

# Controlling Privacy Loss in Sampling Schemes: an Analysis of Stratified and Cluster Sampling

Mark Bun 

Department of Computer Science at Boston University, USA

Jörg Drechsler 

Institute for Employment Research, Germany

Joint Program in Survey Methodology, University of Maryland, USA

Marco Gaboardi 

Department of Computer Science at Boston University, USA

Audra McMillan<sup>1</sup> 

Apple, USA

Jayshree Sarathy 

Harvard John A. Paulson School of Engineering and Applied Sciences, USA

---

## Abstract

Sampling schemes are fundamental tools in statistics, survey design, and algorithm design. A fundamental result in differential privacy is that a differentially private mechanism run on a *simple random* sample of a population provides stronger privacy guarantees than the same algorithm run on the entire population. However, in practice, sampling designs are often more complex than the simple, data-independent sampling schemes that are addressed in prior work. In this work, we extend the study of privacy amplification results to more complex, data-dependent sampling schemes. We find that not only do these sampling schemes often fail to amplify privacy, they can actually result in privacy degradation. We analyze the privacy implications of the pervasive cluster sampling and stratified sampling paradigms, as well as provide some insight into the study of more general sampling designs.

**2012 ACM Subject Classification** Security and privacy → Privacy protections

**Keywords and phrases** privacy, differential privacy, survey design, survey sampling

**Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

## 1 Introduction

Sampling schemes are fundamental tools in statistics, survey design, and algorithm design. For example, they are used in social science research to conduct surveys on a random sample of a target population. They are also used in machine learning to improve the efficiency and accuracy of algorithms on large datasets. In many of these applications, however, the datasets are sensitive and privacy is a concern. Intuition suggests that (sub)sampling a dataset before analysing it provides additional privacy, since it gives individuals plausible deniability about whether their data was included or not. This intuition has been formalized for some types of sampling schemes (such as simple random sampling with and without replacement and Poisson sampling) in a series of papers in the differential privacy literature [23, 33, 11, 31]. Such *privacy amplification by subsampling* results can provide tight privacy accounting when analysing algorithms that incorporate subsampling, e.g. [32, 1, 21, 28, 19]. However, in practice, sampling designs are often more complex than the simple, data independent sampling schemes that are addressed in prior work. In this work, we extend the study of privacy amplification results to more complex and data dependent sampling schemes.

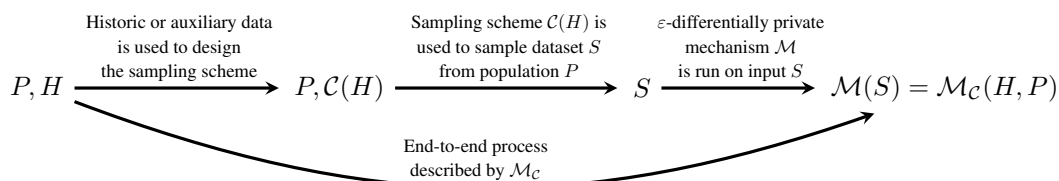
We consider the setting described in Figure 1. We have a *population*  $P$  and a historic or auxiliary data set  $H$  which is used to inform the sampling design. We think about the sampling scheme as a

---

<sup>1</sup> Corresponding author



## 23:2 Controlling Privacy Loss in Sampling Schemes



■ **Figure 1** The structure of using a data-dependent sampling scheme.

42 function  $\mathcal{C}(H)$  of the historic or auxiliary data  $H$ . Using this sampling scheme, we draw a sample  
 43  $S$  from the population  $P$ , on which we run the differentially private mechanism  $\mathcal{M}$ . We can think  
 44 about these multiple steps as comprising a mechanism  $\mathcal{M}_{\mathcal{C}}(H, P)$  working directly on the population  
 45  $P$  and the historic data  $H$  whose privacy depends on the privacy of the mechanism  $\mathcal{M}$  and on the  
 46 properties of the sampling scheme  $\mathcal{C}(H)$ . While this is the general framework for the problem we  
 47 study, we state the technical results in this paper for the simplified case where  $H = P$ ; see Section 2.1  
 48 for further discussion.

