# Causal Inference with Differentially Private (Clustered) Outcomes

Adel Javanmard*†       Vahab Mirrokni †       Jean Pouget-Abadie†

May 1, 2024

## Abstract

Estimating causal effects from randomized experiments is only feasible if participants agree to reveal their potentially sensitive responses. Of the many ways of ensuring privacy, label differential privacy is a widely used measure of an algorithm's privacy guarantee, which might encourage participants to share responses without running the risk of de-anonymization. Many differentially private mechanisms inject noise into the original data-set to achieve this privacy guarantee, which increases the variance of most statistical estimators and makes the precise measurement of causal effects difficult: there exists a fundamental privacy-variance trade-off to performing causal analyses from differentially private data. With the aim of achieving lower variance for stronger privacy guarantees, we suggest a new differential privacy mechanism, CLUSTER-DP, which leverages any given cluster structure of the data while still allowing for the estimation of causal effects. We show that, depending on an intuitive measure of cluster quality, we can improve the variance loss while maintaining our privacy guarantees. We compare its performance, theoretically and empirically, to that of its unclustered version and a more extreme uniform-prior version which does not use any of the original response distribution, both of which are special cases of the CLUSTER-DP algorithm.

## 1 Introduction

Measurement and experimentation are essential tools to improve any user-facing product. Technology companies routinely run randomized experiments (Imbens and Rubin, 2015), also known as A/B tests, to compare the performance of a new product or iteration (the treatment) to some well-chosen baseline (the control). Randomized experiments are also used to evaluate the impact of new drugs, in the form of clinical trials, or to inform public policy. Measuring causal effects from these randomized experiments assumes that participants are willing to share their potentially sensitive or private response to treatment. This assumption is constantly challenged by the rise of privacy concerns and regulations for protecting individuals' online data. Many participants and regulatory guidelines agree with sharing some degree of information, as long as there is so-called plausible deniability, meaning no response can be tracked to any individual user, by sharing only aggregated data for example. However, data aggregation is often not sufficient to entirely prevent the risk of de-anonymization (Sweeney, 2000; Narayanan and Shmatikov, 2008).

---

*Data Sciences and Operations Department, University of Southern California
†Google Research

Differential privacy is one possible framework which has emerged as a solid contender under which user outcomes might be shared while diminishing the risk of deanonymization. It formalizes the notion that two privatized datasets are unlikely to differ in any measurable way if the true responses differ by a single point. Ensuring such a privacy guarantee often comes at the risk of adding additional noise into the original dataset, which increases the variance of statistical estimators. This privacy-variance trade-off is crucial for causal inference applications, since randomized experiments aim to obtain the most precise measurements possible of a causal effect.

Our paper has two main objectives: (a) to mathematically analyze the privacy-variance tradeoff for an intuitive set of algorithms that allow for the estimation of causal effects from differentially private data; and (b) to develop a novel mechanism CLUSTER-DP that enhances this privacy-variance trade-off by leveraging non-sensitive cluster information about the dataset.

Many of the algorithms we consider assume the existence of a central unit that observes all outcomes, and computes and shares a privatized dataset on which causal inference analyses can be run by a third-party. Furthermore, our suggested CLUSTER-DP mechanism assumes that outcomes exhibit some cluster structure, such as geographic regions or broad demographic classes. Notably, these clusters need not satisfy specific cardinality or quality constraints, and can include singleton clusters, random clusters, or even a single cluster containing all units. We do show however that our results improve when the clusters exhibit a measure of cluster quality that we define.

In Section 2, we motivate and define the differential privacy setting and causal objective of our work. In Section 3, we consider several intuitive mechanisms for privatizing a dataset while allowing for unbiased and consistent estimation of the average treatment effect. In particular, we consider the UNIFORM-PRIOR-DP mechanism which samples responses with some probability at random from the space of possible outcomes. In Section 4, we introduce our novel private-and-causal CLUSTER-DP mechanism, and its special case when all units belong to the same cluster, the CLUSTER-FREE DP mechanism. We evaluate their privacy guarantees and their variance gap to their non-differentially-private counterparts. We conclude in Section 5 with numerical experiments on simulated and real graphs to validate our claims and compare the empirical performance of each algorithm.

## 1.1 Related works

There are different approaches to preserving the privacy of user data regardless of any downstream analysis of it. A popular approach involves anonymizing data by removing, aggregating, or randomizing identifying details in some way before releasing the data to the public, in the hope of making 're-identification' of users difficult. Several privacy measures have been proposed under this approach, including $k$-anonymity (Sweeney, 2002, 2000), $\ell$-diversity (Machanavajjhala et al., 2007), and $t$-closeness (Li et al., 2006). While providing a layer of privacy protection, there are well-documented cases where de-identified data has been combined with other data sources to uniquely re-identify large proportion of users (Narayanan and Shmatikov, 2008; Sweeney, 2000).

In this work, we consider the differential privacy measure, which is a property of a data processing algorithm, rather than of a data set (Dwork et al., 2006a,b). Differential privacy is widely used and extensively covered. In this work, we will focus primarily on label differential privacy, introduced by Chaudhuri and Hsu (2011). The literature on label differential privacy is mostly dedicated to classification and regression tasks with the goal of improving excess risk while offering protection for labels (Beimel et al., 2013; Bassily et al., 2018; Wang and Xu, 2019). More recently, there have been several papers improving utility-privacy tradeoffs of label DP algorithms (Esfandiari et al.,

2022; Ghazi et al., 2021, 2022). Our mechanism is inspired by a technique by (Esfandiari et al., 2022), which we adapt to the estimation of causal effects. We further provide privacy-variance tradeoffs and a tighter analysis of the privacy guarantee than the proof methodology of (Esfandiari et al., 2022), which we extend to an $(\varepsilon, \delta)$-type guarantee (See Theorem 4.1).

Panigrahi et al. (2022) study the problem of treatment effect estimation after adjusting for potential confounders from different independent studies. Using a Lasso estimator, a parsimonious model is selected in each study and an unbiased estimator is constructed by aggregating simple summary statistics. While sharing only the summary statistics provides some layer of protection, this work does not provide any differential privacy guarantees.

Closer to our work, Kancharla and Kang (2021) also study the problem of average treatment effect estimation from a randomized control trial, where outcomes have been privatized using a differentially private mechanism. Specifically, they consider a binary outcome space and a mechanism, which for any given true response $y_i$, either returns $\tilde{y}_i = 0$ with probability $r_0$, $\tilde{y}_i = 1$ with probability $r_1$, or the true outcome $\tilde{y}_i = y_i$ with probability $1 - r_0 - r_1$ for some choice of probabilities $r_0, r_1 \in (0, 1)$. Our setting and results differ significantly in that (1) we go beyond binary responses and consider a general discrete outcome space; (2) our mechanism leverages a clustering structure of responses to improve the privacy-variance trade-off, and we quantify the impact of cluster quality on the variance of the estimator. In fact, their proposed procedure does not account for the empirical distribution of responses, such that an extension of their algorithm to a general discrete outcome space would be closer to the UNIFORM-PRIOR DP mechanism, which we consider along with other baselines in Section 3. In the absence of non-compliers, the procedure in (Kancharla and Kang, 2021) can be viewed as a special case of ours, for binary outcomes and assuming no-cluster structure (see Equation 7 for further detail).

Betlei et al. (2021) focus on the randomized control trial set-up and propose a differentially private method, called ADUM, to learn uplift models from data aggregated along a given partition of the feature space. The privacy-utility trade-off is studied by computing the mean-squared error of the estimator and its dependence on the underlying partition size and privacy budget. The analysis is for uni-dimensional feature spaces and makes the assumption that every bin has the same number of treated and controlled units. The ADUM mechanism adds Laplace noise to the count and the sum of responses from treated and controlled groups within each bin, and uses the (noisy) aggregate responses to estimate the conditional average treatment effect (CATE). While the bins can be thought of as clusters based on the features, their work does not use this clustering structure to reduce the variance of the estimator resulting from the added noise to achieve differential privacy. The ADUM mechanism is very similar to the noisy Horvitz-Thompson estimator baseline discussed in Section 3, with the difference that ADUM adds noise to the count and the sum of responses, while the noisy Horvitz-Thompson estimator baseline adds noise directly to the average response of treated and controlled groups in each cluster (cf. Equation (1)).

Niu et al. (2022) propose a meta-algorithm for multi-stage learning under privacy constraints and apply that to CATE estimation. The methodology relies on multiple sample-splittings where different parts of the sample are used to estimate different components of the estimator. The approach uses DP-EBMs (Nori et al., 2021) as the base learner, and conducts a privacy analysis using the sample splitting structure of the algorithm and the parallel composition property of differential privacy. They study the privacy-accuracy trade-off empirically. This framework goes beyond randomized control trials and allows for heterogeneous causal effects. However, they do not leverage the cluster structure of the data to improve the variance-privacy tradeoff, nor do

they provide private 'unit-level' data to the experimenters; instead their work aims to estimate the propensity model and the outcome model in a differentially-private way.

## 2    Causal Objective and Privacy Setting

We are motivated by the real-world scenario of a technology company in the business of selling advertising space to advertisers. Its clients, the advertisers, wish to measure the effectiveness of their advertising campaigns by running A/B tests, but do not want to rely on this technology company to provide their causal estimates. Instead, they would like access to user-level data, such as whether a user clicked on their ad, as well as any meaningful covariates about that user, so that they can conduct these analyses themselves. One reason for this might be that advertisers wish to run their own covariate-adjustment methods, or they would like to investigate proprietary sub-slices of users. On the other hand, this technology company seeks to protect the privacy of its users, and does not wish to share sensitive information about its users. Hence it must act as a central unit which privatizes its datasets before passing them on to advertisers for them to perform their own causal inference analyses on. We now introduce the formal causal objective and privacy setting. To guide the reader through abundant notation, we include a glossary at the beginning of the Appendix.

**Causal Objective.**    We consider a fixed population of $n$ users, henceforth units, where we can assume the Stable Unit Treatment Value Assumption (Imbens and Rubin, 2015). Let $y_i(0)$ be the potential outcome of unit $i$ if it is controlled, and $y_i(1)$ if it is treated. These potential outcomes are sampled from a finite response space $\mathcal{Y}$ of cardinality $K = |\mathcal{Y}|$. While finite response spaces are common in many advertising settings (e.g. number of clicks or impressions), we suggest binning when outcomes are continuous, as illustrated in Section 5.

For our CLUSTER-DP algorithm, we will further assume that there is some known clustered structure of these units into $C = |\mathcal{C}|$ non-overlapping clusters of size $n_c$, and let $c_i \in \mathcal{C}$ be the cluster membership of unit $i$. These clusters may be geographic regions or broad demographic groups the units belong to. We do not make any assumptions on the number of clusters or their size. In particular, our results hold for a single cluster or all singleton clusters. The strength of our results, however, improve with a specific measure of cluster quality, introduced in Section 4.

While we focus on the common finite sample setting, we can easily extend our results to their super-population equivalents. For this purpose, we will sometimes denote by $x_i \in \mathbb{R}^d$ the covariate vector of each unit $i$ such that each $(x_i, y_i(0), y_i(1), c_i)$ is drawn from some joint distribution $\mathcal{P}$. Our causal estimand is the average treatment effect estimand defined in the finite sample regime by

$$\tau = \frac{1}{n} \sum_{i=1}^{n} \left( y_i(1) - y_i(0) \right) .$$

Let $z_i$ correspond to the treatment assignment of unit $i$, $z_i = 1$ if treated and $z_i = 0$ if controlled. Let $n_1$ be the total number of treated units and $n_0$ to be the total number of controlled units across units. Treatment and control are assigned completely at random. When a clustering of the data is available, the treatment assignment is sampled in a completely randomized way over clusters: a fixed number of $n_{1,c}$ (resp $n_{0,c}$) units is chosen uniformly at random to be treated (resp. controlled) within each cluster $\mathcal{C}$, with $n_c = n_{1,c} + n_{0,c}$ the total number of units in cluster $c$.

4

**Privacy Setting.** In the real-world advertising setting presented above, it is assumed that a central unit privatizes the dataset before sharing it externally. Some of the mechanisms we explore also work without the presence of this central unit, a setting known as local differential privacy, but this is not a requirement for our framework.

Furthermore, an important question is which variables (treatment assignment, outcome, cluster assignment and/or covariates if available) are sensitive and should be privatized. In our advertising setting, a unit's treatment assignment is assigned purely at random; it is therefore not sensitive, and can be shared as-is. Outcomes are clearly sensitive and should be privatized. Finally, of the mechanisms we present, only our proposed CLUSTER-DP mechanism uses clusters—or covariates to form these clusters—to improve its privacy-variance tradeoff. While not all covariates may be sensitive (e.g. broad geographic regions), we can decouple the privacy protection of covariates from that of outcomes. In particular, by virtue of the composition property of differential privacy, we can allocate $\varepsilon_1$ privacy loss for the first step (forming clusters from covariates) and $\varepsilon_2$ privacy loss for the second step (computing the estimator). The end-to-end process will be $\varepsilon(= \varepsilon_1 + \varepsilon_2)$-DP. There already exists a rich literature on DP-clustering algorithms, see e.g., Nissim et al. (2007); Feldman et al. (2009, 2017); Stemmer and Kaplan (2018); Cohen et al. (2021). We focus therefore on the second step, namely estimating average treatment effect from differentially private outcomes, assuming a given cluster structure. This more restricted setting is known as label-differential privacy, introduced by Chaudhuri and Hsu (2011) and (re-)defined formally below. In this context, a unit's "label" refers to its observed outcome; we use the words outcome and label interchangeably.

**Definition 2.1.** *(Label Differential Privacy) Consider a randomized mechanism $M : D \to \mathcal{O}$ that takes as input a dataset $D$ and outputs into $\mathcal{O}$. Let $\varepsilon, \delta \in \mathbb{R}_{\geq 0}$. A mechanism $M$ is called $(\varepsilon, \delta)$-label differentially private—or $(\varepsilon, \delta)$-label DP—if for any two datasets $(D, D')$ that differ in the label (outcome) of a single example and any subset $O \subseteq \mathcal{O}$ we have $\mathbb{P}[M(D) \in O] \leq e^\varepsilon \mathbb{P}[M(D') \in O] + \delta$, where $\varepsilon$ is the privacy budget and $\delta$ is the failure probability. If $\delta = 0$, then $M$ is said to be $\varepsilon$-label differentially private, or $\varepsilon$-label DP.*

Achieving label-differential privacy implies that the output of a mechanism does not change too much if a single label in the input dataset is changed. The *privacy loss $\varepsilon$* controls the size of the possible change, and $\delta$ is the *failure probability* in providing such a guarantee. In other words, $(\varepsilon, 0)$-differential privacy ensures that, for *every* run of the mechanism $M$, the observed output is (almost) equally likely to be observed on every other neighboring dataset, simultaneously. The $(\varepsilon, \delta)$-differential privacy property relaxes this constraint and states only that it is unlikely that the observed value $M(D)$ has a much higher or lower chance to be generated under a dataset $D$ compared to a neighboring dataset $D'$. Differential privacy can also be viewed from a statistical hypothesis testing framework, where an attacker aims to distinguish $D$ from $D'$ based on the output of the mechanism. This viewpoint has been put forward by Wasserman and Zhou (2010) and Kairouz et al. (2015), who show that, by using the output of an $(\varepsilon, \delta)$-DP mechanism, the power of any test with significance level $\alpha \in [0, 1]$ is bounded by $e^\varepsilon \alpha + \delta$. For small enough $(\varepsilon, \delta)$, this bound is only slightly larger than $\alpha$, and so any test which aims to distinguishing $D$ from $D'$ is powerless.

# 3   Aggregation-based baselines and the Uniform-Prior-DP mechanism

In this section, we introduce three differentially private algorithms that still allow for the estimation of causal effects. These will serve as important baselines to our proposed CLUSTER-DP algorithm, which will be introduced later in Section 4.

## 3.1   Two aggregation-based baselines

The simplest approach to sharing a differentially private estimate of the average treatment effect is for the central unit to compute some unbiased estimator based on the original responses $y_i$ and add noise to the estimate before sharing it externally. We provide an example in Algorithm 1, written in the broadest generality when a clustering is available. When no clustering is available, we can simply assume that all units belong to the same cluster.

---

**Algorithm 1:** NOISY HORVITZ-THOMPSON mechanism

**Input**: Individual responses $y_1, \ldots, y_n$, (optional) cluster memberships $c_1, \ldots, c_n$
**Output**: Privatized estimate $\hat{\tau}$

$$\text{Return} \quad \hat{\tau} := \sum_{c \in \mathcal{C}} \frac{n_c}{n} \left\{ \sum_{i \in c} \left( \frac{y_i z_i}{n_{1,c}} - \frac{y_i (1 - z_i)}{n_{0,c}} \right) + w_c \right\}, \quad w_c \sim \text{Laplace}(\eta_c). \quad (1)$$

---

The variances of the noise parameters $\eta_c$ determine both the privacy guarantee $\varepsilon$ and additional estimator variance of the Noisy Horvitz-Thompson algorithm. To compute its privacy guarantee, we apply (Dwork et al., 2014, Theorem 3.6) and consider the sensitivity $\Delta_c$ of the inner function $^1/_{n_{1,c}} y_i z_i - ^1/_{n_{0,c}} y_i (1 - z_i)$, defined as the maximum change in its value when changing only one label in the data set. The variance of $\hat{\tau}$ can be expressed easily as a function of the variance of its non-differentially-private equivalent, the Horvitz-Thompson estimator without the Laplace noise:

$$\hat{\tau}_{\text{No-DP}} := \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{i \in c} \left( \frac{y_i z_i}{n_{1,c}} - \frac{y_i (1 - z_i)}{n_{0,c}} \right). \quad (2)$$

**Proposition 3.1.** *The noisy Horvitz-Thompson estimator $\hat{\tau}$ is $\varepsilon$-DP when setting $\eta_c = {}^{\Delta_c}/_\varepsilon$ for every cluster, where $\Delta_c = \min\{n_{0,c}, n_{1,c}\}^{-1} \times \max_{y \in \mathcal{Y}} |y|$. Furthermore, its variance with respect to the treatment assignment $\boldsymbol{z}$ and the Laplace noise (DP) is given by*

$$\text{Var}_{DP, \boldsymbol{z}}[\hat{\tau}] = \text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + 2 \sum_{c \in \mathcal{C}} \left( \frac{n_c}{n} \frac{\Delta_c}{\varepsilon} \right)^2,$$

*where $\hat{\tau}_{\text{No-DP}}$ is the non-differentially-private stratified Horvitz-Thompson estimator defined in Eq. 2, and $\text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$ is its variance with respect to the treatment assignment $\boldsymbol{z}$.*

Because these results hold for any clustering, they also hold when no clustering is available; in that case, we consider all units to be part of the same cluster. We refer the reader to the appendix for the well-known closed-form expression of $\text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$.

A second and slightly more sophisticated approach would be for the central unit to add noise to the frequency of responses in each cluster before sharing the histogram externally, since the estimated treatment effect depends only on the histogram of responses of treated and controlled units in each cluster. We provide an example in Algorithm 2, written in the broadest generality when a clustering is available.

---

**Algorithm 2:** NOISY HISTOGRAM mechanism

---

**Input**: Individual responses $y_1, \ldots, y_n$, (optional) cluster memberships $c_1, \ldots, c_n$
**Output**: Privatized estimate $\hat{\tau}$
Compute the empirical distribution $\hat{p}_a(y|c)$ of treated ($a = 1$) and controlled ($a = 0$) units within cluster $c$.

$$\text{Return} \quad \hat{\tau} := \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{y \in \mathcal{Y}} y \times (\hat{p}_1(y|c) + w_{1,c,y} - \hat{p}_0(y|c) - w_{0,c,y}), \quad w_{a,c,y} \sim \text{Laplace}(\eta_{a,c}). \tag{3}$$

---

Since the $K$ bins corresponding to the $K$ elements of $\mathcal{Y}$ are disjoint, and the sensitivity of the value of each histogram bin is $n_{a,c}^{-1}$, the central unit can share the histogram privately by adding independent draws from $\text{Laplace}((n_{a,c}\varepsilon)^{-1})$ to the frequency of each value. Furthermore, we can compute in closed form the variance gap of the Noisy Histogram mechanism compared to its non-private Horvitz-Thompson counterpart.

