ell$  times in  $S$ , we denote all replicates of this point in the dataset as  $[X_2, \dots, x_\ell]$ . As each point is flipped with probability  $\eta$ , let  $k_i^x = \mathbb{I}\{x_i \text{ is not flipped}\}$ ,

then  $\sum_{i=1}^{\ell} k_i^x$  is a binomial random variable with  $\ell$  trials and success probability  $1 - \eta$ . The algorithm  $A_\eta$  would correctly classify  $x$  if  $\sum_{i=1}^{\ell} k_i^x > \ell/2$ .

The upper bound of the error probability is shown in Sanyal et al. [46], but we reproduce it here for completeness.

$$\begin{aligned}
\mathbb{P} \left[ \sum_{i=1}^{\ell} k_i^x \leq \frac{\ell}{2} \right] &= \mathbb{P} \left[ \sum_{i=1}^{\ell} k_i^x \leq \frac{\ell}{2\mu} \mu + \mu - \mu \right] \\
&= \mathbb{P} \left[ \sum_{i=1}^{\ell} \left( 1 - \left( 1 - \frac{\ell}{2\mu} \right) \right) \mu \right] \\
&\stackrel{(a)}{\leq} e^{-\frac{(1-2\eta)^2}{8(1-\eta)^2} \mu} \\
&\stackrel{(b)}{=} e^{-\frac{(1-2\eta)^2}{8(1-\eta)} \ell}
\end{aligned} \tag{38}$$

where step (a) follows from Chernoff's concentration bound on the lower tail and step (b) due to  $\mu = (1 - \eta)\ell$ .

Then, we show a lower bound for the error probability using an anti-concentration bound on the lower tail. Define  $D(a||p)$  to be the Kullback-Leibler divergence between two Bernoulli random variables with parameter  $a$  and  $p$ , i.e.  $D(a||p) = a \log \left( \frac{a}{p} \right) + (1 - a) \log \left( \frac{1-a}{1-p} \right)$ .

$$\begin{aligned}
\mathbb{P} \left[ \sum_{i=1}^{\ell} k_i^x \leq \frac{\ell}{2} \right] &\stackrel{(a)}{\geq} \frac{1}{\sqrt{2\ell}} e^{-\ell D(\frac{\ell}{2\ell} || 1-\eta)} \\
&= \frac{1}{\sqrt{2\ell}} e^{-\ell D(\frac{\ell}{2\ell} || 1-\eta)} \\
&= \frac{1}{\sqrt{2\ell}} e^{-\frac{\ell}{2} (\log(1/2(1-\eta)) + \log(1/2\eta))} \\
&= \frac{1}{\sqrt{2\ell}} [4\eta(1-\eta)]^{\ell/2}
\end{aligned} \tag{39}$$

where step (a) follows from the binomial anti-concentration bound

**Inequality 2** (From Ash [4]). *Let  $X$  be a binomial random variable with parameters  $(m, p)$ , then for any  $k < mp$ , we can lower bound the left tail of  $X$  as*

$$\mathbb{P}[X \leq k] \geq \frac{1}{\sqrt{2m}} e^{-mD(\frac{k}{m} || p)}.$$

Given the upper bound and lower bound of the error probability, we can show that for any  $\beta_1 \in (0, 1)$ , if a point appears less than  $\ell \leq \frac{-\log(\beta_1)}{1 - \frac{1}{2} \log[4\eta(1-\eta)]}$  times,  $\mathbb{P}_{h \sim \mathcal{A}}[h(x) \neq f(x)] \geq \beta_1$  and for any  $\beta_2 \in (0, 1)$ , if a point appears more than  $\ell \geq -\frac{8(1-\eta)}{(1-2\eta)^2} \log \beta_2$  times in the dataset, then  $\mathbb{P}_{h \sim \mathcal{A}}[h(x) \neq f(x)] \leq \beta_2$ .