### 49 1.1 Our contributions

50 We primarily focus on two classes of sampling schemes that are common in practice: *cluster sampling*  
 51 and *stratified sampling*. In (single-stage) cluster sampling, the population arrives partitioned into  
 52 disjoint clusters. A sample is obtained by selecting a small number of clusters at random, and then  
 53 including all of the individuals from those chosen clusters. In stratified sampling, the population is  
 54 partitioned into “strata.” Individuals are then sampled at different rates according to which stratum  
 55 they belong to.

56 For these more complex schemes, we find that privacy amplification can be negligible even when  
 57 only a small fraction of the population is included in the final sample. Moreover, in settings where  
 58 the sampling design is data dependent, privacy degradation can occur – some sampling designs can  
 59 actually make privacy guarantees worse. Intuitively, this is because the sample design itself can reveal  
 60 sensitive information. Our goal in this paper is to explain how and why these phenomena occur and  
 61 introduce technical tools for understanding the privacy implications of concrete sampling designs.

62 **Understanding randomised and data-dependent sampling.** It is simple to show that deterministic,  
 63 data-dependent sampling designs do not achieve privacy amplification, and can suffer privacy degrad-  
 64 ation. Motivated by this observation, we start by studying the privacy implications of randomised and  
 65 data-dependent sampling, attempting to isolate their effects in the simplest possible setting.

66 Specifically, we aim to understand sampling schemes of the following form: For a possibly  
 67 randomised function  $f$  (an “allocation rule”), sample  $f(P)$  individuals uniformly from  $P$  without  
 68 replacement. In Section 3, we study the case where  $f$  is randomised but data-independent, i.e., the  
 69 number of individuals samples is drawn from a distribution that does not depend on  $P$ . We give  
 70 an essentially complete characterization of what level of amplification is possible in terms of this  
 71 distribution.

72 In Section 4, we turn our attention to data-dependent sampling. We identify necessary conditions  
 73 for allocation rules  $f$  to enable privacy amplification by way of a hypothesis testing perspective;  
 74 intuitively, for  $f$  to be a good amplifier, every differentially private algorithm must fail to distinguish  
 75 the distributions of  $f(P)$  and  $f(P')$  for neighboring  $P, P'$ . We also study a specific natural allocation  
 76 rule called *proportional allocation* that is commonly applied in stratified sampling. We design a  
 77 simple randomised rounding method that offers a minor change to the way proportional allocation is  
 78 generally implemented in practice, but that offers substantially better privacy amplification.

79 **Cluster sampling.** In Section 5, we study cluster sampling where a population partitioned into  $k$   
80 clusters is sampled by selecting  $m$  clusters uniformly at random without replacement. Our results  
81 give tradeoffs between the privacy amplification achievable and the sizes of the clusters. In particular,  
82 privacy amplification is possible when all of the clusters are small. As the cluster sizes grow, the best  
83 achievable privacy loss rapidly approaches the baseline of the privacy guarantee of  $\mathcal{M}$ . We provide  
84 some insight into these results by connecting the privacy loss to the ability of a hypothesis test to  
85 determine from a differentially private output which clusters were included in the sample.

86 **Stratified sampling.** Building on our randomised rounding method for the “single-stratum” case, we  
87 show that stratified sampling with the proportional allocation rule amplifies privacy. Unfortunately, as  
88 in the single stratum case, there are natural lower bounds which limit extending this approach to other  
89 common allocation rules.

90 A common goal when choosing an allocation function  $f$  (a function which decides how many  
91 samples to draw from each stratum) is to minimise the variance of a particular statistic. For example,  
92 the popular Neyman allocation is the optimal allocation for computing the population mean. A natural  
93 question then is how to define and compute the optimal allocation when privacy is a concern? In  
94 this work, we will formulate the notion of an optimal allocation under privacy constraints. This  
95 formulation is somewhat subtle since the privacy implications of different allocation methods need to  
96 be properly accounted for. Our goal is to initiate the study of alternative allocation functions that may  
97 prove useful when privacy is a concern.

## 98 1.2 Related work

99 Several works have studied the privacy amplification of simple sampling schemes. Kasiviswanathan  
100 et al. [23] and Beimel et al. [9] showed that applying Poisson sampling before running a differentially  
101 private mechanism improves its end-to-end privacy guarantee. Subsequently, Bun et al. [11] analyzed  
102 simple random sampling with replacement in a similar way. Beimel et al. [10], Bassily et al. [7], and  
103 Wang et al. [34] analyzed simple random sampling without replacement. Imola and Chaudhuri [20]  
104 provide lower and upper bounds on privacy amplification when sampling from a multidimensional  
105 Bernoulli family, a task which has direct applications to Bayesian inference. Balle et al. [5] unified  
106 the analyses of privacy amplification of these mechanisms using the lenses of *probabilistic couplings*,  
107 an approach that we also use in this paper. The effects that sampling can have on differentially private  
108 mechanisms is also studied from a different perspective in [13]. However, none of the prior works  
109 consider the privacy amplification of more complex, data-dependent sampling schemes commonly  
110 used in practice. To the best of our knowledge, this paper is the first to do so.