**Proposition 3.2.** *The noisy Histogram mechanism $\hat{\tau}$ is $\varepsilon$-DP when setting $\eta_{a,c} = (n_{a,c}\varepsilon)^{-1}$ for every cluster. Furthermore, its variance with respect to the treatment assignment $\boldsymbol{z}$ and the Laplace noise (DP) is given by*

$$\text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}] = \text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \frac{2}{\varepsilon^2} \left( \sum_{y \in \mathcal{Y}} y^2 \right) \sum_{c \in \mathcal{C}} \left( \frac{n_c}{n} \right)^2 \left( \frac{1}{n_{0,c}^2} + \frac{1}{n_{1,c}^2} \right),$$

*where $\hat{\tau}_{\text{No-DP}}$ is the non-differentially-private stratified Horvitz-Thompson estimator defined in Eq. 2, and $\text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$ is its variance with respect to the treatment assignment $\boldsymbol{z}$.*

For the same privacy guarantee, the NOISY-HORVITZ-THOMPSON mechanism has a smaller variance gap than the NOISY-HISTOGRAM mechanism, since $\|y\|_\infty \leq \|y\|_2$ and $\min(n_{0,c}, n_{1,c})^{-2} \leq \min(n_{0,c}, n_{1,c})^{-2} + \max(n_{0,c}, n_{1,c})^{-2} = n_{0,c}^{-2} + n_{1,c}^{-2}$.

**Limitations** These aggregation-based approaches have two drawbacks in the real-world setting of Section 2. First, advertisers expect user-level data even if privatized. Perhaps this is because they wish to analyse their own slices of the user population or perhaps they wish to apply their own proprietary covariate-adjustment method to measure campaign effectiveness. Second, the level of noise in Algorithm 1 is averaged over the number of clusters $C$; in Algorithm 2, it is averaged over the number of possible outcomes $K$. The next mechanisms that we introduce will privatize user-level responses, and the noise level therein is averaged over $n$ samples, the number of users. This becomes an important competitive advantage of user-level methods in the case of one-shot communication between the central unit and the advertisers. In particular, a user-level privatizing

scheme achieves lower finite-sample conditional bias than the prior two aggregation baselines when $n \gg K$ and $n \gg C$, when we condition on the randomness of the DP mechanism and consider the bias with respect to the randomization in the sub-population. We will illustrate this point further in Experiment 5 of Section 5.

## 3.2   The Uniform-prior DP mechanism

Unlike the two prior mechanisms, the next mechanism we consider provides user-level outcomes. Formalized in Algorithm 3, it reports the true outcome with some probability, and otherwise reports an outcome sampled uniformly at random from the space of possible outcomes. We refer to it as the UNIFORM-PRIOR-DP mechanism because it does not leverage any information about the empirical distribution of outcomes beyond its support.

---

**Algorithm 3:** UNIFORM-PRIOR-DP mechanism

---

**Input**: Individual responses $y_1, \ldots, y_n$
**Output**: Privatized responses $\tilde{y}_1, \ldots, \tilde{y}_n$
**for** $i \in \{1, \ldots n\}$ **do**
$$\tilde{y}_i \leftarrow \begin{cases} y_i^0 \sim \mathcal{U}(\mathcal{Y}) & \textit{with probability } \lambda \quad // \ \mathcal{U} \textit{ is the uniform distribution} \\ y_i & \textit{with probability } 1 - \lambda \end{cases}$$
*Return privatized responses* $\{\tilde{y}_1, \ldots, \tilde{y}_n\}$.

---

The UNIFORM-PRIOR-DP mechanism is a generalization of the mechanism proposed by Kancharla and Kang (2021) in the binary-outcome setting, when there are no non-compliers. In the broadest generality when a clustering is available, the following stratified estimator is unbiased for the average treatment effect:

$$\hat{\tau}_\lambda = \frac{1}{1-\lambda} \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{i \in c} \left( \frac{\tilde{y}_i z_i}{n_{1,c}} - \frac{\tilde{y}_i(1-z_i)}{n_{0,c}} \right) . \tag{4}$$

**Proposition 3.3.** *The conditional expectation of the estimator $\hat{\tau}_\lambda$ defined in Eq. (4), with respect to the DP mechanism, is equal to the non-differentially private Horvitz-Thompson estimator. It is therefore unbiased for the average treatment effect $\tau$ over $z$ and the DP mechanism.*

$$\mathbb{E}_{DP}[\hat{\tau}_\lambda | z] = \hat{\tau}_{\text{No-DP}} \quad \text{and} \quad \mathbb{E}_{DP,z}[\hat{\tau}_\lambda] = \tau$$

Having an unbiased estimator for causal inference is an important but not entirely surprising result. In fact, many differentially private mechanisms can recover true labels in expectation; Kancharla and Kang (2021) also propose an unbiased differentially private estimator in the setting of binary potential outcomes $y_i \in \{0, 1\}$. Instead, the main difficulty is to minimize the variance gap with non-differentially-private estimators. To state the variance of $\hat{\tau}$ under the UNIFORM-PRIOR-DP mechanism, we consider the following notation: $\bar{y} := 1/|\mathcal{Y}| \sum_{y \in \mathcal{Y}} y$ and $\overline{y^2} := 1/|\mathcal{Y}| \sum_{y \in \mathcal{Y}} y^2$ over all possible outcomes. For $a \in \{0, 1\}$, we also define $\overline{y_c(a)} := 1/n_c \sum_{i \in c} y_i(a)$ and $\overline{y_c^2(a)} := 1/n_c \sum_{i \in c} y_i^2(a)$ over the units of cluster $c$.

**Theorem 3.4.** *For any $\tilde{\varepsilon} > 0$, the UNIFORM-PRIOR-DP mechanism is $(\tilde{\varepsilon}, \delta)$-label DP when we set $\delta = \max(0, 1 - \lambda + \frac{\lambda}{K}(1 - e^{\tilde{\varepsilon}}))$. In particular, it is $\varepsilon$-label DP with $\varepsilon = \log\left(1 + \frac{(1-\lambda)K}{\lambda}\right)$.*

*Furthermore, the variance of estimator $\hat{\tau}_\lambda$ in* (4) *under the* UNIFORM-PRIOR-DP *mechanism and the treatment assignment* $\boldsymbol{z}$ *is given by*

$$\mathrm{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_\lambda] = \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \sum_{c \in \mathcal{C}} \frac{n_c^2}{n^2} \left( \frac{1}{n_{0,c}} + \frac{1}{n_{1,c}} \right) \frac{\lambda \overline{\mathrm{y}^2} - \lambda^2 \bar{\mathrm{y}}^2}{(1-\lambda)^2}$$

$$+ \sum_{c \in \mathcal{C}} \frac{n_c^2}{n^2} \left[ \frac{\lambda}{1-\lambda} \left( \frac{\overline{y_c^2(0)}}{n_{0,c}} + \frac{\overline{y_c^2(1)}}{n_{1,c}} \right) - \frac{2\lambda\bar{\mathrm{y}}}{1-\lambda} \left( \frac{\overline{y_c(0)}}{n_{0,c}} + \frac{\overline{y_c(1)}}{n_{1,c}} \right) \right], \tag{5}$$

As the sampling probability grows small $\lambda \to 0$, we recover the non-private variance formula $\mathrm{Var}_{DP,\boldsymbol{z}}(\hat{\tau}) \to \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$, but the $\varepsilon$-DP guarantee goes to infinity. Since the UNIFORM-PRIOR-DP mechanism itself does not depend on the clusters, the privacy guarantee does not depend on the clustering properties of the data, if any. The dependence on the clustering in Equation (5) is only due to the definition of the stratified estimator. When a good clustering is not available, the above estimator can be simplified to its unstratified version $\hat{\tau}^u$ by considering that all units belong to the same cluster:

$$\hat{\tau}^u = \frac{1}{1-\lambda} \sum_{i=1}^{n} \left( \frac{\tilde{y}_i z_i}{n_1} - \frac{\tilde{y}_i(1-z_i)}{n_0} \right). \tag{6}$$

The following variance result is a direct corollary of Theorem 3.4.

**Corollary 3.5.** *Under the* UNIFORM-PRIOR-DP *mechanism, the variance of the unstratified estimator $\hat{\tau}^u$ defined in Eq.* 6 *is given by*

$$\mathrm{Var}_{DP,\boldsymbol{z}}[\hat{\tau}^u] = \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}^u_{\text{No-DP}}] + \frac{n}{n_1 n_0} \frac{\lambda \overline{\mathrm{y}^2} - \lambda^2 \bar{\mathrm{y}}^2}{(1-\lambda)^2} + \frac{\lambda}{1-\lambda} \left( \frac{\overline{y^2(0)}}{n_0} + \frac{\overline{y^2(1)}}{n_1} \right) - \frac{2\lambda\bar{\mathrm{y}}}{1-\lambda} \left( \frac{\overline{y(0)}}{n_0} + \frac{\overline{y(1)}}{n_1} \right)$$

*where* $\mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}^u_{\text{No-DP}}]$ *denotes the variance of its non-private equivalent $\hat{\tau}^u_{\text{No-DP}}$. Its differential privacy guarantees are the same as those in Theorem* 3.4.

We refer the reader to the appendix for the well-known closed form formula of $\mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}^u_{\text{No-DP}}]$. The special case of the unstratified estimator $\hat{\tau}^u$ in (6) for binary outcomes $\mathcal{Y} = \{0,1\}$ was previously proposed by Kancharla and Kang (2021), in which case the variance of the estimator can be further simplified:

$$\mathrm{Var}_{DP,\boldsymbol{z}}[\hat{\tau}^u] = \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}^u_{\text{No-DP}}] + \frac{n}{n_0 n_1} \frac{\frac{\lambda}{2}(1 - \frac{\lambda}{2})}{(1-\lambda)^2}. \tag{7}$$

The first two aggregation-based mechanisms in Section 3.1 assumed that a trusted data curator (e.g. a technology company, in the motivating example in Section 2) has access to the true outcomes and computes a differentially private estimate or empirical distribution of these responses. In contrast, the UNIFORM-PRIOR-DP mechanism can be implemented without such a curator: each user can privatize their response before sharing it with the experimenter. In other words, the UNIFORM-PRIOR-DP mechanism provides a local DP guarantee, defined by Kasiviswanathan et al. (2011), which is stronger than a DP guarantee. That said, in our motivating example, assuming the existence of a trusted curator—the technology company—is more natural than putting the burden of privatizing responses on each individual user.

# 4 The Cluster-DP and Cluster-Free-DP mechanisms

We now introduce our main differentially private mechanism, CLUSTER-DP, which not only provides user-level privatized outcomes, but also leverages information about the empirical distribution of outcomes within each cluster to improve its privacy-tradeoff. When no good clustering is available, we consider its special case when all units can be considered part of the same cluster, the CLUSTER-FREE-DP mechanism.

---

**Algorithm 4:** Our suggested differential privacy mechanism: CLUSTER-DP

---

> **Parameters**: threshold $\gamma \in [0, 1/K]$; noise scale $\sigma \geq 0$; re-sampling probability $\lambda \in [0, 1]$
> **Input**: Individual responses $y_1, \ldots, y_n$, treatment assignments $z_1, \ldots, z_n$.
> **Output**: Privatized responses $\tilde{y}_1, \ldots, \tilde{y}_n$
> // *Compute noisy response distribution per cluster and treatment group*
> **for** $c \in \mathcal{C}$ **do**
> > **for** $a \in \{0, 1\}$ **do**
> > > // *Add noise to each empirical probability distribution $\hat{p}_a(y|c)$ and truncate*
> > > **for** $y \in \mathcal{Y}$ **do**
> > > > $q_a(y|c) \leftarrow \max\{\gamma, \min\{1, \hat{p}_a(y|c) + w\}\}, \quad$ where $w \sim \text{Laplace}(\sigma/n_{a,c})$
> >
> > **for** $a \in \{0, 1\}$ **do**
> > > // *Renormalize each distribution*
> > > **for** $y \in \mathcal{Y}$ **do**
> > > > **if** $\sum_y q_a(y|c) > 1$ **then** $\zeta_y \leftarrow q_a(y|c) - \gamma$;
> > > > **else** $\zeta_y \leftarrow 1 - q_a(y|c)$
> > > **for** $y \in \mathcal{Y}$ **do**
> > > > $\tilde{q}_a(y|c) \leftarrow q_a(y|c) + \frac{\zeta_y}{\sum_{y'} \zeta_{y'}} \left(1 - \sum_y q_a(y|c)\right)$
>
> // *Randomize responses*
> **for** $i \in \{1, \ldots n\}$ **do**
> > $\tilde{y}_i \leftarrow \begin{cases} y_i^0 \sim \tilde{q}_{z_i}(\cdot|c_i) & \text{with probability } \lambda \\ y_i & \text{with probability } 1 - \lambda \end{cases}$
>
> *Return privatized responses* $\{\tilde{y}_1, \ldots, \tilde{y}_n\}$.

---

Formalized in Algorithm 4, our proposed mechanism deals with each cluster individually and independently of other clusters, handling treated and controlled groups separately. It returns a privatized potential outcome $\tilde{y}_i$ for each unit, which is either the true outcome with some probability or sampled from a transformed empirical distribution of responses from units in the same cluster. The transformation is inspired from a mechanism in Esfandiari et al. (2022). For the sake of exposition, we focus on the controlled units of a cluster $c \in \mathcal{C}$.

(1) Compute the empirical response distribution of the controlled units in the cluster $\hat{p}_0(y|c)$.

(2) Add noise drawn from a Laplace distribution with parameter $(\sigma/n_{0,c})$ to each response probability. Recall that $n_{0,c}$ is the number of controlled units in cluster $c$.

(3) Truncate the response probabilities to be within the interval $[\gamma, 1]$, with $\gamma \leq 1/K$.

(4) Renormalize the response probabilities to form a distribution. We follow a specific renormalization so that the resulting response probabilities remain in $[\gamma, 1]$, and add up to one.

(5) With probability $\lambda$, each original response is replaced by a random sample from the distribution constructed in the previous step.

## 4.1 Privacy guarantees

The following theorem and its corollary state the differential privacy guarantee of our proposed CLUSTER-DP mechanism.

**Theorem 4.1.** *Let $\tilde{\varepsilon} > 0$ and $\delta := \max(0, 1 - \lambda + \lambda\gamma(1 - e^{\tilde{\varepsilon}}))$. The CLUSTER-DP mechanism described in Algorithm 4 is $(\varepsilon, \delta)$-label DP with $\varepsilon = \min\left(\frac{1}{\sigma}, \frac{2}{\gamma}\right) + \tilde{\varepsilon}$. By setting $\tilde{\varepsilon} = \log(1 + \frac{1-\lambda}{\lambda\gamma})$, we have $\delta = 0$, and therefore the CLUSTER-DP mechanism is also $\varepsilon$-label DP, with $\varepsilon = \min\left(\frac{1}{\sigma}, \frac{2}{\gamma}\right) + \log\left(1 + \frac{1-\lambda}{\lambda\gamma}\right)$.*

We refer the reader to the Appendix for a full proof of Theorem 4.1 and provide here some intuition for its stated privacy loss $\varepsilon$. The first term $\min\left(1/\sigma, 2/\gamma\right)$ is the privacy budget used to privately estimate the empirical response distribution $\tilde{q}_a(\cdot|c)$ for each cluster. Fixing the transformed empirical distributions $\tilde{q}_a(\cdot|c)$, the log term is the privacy budget used to generate the privatized responses $\tilde{y}_i$. By the composition theorem for differential privacy (Dwork et al., 2014, Theorem B.1), the total privacy loss is given by the sum of these two losses. As expected, when the resampling probability goes to zero ($\lambda \to 0$), the privacy loss grows large ($\varepsilon \to +\infty$). Similarly, as the Laplace noise $\sigma$ and truncation parameter $\gamma$ grow large, the privacy guarantee improves ($\varepsilon \to 0$).

Because these privacy guarantees do not depend on the size, cardinality, or quality of the clusters, Theorem 4.1 also holds for the special case where there is no cluster structure to the data, in which case we can repeat the same mechanism as if all units belong to the same large cluster. Known as the CLUSTER-FREE-DP mechanism, it has the same privacy guarantee as the CLUSTER-DP mechanism. We will show the benefit that clusters may have in the following section on the variance gap.

**Similarity with uniform-prior-DP**    The careful reader may have noticed similarities between the privacy guarantees in Theorem 3.4 and Theorem 4.1. In fact, our CLUSTER-DP mechanism is a generalization of the UNIFORM-PRIOR-DP mechanism with the right choice of parameters, and by extension also a generalization of the mechanism proposed by Kancharla and Kang (2021) with no non-compliers. To prove this, we observe that the distributions $\tilde{q}_a(y|c)$ constructed in the CLUSTER-DP mechanism obey the following properties: $\tilde{q}_a(y|c) \geq \gamma$ for all $y \in \mathcal{Y}$, and $\sum_y \tilde{q}_a(y|c) = 1$. When setting the truncation parameter $\gamma = 1/K$, these distributions reduce to uniform distributions over the space of all outcomes, in which case the cluster-DP mechanism amounts to the simpler UNIFORM-PRIOR-DP mechanism, regardless of the value of the Laplace noise variance $\sigma^2$.

## 4.2 Estimation and variance guarantees

We now consider estimating causal effects from the privatized outcomes provided by our suggested CLUSTER-DP mechanism. For each cluster $c \in \mathcal{C}$ and each value $a \in \{0, 1\}$ of treatment, we
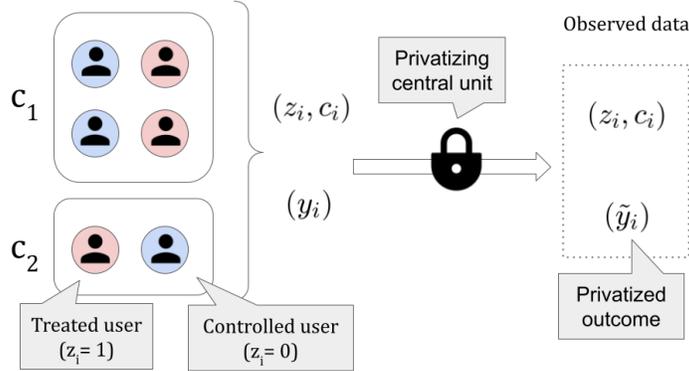
Figure 1: Illustration of CLUSTER-DP mechanism with a central unit computing the (clustered) privatized outcomes for valid causal inference.

construct the response randomization matrix $Q_{c,a} \in \mathbb{R}^{K \times K}$:

$$Q_{c,a}[y', y] := (1 - \lambda)\mathbb{I}(y' = y) + \lambda \tilde{q}_a(y'|c). \tag{8}$$

Conditional on its true outcome $y_i$, treatment assignment $z_i$, and cluster assignment $c_i$, the privatized response $\tilde{y}_i$ of unit $i$ is distributed according to $Q_{c_i,z_i}[\tilde{y}_i, y_i]$; in other words, $\forall y'$, $P(\tilde{y}_i = y'|c_i, z_i, y_i) = Q_{c_i,z_i}[y', y_i]$.

We use the inverse of the response randomization matrix to debias the privatized responses. Recall the notation $\mathsf{y}$ to represent in vector form the space of all possible potential outcomes, with similar ordering of rows and columns as $Q_{c_i,z_i}$. With a small abuse of notation, we write the index $\tilde{y}_i$ of the vector $\mathsf{y}^T Q_{c,z_i}^{-1}$ as $\mathsf{y}^T Q_{c,z_i}^{-1}[\tilde{y}_i]$ and show that it is an unbiased estimate for $y_i$ over the randomness of Algorithm 4. As a result, by reweighting each privatized outcome by the inverse of its conditional probability of occurring $Q_{c_i,z_i}[\tilde{y}_i, y_i]$, we propose the following Horvitz-Thompson-like estimator for the average treatment effect:

$$\hat{\tau}_Q := \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{i \in c} \left( \mathsf{y}^T Q_{c,z_i}^{-1}[\tilde{y}_i] \frac{z_i}{n_{1,c}} - \mathsf{y}^T Q_{c,z_i}^{-1}[\tilde{y}_i] \frac{1 - z_i}{n_{0,c}} \right). \tag{9}$$

A proof of the following theorem can be found in the Appendix.