- Consider the case  $\ell \leq -\frac{\log(\beta_1)}{1 - \frac{1}{2} \log[4\eta(1-\eta)]}$ . We can rewrite this as

$$\begin{aligned}
-\ell + \frac{\ell}{2} \log[4\eta(1-\eta)] &\geq \log(\beta_1) \\
-\frac{1}{2} \log(2\ell) + \frac{\ell}{2} \log[4\eta(1-\eta)] &\geq \log \beta_1
\end{aligned} \tag{40}$$

where the first inequality is obtained by using the identity  $-1 + \frac{1}{2} \log[4\eta(1-\eta)] < 0$  when  $\eta \in (0, \frac{1}{2})$  and the second inequality is because  $\log(2\ell) \leq 2\ell$ . Substituting this into Equation (39), we obtain the desired result

$$\mathbb{P}_{h \sim \mathcal{A}}[h(x) \neq f(x)] \geq \frac{1}{\sqrt{2\ell}} [4\eta(1-\eta)]^{\frac{\ell}{2}} \geq \beta_1.$$

- Next, we consider the case  $\ell \geq -\frac{8(1-\eta)}{(1-2\eta)^2} \log(\beta_2)$ . We can write this as

$$-\frac{(1-2\eta)^2 \ell}{8(1-\eta)} \leq \log \beta_2$$

Substituting this into the upper bound on error from Equation (38), we get the desired result

$$\mathbb{P}_{h \sim \mathcal{A}} [h(x) \neq f(x)] \leq e^{-\frac{(1-2\eta)^2 \ell}{8(1-\eta)}} \leq \beta_2$$

This completes the proof of Lemma 2. □

## B Experimental Details and Additional Experiments

### B.1 Synthetic data

**Constructing the synthetic data distribution** Given  $m \in \mathbb{N}, p \in (0, 0.5), C > 0$ , and  $k \in \mathbb{N}$ , we first create two boolean hypercubes of dimension  $d_{\min}$  and  $d_{\max}$  respectively, where  $d_{\min} = O(\lceil \log Cm \rceil)$  and  $d_{\max} = O(\lceil \log(k(1-p)) \rceil)$  respectively. For the minority subpopulation, the data is distributed uniformly across a mixture of  $N$  Gaussian clusters positioned on a randomly chosen set of  $d_{\min}$  vertices of this hyper-cube. Similarly, for the majority class the data is distributed across a uniform mixture of  $d_{\max}$  Gaussian clusters situated on randomly chosen vertices of the second hypercube. Each of the clusters is randomly assigned a binary label. Finally,  $m$  samples are sampled from the minority distribution, along with their labels, and each data covariate is appended with a  $d_{\max}$ -dimensional vector of all  $10^{-4}$  to create a  $d_{\min} + d_{\max}$  dimensional dataset of size  $m$ . Similarly,  $m$  samples are sampled from the majority distribution and the covariates are appended with a  $d_{\min}$  vector of all  $10^{-4}$  to create a  $d_{\min} + d_{\max}$  dimensional dataset of size  $m$ . Finally,  $m_1$  points are chosen randomly from the minority dataset and  $m - m_1$  points are chosen randomly from the majority dataset where  $m_1$  is sampled from a binomial distribution with parameters  $m$  and  $p$ .

**DP Algorithm for synthetic data** We use PyTorch Opacus [57] for all our experiments. We use multiple values of  $\epsilon$  in the range  $(0.1, 10)$ . We clip the maximum gradient norm to 1, and then train the algorithm for  $10^3$  epochs. Note that the clipping is necessary to control the lipschitzness of the gradient descent algorithm update step. The DP-SGD algorithm uses the value of  $\delta$ , the maximum gradient norm, and the number of epochs mentioned above to compute the noise-multiplier of the algorithm that is necessary to attain the required privacy guarantee.

**Additional weight of minority group** In this section, we plot additional experiments with the same setup as Section 3.1. In particular, we show results from the same experiment as Figure 2 but with  $p = 0.5$  instead of  $p = 0.2$ . Our experiments show that the same trend is held in particular for how the accuracy discrepancy changes with  $c$  and with  $\epsilon$ .