## 111 2 Background

### 112 2.1 Data-dependent sampling schemes

113 In the data-driven sciences, data is often obtained by sampling a fraction of the population of interest.  
114 This sample can be created in a wide variety of ways, referred to as the sample design. Sample  
115 designs can vary from simple designs such as taking a uniformly random subset of a fixed size, to  
116 more complex data-dependent sampling designs like cluster or stratified sampling. Data-dependent  
117 sampling designs achieve accuracy and meet budgeting goals by using historic or auxiliary data to  
118 exploit structure in the population. The privacy implications of simple random sampling are quite well  
119 understood from prior work. In this work, we will move beyond simple random sampling to analyse  
120 the privacy implications of more complex sampling designs, including data-dependent sampling.

121 An outline of the schema for data dependent sampling designs is given in Figure 1. There are  
122 ostensibly two datasets:  $H$ , the historic or auxiliary data that is used to design the sampling scheme

## 23:4 Controlling Privacy Loss in Sampling Schemes

123  $\mathcal{C}(H)$ , and  $P$ , the current population that is sampled from. For the remainder of this paper, we  
124 make the simplifying assumption that  $H = P$ . That is, we will not distinguish between the historic  
125 or auxiliary data and the “current” data. Even if we only care about maintaining the privacy of  
126 the individuals in population  $P$ , this assumption is required if we have no information about the  
127 relationship between  $H$  and  $P$ . Thus, we view the function  $\mathcal{M}_{\mathcal{C}}(P, H)$  as simply a function of  $P$ . We  
128 will refer to the size of the sample  $S$  as the *sample size*, and the fraction  $|S|/|P|$  as the *sampling rate*.

129 More refined models can be obtained by imposing specific assumptions on the relationship  
130 between  $H$  and  $P$ , for example, by modeling the temporal correlation between historic and current  
131 data. We leave this for future work.

### 132 2.2 Differential privacy

133 Differential privacy (DP) is a measure of stability for randomised algorithms. It bounds the change in  
134 the distribution of the outputs of a randomised algorithm when provided with two datasets differing  
135 on the data of a single individual. We will call such datasets neighboring. In order to formalise what  
136 a “bounded change” means, we define  $(\epsilon, \delta)$ -indistinguishability. Two random variables  $P$  and  $Q$   
137 over the same probability space are  $(\epsilon, \delta)$ -indistinguishable if for all sets of outcomes  $E$  over that  
138 probability space,

$$139 \quad e^{-\epsilon} (\Pr(Q \in E) - \delta) \leq \Pr(P \in E) \leq e^{\epsilon} \Pr(Q \in E) + \delta.$$

140 If  $\delta = 0$  then we will say that  $P$  and  $Q$  are  $\epsilon$ -indistinguishable. For any  $n \in \mathbb{N}$ , let  $\mathcal{U}^n$  be the set of all  
141 datasets of size  $n$  over elements of the data universe  $\mathcal{U}$ . Let  $\mathcal{U}^* = \cup_{n \in \mathbb{N}} \mathcal{U}^n$  be the set of all possible  
142 datasets. We discuss two privacy definitions in this work corresponding to two different neighboring  
143 relations: *unbounded* differential privacy and *bounded* differential privacy. We will say two datasets  
144 are *unbounded neighbors* if one can be obtained from the other by adding or removing a single data  
145 point, and *bounded neighbors* if they have the same size, and one can be obtained from the other by  
146 changing the data of a single individual.

147 ► **Definition 1.** A mechanism  $\mathcal{M} : \mathcal{U}^* \rightarrow \mathcal{O}$  is  $(\epsilon, \delta)$ -unbounded (resp. bounded) differentially  
148 private (DP) if for all pairs of unbounded (resp. bounded) neighboring datasets  $P$  and  $P'$ ,  $\mathcal{M}(P)$   
149 and  $\mathcal{M}(P')$  are  $(\epsilon, \delta)$ -indistinguishable.