**Theorem 4.2.** *Conditionally on the randomness of the treatment assignment, $\hat{\tau}_Q$ is equal in expectation over the randomness of the DP mechanism to the stratified difference-in-means estimator, such that $\hat{\tau}_Q$ is an unbiased and consistent estimator of $\tau$.*

$$\mathbb{E}_{DP}[\hat{\tau}_Q|\boldsymbol{z}] = \sum_{c \in \mathcal{C}} \frac{n_c}{n} \left( \sum_{i=1}^{n} y_i(1) \frac{z_i}{n_{1,c}} - \sum_{i=1}^{n} y_i(0) \frac{1 - z_i}{n_{0,c}} \right) = \hat{\tau}_{\text{No-DP}} \quad and \quad \mathbb{E}_{DP,\boldsymbol{z}}[\hat{\tau}_Q] = \tau$$

For the third parties to compute this estimator themselves, the central unit must pass along the cluster assignment, the treatment assignment, the privatized response $\tilde{y}_i$, as well as the vector of probabilities $\mathsf{y}^T Q_{c,z_i}^{-1}$, as illustrated in Figure 1. Since $\tilde{q}_a(\cdot|c)$ and the responses $\tilde{y}_i$ are $\varepsilon$-DP, by the post-processing property of differential privacy (Dwork et al., 2014, Proposition 2.1), all the

information passed to the third-party, as well as any estimation based on this information, is also $\varepsilon$-DP.

Our goal for Algorithm 4 is to make the gap between the variance of our differentially-private estimator $\hat{\tau}$ and its non-differentially private counterpart $\hat{\tau}_{\text{No-DP}}$ as small as possible for a given privacy guarantee. While all our results hold for any given clustering, they are greatly improved when clusters are homogeneous, as defined below.

**Definition 4.3** (Cluster homogeneity). *For $a \in \{0,1\}$, define a clustering's homogeneity as the average intra-cluster variance of outcomes $\phi_a \geq 0$, defined as*

$$\phi_a := \sum_{c \in \mathcal{C}} \frac{n_c^2}{n^2} \frac{S^2(\vec{y}_c(a))}{n_{a,c}} \,,$$

*where for any vector $\vec{u} \in \mathbb{R}^d$, $S^2(\vec{u}) := \frac{1}{d-1} \sum_{u \in \vec{u}} (u - \bar{u})^2$ and $\bar{u} := \frac{1}{d} \sum_{u \in \vec{u}} u$.*

The quantity $\phi_a$ has a natural super-population interpretation when taking its expectation of over the distribution $\mathcal{P}$: $\phi_a = \mathbb{E}[\text{Var}(y(a)|c)] = \text{Var}(y(a)) - \text{Var}(\mathbb{E}[y(a)|c]) > 0$. Holding $\text{Var}(y(a))$ constant, lower values of $\phi_a$ implies that clusters are better separated. For $\phi_a = 0$, the outcome values of each clusters are contained within a singleton set. On the other hand, if $\phi_a$ is high, clusters contain a wide range of responses, up to the variation of outcomes of the entire population. The following theorem provides a bound on the variance of our estimator $\hat{\tau}_Q$, with respect to the randomness of Algorithm 4 and the random assignment $\mathbf{z}$, as a function of $\text{Var}_{\mathbf{z}}[\hat{\tau}_{\text{No-DP}}]$ and $\phi_a$.

**Theorem 4.4.** *The variance of the estimator $\hat{\tau}$ defined in (9) is bounded by*

$$0 \leq \text{Var}_{DP,\mathbf{z}}[\hat{\tau}_Q] - \text{Var}_{\mathbf{z}}[\hat{\tau}_{\text{No-DP}}] \leq \left( \frac{1}{(1-\lambda)^2} - 1 \right) \sum_{a \in \{0,1\}} \phi_a + \sum_{a \in \{0,1\}} \sum_{c \in \mathcal{C}} \frac{n_c^2}{n^2} \frac{A(n_{a,c})}{n_{a,c}},$$

*where $\phi_a$ is the measure of cluster homogeneity defined in Definition 4.3, and for any $x$,*

$$A(x) := 2K \left[ \frac{3\|\mathsf{y}\|_\infty^2 + (\lambda\sqrt{K} + 1)^2 + \|\mathsf{y}\|_2^2(1 - \lambda(K-1)\gamma)}{(1-\lambda)^2} + 2\|\mathsf{y}\|_\infty^2 \right] \left[ \gamma + \frac{\sigma}{x} \left( e^{-\gamma x/\sigma} - e^{-x/\sigma} \right) \right] \,,$$

*with $K$ the number of possible potential outcomes and $\mathsf{y} \in \{\mathcal{Y}\}^K$ the vector of all possible outcomes.*

Theorem 4.1 and Theorem 4.4 together allow us to capture the privacy-variance trade-off of our proposed mechanism. Recall that the privacy guarantee of Theorem 4.1 is agnostic to the clustering. On the other hand, the variance gap in Theorem 4.4 depends first on the homogeneity of clusters, as defined in Definition 4.3, and a second term that is agnostic to the clustering. As a result, more homogeneous clusters—those with low $\phi_a$—result in a smaller variance gap with equal privacy guarantees, leading to a better privacy-variance trade-off than less homogeneous clusters, all else being equal.

We now provide some intuition for the second term $A(x)$. By choosing $\gamma$ and $\sigma$ to be arbitrarily small, we can make this second term arbitrarily small. As expected, the privacy guarantees of Theorem 4.1 suffer in that regime. When setting $\lambda = 0$, our CLUSTER-DP mechanism always outputs the true outcome, and we no longer produce privatized outcomes. In that case, we can set the truncation parameter $\gamma$ and the Laplace noise $\sigma$ to be zero with no consequence to recover the

trivial equality $\text{Var}_{DP,\boldsymbol{z}}(\hat{\tau}) = \text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$ from our bound above. Naturally, the more interesting setting from a privacy perspective is $\lambda \in (0, 1)$.

As discussed previously, the privacy guarantee in Theorem 4.1 for the CLUSTER-DP and CLUSTER FREE-DP mechanisms reduces to the guarantee of the UNIFORM-PRIOR-DP mechanism in Theorem 3.4 when setting the truncation parameter $\gamma = 1/K$ and $\sigma = \infty$. Yet, because both CLUSTER-DP and CLUSTER-FREE-DP mechanisms use data-dependent priors, there may exist choices of $(\sigma, \gamma, \lambda)$ which result in better privacy-variance trade-offs than the latter for certain outcome distributions. Rather than computing the variance gaps of each mechanism in closed-form, we encourage practitioners to compute their performance empirically for different values of each mechanism's parameters, and to keep track of the resulting privacy guarantee. In the following section, we conduct empirical evaluations of the privacy-variance trade-off of the different mechanisms.

## 5 Numerical experiments

In this section, we perform a series of simulated experiments to validate the theoretical claims we make in the paper and to illustrate their usefulness. Except for Experiment 5, we focus our attention on the UNIFORM-PRIOR DP, CLUSTER-FREE-DP, and CLUSTER-DP mechanisms, due to the strong limitations of the aggregation-based baselines in our real-world setting compared to these three user-level privatization schemes.

We start by considering a Gaussian Mixture Model setting where for every unit $i$ in cluster $c$, a continuous quantity $y_i'$ is given by

$$\forall i \in c, \; y_i' = \sqrt{\beta}\mu_c + \sqrt{v - \beta}w_i, \tag{10}$$

where $\mu_c$ and $w_i$ are drawn from the standard normal distribution. The coefficient $\beta \in [0, v]$ measures the dependence of the response on the cluster center. This specific parameterization is chosen to fix the variance of the response, equal to $v$, as $\beta$ varies. Since the proposed mechanism is for discrete outcome spaces, we quantize the response in the following way:

$$y_i(1) = y_i(0) + \tau \quad \text{and} \quad y_i(0) = \begin{cases} K' \text{ if } y_i' > 2\sqrt{v} \\ -K' \text{ if } y_i' < -2\sqrt{v} \\ [y/\Delta] \text{ otherwise} \end{cases}$$

where $\Delta := 2\sqrt{v}/K'$ and $[x]$ denotes the rounding of $x$ to the nearest integer. The treatment effect is an additive $\tau$ term on the potential outcome under control. We fix $\tau = 1$, such that the outcomes take values in the set $\mathcal{Y} = \{-K', \ldots, 0, \ldots, K', K' + 1\}$. We denote by $K := 2(K' + 1)$ the size of outcome space.

Unless otherwise specified, and with no particular reason to fix parameters one way or another, we take $K' = 5$, $v = 5$, and $\beta = 4.5$. We consider $C = 3$ clusters of sizes $500$, $10^3$, $2 \times 10^3$ with an equal number of controlled and treated units in each cluster. To display confidence intervals around certain results, we consider a super-population of three clusters of sizes $2.5 \times 10^3$, $5 \times 10^3$, and $10^4$ units, and repeatedly draw uniformly at random sub-populations of three clusters from these original clusters.

For any given sub-population, we compute the variance $\text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_Q]$ by empirically computing the variance (or histogram) of $\hat{\tau}_Q$ empirically over 500 realizations of the randomness in the corresponding DP mechanism (e.g. Laplace noise and response randomization), as well as the treatment
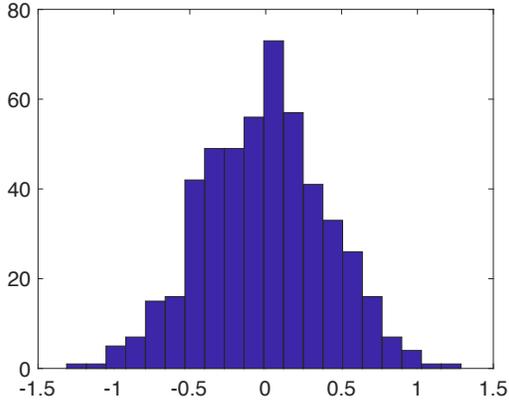
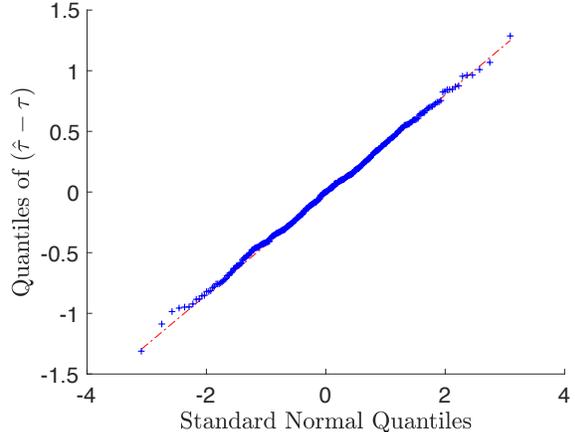Figure 2: histogram of $\hat{\tau} - \tau$



Figure 3: qq-plot of $\hat{\tau} - \tau$

assignments, which are done by choosing balanced set of treated and controlled units uniformly at random within each cluster. Unless otherwise specified, for the CLUSTER-DP mechanism, we set the truncation parameter $\gamma = 0.02$, the Laplace noise $\sigma = 10$, and the resampling probability $\lambda = 0.8$.

**Experiment 1. (Bias and Gaussianity)** We first verify that our CLUSTER-DP estimator $\hat{\tau}_Q$, given by (9), is unbiased and admits an asymptotically Gaussian distribution by plotting the histogram and the qq-plot of $\hat{\tau} - \tau$ in Figures 2 and 3.

**Experiment 2. (Privacy-variance trade-off)** We next compare the privacy-variance trade-off of our suggested CLUSTER-DP mechanism with that of the CLUSTER-FREE-DP mechanism, as well as the stratified and unstratified versions of the UNIFORM-PRIOR-DP mechanism. We observe that the CLUSTER-DP can have significantly lower variance for its estimator, compared to the other mechanisms, for the same privacy loss $(\varepsilon, \delta)$.

In Figure 4, we aim to fix the privacy loss to $\varepsilon = 0.2$ and $\delta = 10^{-4}$ for all three mechanisms. For the CLUSTER-DP and CLUSTER-FREE-DP, we set the Laplace parameter to $\sigma = 10$, and vary the truncation parameter $\gamma \in [0.1/K, 1/K]$. Following Theorem 4.1, we first choose $\tilde{\varepsilon}$ so that the corresponding privacy $\varepsilon$, is equal to its target $\varepsilon = 0.2$, and then choose the re-sampling probability $\lambda$ to obtain the failure probability $\delta = 10^{-4}$. Likewise, for the UNIFORM-PRIOR-DP, we set the re-sampling probability $\lambda$ according to Theorem 3.4, such that $\varepsilon = 0.2$ and $\delta = 10^{-4}$. In summary, as the truncation parameter $\gamma$ varies, we compare the three mechanisms at the same privacy loss. As we observe in Figure 4, for small values of $\gamma$, the CLUSTER-DP achieves significantly lower variance compared to to the other mechanisms. When $\gamma = 1/K$ and $\sigma = \infty$, the theory tells us that CLUSTER-DP reduces to UNIFORM-PRIOR-DP (stratified) and the CLUSTER FREE-DP reduces to UNIFORM-PRIOR-DP (unstratified). However, since we have set $\sigma = 10$, we observe that the variance for the UNIFORM-PRIOR-DP becomes lower than the other mechanisms for $\gamma = 1/K$. The error-bars here correspond to 50 independent draws of the sub-population.

In Figure 5, we plot the variance of the estimators versus the privacy loss $\varepsilon$, as we fix $\delta = 10^{-4}$. Here, we optimize the choice of Laplace parameter $\sigma \in \{10, 20, \infty\}$ and the truncation parameter $\gamma \in \{0.01/K, 0.1/K, 1/K\}$. We observe that both CLUSTER-DP and CLUSTER FREE-DP
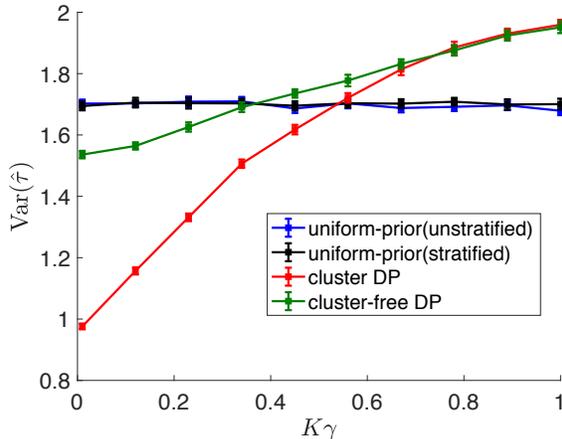
15

Figure 4: Variances for the DP mechanisms, as we vary the truncation level $\gamma \in [0.1/K, 1/K]$, under the setting of Experiment 2. The privacy loss is fixed at $\varepsilon = 0.2$ and $\delta = 10^{-4}$.
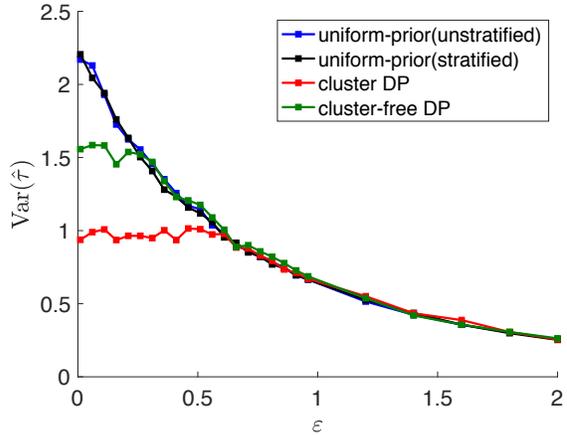


Figure 5: Privacy-variance trade-off for the four DP mechanisms under the setting of Experiment 2. We fix the failure probability in DP guarantee to $\delta = 10^{-4}$, and optimize the choice of $\sigma$ and $\gamma$ in the sets $\sigma \in \{10, 20, \infty\}$ and $\gamma \in \{0.01/K, 0.1/K, 1/K\}$.

estimators, which use data-dependent priors achieve a better trade-off than either version of the UNIFORM-PRIOR-DP mechanism. Furthermore, the CLUSTER-DP mechanism, which also leverages the clustering structure, showcases an even better trade-off compared to the CLUSTER FREE-DP mechanism.

**Experiment 3. (Role of clustering quality)**  In this experiment we show that, as the clustering quality improves, the variance of the estimator for the CLUSTER-DP mechanism decreases when compared to the variance of the estimator for the CLUSTER FREE-DP mechanism, without affecting their privacy guarantees, since these are agnostic to the clustering according to Theorem 4.1. Under our specified potential outcome model (10), the cluster homogeneity $\phi_a$, as defined in Definition 4.3, is given by $\phi_0 = \mathbb{E}(\text{Var}(y_i(0)|c)) \propto v - \beta = \phi_1$, hence our clusters become more homogeneous as $\beta$ increases. From Theorem 4.4, the clustering structure reduces the variance of the estimator at more homogeneous clusters, i.e. lower values of $\phi_0, \phi_1$, and $\lambda$. We verify this in Figure 6, which plots the ratio of the variances for two values of $\lambda \in \{0.5, 0.8\}$ as we vary $\beta$. As $\beta$ grows, we observe a stronger reduction in the variance using the clustering structure of data. This effect is stronger at smaller values of $\lambda$.

**Experiment 4. (Validation of theoretical bound)**  In Theorem 4.4, we bounded the excess variance of the private estimator (9) compared to the non-private estimator (2). The bound had two additive terms. The first one depends on the cluster structure of data, namely the cluster homogeneity quantities $\phi_0$, $\phi_1$, and the second term did not depend on the clusters, capturing instead an increase in the variance due to the randomness of the CLUSTER-DP mechanism. In Figure 7, we compute the gap $\text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}] - \text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$ empirically, by averaging over 500 different realizations of the randomness in the DP mechanism and the treatment assignments in the same setting as the previous experiment. We plot this gap as we vary $\beta$, along with a shaded region whose upper boundary corresponds to the upper bound given in Theorem 4.4 and its lower boundary
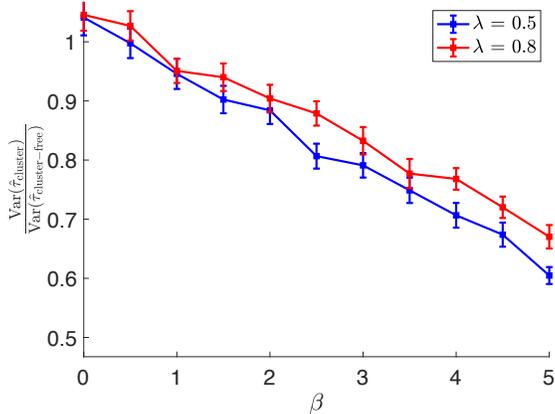
Figure 6: Ratio of the variance of the estimators under the CLUSTER-DP and CLUSTER FREE-DP mechanisms in Experiment 3. The benefit of CLUSTER-DP mechanism is stronger at larger $\beta$ and smaller value of $\lambda$.
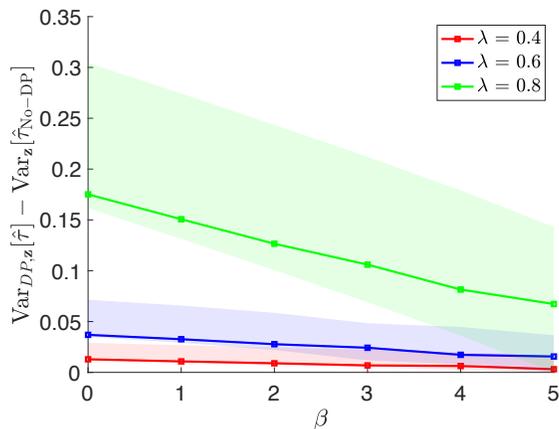


Figure 7: The variance gap between the private estimator $\hat{\tau}_Q$, given by (9), and the non-private estimator $\hat{\tau}_{\text{No-DP}}$ in the setting of Experiment 4. The upper boundary of the shaded area corresponds to the upper bound derived in Theorem 4.4, and it lower boundary corresponds to the the first term in that bound. As we see the gap remains between the two boundaries.

corresponds to only the first term in that bound. We observe that the variance gap remains in the shaded area which validates the theoretical upper bound given by Theorem 4.4, and shows that the derived bound is tight, up to the second term.