### B.2 Additional sizes of subpopulations in CelebA

To construct the subpopulations in CelebA, we use the following attributes 1. Pointy Nose, 2. Wearing Earrings, 3. Wavy Hair, 4. Wearing Lipstick, 5. Heavy Makeup, 6. Attractive, 7. Receding Hairline, 8. Blurry, 9. Bangs, 10. Wearing Hat, 11. Eyeglasses.. This also corresponds to Set A used in Figure 7.

We use the following features for Set B in our experiments in Figure 7. 1. Arched Eyebrows 2. Bags Under Eyes 3. Blond Hair 4. Double Chin 5. High Cheekbones 6. Pale Skin 7. Rosy Cheeks 8. Straight Hair 9. Wearing Necklace 10. Wearing Necktie 11. Young

In Figure 4, we reported results with the subpopulation size set to 40. In this section, we report results where the size belongs to  $\{5, 10, 20, 60, 80, 100\}$ . In Figure 9, we show how disparate accuracy changes with  $\epsilon$  for various subpopulation sizes. In Figure 10, we show how the minority group accuracy and the overall

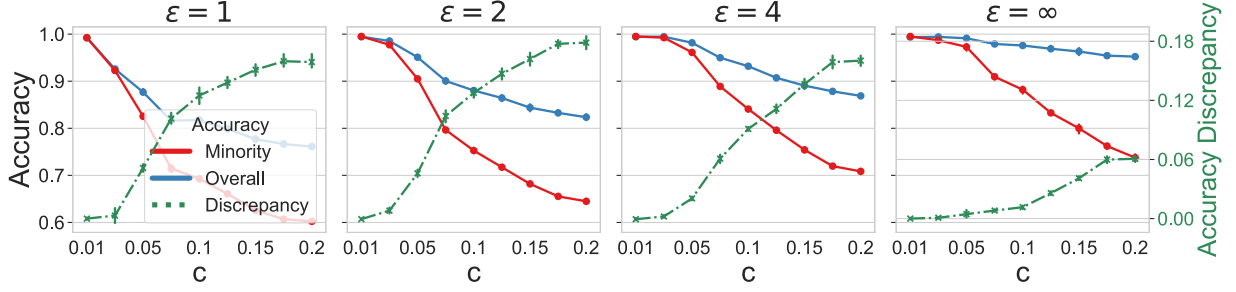


Figure 8: Here  $p = 0.5$ . Each figure plots the disparate accuracy (higher is less fair) in green dashed line, the accuracy of the minority group with red, and the overall accuracy with blue on the y-axis and the parameter  $c$  in the X-axis. The left most ( $\epsilon = 1$ ) achieves the strictest level of privacy and the right most ( $\epsilon = \infty$ ) is vanilla training without any privacy constraints.