150 We will use both bounded and unbounded DP throughout the paper as they are appropriate in  
151 different settings. When considering which notion to choose, it is important to consider which  
152 guarantees are meaningful in context. For example, it will be common in the sample designs we cover  
153 for the size of the sample  $S$  (see Figure 1) to be data-dependent. When considering these sampling  
154 designs, we will focus on mechanisms  $\mathcal{M}$  that satisfy unbounded DP since bounded DP does not  
155 protect the sample size. However, bounded DP may be more appropriate for the privacy guarantee on  
156  $\mathcal{M}_{\mathcal{C}}$  in applications where it is unrealistic to assume that an individual can choose not to be part of  
157 the auxiliary dataset or the population. For example, the auxiliary data may be administrative data,  
158 data from a mandatory census, or data from a monopolistic service provider. Results and intuition are  
159 often similar between unbounded and bounded DP, although care should be taken when translating  
160 between the two notions. We note in particular that any  $\epsilon$ -unbounded DP mechanism is  $2\epsilon$ -bounded  
161 DP.

### 162 2.3 Privacy amplification with uniform random sampling

163 Sampling does not provide strong differential privacy guarantees on its own. But when employed as a  
164 pre-processing step in a differentially private algorithm, it can amplify existing privacy guarantees.  
165 Intuitively, this is because if the choice of individuals is kept secret, sampling provides data subjects

166 the plausible deniability to claim that their data was or was not in the final data set. This effect was  
 167 first explicitly articulated in [29], and a formal treatment of the phenomenon was given in [5]. Three  
 168 types of sampling are analysed in [5]: simple random sampling with replacement, simple random  
 169 sampling without replacement, and Poisson sampling. In all three settings the privacy amplification is  
 170 proportional to the probability of an individual not being included in the final computation. To gain  
 171 some intuition before we move into the more complicated sampling schemes that are the focus on  
 172 this paper, let us state and discuss the results from [5].

173 ► **Theorem 2.** [5] *Let  $\mathcal{C}$  be a sampling scheme that samples  $m$  values out of  $n$  possible values*  
 174 *without replacement. Given an  $(\varepsilon, \delta)$ -bounded differentially private mechanism  $\mathcal{M}$ , we have that*  
 175  *$\mathcal{M}_{\mathcal{C}}$  is  $(\varepsilon', \delta')$ -bounded differentially private for  $\varepsilon' = \log(1 + \frac{m}{n}(e^{\varepsilon} - 1))$  and  $\delta' = \frac{m}{n}\delta$ .*

176 To consider the implications of this result, notice that  $\varepsilon' \leq \varepsilon$  for all values of  $m \leq n$  so the  
 177 sampled mechanism  $\mathcal{M}_{\mathcal{C}}$  is strictly more private than the original mechanism  $\mathcal{M}$ . Further, taking  
 178 into account the following two approximations which hold for small  $x$ ,

$$179 \quad e^x - 1 \approx x \tag{1}$$

$$180 \quad \log(1 + x) \approx x, \tag{2}$$

182 we have that for small  $\varepsilon$ ,  $\varepsilon' \approx \frac{m}{n}\varepsilon$ . So the degree of amplification in both parameters is roughly  
 183 proportional to the sampling rate  $m/n$ .

## 184 2.4 How do people use subsampling amplification results?

185 Suppose we have a dataset that contains  $n$  records, and we want to estimate the proportion of  
 186 individuals that satisfy some attribute in an  $\varepsilon$ -DP manner. Let us set our target privacy guarantee to be  
 187  $\varepsilon = 1$ . To do this, we can simply compute the proportion non-privately and add Laplace noise with  
 188 scale  $1/n$ . But, if we know that the dataset is a secret and simple random sample from a population  
 189 of  $100n$  individuals, then adding Laplace noise with scale  $1/n$  as before will actually yield a stronger  
 190 privacy guarantee of  $\varepsilon' = 0.01$  for the underlying population. To get  $\varepsilon' = 1$ , we will need to add  
 191 noise with scale only  $1/(100n)$ . In other words, the secrecy of the sample means that the computation  
 192 has more privacy inherently, and therefore, we can add less noise in order to achieve the desired  
 193 privacy guarantee.