**Experiment 5. (Comparisons with aggregation-based baselines)** We next compare the privacy-variance trade-off of the estimator based on the CLUSTER-DP mechanism with the other baselines discussed in Section 3, namely the noisy Horvitz-Thompson estimator and the noisy histogram estimator. The goal of this experiment is to show that in the case of one-shot communication between the central unit and the advertisers, the CLUSTER-DP estimator achieves lower finite-sample conditional bias than the other two baselines. To demonstrate this point, we fix the noise and randomization in each DP mechanisms for the super-population and compute the bias of each estimator with respect to random draws from the super-population and of the treatment
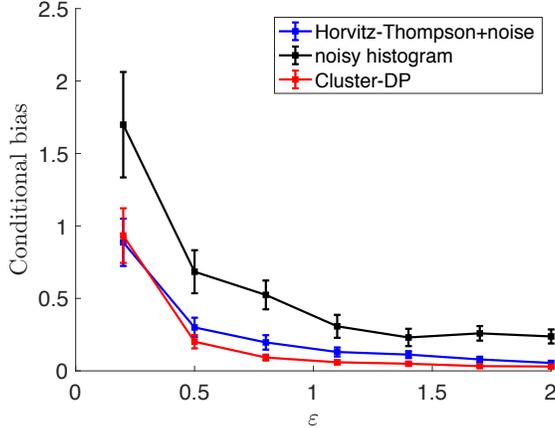
17

Figure 8: Bias of the CLUSTER-DP, noisy Horvitz-Thompson and noisy histogram estimators under one shot communication between the central unit and the advertisers in the setting of Experiment 5.

assignments. Specifically we compute the expectation of the treatment effect estimator over 500 sub-populations, each consisting of 500, 1000, 2000 units from each cluster, uniformly at random with a balanced number of treated and controlled units in each cluster. The bias is then computed as the difference between the expectation of the estimator and the true treatment effect. As we see in Figure 8, CLUSTER-DP estimator achieves a lower conditional bias compared to the other two baselines, as we vary the privacy loss $\varepsilon$. The error bars are obtained by considering 50 different realizations of the noise/randomization in the DP mechanisms.

## 5.1 Simulation on the Youtube social network

We now use a subset of the Youtube social network to replicate two experiment results in a setting with natural clusters. First, we demonstrate that the proposed stratified estimator combined with the CLUSTER-DP mechanism is unbiased and admits a Gaussian distribution, replicating the results of Experiment 1. We then compare the variance of our suggested estimator for the CLUSTER-DP mechanism with its variance when using the CLUSTER FREE-DP mechanism to show the benefit of leveraging the clustering structure, replicating the results of Experiment 2.

The Youtube social network dataset (Leskovec and Krevl, 2014) contains the friendship links of a set of users on Youtube, and the ground-truth clusters correspond to groups created by users. We form a smaller dataset, by considering only the 50 largest communities, which includes a total of 22,179 users with a minimum cluster size of 199. We generate the potential outcomes for the users as follows:

$$y_i(0) = x_i^\mathsf{T}\beta + w_i, \quad y_i(1) = y_i(0) + \tau,$$

with $w_i \sim \mathsf{N}(0, v^2)$ capturing individual $i$'s effect and the $x_i^\mathsf{T}\beta$ term capturing the cluster-level effect. We follow a similar model as in (Zhou et al., 2020) and consider a four-dimensional feature vector $x_i$, with $x_{i1}$ being the number of nodes in cluster $c_i$ (the cluster of user $i$), $x_{i2}$ the number of edges in $c_i$, $x_{i3}$ the number of edges in $c_i$ with other clusters, and $x_{i4}$ the density of cluster $c_i$. Recall that for a cluster with $n$ nodes and $e$ edges, its density is defined as $e \times \binom{n}{2}^{-1}$.

Since the proposed mechanism is for discrete outcome spaces, we quantize the responses into $K = 8$ levels. We standardize the features by making each of the four features zero mean and unit
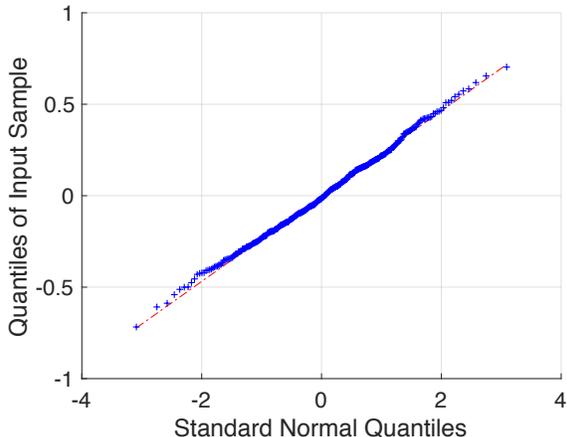
18

Figure 9: qqplot of $\hat{\tau} - \tau$, with $\hat{\tau}$ the CLUSTER-DP estimator using 500 realizations of randomness in the outcomes and the DP mechanism.
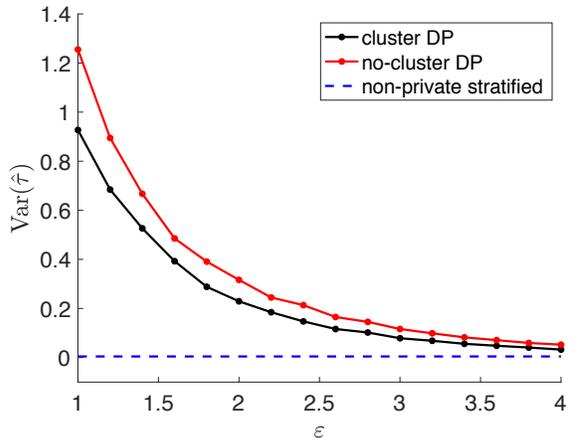


Figure 10: Privacy-variance trade-off of the CLUSTER-DP and CLUSTER FREE-DP stratified estimators. The dotted line represents the variance of the non-private stratified estimator.

norm across clusters, and setting the standard deviation of the Gaussian noise $w_i$ to $v = 0.1$. In our experiments, we set $\beta = (1, 1, 1, 1)^{\mathsf{T}}$ and $\tau = 1$. In the CLUSTER-DP mechanism, we set the truncation threshold to $\gamma = 0.1/K$ and the Laplace noise level to $\sigma = 5$.

Figure 9 shows the qqplot of $\hat{\tau} - \tau$ with $\hat{\tau}$ being the CLUSTER-DP mechanism, using 500 realizations of the randomness in the outcomes and the DP mechanism. As the plot demonstrates $\hat{\tau}$ is an unbiased and Gaussian estimator. In Figure 10, we plot the privacy-variance trade-off for the CLUSTER-DP and the CLUSTER FREE-DP mechanisms, along with the variance of the non-private stratified estimator, finding once again that the CLUSTER-DP mechanism achieves a better trade-off by leveraging the natural cluster structure of the Youtube users.

## Acknowledgement

We would like to thank Nick Doudchenko, Ian Waudby-Smith, and many others for helpful discussions on this work.

## References

Bassily, R., Thakurta, A. G., and Thakkar, O. D. (2018). Model-agnostic private learning. *Advances in Neural Information Processing Systems*. 2

Beimel, A., Nissim, K., and Stemmer, U. (2013). Private learning and sanitization: Pure vs. approximate differential privacy. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 363–378. Springer. 2

Betlei, A., Gregoir, T., Rahier, T., Bissuel, A., Diemert, E., and Amini, M.-R. (2021). Differentially private individual treatment effect estimation from aggregated data. *hal-03339723*. 3

Chaudhuri, K. and Hsu, D. (2011). Sample complexity bounds for differentially private learning. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 155–186. JMLR Workshop and Conference Proceedings. 2, 5

Cohen, E., Kaplan, H., Mansour, Y., Stemmer, U., and Tsfadia, E. (2021). Differentially-private clustering of easy instances. In *International Conference on Machine Learning*, pages 2049–2059. PMLR. 5

Dwork, C., Kenthapadi, K., McSherry, F., Mironov, I., and Naor, M. (2006a). Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology-EUROCRYPT 2006: 24th Annual International Conference on the Theory and Applications of Cryptographic Techniques, St. Petersburg, Russia, May 28-June 1, 2006. Proceedings 25*, pages 486–503. Springer. 2

Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006b). Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer. 2

Dwork, C., Roth, A., et al. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4):211–407. 6, 11, 12, 24, 27

Esfandiari, H., Mirrokni, V., Syed, U., and Vassilvitskii, S. (2022). Label differential privacy via clustering. In *International Conference on Artificial Intelligence and Statistics*, pages 7055–7075. PMLR. 2, 3, 10, 39, 40, 43

Feldman, D., Fiat, A., Kaplan, H., and Nissim, K. (2009). Private coresets. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pages 361–370. 5

Feldman, D., Xiang, C., Zhu, R., and Rus, D. (2017). Coresets for differentially private k-means clustering and applications to privacy in mobile sensor networks. In *Proceedings of the 16th ACM/IEEE International Conference on Information Processing in Sensor Networks*, pages 3–15. 5

Ghazi, B., Golowich, N., Kumar, R., Manurangsi, P., and Zhang, C. (2021). Deep learning with label differential privacy. *Advances in neural information processing systems*, 34:27131–27145. 3

Ghazi, B., Kamath, P., Kumar, R., Leeman, E., Manurangsi, P., Varadarajan, A., and Zhang, C. (2022). Regression with label differential privacy. *arXiv preprint arXiv:2212.06074*. 3

Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press. 1, 4

Kairouz, P., Oh, S., and Viswanath, P. (2015). The composition theorem for differential privacy. In *International conference on machine learning*, pages 1376–1385. PMLR. 5

Kancharla, M. and Kang, H. (2021). A robust, differentially private randomized experiment for evaluating online educational programs with sensitive student data. *arXiv preprint arXiv:2112.02452*. 3, 8, 9, 11

Kasiviswanathan, S. P., Lee, H. K., Nissim, K., Raskhodnikova, S., and Smith, A. (2011). What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826. 9

Leskovec, J. and Krevl, A. (2014). SNAP Datasets: Stanford large network dataset collection. http://snap.stanford.edu/data. 18

Li, N., Li, T., and Venkatasubramanian, S. (2006). t-closeness: Privacy beyond k-anonymity and l-diversity. In *2007 IEEE 23rd international conference on data engineering*, pages 106–115. IEEE. 2

Machanavajjhala, A., Kifer, D., Gehrke, J., and Venkitasubramaniam, M. (2007). l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):3–es. 2

Narayanan, A. and Shmatikov, V. (2008). Robust de-anonymization of large sparse datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pages 111–125. IEEE. 1, 2

Nissim, K., Raskhodnikova, S., and Smith, A. (2007). Smooth sensitivity and sampling in private data analysis. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 75–84. 5

Niu, F., Nori, H., Quistorff, B., Caruana, R., Ngwe, D., and Kannan, A. (2022). Differentially private estimation of heterogeneous causal effects. In *Conference on Causal Learning and Reasoning*, pages 618–633. PMLR. 3

Nori, H., Caruana, R., Bu, Z., Shen, J. H., and Kulkarni, J. (2021). Accuracy, interpretability, and differential privacy via explainable boosting. In *International Conference on Machine Learning*, pages 8227–8237. PMLR. 3

Panigrahi, S., Wang, J., and He, X. (2022). Treatment effect estimation with efficient data aggregation. *arXiv preprint arXiv:2203.12726*. 3

Stemmer, U. and Kaplan, H. (2018). Differentially private k-means with constant multiplicative error. *Advances in Neural Information Processing Systems*, 31. 5

Sweeney, L. (2000). Simple demographics often identify people uniquely. *Health (San Francisco)*, 671(2000):1–34. 1, 2

Sweeney, L. (2002). k-anonymity: A model for protecting privacy. *International journal of uncertainty, fuzziness and knowledge-based systems*, 10(05):557–570. 2

Wang, D. and Xu, J. (2019). On sparse linear regression in the local differential privacy model. In *International Conference on Machine Learning*, pages 6628–6637. PMLR. 2

Wasserman, L. and Zhou, S. (2010). A statistical framework for differential privacy. *Journal of the American Statistical Association*, 105(489):375–389. 5

Zhou, Y., Liu, Y., Li, P., and Hu, F. (2020). Cluster-adaptive network a/b testing: From randomization to estimation. *arXiv preprint arXiv:2008.08648*. 18

# A Appendix

## A.1 Notation and common formulas

We recall here all the notations used in the paper and in the proofs:

- $n$: total number of units.

- $n_0$ (resp $n_1$): total number of controlled (resp. treated) units.

- $z_i \in \{0, 1\}$: the treatment assignment of unit $i$. $z_i = 1$ (resp. 0) implies unit $i$ is treated (resp. controlled).

- $\mathcal{Y}$: response space of cardinality $K = |\mathcal{Y}| < \infty$.

- $y_i(z) \in \mathcal{Y}$: the potential treatment outcome of unit $i$ under treatment assignment value $z$.

- $\tilde{y}_i \in \mathcal{Y}$: privatized outcome of unit $i$ returned by the DP mechanism.

- $\mathsf{y} = (y)_{y \in \mathcal{Y}}$ : vector notation of the entire response space.

- $\mathcal{C}$: set of all clusters of cardinality $C = |\mathcal{C}| < \infty$.

- $c_i$: cluster of unit $i$.

- $n_c = |\{i \in [n] : c_i = c\}|$: number of units in cluster $c$.

- $\mathcal{O}_{a,c} := \{i : c_{x_i} = c, z_i = a\}$ : units belonging to cluster $c$ and treatment assignment $z_i = a$.

- $\mathcal{O}_c := \{i : c_i = c\}$ : units belonging to cluster $c$

- $n_{a,c} = |\mathcal{O}_{a,c}|$ for $a \in \{0, 1\}$.

- $\gamma \in [0, 1/K]$: minimum of clipped empirical distribution

- $\sigma$: noise scale

- $1 - \lambda \in [0, 1]$: true response sampling probability

- $p_a(y|x) = \mathbb{P}(Y(a) = y|X = x)$: true distribution for treated ($a = 1$) and controlled ($a = 0$) units.

- $p_a(y|c) = \mathbb{P}(Y(a) = y|c_X = c)$: true distribution for treated ($a = 1$) and controlled ($a = 0$) units within cluster $c$.

- $\hat{p}_a(y|c) = \dfrac{|\{i : c_{x_i} = c, z_i = a, y_i = y\}|}{|\{i : c_{x_i} = c, z_i = a\}|}$: empirical distribution for treated ($a = 1$) and controlled ($a = 0$) units within cluster $c$.

- When the context is clear, we sometimes adopt this lightened notation:

  - $\hat{p}_\ell = \frac{1}{n_{0,c}}|\{i \in \mathcal{O}_{0,c} : y_i = \ell\}|$ : empirical probability of outcome $\ell$ for the controlled units in cluster $c$

- $\hat{\boldsymbol{p}} = [\hat{p}_l]_{l \in \mathcal{Y}} \in \mathbb{R}^{K \times 1}$: empirical distribution of outcomes of the controlled units in cluster $c$, arranged into a $K$-dimensional vector, with coordinates $\hat{p}_\ell$.

- $\phi_a = \mathrm{Var}(\mathbb{E}[Y(a)|c_X])$: clustering quality

- $Q_{c,a}[y', y] = (1 - \lambda)\mathbb{I}(y' = y) + \lambda\tilde{q}_a(y'|c)$ : response randomization matrix

- When the context is clear, we sometimes adopt this lightened notation

  - $Q = Q_{c,0}$
  - $Q_{a,b} = Q_{c,0}[a, b]$,
  - $Q_{a,b}^{-1} = Q_{c,0}^{-1}[a, b]$,
  - $Q^{-\mathsf{T}} = (Q^{-1})^{\mathsf{T}}$.
  - $\tilde{\boldsymbol{q}} = [\tilde{q}_0(y|c)]_{y \in \mathcal{Y}} \in \mathbb{R}^{K \times 1}$ the distribution constructed in the DP mechanism (after adding noise to empirical distribution $\hat{p}_{0,D}$, truncation and normalization)
  - $\tilde{q}_l$: coordinate $l$ of vector $\tilde{\boldsymbol{q}}$.

- $u_{a,c} := \sum_{i \in \mathcal{O}_{a,c}} \sum_{y' \in \mathcal{Y}} Q_{c_i,a}^{-1}[y', \tilde{y}_i]y'$ .

- $\vec{y}_{a,c} := \{y_i(a) : i \in \mathcal{O}_{a,c}\}$, for $a \in \{0, 1\}$. Note that $\vec{y}_{a,c}$ is observed by the experimenter.

- $\vec{y}_c(a) := \{y_i(a) : c_i = c\}$. Note that $\vec{y}_c(a)$ contains unobserved values.

- $e_\ell \in \mathbb{R}^{K \times 1}$ with 1 at the $\ell$-th position and zero everywhere else

- $A(x) = 2K\left\{B^2\left(\frac{3}{(1-\lambda)^2} + 2\right) + \frac{(\lambda\sqrt{K}+1)^2}{(1-\lambda)^2}\|\mathcal{Y}\|^2(1 - \lambda(K - 1)\gamma)\right\}\left[\gamma + \frac{\sigma}{x}\left(e^{-\gamma x/\sigma} - e^{-x/\sigma}\right)\right]$: recurring expression in variance bounds.

We now state the well-known expression for the variance of the stratified Horvitz-Thompson estimator, which we recall below:

$$\hat{\tau}_{\text{No-DP}} := \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{i \in c} \left(\frac{y_i z_i}{n_{1,c}} - \frac{y_i(1 - z_i)}{n_{0,c}}\right) .$$

Let $\vec{y}_c := \{y_i : c_i = c\} \in \mathcal{Y}^{n_c}$ be the vector of outcomes of units in cluster $c$, and $\vec{\tau}_c := \vec{y}_c(1) - \vec{y}_c(0)$ be the vector of the differences between each unit's potential outcome in treatment and in control. The variance of $\hat{\tau}_{\text{No-DP}}$ is given by

$$\mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] = \sum_{c \in \mathcal{C}} \frac{n_c^2}{n^2} \left(\frac{S^2(\vec{y}_c(1))}{n_{1,c}} + \frac{S^2(\vec{y}_c(0))}{n_{0,c}} - \frac{S^2(\vec{\tau}_c)}{n_c}\right) ,$$

where, for any vector $u \in \mathbb{R}^d$, $S^2(\vec{u}) := \frac{1}{d-1} \sum_{u \in \vec{u}} (u - \bar{u})^2$ and $\bar{u} := \frac{1}{d} \sum_{u \in \vec{u}} u$ . The formula for $\mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}^u]$ can be obtained from the formula above where all units belong to a single cluster $|\mathcal{C}| = 1$ and $n_c = n$.

## A.2 Proof of Proposition 3.1 and 3.2

We start by proving Proposition 3.1. Recall the noisy Horvitz-Thompson estimator $\hat{\tau}$ given by (1):

$$\hat{\tau} := \sum_{c \in \mathcal{C}} \frac{n_c}{n} \left\{ \sum_{i \in c} \left( \frac{y_i z_i}{n_{1,c}} - \frac{y_i(1-z_i)}{n_{0,c}} \right) + w_c \right\}, w_c \sim \text{Laplace}(\eta_c).$$

To show its privacy guarantee, we apply (Dwork et al., 2014, Theorem 3.6). Consider the sensitivity $\Delta_c$ of the inner function $1/n_{1,c} y_i z_i - 1/n_{0,c} y_i(1-z_i)$, defined as the maximum change in its value when changing only one label in the data set. Since the assignments are not private, we keep them intact in computing the sensitivity. Therefore, changing only on label will change the inner function by at most $\Delta_c = \min\{n_{0,c}, n_{1,c}\}^{-1} \times \max_{y \in \mathcal{Y}} |y|$. By using (Dwork et al., 2014, Theorem 3.6), adding Laplace noise with parameter $\Delta_c/\varepsilon$ will make each of the inner terms $\varepsilon$-DP and by the post-processing property (Dwork et al., 2014, Proposition 2.1 ), $\hat{\tau}$ is also $\varepsilon$-DP.