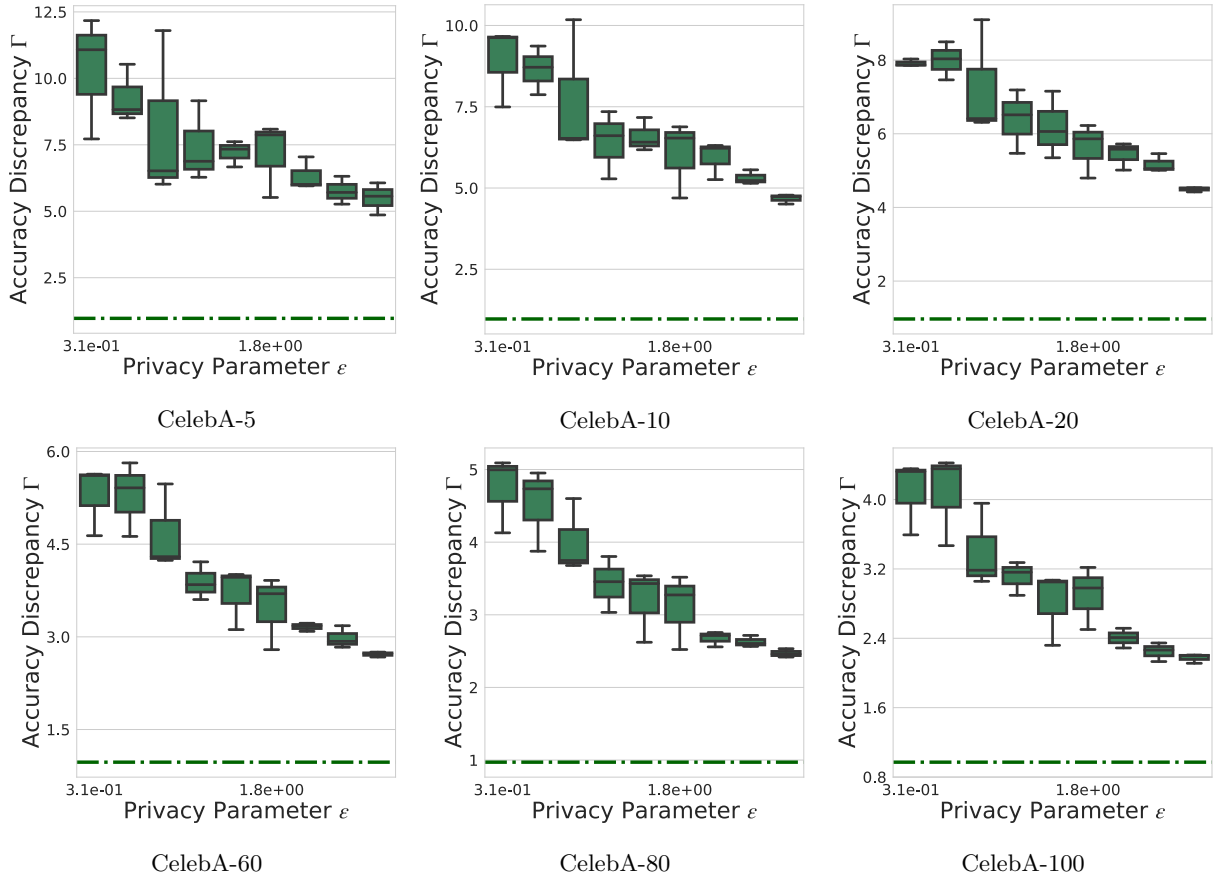


Figure 9: CelebA- $x$  represents a minority majority division where the minority group is comprised of subpopulations whose size is atmost  $x$ . Disparate accuracy decreases (low is fairer) with increasing  $\epsilon$  (high is less privacy) i.e. fairness is worse for stricter privacy for all sub-population sizes. Dashed lines show the corresponding accuracy discrepancy of a vanilla model (no privacy criterion).

accuracy changes with  $\epsilon$  for various subpopulation sizes. Overall the results, reflect the same trend as Figure 4 but the worst case accuracy discrepancy decreases for increasing subpopulation sizes.