194 Existing DP data analysis tools such as DP Creator [18, 17] employ privacy amplification results  
 195 to provide better statistical utility. For example, the DP Creator interface prompts the user to input the  
 196 population size if the data is a secret and random sample from a larger population of known size and  
 197 take advantage of the resulting boost in accuracy without changing the privacy guarantee.

198 As we discussed before, privacy amplification results are also used to analyse algorithms that  
 199 incorporate subsampling as one of their components. Privacy amplification results permit a tighter  
 200 analysis of the privacy that these algorithm can guarantee. In particular, these algorithms are quite  
 201 common in learning tasks, e.g. [32, 1, 21, 28, 19].

## 202 3 Randomised data-independent sampling rates

203 While we are ultimately interested in data-dependent sampling designs, we begin with an extension of  
 204 Theorem 2 to non-constant but data-independent sampling rates. Prior results on privacy amplification  
 205 by subsampling [23, 33, 11, 31, 6] all focus on constant sampling rates where the sampling rate (the  
 206 fraction of the data set sampled) is fixed in advance. However, we will eventually see that randomising  
 207 the sample rate is essential to privacy amplification when the target rate is data dependent. To work  
 208 toward this eventual discussion, we first study the data-independent case to gain intuition for what

## 23:6 Controlling Privacy Loss in Sampling Schemes

209 properties of the distribution on sampling rates characterize how much privacy amplification is  
210 possible.

211 Suppose that there is a random variable  $t$  on  $[n]$  and the sampling scheme is as follows: given a  
212 dataset  $P$ , a sample  $m$  is drawn from  $t$ , and then  $m$  subjects are drawn without replacement from  
213  $P$  to form the sample  $S$ . In this section we consider unbounded differential privacy<sup>2</sup> for  $\mathcal{M}$  and  
214 bounded differential privacy for  $\mathcal{M}_C$ , where the total number of cases,  $n$ , is known and fixed. A  
215 simple generalisation of Theorem 2 immediately implies that the privacy loss of this randomised  
216 scheme is no worse than if  $t$  was concentrated on the maximum value in its support. However, prior  
217 work does not give insight into what happens when  $t$  is concentrated below its maximum or is evenly  
218 spread. What property of the distribution characterises its potential for privacy amplification? The  
219 following theorem characterizes the privacy amplification of sampling without replacement with  
220 data-independent randomised sampling rates.

221 **► Theorem 3.** *Let  $P$  be a dataset of size  $n$ , let  $t$  be a distribution over  $\{0, 1, \dots, n\}$ , and let  
222  $\mathcal{C} : \mathcal{X} \rightarrow \mathcal{U}^*$  be the randomised, dataset-independent sampling scheme that randomly draws  
223  $m \sim t$  and samples  $m$  records from  $P$  without replacement. Define the distribution  $\tilde{t}$  on  $[n]$  where  
224  $\tilde{t}(m) \propto e^{\varepsilon m} \cdot t(m)$  for all  $m \in [n]$ .*

225 **Upper bound:** *Let  $\mathcal{M} : \mathcal{U}^* \rightarrow \mathcal{O}$  be an  $\varepsilon$ -unbounded DP algorithm. Then,  $\mathcal{M}_C$  is  $\varepsilon'$ -bounded DP,  
226 where*

$$227 \quad \varepsilon' = \log \left( 1 + \frac{1}{n} \cdot \mathbb{E}_{m \sim \tilde{t}}[m] \cdot (e^\varepsilon - 1) \right).$$

228 **Lower bound:** *There exists neighboring datasets  $P$  and  $P'$  of size  $n$ , and an  $\varepsilon$ -unbounded DP  
229 mechanism  $\mathcal{M}$  such that if  $\mathcal{M}_C(P)$  and  $\mathcal{M}_C(P')$  are  $\varepsilon'$ -indistinguishable then*

$$230 \quad \varepsilon' \geq -\log \left( 1 - \frac{1}{n} \cdot \mathbb{E}_{m \sim \tilde{t}}[m] \cdot (1 - e^{-\varepsilon}) \right)$$