For the variance, recall the non-differentially-private Horvitz-Thompson estimator $\hat{\tau}_{\text{No-DP}}$ from (2), by which we can write

$$\hat{\tau} = \hat{\tau}_{\text{No-DP}} + \sum_{c \in \mathcal{C}} \frac{n_c}{n} w_c.$$

Since $w_c$ are drawn independently from each other and also independent from the assignments $z_i$, we have

$$\text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}] = \text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \sum_{c \in \mathcal{C}} \left( \frac{n_c}{n} \right)^2 \text{Var}[w_c]$$

$$= \text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + 2 \sum_{c \in \mathcal{C}} \left( \frac{n_c}{n} \frac{\Delta_c}{\varepsilon} \right)^2,$$

where the last step holds because $w_c \sim \text{Laplace}(\eta_c)$ with $\eta_c = \Delta_c/\varepsilon$.

We next proceed with proving Proposition 3.2. Its privacy guarantee follows easily from the fact that $\hat{p}_a(y|c)$ has sensitivity $1/n_{a,c}$ (histogram queries) and therefore adding independent draws from $\text{Laplace}((n_{a,c}\varepsilon)^{-1})$ to the frequency of each value will make the histogram $\varepsilon$-DP. To prove the claim on its variance, we note that the non-private Horvitz-Thompson estimator can be written as

$$\hat{\tau}_{\text{No-DP}} = \sum_{c \in \mathcal{C}} \frac{n_c}{n} \left( \sum_{y \in \mathcal{Y}} y \hat{p}_1(y|c) - \sum_{y \in \mathcal{Y}} y \hat{p}_0(y|c) \right).$$

Therefore, we can write the noisy Histogram estimator (3) as

$$\hat{\tau} = \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{y \in \mathcal{Y}} y(\hat{p}_1(y|c) + w_{1,c,y} - p_0(y|c) - w_{0,c,y})$$

$$= \hat{\tau}_{\text{No-DP}} + \sum_{c \in \mathcal{C}} \frac{n_c}{n} \sum_{y \in \mathcal{Y}} y(w_{1,c,y} - w_{0,c,y}).$$

Since $w_{a,c,y}$ are independent from each other and $w_{a,c,y} \sim \text{Laplace}(\eta_{a,c})$, we get

$$\text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}] = \text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \sum_{c \in \mathcal{C}} \sum_{y \in \mathcal{Y}} \left(\frac{n_c}{n} y\right)^2 \left(\text{Var}[w_{1,c,y}] + \text{Var}[w_{0,c,y}]\right)$$

$$= \text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \left(\sum_{y \in \mathcal{Y}} y^2\right) \sum_{c \in \mathcal{C}} \left(\frac{n_c}{n}\right)^2 (2\eta_{1,c}^2 + 2\eta_{0,c}^2)$$

$$= \text{Var}_{DP,\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \frac{2}{\varepsilon^2} \left(\sum_{y \in \mathcal{Y}} y^2\right) \sum_{c \in \mathcal{C}} \left(\frac{n_c}{n}\right)^2 \left(\frac{1}{n_{1,c}^2} + \frac{1}{n_{0,c}^2}\right).$$

This completes the proof of Proposition 3.2.

## A.3 Proof of Theorem 3.4

Because the UNIFORM-PRIOR DP mechanism is a special case of the CLUSTER-DP mechanism, we follow the proof of Theorem 4.4, which is detailed below and which the reader might prefer reading first. In this special case, we can obtain an exact form for the variance gap. Hence, we continue from (27), which in the case that there is no Laplace noise added, reads as

$$\text{Var}_{DP}(u_{0,c}|\boldsymbol{z}, \mathcal{P}) = n_{0,c} \mathsf{y}^\mathsf{T} Q^{-1} \text{diag}(Q\hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \mathsf{y} - \sum_{i \in \mathcal{O}_{0,c}} y_i^2(0). \tag{11}$$

Note that we are using the lightened notation $\hat{\boldsymbol{p}}$ to indicate the empirical distribution of outcomes of the controlled units in cluster $c$.

In the mechanism described by Algorithm 4, $\tilde{\boldsymbol{q}}$ is data-dependent and so correlated to $\hat{\boldsymbol{p}}$. In that case, we analyzed the first term via the decomposition $\text{diag}(Q\hat{\boldsymbol{p}}) = \text{diag}(Q\tilde{\boldsymbol{q}}) + \text{diag}(Q(\hat{\boldsymbol{p}} - \tilde{\boldsymbol{q}}))$ and bounding $\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1$. In the current case that $\tilde{\boldsymbol{q}}$ is the uniform distribution, this approach is not tight as $\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1$ would be large. However, since $\tilde{\boldsymbol{q}}$ (and therefore $Q$) is data-independent we can directly analyze the first term as follows:

$$\mathsf{y}^\mathsf{T} Q^{-1} \text{diag}(Q\hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \mathsf{y}$$
$$= \frac{1}{(1-\lambda)^2} \mathsf{y}^\mathsf{T} \left(I - \frac{\lambda}{K}\mathbf{1}\mathbf{1}^\mathsf{T}\right) \text{diag}\left((1-\lambda)\hat{\boldsymbol{p}} + \frac{\lambda}{K}\mathbf{1}\right) \left(I - \frac{\lambda}{K}\mathbf{1}\mathbf{1}^\mathsf{T}\right) \mathsf{y}$$
$$= \mathsf{y}^\mathsf{T} \left\{\frac{1}{(1-\lambda)}\text{diag}(\hat{\boldsymbol{p}}) + \frac{\lambda}{K(1-\lambda)^2}\text{diag}(\mathbf{1}) - \frac{2\lambda}{K(1-\lambda)}\mathbf{1}\hat{\boldsymbol{p}}^\mathsf{T} - \frac{\lambda^2}{K^2(1-\lambda)^2}\mathbf{1}\mathbf{1}^\mathsf{T}\right\} \mathsf{y}$$
$$= \frac{\overline{\boldsymbol{y}_{0,c}^2}}{(1-\lambda)} + \frac{\lambda \overline{y^2}}{(1-\lambda)^2} - \frac{2\lambda \bar{y}}{(1-\lambda)}\overline{\boldsymbol{y}_{0,c}} - \frac{\lambda^2 \bar{y}^2}{(1-\lambda)^2}$$
$$= \frac{\overline{\boldsymbol{y}_{0,c}^2}}{(1-\lambda)} + \frac{\lambda \overline{y^2} - \lambda^2 \bar{y}^2}{(1-\lambda)^2} - \frac{2\lambda \bar{y}}{(1-\lambda)}\overline{\boldsymbol{y}_{0,c}}, \tag{12}$$

where we use the shorthand

$$\overline{\boldsymbol{y}_{0,c}^2} = \frac{1}{n_{0,c}} \sum_{i \in \mathcal{O}_{0,c}} y_i^2(0), \quad \overline{\boldsymbol{y}_{0,c}} = \frac{1}{n_{0,c}} \sum_{i \in \mathcal{O}_{0,c}} y_i(0).$$

25

Using (12) in (11), we arrive at

$$\text{Var}_{DP}(u_{0,c}|\boldsymbol{z},\mathcal{P}) = n_{0,c}\left[\frac{\overline{\boldsymbol{y}_{0,c}^2}}{(1-\lambda)} + \frac{\lambda\overline{y^2} - \lambda^2\bar{y}^2}{(1-\lambda)^2} - \frac{2\lambda\bar{y}}{(1-\lambda)}\overline{\boldsymbol{y}_{0,c}}\right] - n_{0,c}\overline{\bar{y}_{0,c}^2}$$

$$= n_{0,c}\left[\frac{\lambda}{1-\lambda}\overline{\boldsymbol{y}_{0,c}^2} + \frac{\lambda\overline{y^2} - \lambda^2\bar{y}^2}{(1-\lambda)^2} - \frac{2\lambda\bar{y}}{(1-\lambda)}\overline{\boldsymbol{y}_{0,c}}\right].$$

Invoking (23), the above characterization yields the following:

$$\text{Var}_{DP}(\hat{\tau}|\boldsymbol{z},\mathcal{P}) = \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left[\frac{\lambda}{1-\lambda}\left(\frac{\overline{\boldsymbol{y}_{0,c}^2}}{n_{0,c}} + \frac{\overline{\boldsymbol{y}_{1,c}^2}}{n_{1,c}}\right) - \frac{2\lambda\bar{y}}{(1-\lambda)}\left(\frac{\overline{\boldsymbol{y}_{0,c}}}{n_{0,c}} + \frac{\overline{\boldsymbol{y}_{1,c}}}{n_{1,c}}\right)\right]$$

$$+ \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left(\frac{1}{n_{0,c}} + \frac{1}{n_{1,c}}\right)\frac{\lambda\overline{y^2} - \lambda^2\bar{y}^2}{(1-\lambda)^2}. \tag{13}$$

We next compute $\mathbb{E}_{\boldsymbol{z}}[\text{Var}_{DP}(\hat{\tau}|\boldsymbol{z},\mathcal{P})]$. Since we are fixing $n_{a,c}$ for each cluster, we have $\mathbb{P}(z_i = a) = \frac{n_{a,c}}{n_c}$ for $i \in \mathcal{O}_c$ and $a \in \{0,1\}$. We therefore have

$$\mathbb{E}_{\boldsymbol{z}}[\overline{\boldsymbol{y}_{a,c}}] = \mathbb{E}_{\boldsymbol{z}}\left[\frac{1}{n_{a,c}}\sum_{i\in\mathcal{O}_c}\mathbb{I}(z_i = a)y_i(a)\right] = \frac{1}{n_c}\sum_{i\in\mathcal{O}_c}y_i(a) = \overline{\vec{y}_c(a)}.$$

Likewise we have $\mathbb{E}_{\boldsymbol{z}}[\overline{\boldsymbol{y}_{a,c}^2}] = \overline{\vec{y}_c^2(a)}$. Using this identities in (13), we obtain

$$\mathbb{E}_{\boldsymbol{z}}[\text{Var}_{DP}(\hat{\tau}|\boldsymbol{z},\mathcal{P})] = \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left[\frac{\lambda}{1-\lambda}\left(\frac{\overline{\boldsymbol{y}_c^2(0)}}{n_{0,c}} + \frac{\overline{\boldsymbol{y}_c^2(1)}}{n_{1,c}}\right) - \frac{2\lambda\bar{y}}{(1-\lambda)}\left(\frac{\overline{\boldsymbol{y}_c(0)}}{n_{0,c}} + \frac{\overline{\boldsymbol{y}_{c(1)}}}{n_{1,c}}\right)\right]$$

$$+ \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left(\frac{1}{n_{0,c}} + \frac{1}{n_{1,c}}\right)\frac{\lambda\overline{y^2} - \lambda^2\bar{y}^2}{(1-\lambda)^2}.$$

We next recall (22):

$$\text{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z},\mathcal{P})] = \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left(\frac{S^2(\vec{y}_c(1))}{n_{1,c}} + \frac{S^2(\vec{y}_c(0))}{n_{0,c}} - \frac{S^2(\vec{\tau}_c)}{n_c}\right),$$

which is the variance of the typical estimator with no-differential-privacy and so was written as $\text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$. Finally, from the law of total variance, we have:

$$\text{Var}(\hat{\tau}|n_0, n_1, \mathcal{P}) = \mathbb{E}_{\boldsymbol{z}}[\text{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, n_0, n_1, \mathcal{P})] + \text{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z}, n_0, n_1, \mathcal{P})]$$

$$= \mathbb{E}_{\boldsymbol{z}}[\text{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})] + \text{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})]$$

$$= \text{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left(\frac{1}{n_{0,c}} + \frac{1}{n_{1,c}}\right)\frac{\lambda\overline{y^2} - \lambda^2\bar{y}^2}{(1-\lambda)^2}$$

$$+ \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left[\frac{\lambda}{1-\lambda}\left(\frac{\overline{\boldsymbol{y}_c^2(0)}}{n_{0,c}} + \frac{\overline{\boldsymbol{y}_c^2(1)}}{n_{1,c}}\right) - \frac{2\lambda\bar{y}}{(1-\lambda)}\left(\frac{\overline{\boldsymbol{y}_c(0)}}{n_{0,c}} + \frac{\overline{\boldsymbol{y}_c(1)}}{n_{1,c}}\right)\right].$$

## A.4 Proof of Theorem 4.1

The CLUSTER-DP mechanism randomizes the labels using the empirical probability of units with the same treatment status (treated or controlled) within the same cluster, so we can focus on the controlled units within one cluster, and drop the index $a, c$ from our notation, unless needed for clarification. With slight abuse of notation, suppose that there are $n$ controlled units in the cluster and denote by $M$ the mechanism described in Algorithm 4.

We can think of $M$ as composition of two mechanisms $M_1$ and $M_2$ with $M(D) = M_2(D, M_1(D))$, where $M_1(D)$ represents the mechanism that returns the noisy cluster label distribution $\tilde{q}$, and $M_2(D, \tilde{q})$ represents the mechanism which uses $\tilde{q}$ to re-sample the labels and use them to form the average treatment effect estimator $\hat{\tau}$. By composition theorem for $(\varepsilon, \delta)$-DP (see e.g. (Dwork et al., 2014, Theorem B.1)), if $M_1$ is $(\varepsilon_1, \delta_1)$-DP and $M_2$ is $(\varepsilon_2, \delta_2)$-DP, then $M$ is $(\varepsilon_1 + \varepsilon_2, \delta_1 + \delta_2)$-DP.

After adding noise terms $w_{y,c}$ to empirical distributions $\hat{p}_{0,D}(y|c)$ and $\hat{p}_{1,D}(y|c)$, the dataset $D$ is not accessed anymore. Furthermore, the empirical distributions have sensitivity $1/n_c$ and the Laplace noise used in $M_1$ is of scale $\sigma/n_c$, which imply that $M_1$ is $(1/\sigma, 0)$-DP (see e.g. (Dwork et al., 2014, Theorem 3.6) for an argument).

For mechanism $M_2$, note that it is a randomization per label mechanism (using perturbed distribution $\tilde{q}$), followed by post-processing (computing average treatment effect estimator). We next show that $M_2$ is $(\tilde{\varepsilon}, \delta)$-DP. Note that for all $y \in \mathcal{Y}$, we have

$$\mathbb{P}(\tilde{y}_i = y | y_i = y) = 1 - \lambda + \lambda\tilde{q}(y), \quad \mathbb{P}(\tilde{y}_i = y | y_i \neq y) = \lambda\tilde{q}(y).$$

Since $\mathbb{P}(\tilde{y}_i = y | y_i \neq y)$ is independent of $y_i$ and $\mathbb{P}(\tilde{y}_i = y | y_i \neq y) < \mathbb{P}(\tilde{y}_i = y | y_i = y)$, the only condition we need to verify is the following:

$$\mathbb{P}(\tilde{y}_i = y | y_i = y) \leq e^{\tilde{\varepsilon}} \mathbb{P}(\tilde{y}_i = y | y_i \neq y) + \delta.$$

By substituting for the events probabilities and $\delta$, the above condition becomes equivalent to

$$0 \leq \lambda(\gamma - \tilde{q}(y))(1 - e^{\tilde{\varepsilon}}),$$

which holds since $\tilde{\varepsilon} > 0$ and $\gamma \leq \tilde{q}(y)$. To summarize, by applying composition theorem for $(\varepsilon, \delta)$-DP, we obtain that $M$ is $(\varepsilon', \delta)$-label DP, with $\varepsilon' = 1/\sigma + \tilde{\varepsilon}$.

We next show that $M$ is $(\varepsilon'', \delta)$-label DP, with $\varepsilon'' = 2/\gamma + \sigma$, which along with the previous result gives the claim of Theorem 4.1. Let $Y_{1:n}$ be the random vector denoting the labels of the units. We need to show that for any two neighboring data sets $D = (y_{1:n}, x_{1:n})$ and $D' = (y'_{1:n}, x_{1:n})$ (where $y_{1:n}$ and $y'_{1:n}$ differ only in one entry) we have

$$\mathbb{P}(M(y_{1:n}) \in O) \leq e^{\varepsilon''} \mathbb{P}(M(y'_{1:n}) \in O) + \delta, \tag{14}$$

for any set $O \in \mathcal{Y}^n$. Proof of this part requires more effort. Let $W_{1:K}$ be the random vector representing the noise values added to the set of possible labels in the data set. We then have

$$\begin{aligned}
\mathbb{P}(M(y_{1:n}) \in O) &= \sum_{w_{1:K}} \mathbb{P}(M(y_{1:n}) \in O) | Y_{1:n} = y_{1:n}, W_{1:K} = w_{1:K}) \, \mathbb{P}(W_{1:K} = w_{1:K}), \\
\mathbb{P}(M(y'_{1:n}) \in O) &= \sum_{w_{1:K}} \mathbb{P}(M(y'_{1:n}) \in O) | Y_{1:n} = y'_{1:n}, W_{1:K} = w_{1:K}) \, \mathbb{P}(W_{1:K} = w_{1:K}).
\end{aligned} \tag{15}$$

It suffices to show that for any value of $w_{1:K}$, we have

$$\mathbb{P}(M(y_{1:n}) \in O)|Y_{1:n} = y_{1:n}, W_{1:K} = w_{1:K}) \leq e^{\varepsilon''}\mathbb{P}(M(y'_{1:n}) \in O)|Y_{1:n} = y'_{1:n}, W_{1:K} = w_{1:K}) + \delta. \tag{16}$$

By multiplying both sides of the above equation with $\mathbb{P}(W_{1:K} = w_{1:K})$, and summing over $w_{1:K}$, and using that $\sum_{w_{1:K}} \mathbb{P}(W_{1:K} = w_{1:K}) = 1$, we get the desired bound in (14).

Let $\tilde{q}$ and $\tilde{q}'$ be the empirical distributions of the DP mechanism, as defined in Algorithm 4. we continue by establishing a lemma on $\|\tilde{q} - \tilde{q}'\|_\infty$, proven in the next section.

**Lemma A.1.** *For all $y \in \mathcal{Y}$, $|\tilde{q}(y) - \tilde{q}'(y)| \leq \frac{2}{n}$.*

Define the shorthand $R := M(y_{1:n})$ and $R' := M(y'_{1:n})$. In order to prove (16), it suffices to show that for all $o_{1:n} \in \mathcal{Y}^n$, we have

$$\mathbb{P}(R = o_{1:n}|Y_{1:n} = y_{1:n}, W_{1:K} = w_{1:K}) \leq e^{\varepsilon''}\mathbb{P}(R' = o_{1:n}|Y_{1:n} = y'_{1:n}, W_{1:K} = w_{1:K}) + \delta. \tag{17}$$

By the definition of the mechanism $M$ we have

$$\mathbb{P}(R = o_{1:n}|Y_{1:n} = y_{1:n}, W_{1:K} = w_{1:K}) = \mathbb{P}(R = o_{1:n}|Y_{1:n} = y_{1:n}, \tilde{q}(\cdot)) = \prod_{i=1}^{n} \mathbb{P}(R_i = o_i|Y_i = y_i, \tilde{q}(\cdot))$$

$$\mathbb{P}(R' = o_{1:n}|Y_{1:n} = y'_{1:n}, W_{1:K} = w_{1:K}) = \mathbb{P}(R = o_{1:n}|Y_{1:n} = y'_{1:n}, \tilde{q}'(\cdot)) = \prod_{i=1}^{n} \mathbb{P}(R'_i = o_i|Y'_i = y'_i, \tilde{q}'(\cdot))$$

For ease in presentation, we adopt the shorthand

$$A_i := \mathbb{P}(R_i = o_i|Y_i = y_i, \tilde{q}(\cdot)), \quad B_i := \mathbb{P}(R'_i = o_i|Y'_i = y'_i, \tilde{q}'(\cdot)),$$

for $i = 1, \ldots, n$. Our next lemma bounds the event probability $A_i$ in terms of the event probability $B_i$. Proof of Lemma A.2 is deferred to Section A.4.