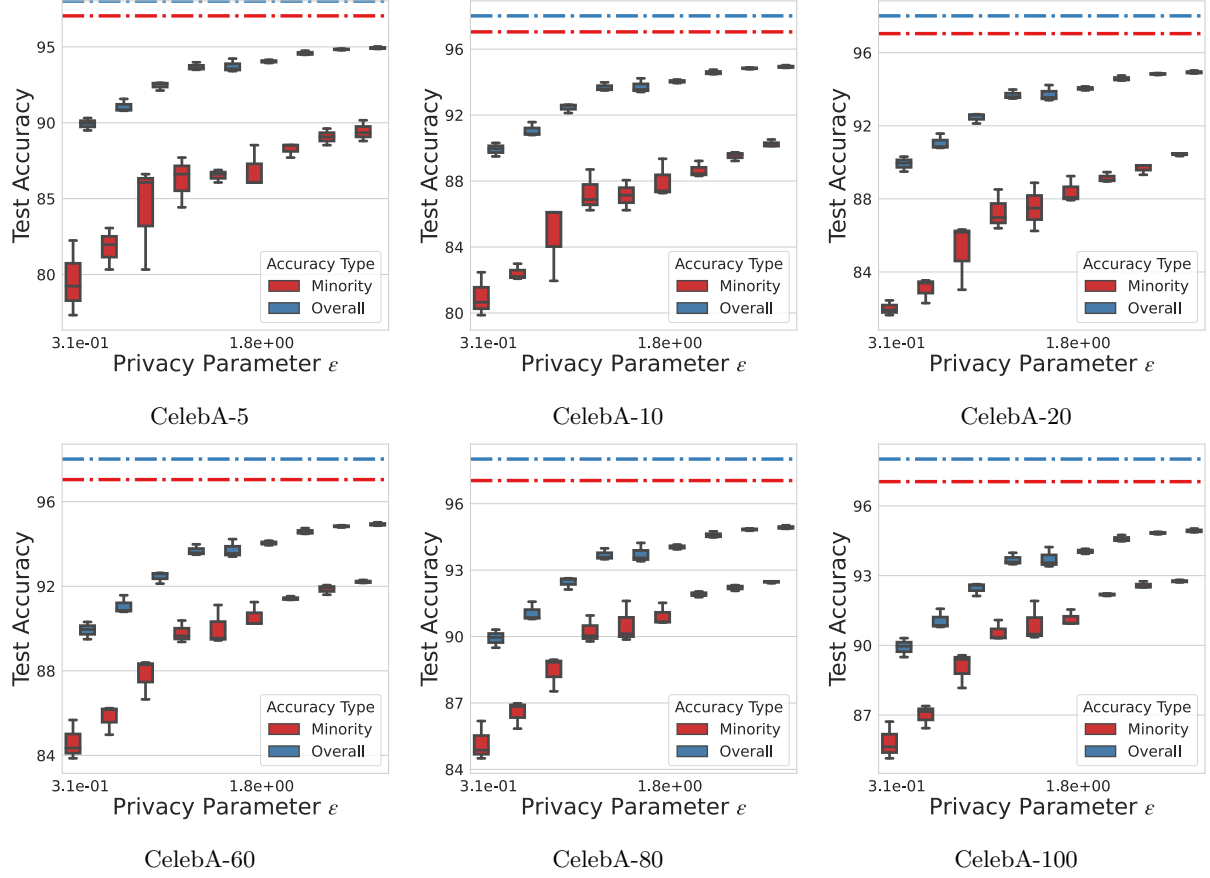


Figure 10: CelebA- $x$  represents a minority majority division where the minority group is comprised of subpopulations whose size is at most  $x$ . Both minority and overall accuracy increases with increasing  $\epsilon$  (high is less privacy) but overall accuracy increases faster. Dashed lines show the corresponding accuracy of a vanilla model (no privacy criterion).

### B.3 Additional details and experiments on CIFAR-10

In Figure 11, we plot the same results as Figure 5 but with the threshold  $\rho$  set to 0.01.



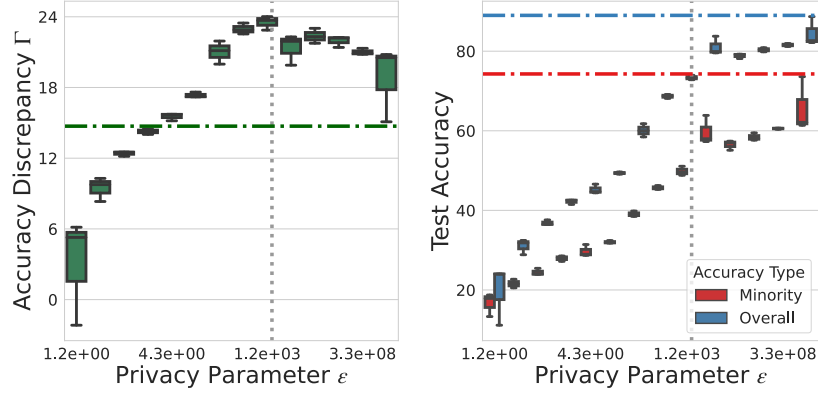


Figure 11: CIFAR-10—**Left** Disparate accuracy decreases (low is fairer) with increasing  $\epsilon$  (high is less privacy) i.e. fairness is worse for stricter privacy. **Right** figure shows that both minority and overall accuracy increases with  $\epsilon$ . Dashed lines show the corresponding metric of a vanilla model (no privacy criterion).