231 First notice that Theorem 3 comports with the generalization of Theorem 2; as expected, if the  
232 support of  $t$  is contained within  $[0, m']$  then  $\mathbb{E}_{m \sim \tilde{t}}[m] \leq m'$ , so the randomised scheme is at least as  
233 private as if  $t$  was concentrated on  $m'$ . It also determines that the property of  $t$  that determines the  
234 privacy amplification is  $\mathbb{E}_{m \sim \tilde{t}}[m]$ , the expectation of an exponential re-weighting of the distribution  
235 that gives more weight to larger sample sizes. When  $\varepsilon$  is small, the simple approximations  $e^x - 1 \approx x$ ,  
236  $1 - e^{-x} \approx x$ , and  $\log(1 + x) \approx x$  mean that both the upper and lower bounds amount to

$$237 \quad \varepsilon' \approx \frac{\mathbb{E}_{m \sim \tilde{t}}[m]}{n} \cdot \varepsilon.$$

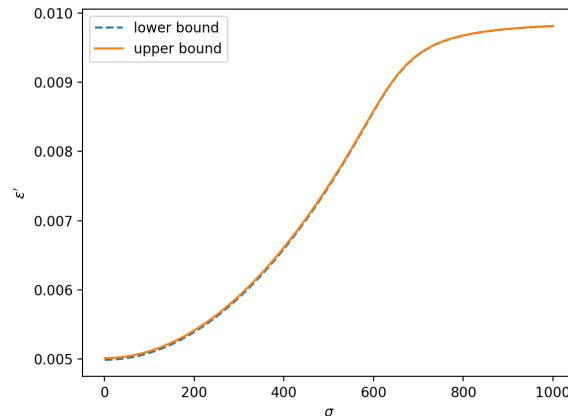
238 Due to the exponential re-weighting,

$$239 \quad \mathbb{E}_{m \sim \tilde{t}}[m] = \frac{\sum_{m=0}^n e^{\varepsilon m} \Pr(t = m)m}{\sum_{m=0}^n e^{\varepsilon m} \Pr(t = m)}$$

240 rapidly approaches  $n$  as the weight of  $t$  on values close to  $n$  increases. Intuitively, this means that  
241 even a small probability of sampling the entire dataset can be enough to ensure that there is no privacy  
242 amplification, even if the mode of  $t$  is much smaller than  $n$ . Conversely, if  $t$  is a light tailed distribution  
243 (say, subgaussian) concentrated on a value much smaller than  $n$ , then privacy amplification is possible.

244 For example, suppose that  $t$  is a truncated Gaussian on  $[0, n]$  with mean  $n/2$  and standard deviation  
245  $\sigma$ . If  $t$  is highly concentrated then we expect the privacy guarantee of  $\mathcal{M}_C$  to be  $\approx \varepsilon/2$ . As  $\sigma$  grows

<sup>2</sup> Note that we must use the unbounded differential privacy definition for  $\mathcal{M}$  in this setting; otherwise, the sample size  $m$  would be fixed.



■ **Figure 2** Numerical computation of the upper and lower bounds from Theorem 3 when  $t$  is truncated Gaussian supported on  $[0, n]$  with mean  $n/2$ , where  $n = 10^4$  and standard deviation  $\sigma$  varies from 1 to  $10^3$ . The privacy parameter of the mechanism  $\mathcal{M}$  is 0.01.

246 we expect the privacy guarantee to tend towards  $\varepsilon$  as more weight is placed near  $n$ . In Figure 2, we  
 247 illustrate the bounds of Theorem 3 numerically with this Gaussian example. We can see that when  
 248  $n = 10,000$  and  $\sigma \approx 800$ , the privacy guarantee of  $\mathcal{M}_{\mathcal{C}}$  is already close to  $\varepsilon = 0.01$ , the privacy  
 249 guarantee of  $\mathcal{M}$ .

## 250 4 Data-dependent sampling rates

251 We now turn our attention to sampling schemes where sampling rates may depend on the data. The  
 252 results in this section are motivated by *stratified sampling*, where the population is stratified into  
 253  $k$  disjoint sub-populations called strata, and an allocation function is used to determine how many  
 254 samples to draw from each stratum. We will discuss stratified sampling with  $k > 1$  in Section 6, but  
 255 for simplicity and clarity, we first focus on the “single stratum” case. In this section, we develop tools  
 256 and statements that we expect to be more broadly useful in understanding complex sampling designs.