**Lemma A.2.** *Let $\tilde{\varepsilon} > 0$ and define $\delta := (1 - \lambda + \lambda\gamma(1 - e^{\tilde{\varepsilon}}))_+ < 1$. Without loss of generality suppose that the neighboring label sets $y_{1:n}$ and $y'_{1:n}$ differs in the first coordinate. We then have*

$$A_1 \leq e^{\tilde{\varepsilon}}\left(1 + \frac{2}{\gamma n}\right)B_1 + \delta, \tag{18}$$

$$A_i \leq B_i\left(1 + \frac{2}{\gamma n}\right), \quad for \ i = 2, \ldots, n. \tag{19}$$

We are now ready to prove inequality (17). Using Lemma A.2, we write

$$\mathbb{P}(R = o_{1:n}|Y_{1:n} = y_{1:n}, W_{1:K} = w_{1:K}) = \prod_{i=1}^{n} A_i$$

$$= A_1 \min\left\{1, \prod_{i=2}^{n} A_i\right\}$$

$$\leq \left[e^{\tilde{\varepsilon}}\left(1 + \frac{2}{\gamma n}\right)B_1 + \delta\right]\min\left\{1, \left(1 + \frac{2}{\gamma n}\right)^{n-1}\prod_{i=2}^{n} B_i\right\}$$

$$\leq e^{\tilde{\varepsilon}}\left(1 + \frac{2}{\gamma n}\right)^{n}\prod_{i=1}^{n} B_i + \delta$$

$$\leq e^{\tilde{\varepsilon}+2/\gamma}\prod_{i=1}^{n} B_i + \delta$$

$$= e^{\varepsilon''}\mathbb{P}(R' = o_{1:n}|Y_{1:n} = y'_{1:n}, W_{1:K} = w_{1:K}) + \delta .$$

where the second equality holds since $A_i \leq 1$, for all $i$.

**Proof of Lemma A.1**

Recall the notation of Theorem 4.1. We consider two neighboring datasets $D = (y_{1:n}, x_{1:n})$ and $D' = (y'_{1:n}, x_{1:n})$, where $y_{1:n}$ and $y'_{1:n}$ differ only in one entry. Define the function $f_\gamma$ as follows:

$$f_\gamma(x) = \max\{\gamma, \min\{1, x\}\} = \begin{cases} \gamma, & x \leq \gamma \\ x, & \gamma \leq x \leq 1 \\ 1, & x > 1 \end{cases}$$

We consider $q(y) := f_\gamma(\hat{p}(y) + w_y)$ and $q'(y) := f_\gamma(\hat{p}'(y) + w_y)$, where $\hat{p}$ and $\hat{p}'$ respectively denote the empirical distribution of $y_{1:n}$ and $y'_{1:n}$ and $w_y$ indicates the component of $w_{1:K}$ corresponding to label $y$. We wish to bound the difference between distributions $\tilde{q}(y)$ and $\tilde{q}'(y)$, defined in Algorithm 4, and recalled below:

$$\tilde{q}(y) = q(y) + \frac{\zeta_y}{\sum_{y'}\zeta_{y'}}\Delta, \quad \tilde{q}'(y) = q'(y) + \frac{\zeta'_y}{\sum_{y'}\zeta'_{y'}}\Delta',$$

where $\Delta = 1 - \sum_y q(y)$ and $\Delta' = 1 - \sum_y q'(y)$. To achieve this, we will need a bound on $|q(y) - q'(y)|$ and on a bound on $|\Delta - \Delta'|$.

- Since $f_\gamma$ is 1-Lipschitz, we have for any $y \in \mathcal{Y}$

$$|q(y) - q'(y)| = |f_\gamma(\hat{p}(y) + w_y) - f_\gamma(\hat{p}'(y) + w_y)| \leq |\hat{p}(y) - \hat{p}'(y)| \leq \frac{1}{n},$$

where the last inequality holds because the datasets $D$ and $D'$ differ in only one label.

29

- We now show that $|\Delta - \Delta'| \leq 1/n$. Without loss of generality, we can assume that the neighboring label sets $y_{1:n}$ and $y'_{1:n}$ differ in the first coordinate, with $y_1 = \ell$, $y'_1 = \ell'$ for $\ell, \ell' \in \mathcal{Y}$, such that

$$\hat{p}(\ell) = \hat{p}'(\ell) + \frac{1}{n}, \quad \hat{p}(\ell') = \hat{p}'(\ell') - \frac{1}{n}.$$

It follows that

$$
\begin{aligned}
\Delta' - \Delta &= \sum_y q(y) - \sum_y q'(y) \\
&= f_\gamma(\hat{p}(\ell) + w_\ell) + f_\gamma(\hat{p}(\ell') + w_{\ell'}) - f_\gamma(\hat{p}'(\ell) + w_\ell) - f_\gamma(\hat{p}'(\ell') + w_{\ell'}) \\
&= f_\gamma\left(\hat{p}'(\ell) + \frac{1}{n} + w_\ell\right) - f_\gamma(\hat{p}'(\ell) + w_\ell) + f_\gamma\left(\hat{p}'(\ell') - \frac{1}{n} + w_{\ell'}\right) - f_\gamma(\hat{p}'(\ell') + w_{\ell'}) \\
&\leq f_\gamma\left(\hat{p}'(\ell) + \frac{1}{n} + w_\ell\right) - f_\gamma(\hat{p}'(\ell) + w_\ell) \\
&\leq \frac{1}{n},
\end{aligned}
$$

where the second to last inequality holds since $f_\gamma$ is a non-decreasing function, and the last step follows from 1-Lipschitzness of $f_\gamma$. Likewise, we can show $\Delta - \Delta' \leq 1/n$ in order to obtain $|\Delta - \Delta'| \leq 1/n$.

With this, we next bound the difference between distributions $\tilde{q}(y)$ and $\tilde{q}'(y)$, defined above. Consider three different cases:

- $\Delta > 0, \Delta' < 0$. We have

$$
|\tilde{q}(y) - \tilde{q}'(y)| \leq |q(y) - q'(y)| + \left| \frac{\zeta_y}{\sum_{y'} \zeta_{y'}} \Delta - \frac{\zeta'_y}{\sum_{y'} \zeta'_{y'}} \Delta' \right|
$$

$$
\leq |q(y) - q'(y)| + |\Delta - \Delta'| \leq \frac{2}{n}.
$$

The case of $\Delta < 0, \Delta' > 0$ can be handled similarly.

- $\Delta, \Delta' < 0$. We have

$$
\begin{aligned}
\tilde{q}(y) &= q(y) + (q(y) - \gamma) \frac{\Delta}{\sum_{y'}(q(y') - \gamma)} \\
&= q(y) + (q(y) - \gamma) \frac{\Delta}{1 - \Delta - K\gamma} \\
&= \gamma + (q(y) - \gamma) + (q(y) - \gamma) \frac{\Delta}{1 - \Delta - K\gamma} \\
&= \gamma + (q(y) - \gamma) \frac{1 - K\gamma}{1 - \Delta - K\gamma}.
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
|\tilde{q}(y) - \tilde{q}'(y)| &\leq (q(y) - \gamma)\left|\frac{1 - K\gamma}{1 - \Delta - K\gamma} - \frac{1 - K\gamma}{1 - \Delta' - K\gamma}\right| + |q'(y) - q(y)|\frac{1 - K\gamma}{1 - \Delta' - K\gamma} \\
&= (q(y) - \gamma)\frac{(1 - K\gamma)|\Delta - \Delta'|}{(1 - \Delta - K\gamma)(1 - \Delta' - K\gamma)} + \frac{1}{n}\frac{1 - K\gamma}{1 - \Delta' - K\gamma} \\
&= \frac{1 - K\gamma}{1 - \Delta' - K\gamma}\left[\frac{(q(y) - \gamma)|\Delta - \Delta'|}{(1 - \Delta - K\gamma)} + \frac{1}{n}\right] \\
&\overset{(a)}{\leq} \frac{1}{n}\frac{1 - K\gamma}{1 - \Delta' - K\gamma}\left(\frac{q(y) - \gamma}{1 - \Delta - K\gamma} + 1\right) \\
&\overset{(b)}{\leq} \frac{2}{n}\frac{1 - K\gamma}{1 - \Delta' - K\gamma} \leq \frac{2}{n}.
\end{aligned}
$$

$(a)$ holds since $|\Delta - \Delta'| \leq 1/n$. $(b)$ follows from the fact that, since $q(y) \geq \gamma$ for all $y$,

$$
q(y) + (K - 1)\gamma \leq q(y) + \sum_{y' \neq y} q(y') = 1 - \Delta,
$$

such that $q(y) - \gamma \leq 1 - \Delta - K\gamma$.

- $\Delta, \Delta' > 0$. We have

$$
\begin{aligned}
\tilde{q}(y) &= q(y) + (1 - q(y))\frac{\Delta}{\sum_{y'}(1 - q(y'))} \\
&= 1 + (1 - q(y))\left(\frac{\Delta}{\Delta + K - 1} - 1\right) \\
&= 1 + (1 - q(y))\frac{1 - K}{\Delta + K - 1}.
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
|\tilde{q}(y) - \tilde{q}'(y)| &\leq (1 - q(y))\left|\frac{1 - K}{\Delta + K - 1} - \frac{1 - K}{\Delta' + K - 1}\right| + |q(y) - q'(y)|\frac{K - 1}{\Delta' + K - 1} \\
&\leq (1 - q(y))\left[\frac{|\Delta' - \Delta|}{\Delta + K - 1} + \frac{1}{n}\right]\frac{K - 1}{\Delta' + K - 1} \\
&\leq (1 - q(y))\frac{1}{n(K - 1)} + \frac{1}{n} \leq \frac{1}{n}\frac{K}{K - 1} \leq \frac{2}{n}.
\end{aligned}
$$

Combining the above three cases together, we obtain our stated lemma.

## Proof of Lemma A.2

We start by proving (18). Consider three different cases:

- $o_1 \neq y_1, y_1'$: In this case, we have

$$
A_1 = \lambda\tilde{q}(o_1), \quad B_1 = \lambda\tilde{q}'(o_1).
$$

Therefore, we can write

$$
\begin{aligned}
A_1 &\leq \lambda \tilde{q}'(o_1) + \lambda \|\tilde{q} - \tilde{q}'\|_\infty \\
&\leq \lambda \tilde{q}'(o_1) + \frac{2\lambda}{n} \\
&\leq \lambda \tilde{q}'(o_1) \left(1 + \frac{2}{n\gamma}\right) \\
&= B_1 \left(1 + \frac{2}{n\gamma}\right) \\
&\leq e^{\tilde{\varepsilon}} B_1 \left(1 + \frac{2}{n\gamma}\right) + \delta,
\end{aligned}
$$

where the second step follows from Lemma A.1, third step holds since $\gamma \leq \tilde{q}'(o_1)$, and the last step holds since $\tilde{\varepsilon}, \delta > 0$. So the claim (18) is proved in this case.

- $o_1 = y_1$: In this case,

$$
A_1 = 1 - \lambda + \lambda \tilde{q}(o_1), \quad B_1 = \lambda \tilde{q}'(o_1).
$$

We then have

$$
\begin{aligned}
A_1 &\leq 1 - \lambda + \lambda \tilde{q}'(o_1) + \lambda \|\tilde{q} - \tilde{q}'\|_\infty \\
&\leq 1 - \lambda + \lambda \tilde{q}'(o_1) + \lambda \frac{2}{n\gamma} \tilde{q}'(o_1) e^{\tilde{\varepsilon}},
\end{aligned} \tag{20}
$$

where we used Lemma A.1 along with the facts that $\tilde{q}'(o_1) \geq \gamma$ and $\tilde{\varepsilon} > 0$.

We next recall the definition $\delta := (1 - \lambda + \lambda\gamma(1 - e^{\tilde{\varepsilon}}))_+ < 1$. By a simple rearrangement of the terms and using that $\tilde{q}'(o_1) \geq \gamma$ and $\tilde{\varepsilon} > 0$, we can verify the following,

$$
1 - \lambda + \lambda \tilde{q}'(o_1) \leq e^{\tilde{\varepsilon}} \lambda \tilde{q}'(o_1) + \delta. \tag{21}
$$

Therefore, by combining equations (20) and (21), we get

$$
A_1 \leq e^{\tilde{\varepsilon}} \lambda \tilde{q}'(o_1) + \delta + \lambda \frac{2}{n\gamma} \tilde{q}'(o_1) e^{\tilde{\varepsilon}} = e^{\tilde{\varepsilon}} = \left(1 + \frac{2}{n\gamma}\right) B_1 + \delta,
$$

which completes the proof of claim (18) in this case.

- $o_1 = y_1'$: In this case,

$$
A_1 = \lambda \tilde{q}(o_1), \quad B_1 = 1 - \lambda + \lambda \tilde{q}'(o_1).
$$

The proof of claim (18) in this case follows readily from case 1, because $A_1$ is the same as in there, while $B_1$ is larger.

This concludes the proof of Claim 18. We next prove Claim 19. Note that for $i = 2, \ldots, n$, we have $y_i = y_i'$. Consider the following two cases:

- $o_i = y_i = y_i'$: In this case we have

$$
\frac{A_i}{B_i} = \frac{1 - \lambda + \lambda \tilde{q}(o_i)}{1 - \lambda + \lambda \tilde{q}'(o_i)}.
$$

32

- $o_i \neq y_i$: Since $y_i = y'_i$, we also have $o_i \neq y'_i$. In this case,

$$\frac{A_i}{B_i} = \frac{\lambda \tilde{q}(o_i)}{\lambda \tilde{q}'(o_i)}.$$

By symmetry, we can assume $\tilde{q}(o_i) \leq \tilde{q}'(o_i)$, without loss of generality, and therefore, the maximum value of the ratio $A_1/B1$ is achieved in the second case, for which we have

$$\frac{A_i}{B_i} = \frac{\tilde{q}(o_i)}{\tilde{q}'(o_i)} \leq 1 + \frac{\|\tilde{q} - \tilde{q}'\|_\infty}{\tilde{q}'(o_i)} \leq 1 + \frac{2}{n\gamma}.$$

This completes the proof of (19).

## A.5 Proof of Theorem 4.2

We would like to express the expectation $\mathbb{E}(\hat{\tau}|n_0, n_1)$. Recall that there are three sources of randomness:

- the differential privacy mechanism $DP$: determines the Laplace noise $\boldsymbol{w}$ and the $\lambda$ probability of reporting the true outcome.

- the randomized assignment $\boldsymbol{z}$: determines which units get assigned to treatment and which units get assigned to control.

- the super-population $\mathcal{P}$: determines the potential outcomes as well as the cluster assignments.

For a given unit $i$ with $(y_i(0), y_i(1), c_i) \sim \mathcal{P}$ and $z_i = a$,

$$\mathbb{E}_{DP}\left[\sum_{y'\in\mathcal{Y}} Q^{-1}_{c_i,z_i}[y', \tilde{y}_i] y' z_i \,\bigg|\, \boldsymbol{z}, \mathcal{P}\right] = \mathbb{E}_{DP}\left[\sum_{y'\in\mathcal{Y}}\sum_{y\in\mathcal{Y}} \mathbb{I}(\tilde{y}_i = y) Q^{-1}_{c_i,z_i}[y', y] y' z_i \,\bigg|\, \boldsymbol{z}, \mathcal{P}\right]$$

$$\overset{(a)}{=} \sum_{y'\in\mathcal{Y}} \mathbb{E}_{DP}\left[\sum_{y\in\mathcal{Y}} \mathbb{I}(\tilde{y}_i = y) Q^{-1}_{c_i,z_i}[y', y] \,\bigg|\, \boldsymbol{z}, \mathcal{P}\right] y' z_i$$

$$\overset{(b)}{=} \sum_{y'\in\mathcal{Y}} \mathbb{E}_{\boldsymbol{w}}\left[\sum_{y\in\mathcal{Y}} \mathbb{E}_\lambda\left[\mathbb{I}(\tilde{y}_i = y)\right] Q^{-1}_{c_i,z_i}[y', y] \,\bigg|\, \boldsymbol{z}, \mathcal{P}\right] y' z_i$$

$$\overset{(c)}{=} \sum_{y'\in\mathcal{Y}} \mathbb{E}_{\boldsymbol{w}}\left[\sum_{y\in\mathcal{Y}} Q_{c_i,z_i}[y, y_i] Q^{-1}_{c_i,z_i}[y', y] \,\bigg|\, \boldsymbol{z}, \mathcal{P}\right] y' z_i$$

$$\overset{(d)}{=} \sum_{y'\in\mathcal{Y}} \mathbb{E}_{\boldsymbol{w}}\left[\mathbb{I}(y_i = y') \,\big|\, \boldsymbol{z}, \mathcal{P}\right] y' z_i$$

$$\overset{(e)}{=} \sum_{y'\in\mathcal{Y}} \mathbb{I}(y_i = y') y' z_i$$

$$= y_i(1) z_i.$$

$(a)$ holds since assignments $z_i$ is independent from $\{y_i(0), y_i(1), c_i\}$; $(b)$ holds from the law of iterated expectation and the fact that there are two sources of randomness in the differential privacy

mechanism: $(\lambda, \boldsymbol{w})$ with $Q_{c,a}$ independent of the Bernoulli $\lambda$; $(c)$ follows from the definition of $Q_{c,a}$: $Q_{c,a}[y', y] = (1 - \lambda)\mathbb{I}(y' = y) + \lambda \tilde{q}_a(y'|c)$; and $(d)$ follows from the fact that $I = Q_{c_i}^{-1} Q_{c_i}$ therefore, for any $a, b \in [K]$,

$$\sum_y Q_{c_i}^{-1}[a, y] Q_{c_i}[y, b] = I_{a,b} = \mathbb{I}(a = b).$$

Finally, $(e)$ follows from the fact that $\boldsymbol{w}$ is independent from $\{y_i(0), y_i(1), c_i\}$. Similarly,

$$\mathbb{E}_{DP}\left[\sum_{y' \in \mathcal{Y}} Q_{c_i, z_i}^{-1}[y', \tilde{y}_i] y'(1 - z_i) \middle| \boldsymbol{z}, \mathcal{P}\right] = y_i(0)(1 - z_i)$$

As a result, with $n_{0,c}$ (resp. $n_{1,c}$) the total number of controlled (resp. treated) units in cluster $c$ and $n_c := n_{0,c} + n_{1,c}$,

$$\mathbb{E}_{DP}[\hat{\tau}|\boldsymbol{z}, \mathcal{P}] = \sum_{c \in \mathcal{C}} \frac{n_c}{n} \left(\sum_{i=1}^n y_i(1)\frac{z_i}{n_{1,c}} - \sum_{i=1}^n y_i(0)\frac{1 - z_i}{n_{0,c}}\right)$$

We recover the standard form of the difference-in-means estimator. From the law of iterated expectations, we have

$$\mathbb{E}_{DP,\boldsymbol{z}}[\hat{\tau}] = \mathbb{E}_{\boldsymbol{z}}\left[\mathbb{E}_{DP}[\hat{\tau}|\boldsymbol{z}]|n_0, n_1, \mathcal{P}\right] = \tau.$$

## A.6 Proof of Theorem 4.4

We would like to express the variance $\operatorname{Var}_{DP,\boldsymbol{z}}(\hat{\tau})$. We begin by expressing the variance with respect to the first two, considering the third fixed. From the law of total variance, we have:

$$\operatorname{Var}_{DP,\boldsymbol{z}}(\hat{\tau}) = \mathbb{E}_{\boldsymbol{z}}[\operatorname{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, n_0, n_1, \mathcal{P})] + \operatorname{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z}, n_0, n_1, \mathcal{P})]$$
$$= \mathbb{E}_{\boldsymbol{z}}[\operatorname{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})] + \operatorname{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})]$$

We bound the term $\operatorname{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})$ in a separate proposition

**Proposition A.3.** *For the average treatment effect estimator $\hat{\tau}$ given by* (9) *we have*

$$\operatorname{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P}) \leq \sum_{a \in \{0,1\}} \sum_{c \in \mathcal{C}} \frac{n_c^2}{n^2}\left[\left(\frac{1}{(1 - \lambda)^2} - 1\right)\frac{S^2(\vec{y}_{a,c})}{n_{a,c}} + \frac{A(n_{a,c})}{n_{a,c}}\right],$$

*where $\vec{y}_{a,c} := \{y_i(a) : i \in \mathcal{O}_{a,c}\}$, and*

$$A(x) = 2KB^2\left(\left(\frac{3}{(1 - \lambda)^2} + 2\right) + \frac{(\lambda\sqrt{K} + 1)^2}{(1 - \lambda)^2}\|\mathbf{y}\|^2(1 - \lambda(K - 1)\gamma)\right)\left[\gamma + \frac{\sigma}{x}\left(e^{-\gamma x/\sigma} - e^{-x/\sigma}\right)\right]$$

We now take its expectation with respect to $\boldsymbol{z}$. We assume that, for each cluster c, there is a fixed number of units $(n_{1,c})$ assigned to treatment and a fixed number of units $(n_{0,c})$ assigned to

control, regardless of the cluster assignment. We compute the expectation with $\mathbb{E}_{\boldsymbol{z}}\left[S^2(\vec{y}_{a,c})\right]$,

$$(n_{a,c}-1)\mathbb{E}_{\boldsymbol{z}}\left[S^2(\vec{y}_{a,c})\right]$$

$$=\mathbb{E}_{\boldsymbol{z}}\left[\sum_{i\in\mathcal{O}_{a,c}}\left(y_i(a)-\overline{\{y_i(a)\}_{i\in\mathcal{O}_{a,c}}}\right)^2\right]$$

$$=\mathbb{E}_{\boldsymbol{z}}\left[\sum_{i\in\mathcal{O}_{a,c}}y_i^2(a)-n_{a,c}\left(\frac{1}{n_{a,c}}\sum_{i\in\mathcal{O}_{a,c}}y_i(a)\right)^2\right]$$

$$=\sum_{i\in\mathcal{O}_c}\mathbb{P}\left(z_i=a\right)y_i^2(a)-n_{a,c}\mathbb{E}_{\boldsymbol{z}}\left[\left(\frac{1}{n_{a,c}}\sum_{i\in\mathcal{O}_c}\mathbb{I}(z_i=a)y_i(a)\right)^2\right]$$

$$=\sum_{i\in\mathcal{O}_c}\frac{n_{a,c}}{n_c}y_i^2(a)-\frac{1}{n_{a,c}}\sum_{i\in\mathcal{O}_c}\sum_{j\in\mathcal{O}_c}\mathbb{P}\left(z_i=a,z_j=a\right)y_i(a)y_j(a)$$

$$=\sum_{i\in\mathcal{O}_c}\frac{n_{a,c}}{n_c}y_i^2(a)-\frac{1}{n_{a,c}}\sum_{i\in\mathcal{O}_c}\mathbb{P}\left(z_i=a\right)y_i^2(a)-\frac{1}{n_{a,c}}\sum_{j\neq i\in\mathcal{O}_c}\mathbb{P}\left(z_i=a,z_j=a\right)y_i(a)y_j(a)$$

$$=\sum_{i\in\mathcal{O}_c}\frac{n_{a,c}}{n_c}y_i^2(a)-\frac{1}{n_c}\sum_{i\in\mathcal{O}_c}y_i^2(a)-\frac{1}{n_{a,c}}\sum_{j\neq i\in\mathcal{O}_c}\frac{n_{a,c}(n_{a,c}-1)}{n_c(n_c-1)}y_i(a)y_j(a)\,.$$

Adding and subtracting $\frac{n_{a,c}-1}{n_c(n_c-1)}\sum_{i\in\mathcal{O}_c}y_i^2(a)$, we get:

$$(n_{a,c}-1)\mathbb{E}_{\boldsymbol{z}}\left[S^2(\vec{y}_{a,c})\right]$$

$$=\sum_{i\in\mathcal{O}_c}\frac{n_{a,c}}{n_c}y_i^2(a)-\frac{1}{n_c}\sum_{i\in\mathcal{O}_c}y_i^2(a)+\frac{n_{a,c}-1}{n_c(n_c-1)}\sum_{i\in\mathcal{O}_c}y_i^2(a)-\frac{n_{a,c}-1}{n_c(n_c-1)}\left(\sum_{i\in\mathcal{O}_c}y_i(a)\right)^2$$

$$=\left(\frac{n_{a,c}}{n_c}-\frac{1}{n_c}+\frac{n_{a,c}-1}{n_c(n_c-1)}\right)\sum_{i\in\mathcal{O}_c}y_i^2(a)-(n_{a,c}-1)\frac{n_c}{n_c-1}\left(\frac{1}{n_c}\sum_{i\in\mathcal{O}_c}y_i(a)\right)^2$$

$$=\frac{n_{a,c}-1}{n_c-1}\sum_{i\in\mathcal{O}_c}y_i^2(a)-(n_{a,c}-1)\frac{n_c}{n_c-1}\left(\frac{1}{n_c}\sum_{i\in\mathcal{O}_c}y_i(a)\right)^2$$

$$=(n_{a,c}-1)S^2(\vec{y}_c(a))\,.$$

For the second term, we again make the assumption that the number of treated units is fixed at the cluster level. For the second term, from the proof of Theorem 4.2, we have:

$$\mathbb{E}_{DP}[\hat{\tau}|\boldsymbol{z},\mathcal{P}]=\sum_{c\in\mathcal{C}}\frac{n_c}{n}\left(\sum_{i=1}^{n}y_i(1)\frac{z_i}{n_{1,c}}-\sum_{i=1}^{n}y_i(0)\frac{1-z_i}{n_{0,c}}\right)$$

As a result, the second term is given by the usual formula for the variance of the stratified estimator:

$$\text{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z},\mathcal{P})]=\sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left(\frac{S^2(\vec{y}_c(1))}{n_{1,c}}+\frac{S^2(\vec{y}_c(0))}{n_{0,c}}-\frac{S^2(\vec{\tau}_c)}{n_c}\right) \tag{22}$$

35

where, for any vector $\vec{u}$ of length $n$, $S^2(u) = \frac{1}{n-1}\sum_{i=1}^{n}(u_i - \bar{u})^2$ and $\bar{u} = \frac{1}{n}\sum_{i=1}^{n}u_i$. Recall that $\vec{\tau}_c = \{y_i(1) - y_i(0)\}_{i:c_i=c} = \vec{y}_c(1) - \vec{y}_c(0)$. Since this is the variance of the typical estimator with no-differential-privacy, we write this term:

$$\mathrm{Var}_{\boldsymbol{z}}[\mathbb{E}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})] = \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}]$$

As a result, we obtain

$$\mathrm{Var}_{DP,\boldsymbol{z}}(\hat{\tau}|n_{0,c}, n_{1,c}, \mathcal{P})$$

$$\leq \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \sum_{a\in\{0,1\}}\sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\left[\left(\frac{1}{(1-\lambda)^2} - 1\right)\frac{S^2(\vec{y}_c(a))}{n_{a,c}} + \frac{A(n_{a,c})}{n_{a,c}}\right]$$

$$\leq \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \left(\frac{1}{(1-\lambda)^2} - 1\right)\sum_{a\in\{0,1\}}\sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\frac{S^2(\vec{y}_c(a))}{n_{a,c}} + \sum_{a\in\{0,1\}}\sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\frac{A(n_{a,c})}{n_{a,c}}$$

which we can rewrite as:

$$\mathrm{Var}_{DP,\boldsymbol{z}}(\hat{\tau}|n_{0,c}, n_{1,c}, \mathcal{P}) \leq \mathrm{Var}_{\boldsymbol{z}}[\hat{\tau}_{\text{No-DP}}] + \left(\frac{1}{(1-\lambda)^2} - 1\right)\sum_{a\in\{0,1\}}\phi_a + \sum_{a\in\{0,1\}}\sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\frac{A(n_{a,c})}{n_{a,c}}$$

where we have defined

$$\phi_a := \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\frac{S^2(\vec{y}_c(a))}{n_{a,c}} \geq 0$$

## Proof of Proposition A.3

We seek to compute $\mathrm{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P})$. We can rewrite $\hat{\tau}$ as

$$\hat{\tau} = \sum_{c\in C}\frac{n_c}{n}\sum_{a\in\{0,1\}}\frac{u_{a,c}}{n_{a,c}},$$

where $\mathcal{O}_{a,c} := \{i \in [n] : c_i = c, z_i = a\}$ and $u_{a,c} := \sum_{i\in\mathcal{O}_{a,c}}\sum_{y'\in\mathcal{Y}}Q_{c_i,a}^{-1}[y', \tilde{y}_i]y'$. Since $(y_i(0), y_i(1))$ are i.i.d across units, and the DP mechanism is applied to each clusters separately, such that the privatized outcomes $\tilde{y}_i$ are independent across clusters, we have that $u_{0,c}$ and $u_{1,c}$ are independent across clusters.

$$\mathrm{Var}_{DP}(\hat{\tau}|\boldsymbol{z}, \mathcal{P}) = \sum_{c\in\mathcal{C}}\frac{n_c^2}{n^2}\sum_{a\in\{0,1\}}\frac{1}{n_{a,c}^2}\mathrm{Var}_{DP}(u_{a,c}|\boldsymbol{z}, \mathcal{P}). \tag{23}$$

We proceed by calculating $\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z}, \mathcal{P})$. The computation for $\mathrm{Var}_{DP}(u_{1,c}|\boldsymbol{z}, \mathcal{P})$ is identical.

● **Computing** $\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z}, \mathcal{P})$

We have

$$\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z}, \mathcal{P}) = \mathbb{E}_{DP}[u_{0,c}^2|\boldsymbol{z}, \mathcal{P}] - \mathbb{E}_{DP}[u_{0,c}|\boldsymbol{z}, \mathcal{P}]^2$$

We begin by computing $\mathbb{E}_{DP}[u_{0,c}|\boldsymbol{z}, \mathcal{P}]$. Fixing cluster $c$, we lighten the notation by using the shorthand $Q = Q_{c,0}$, $Q_{a,b} = Q_{c,0}[a, b]$, $Q_{a,b}^{-1} = Q_{c,0}^{-1}[a, b]$, and $Q^{-\mathsf{T}} = (Q^{-1})^{\mathsf{T}}$. Finally, recall that $\mathsf{y} = (y)_{y\in\mathcal{Y}}$ is the set of possible outcomes arranged into a vector with the same ordering as the

columns of $Q$, and $e_\ell \in \mathbb{R}^{K\times 1}$ is the vector with 1 at the $\ell$-th position and zero everywhere else. Writing in matrix form, we have

$$u_{0,c} = \sum_{i\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T} Q^{-1}_{\cdot,\tilde{y}_i} .$$

Let $\tilde{\boldsymbol{y}}$, $\boldsymbol{y}$, $\boldsymbol{z}$, $\boldsymbol{w}$ be the vectors of variables $\tilde{y}_i$, $y_i$, $z_i$, $(w)_{y,c}$ respectively. We then have

$$\mathbb{E}_{DP}[u_{0,c}|\boldsymbol{z},\mathcal{P}] = \sum_{i\in\mathcal{O}_{0,c}} \mathbb{E}_{DP}\left[\mathsf{y}^\mathsf{T} Q^{-1}_{\cdot,\tilde{y}_i}\Big|\boldsymbol{z},\mathcal{P}\right]$$

$$= \sum_{i\in\mathcal{O}_{0,c}} \mathbb{E}_{DP}\left[\mathsf{y}^\mathsf{T} \sum_{\ell\in\mathcal{Y}} Q^{-1}_{\cdot,\ell}\mathbb{I}(l=\tilde{y}_i)\Big|\boldsymbol{z},\mathcal{P}\right]$$

$$= \sum_{i\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T} \sum_{\ell\in\mathcal{Y}} \mathbb{E}_{DP}\left[Q^{-1}_{\cdot,\ell}\mathbb{I}(l=\tilde{y}_i)\Big|\boldsymbol{z},\mathcal{P}\right]$$

Following similar steps to the proof of Theorem 4.2, we have

$$\mathbb{E}_{DP}\left[Q^{-1}_{\cdot,\ell}\mathbb{I}(l=\tilde{y}_i)\Big|\boldsymbol{z},\mathcal{P}\right] = \mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}_{\cdot,\ell}\mathbb{E}_\lambda\left[\mathbb{I}(l=\tilde{y}_i)|\boldsymbol{w}\right]\Big|\boldsymbol{z},\mathcal{P}\right] = \mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}_{\cdot,\ell}Q_{l,y_i}\Big|\boldsymbol{z},\mathcal{P}\right] = Ie_i .$$

It follows

$$\mathbb{E}_{DP}[u_{0,c}|\boldsymbol{z},\mathcal{P}] = \sum_{i\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T} Ie_i = \sum_{i\in\mathcal{O}_{0,c}} y_i(0) . \tag{24}$$

We next calculate $\mathbb{E}[u^2_{0,c}|\boldsymbol{z},\mathcal{P}]$.

$$\mathbb{E}_{DP}\left[u^2_{0,c}|\boldsymbol{z},\mathcal{P}\right] = \sum_{i,j\in\mathcal{O}_{0,c}} \mathbb{E}_{DP}\left[\mathsf{y}^\mathsf{T} Q^{-1}_{\cdot,\tilde{y}_j}(Q^{-1}_{\cdot,\tilde{y}_i})^\mathsf{T}\mathsf{y}\Big|\boldsymbol{z},\mathcal{P}\right]$$

$$= \sum_{i,j\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T}\mathbb{E}_{DP}\left[\sum_{l,l'\in\mathcal{Y}} Q^{-1}_{\cdot,l'}(Q^{-1}_{\cdot,l})^\mathsf{T}\mathbb{I}(l=\tilde{y}_i)\mathbb{I}(l'=\tilde{y}_j)\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}$$

$$= \sum_{i,j\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[\sum_{l,l'\in\mathcal{Y}} Q^{-1}_{\cdot,l'}(Q^{-1}_{\cdot,l})^\mathsf{T}\mathbb{E}_\lambda\left[\mathbb{I}(l=\tilde{y}_i)\mathbb{I}(l'=\tilde{y}_j)|\boldsymbol{w}\right]\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}$$

$$= \sum_{i\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[\sum_{\ell\in\mathcal{Y}} Q^{-1}_{\cdot,\ell}(Q^{-1}_{\cdot,\ell})^\mathsf{T}Q_{\ell,y_i}\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}$$

$$+ \sum_{i\neq j\in\mathcal{O}_{0,c}} \mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[\sum_{\ell,\ell'\in\mathcal{Y}} Q^{-1}_{\cdot,\ell'}(Q^{-1}_{\cdot,\ell})^\mathsf{T}Q_{\ell,y_i}Q_{\ell',y_j}\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y} \tag{25}$$

For the first term, we write

$$\sum_{\ell\in\mathcal{Y}} Q^{-1}_{\cdot,\ell}(Q^{-1}_{\cdot,\ell})^\mathsf{T}Q_{\ell,y_i} = \sum_{\ell\in\mathcal{Y}} Q^{-1}_{\cdot,\ell}(Q^{-1}_{\cdot,\ell})^\mathsf{T}(Qe_{y_i})_\ell = Q^{-1}\mathrm{diag}(Qe_{y_i})Q^{-\mathsf{T}} , \tag{26}$$

Let $\hat{\boldsymbol{p}} = [\hat{p}_0(y|c)]_{y\in\mathcal{Y}} \in \mathbb{R}^{K\times 1}$ be the empirical distributions of outcomes of the controlled units in cluster $c$, arranged into a $K$-dimensional vector, such that $\hat{p}_\ell = \frac{1}{n_{0,c}}|\{i\in\mathcal{O}_{0,c}: y_i=\ell\}|$. In vector

form,

$$\hat{\boldsymbol{p}} = \frac{1}{n_{0,c}} \sum_{i \in \mathcal{O}_{0,c}} e_{y_i} \, .$$

Taking the expectation of both sides in (26) and summing over $i \in \mathcal{O}_{0,c}$, we get

$$\sum_{i \in \mathcal{O}_{0,c}} \mathbb{E}_{\boldsymbol{w}} \left[ \sum_{\ell \in \mathcal{Y}} Q_{\cdot,\ell}^{-1} (Q_{\cdot,\ell}^{-1})^{\mathsf{T}} Q_{\ell,y_i} \Big| \boldsymbol{z}, \mathcal{P} \right] = \mathbb{E}_{\boldsymbol{w}} \left[ Q^{-1} \text{diag} \left( Q \sum_{i \in \mathcal{O}_{0,c}} e_{y_i} \right) Q^{-\mathsf{T}} \Big| \boldsymbol{z}, \mathcal{P} \right]$$

$$= \mathbb{E}_{\boldsymbol{w}} \left[ Q^{-1} \text{diag}(n_{0,c} Q \hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \Big| \boldsymbol{z}, \mathcal{P} \right]$$

$$= n_{0,c} \mathbb{E}_{\boldsymbol{w}} \left[ Q^{-1} \text{diag}(Q \hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \Big| \boldsymbol{z}, \mathcal{P} \right]$$

We next proceed with the second term on the right-hand side of (25). We have

$$\sum_{\ell,\ell' \in \mathcal{Y}} Q_{\cdot,\ell'}^{-1} Q_{\cdot,\ell}^{-\mathsf{T}} Q_{\ell,y_i} Q_{\ell',y_j} = \sum_{\ell,\ell' \in \mathcal{Y}} Q_{\cdot,\ell'}^{-1} Q_{\cdot,\ell}^{-\mathsf{T}} (Q e_{y_j})_{\ell'} (Q e_{y_i})_{\ell}$$

$$= \sum_{\ell,\ell' \in \mathcal{Y}} Q_{\cdot,\ell'}^{-1} [(Q e_{y_j})_{\ell'} (Q e_{y_i})_{\ell}] Q_{\cdot,\ell}^{-\mathsf{T}}$$

$$= \sum_{\ell,\ell' \in \mathcal{Y}} Q_{\cdot,\ell'}^{-1} (Q e_{y_j} e_{y_i}^{\mathsf{T}} Q^{\mathsf{T}})_{\ell',\ell} Q_{\cdot,\ell}^{-\mathsf{T}}$$

$$= Q^{-1} Q e_{y_j} e_{y_i}^{\mathsf{T}} Q^{\mathsf{T}} Q^{-T} = e_{y_j} e_{y_i}^{\mathsf{T}} \, .$$

Taking the expectation of both sides of the above equation, we arrive at

$$\mathsf{y}^{\mathsf{T}} \mathbb{E}_{\boldsymbol{w}} \left[ \sum_{\ell,\ell' \in \mathcal{Y}} Q_{\cdot,\ell'}^{-1} (Q_{\cdot,\ell}^{-1})^{\mathsf{T}} Q_{\ell,y_i} Q_{\ell',y_j} \Big| \boldsymbol{z}, \mathcal{P} \right] \mathsf{y} = \mathbb{E}_{\boldsymbol{w}} \left[ \mathsf{y}^{\mathsf{T}} e_{y_j} e_{y_i}^{\mathsf{T}} \mathsf{y} \Big| \boldsymbol{z}, \mathcal{P} \right]$$

$$= y_i(0) y_j(0) \, ,$$

where the second equality holds since $i \neq j \in \mathcal{O}_{0,c}$. Putting these pieces together, we obtain

$$\mathbb{E}_{DP}[u_{0,c}^2 | \boldsymbol{z}, \mathcal{P}] = n_{0,c} \mathsf{y}^{\mathsf{T}} \mathbb{E}_{\boldsymbol{w}} \left[ Q^{-1} \text{diag}(Q \hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \Big| \boldsymbol{z}, \mathcal{P} \right] \mathsf{y} + \sum_{i \neq j \in \mathcal{O}_{0,c}} y_i(0) y_j(0) \, ,$$

which along with (24) gives us

$$\text{Var}_{DP}(u_{0,c} | \boldsymbol{z}, \mathcal{P})$$

$$= \mathbb{E}_{DP}[u_{0,c}^2 | \boldsymbol{z}, \mathcal{P}] - \mathbb{E}_{DP}[u_{0,c} | \boldsymbol{z}, \mathcal{P}]^2$$

$$= n_{0,c} \mathsf{y}^{\mathsf{T}} \mathbb{E}_{\boldsymbol{w}} \left[ Q^{-1} \text{diag}(Q \hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \Big| \boldsymbol{z}, \mathcal{P} \right] \mathsf{y} + \sum_{i \neq j \in \mathcal{O}_{o,c}} y_i(0) y_j(0) - \left( \sum_{i \in \mathcal{O}_{0,c}} y_i(0) \right)^2$$

$$= n_{0,c} \mathsf{y}^{\mathsf{T}} \mathbb{E}_{\boldsymbol{w}} \left[ Q^{-1} \text{diag}(Q \hat{\boldsymbol{p}}) Q^{-\mathsf{T}} \Big| \boldsymbol{z}, \mathcal{P} \right] \mathsf{y} - \sum_{i \in \mathcal{O}_{0,c}} y_i^2(0) \, . \tag{27}$$

- **Decomposing** $\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z},\mathcal{P})$

We wish to bound $\mathrm{Var}(u_{0,c}|\boldsymbol{z},\mathcal{P})$. We begin by decomposing it into two distinct terms. Let $\tilde{\boldsymbol{q}} = [\tilde{q}_0(y|c)]_{y\in\mathcal{Y}} \in \mathbb{R}^{K\times 1}$ be the distribution constructed in the DP mechanism after adding noise to the empirical distribution $\hat{\boldsymbol{p}}$, truncation, and normalization. We consider the following decomposition:

$$Q^{-1}\mathrm{diag}(Q\hat{\boldsymbol{p}})Q^{-\mathsf{T}} = Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}} + Q^{-1}\mathrm{diag}(Q(\hat{\boldsymbol{p}}-\tilde{\boldsymbol{q}}))Q^{-\mathsf{T}}.$$

Plugging into (27),

$$\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z},\mathcal{P}) =$$

$$n_{0,c}\left(\underbrace{\mathsf{y}^{\mathsf{T}}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}}|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}}_{\text{term I}} + \underbrace{\mathsf{y}^{\mathsf{T}}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\mathrm{diag}(Q(\hat{\boldsymbol{p}}-\tilde{\boldsymbol{q}}))Q^{-\mathsf{T}}|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}}_{\text{term II}}\right) - \sum_{i\in\mathcal{O}_{0,c}}y_i^2(0).$$

- **Bounding Term I**

By definition, we can write $Q$ as $Q = (1-\lambda)I + \lambda\tilde{\boldsymbol{q}}\mathbf{1}^{\mathsf{T}}$, with $\mathbf{1}\in\mathbb{R}^{K\times 1}$ indicating the all-one vector. Furthermore, $\mathbf{1}^{\mathsf{T}}\tilde{\boldsymbol{q}} = 1$ because $\tilde{\boldsymbol{q}}$ is a probability distribution (see (Esfandiari et al., 2022, Theorem 6)). Using the Sherman–Morrison formula, we obtain

$$Q^{-1} = \frac{1}{1-\lambda}I - \frac{\lambda}{1-\lambda}\tilde{\boldsymbol{q}}\mathbf{1}^{\mathsf{T}}.$$

Plugging for $Q$ and $Q^{-1}$, we have the following chain of identities:

$$Q\tilde{\boldsymbol{q}} = (1-\lambda)\tilde{\boldsymbol{q}} + \lambda\tilde{\boldsymbol{q}} = \tilde{\boldsymbol{q}},$$

$$Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}}) = \frac{1}{1-\lambda}\mathrm{diag}(\tilde{\boldsymbol{q}}) - \frac{\lambda}{1-\lambda}\tilde{\boldsymbol{q}}\tilde{\boldsymbol{q}}^{\mathsf{T}},$$

$$Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}})Q^{-T} = \frac{1}{(1-\lambda)^2}\mathrm{diag}(\tilde{\boldsymbol{q}}) + \frac{\lambda^2-2\lambda}{(1-\lambda)^2}\tilde{\boldsymbol{q}}\tilde{\boldsymbol{q}}^{\mathsf{T}}.$$

Using the last identity, we have

$$\mathsf{y}^{\mathsf{T}}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}}|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}$$

$$= \frac{1}{(1-\lambda)^2}\mathsf{y}^{\mathsf{T}}\mathbb{E}_{\boldsymbol{w}}\left[\mathrm{diag}(\tilde{q})\right]\mathsf{y} + \frac{\lambda^2-2\lambda}{(1-\lambda)^2}\mathsf{y}^{\mathsf{T}}\mathbb{E}_{\boldsymbol{w}}\left[\tilde{\boldsymbol{q}}\tilde{\boldsymbol{q}}^{\mathsf{T}}\right]\mathsf{y}$$

$$= \frac{1}{(1-\lambda)^2}\sum_{y\in\mathcal{Y}}\mathbb{E}_{\boldsymbol{w}}\left[\tilde{q}_y\right]y^2 + \frac{\lambda^2-2\lambda}{(1-\lambda)^2}\mathbb{E}_{\boldsymbol{w}}\left[\left(\sum_{y\in\mathcal{Y}}y\tilde{q}_y\right)^2\right]$$

$$= \mathbb{E}_{\boldsymbol{w}}\left[\frac{1}{(1-\lambda)^2}\mathbb{E}_{\tilde{\boldsymbol{q}}}[y_i^2(0)] + \left(1 - \frac{1}{(1-\lambda)^2}\right)\mathbb{E}_{\tilde{\boldsymbol{q}}}[y_i(0)]^2\,\Big|\,\boldsymbol{z},\mathcal{P}\right]$$

$$= \mathbb{E}_{\boldsymbol{w}}\left[\frac{1}{(1-\lambda)^2}\left(\mathbb{E}_{\tilde{\boldsymbol{q}}}[y_i^2(0)] - \mathbb{E}_{\tilde{\boldsymbol{q}}}[y_i(0)]^2\right) + E_{\tilde{\boldsymbol{q}}}[y_i(0)]^2\,\Big|\,\boldsymbol{z},\mathcal{P}\right]. \qquad (28)$$

which can also be written as:

$$\mathsf{y}^{\mathsf{T}}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}}\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y} = \frac{1}{(1-\lambda)^2}\left(\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i^2(0)|\boldsymbol{z},\mathcal{P}] - \mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)|\boldsymbol{z},\mathcal{P}]^2\right) + E_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)|\boldsymbol{z},\mathcal{P}]^2$$

39

In the following lemma, we relate the expectation of outcomes with respect to $\tilde{\boldsymbol{q}}$ to their expectation with respect to $\hat{\boldsymbol{p}}$.

**Lemma A.4.** *If outcomes are bounded by $B$,*

$$\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i^2(0)]|\boldsymbol{z},\mathcal{P}] \leq B^2 \mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1 + \frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)^2\,.$$

$$\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)|\boldsymbol{z},\mathcal{P}]^2 \geq \left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)\right)^2 - 2B^2 \mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1\,.$$

Using the above lemma, we obtain:

$$\mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\text{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}}\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y} \leq \frac{1}{(1-\lambda)^2}\left[B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1 + \frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)^2\right]$$

$$-\left(\frac{1}{(1-\lambda)^2}-1\right)\left\{\left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)\right)^2 - 2B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1\right\}\,,$$

which can be simplified to

$$\mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\text{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}}\Big|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}$$

$$\leq \frac{1}{(1-\lambda)^2}\left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)^2 - \left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)\right)^2\right)$$

$$+\left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)\right)^2 + B^2\left(\frac{3}{(1-\lambda)^2}+2\right)\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1$$

$$\leq \frac{n_{0,c}-1}{n_{0,c}}\frac{S^2(\vec{y}_{0,c})}{(1-\lambda)^2} + \left(\overline{\vec{y}_c(0)}\right)^2 + B^2\left(\frac{3}{(1-\lambda)^2}+2\right)\mathbb{E}_{\boldsymbol{w}}\left[\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1\right]\,, \tag{29}$$

where $S^2(\vec{u}) = \frac{1}{|\vec{u}-1|}\sum_{a\in\vec{u}}(a-\bar{a})^2$ and $\bar{\bar{u}} = \frac{1}{|\vec{u}|}\sum_{a\in\vec{u}} a$.

• **Bounding Term II**

We begin with the two inequalities

$$\mathsf{y}^\mathsf{T}Q^{-1}\text{diag}(Q(\hat{\boldsymbol{p}} - \tilde{\boldsymbol{q}}))Q^{-\mathsf{T}}\mathsf{y} \leq \|Q^{-\mathsf{T}}\mathsf{y}\|^2\|Q(\hat{\boldsymbol{p}} - \tilde{\boldsymbol{q}})\|_\infty$$

$$\leq \|Q^{-\mathsf{T}}\mathsf{y}\|^2|Q|_\infty\|\hat{\boldsymbol{p}} - \tilde{\boldsymbol{q}}\|_1\,, \tag{30}$$

where $|Q|_\infty = \max_{i,j}|Q_{ij}|$. For every $\ell \in \mathcal{Y}$ $\tilde{q}_\ell \geq \gamma$ (see (Esfandiari et al., 2022, Theorem 6)), and $\sum_{\ell\in\mathcal{Y}}\tilde{q}_l = 1$, which implies that $\forall l \in \mathcal{Y}, \tilde{q}_l \leq 1 - (K-1)\gamma$. Therefore, by definition of $Q$, we have

$$|Q|_\infty \leq 1 - \lambda + \lambda(1 - (K-1)\gamma) = 1 - \lambda(K-1)\gamma\,.$$

The following lemma bounds the maximum singular value of matrix $Q$.

**Lemma A.5.** *The maximum singular value of label randomization matrix $Q^{-1}$ is at most $\frac{\lambda\sqrt{K}+1}{1-\lambda}$.*

Using Lemma A.5, we get

$$\mathsf{y}^\mathsf{T} Q^{-1}\mathrm{diag}(Q(\hat{\boldsymbol{p}}-\tilde{\boldsymbol{q}}))Q^{-\mathsf{T}}\mathsf{y} \leq \frac{(\lambda\sqrt{K}+1)^2}{(1-\lambda)^2}\|\mathsf{y}\|^2(1-\lambda(K-1)\gamma)\|\hat{\boldsymbol{p}}-\tilde{\boldsymbol{q}}\|_1. \tag{31}$$

• **Bounding** $\mathrm{Var}(u_{0,c}|\boldsymbol{c})$

Combining (29) and (31) with the expression of $\mathrm{Var}(u_{0,c}|\boldsymbol{c},\boldsymbol{z})$, we get

$\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z},\mathcal{P})$

$$= n_{0,c}\left(\underbrace{\mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\mathrm{diag}(Q\tilde{\boldsymbol{q}})Q^{-\mathsf{T}}|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}}_{\text{term I}}+\underbrace{\mathsf{y}^\mathsf{T}\mathbb{E}_{\boldsymbol{w}}\left[Q^{-1}\mathrm{diag}(Q(\hat{\boldsymbol{p}}-\tilde{\boldsymbol{q}}))Q^{-\mathsf{T}}|\boldsymbol{z},\mathcal{P}\right]\mathsf{y}}_{\text{term II}}\right) - \sum_{i\in\mathcal{O}_{0,c}}y_i^2(0)$$

$$\leq (n_{0,c}-1)\frac{S^2(\vec{y}_{0,c})}{(1-\lambda)^2} + n_{0,c}\left((\overline{\vec{y}_{0,c}})^2+B^2\left(\frac{3}{(1-\lambda)^2}+2\right)\mathbb{E}_{\boldsymbol{w}}\left[\|\tilde{\boldsymbol{q}}-\hat{\boldsymbol{p}}\|_1\right]+\right.$$

$$\left.\frac{(\lambda\sqrt{K}+1)^2}{(1-\lambda)^2}\|\mathsf{y}\|^2(1-\lambda(K-1)\gamma)\mathbb{E}_{\boldsymbol{w}}\|\hat{\boldsymbol{p}}-\tilde{\boldsymbol{q}}\|_1\right) - \sum_{i\in\mathcal{O}_{0,c}}y_i^2(0)$$

$$= \frac{n_{0,c}-1}{(1-\lambda)^2}S^2(\vec{y}_{0,c}) + n_{0,c}(\overline{\vec{y}_{0,c}})^2 - \sum_{i\in\mathcal{O}_{0,c}}y_i^2(0) + n_{0,c}A'_{0,c}\mathbb{E}_{\boldsymbol{w}}[\|\tilde{\boldsymbol{q}}-\hat{\boldsymbol{p}}\|_1]$$

$$= (n_{0,c}-1)\left(\frac{1}{(1-\lambda)^2}-1\right)S^2(\vec{y}_{0,c}) + n_{0,c}A'_{0,c}\mathbb{E}_{\boldsymbol{w}}[\|\tilde{\boldsymbol{q}}-\hat{\boldsymbol{p}}\|_1], \tag{32}$$

with

$$A'_{0,c} := B^2\left(\frac{3}{(1-\lambda)^2}+2\right) + \frac{(\lambda\sqrt{K}+1)^2}{(1-\lambda)^2}\|\mathsf{y}\|^2(1-\lambda(K-1)\gamma).$$

The final step is bounding the term $\mathbb{E}_{\boldsymbol{w}}[\|\tilde{\boldsymbol{q}}-\hat{\boldsymbol{p}}\|_1]$.

**Lemma A.6.** *Recall the notation $\tilde{\boldsymbol{q}}=[\tilde{q}_0(y|c)]_{y\in\mathcal{Y}}$ and $\hat{\boldsymbol{p}}=[\hat{p}_0(y|c)]$. Then,*

$$\mathbb{E}_{\boldsymbol{w}}[\|\tilde{\boldsymbol{q}}-\hat{\boldsymbol{p}}\|_1] \leq 2K\left[\gamma + \frac{\sigma}{n_{0,c}}\left(e^{-\gamma n_{0,c}/\sigma}-e^{-n_{0,c}/\sigma}\right)\right].$$

By using Lemma A.6 in (32), we obtain

$$\mathrm{Var}_{DP}(u_{0,c}|\boldsymbol{z},\mathcal{P}) \leq (n_{0,c}-1)\left(\frac{1}{(1-\lambda)^2}-1\right)S^2(\vec{y}_c(0)) + n_{0,c}A(n_{0,c}), \tag{33}$$

with

$$A(x) = A'_{0,c}2K\left[\gamma + \frac{\sigma}{x}\left(e^{-\gamma x/\sigma}-e^{-x/\sigma}\right)\right]$$

$$= 2K\left(B^2\left(\frac{3}{(1-\lambda)^2}+2\right) + \frac{(\lambda\sqrt{K}+1)^2}{(1-\lambda)^2}\|\mathsf{y}\|^2(1-\lambda(K-1)\gamma)\right)\left[\gamma + \frac{\sigma}{x}\left(e^{-\gamma x/\sigma}-e^{-x/\sigma}\right)\right]$$

A similar bound can be derived for $\mathrm{Var}(u_{1,c}|\boldsymbol{c})$, which in conjunction with (23) gives the desired result.

## Proof of Lemma A.4

We recall the statement of Lemma A.4 below for convenience.

$$\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i^2(0)]|\boldsymbol{z},\mathcal{P}] \leq B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1 + \frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)^2. \tag{34}$$

$$\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)|\boldsymbol{z},\mathcal{P}]^2 \geq \left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)\right)^2 - 2B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1. \tag{35}$$

*Proof.* Since the outcomes are bounded by $B$, we have

$$\left|\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i^2(0)|\boldsymbol{z},\mathcal{P}] - \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i^2(0)]\right| = \left|\mathbb{E}_{\boldsymbol{w}}\left[\sum_{y\in\mathcal{Y}} y^2(\tilde{\boldsymbol{q}}_y - \hat{\boldsymbol{p}}_y)\Big|\boldsymbol{z},\mathcal{P}\right]\right|$$

$$\leq \mathbb{E}_{\boldsymbol{w}}\left[\left|\sum_{y\in\mathcal{Y}} y^2(\tilde{\boldsymbol{q}}_y - \hat{\boldsymbol{p}}_y)\right|\ \Big|\boldsymbol{z},\mathcal{P}\right] \leq B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1.$$

Therefore,

$$\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i^2(0)]|\boldsymbol{z},\mathcal{P}] \leq \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i^2(0)|\boldsymbol{z},\mathcal{P}] + B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1.$$

We next note that

$$\mathbb{E}_{\hat{\boldsymbol{p}}}[y_i^2(0)|\boldsymbol{z},\mathcal{P}] = \mathbb{E}\left[\sum_{y\in\mathcal{Y}}\sum_{i\in\mathcal{O}_{0,c}} \frac{\mathbb{I}(y_i = y)}{n_{0,c}} y^2\Big|\boldsymbol{z},\mathcal{P}\right] = \frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)^2.$$

This completes the proof of (34). Likewise we have

$$\left|\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)]^2 - \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i(0)]^2\right| = \left|\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)] - \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i(0)]\right| \cdot \left|\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)] + \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i(0)]\right|$$

$$\leq 2B\left|\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i(0)] - \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i(0)]\right|$$

$$\leq 2B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1.$$

Therefore, we obtain:

$$\mathbb{E}_{\tilde{\boldsymbol{q}},\boldsymbol{w}}[y_i^2(0)|\boldsymbol{z},\mathcal{P}] \geq \mathbb{E}_{\hat{\boldsymbol{p}}}[y_i(0)]^2 - 2B^2\mathbb{E}_{\boldsymbol{w}}\|\tilde{\boldsymbol{q}} - \hat{\boldsymbol{p}}\|_1.$$

We next note that

$$\mathbb{E}_{\hat{\boldsymbol{p}}}[y_i(0)]^2 = \left(\sum_{y\in\mathcal{Y}}\sum_{i\in\mathcal{O}_{0,c}} \frac{\mathbb{I}(y_i = y)}{n_{0,c}} y\right)^2 = \left(\frac{1}{n_{0,c}}\sum_{i\in\mathcal{O}_{0,c}} y_i(0)\right)^2.$$

This completes the proof of (35). □

## Proof of Lemma A.5

**Lemma A.5.** *The maximum singular value of label randomization matrix $Q^{-1}$ is at most $\frac{\lambda\sqrt{K}+1}{1-\lambda}$.*

*Proof.* For any unit norm vector $u$ we have

$$u^\mathsf{T}Q^{-1} = \frac{1}{1-\lambda}u^\mathsf{T} - \frac{\lambda}{1-\lambda}u^\mathsf{T}\tilde{q}\mathbf{1}^\mathsf{T}.$$

Therefore, by triangle inequality

$$\|u^\mathsf{T}Q^{-1}\| \leq \frac{1}{1-\lambda} + \frac{\lambda}{1-\lambda}\|\tilde{q}\| \cdot \|\mathbf{1}\| \leq \frac{1+\lambda\sqrt{K}}{1-\lambda},$$

where in the last step we used $\|u\| = 1$ and $\|\tilde{q}\| \leq \|\tilde{q}\|_1 = 1$. $\qquad\square$

## Proof of Lemma A.6

**Lemma A.6.** *Recall the notation $\tilde{q} = [\tilde{q}_0(y|c)]_{y\in\mathcal{Y}}$ and $\hat{p} = [\hat{p}_{0,D}(y|c)]$, where we dropped the subscript $c$ to lighten the notation. Then,*

$$\mathbb{E}_{w}[\|\tilde{q} - \hat{p}\|_1] \leq 2K\left[\gamma + \frac{\sigma}{n_{0,c}}\left(e^{-\gamma n_{0,c}/\sigma} - e^{-n_{0,c}/\sigma}\right)\right].$$

We follow the proof of (Esfandiari et al., 2022, Lemma 5). By a tighter derivation which carries over in a straightforward way, we obtain the following bound analogous to Equation (6) therein:

$$\mathbb{E}_{w}[\|\tilde{q} - \hat{p}\|_1] \leq 2\sum_{y\in\mathcal{Y}}\mathbb{E}_{w}[\max(\gamma, \min(1, |w_{y,c}|))] = 2K\mathbb{E}[\max(\gamma, \min(1, V))],$$

where $V = |w_{y,c}| \sim \mathrm{Exp}(n_{0,c}/\sigma)$, since $w_{y,c} \sim \mathrm{Laplace}(\sigma/n_{0,c})$. For a random variable $V \sim \mathrm{Exp}(\alpha)$, we have

$$\begin{aligned}
\mathbb{E}[\max(\gamma, \min(1, V))] &= \int_0^\gamma \gamma\alpha e^{-\alpha v}\mathrm{d}v + \int_\gamma^1 v\alpha e^{-\alpha v}\mathrm{d}v + \int_1^\infty \alpha e^{-\alpha v}\mathrm{d}v \\
&= -\gamma e^{-\alpha v}\Big|_0^\gamma - (v + \frac{1}{\alpha})e^{-\alpha v}\Big|_\gamma^1 - e^{-\alpha v}|_1^\infty \\
&= \gamma + \frac{1}{\alpha}(e^{-\alpha\gamma} - e^{-\alpha}),
\end{aligned}$$

which after substituting for $u = n_{0,c}/\sigma$ gives the claim.