257 Specifically, we consider the sampling design where one selects a number of cases according to a  
 258 data-dependent function, and then samples that many cases via simple random sampling. That is, let  
 259  $\tilde{f} : \mathcal{U}^* \rightarrow \mathbb{N}$  be a possibly randomised function and let  $\mathcal{C}_f$  be the sampling function that on input  $P$   
 260 samples  $f(P)$  data points uniformly without replacement from  $P$ . If  $\mathcal{M}$  is an  $\varepsilon$ -DP algorithm, then  
 261 how private is  $\mathcal{M}_{\mathcal{C}_f}$ ?

### 262 4.1 Sensitivity and privacy degradation

263 We first observe that if the function  $f$  used to determine sample size is highly sensitive, then privacy  
 264 *degradation* may occur. That is, if the number of cases sampled may change dramatically on  
 265 neighboring populations, then the output of a DP mechanism can immediately be used to distinguish  
 266 between those populations. For example, suppose  $P$  and  $P'$  are neighboring populations, and  $f$  is  
 267 a function where  $f(P) = m$  and  $f(P') = m + \Delta$ . (That is, the local sensitivity of  $f$  at  $P$  is at  
 268 least  $\Delta$ .) Consider the  $\varepsilon$ -DP algorithm  $\mathcal{M}^{\text{count}}$  that, on input a sample  $S$ , outputs the noisy count  
 269  $|S| + \text{Lap}(1/\varepsilon)$  of the number of cases in the sample. Then  $\mathcal{M}_{\mathcal{C}_f}^{\text{count}}(P)$  is distributed as  $m + \text{Lap}(1/\varepsilon)$   
 270 whereas  $\mathcal{M}_{\mathcal{C}_f}^{\text{count}}(P')$  is distributed as  $m + \Delta + \text{Lap}(1/\varepsilon)$ . When  $\Delta \gg 1$ , these distributions are  
 271 far apart; the privacy loss between these two populations is  $\Delta \cdot \varepsilon \gg \varepsilon$ .

272 Thus, a *necessary* condition for achieving privacy amplification (rather than degradation) is that  
 273 the function  $f$  has low sensitivity. In the following sections, we explore other conditions on low  
 274 sensitivity functions that are necessary and sufficient for amplification.

## 275 4.2 Data dependent sampling and hypothesis testing

276 We established in the previous section that using a deterministic function to determine sample size  
 277 results in privacy degradation. This raises the question: how much randomness is necessary to ensure  
 278 privacy control? That is, what can we say about a randomised function  $\tilde{f} : \mathcal{U}^* \rightarrow \mathbb{N}$  with the property  
 279 that  $\mathcal{M}_{\mathcal{C}_{\tilde{f}}}$  is  $\varepsilon'$ -DP for every  $\varepsilon$ -DP mechanism  $\mathcal{M}$ ? In this section we establish a connection between  
 280 the amplification properties of a function  $\tilde{f}$  and hypothesis testing.

281 A simple hypothesis testing problem is specified by two distributions  $X$  and  $Y$ . A hypothesis test  
 282  $H$  for this problem attempts to determine whether the samples given as input are drawn i.i.d from  $X$   
 283 or from  $Y$ . If a hypothesis test is only given a single sample then we define the advantage of  $H$  to be

$$284 \quad \text{adv}(H; X, Y) = \Pr_{m \sim X} [H(m) = X] - \Pr_{m \sim Y} [H(m) = X].$$

285 That is, the advantage is a measure of how likely the hypothesis test  $H$  is to correctly guess which  
 286 distribution the sample was drawn from. The closer the advantage is to 1, the better the test is at  
 287 distinguishing  $X$  from  $Y$ .

288 One common explanation of differential privacy is that an algorithm is differentially private if it is  
 289 impossible to confidently guess from the output which of two neighbouring datasets was the input  
 290 dataset. This interpretation can be formalised, following [35], by noting that if  $\mathcal{M}$  is  $\varepsilon$ -DP and  $P$  and  
 291  $P'$  are neighbouring populations then for every hypothesis test  $H$ ,

$$292 \quad \text{adv}(H; \mathcal{M}(P), \mathcal{M}(P')) \leq e^\varepsilon - 1 \approx \varepsilon.$$

293 We can establish a similar bound and interpretation of what it means for  $\tilde{f}$  to amplify or pre-  
 294 serve privacy. Suppose that  $\tilde{f}$  is such that  $\mathcal{M}_{\mathcal{C}_{\tilde{f}}}$  is  $\varepsilon'$ -DP for every  $\varepsilon$ -DP mechanism  $\mathcal{M}$ . Then  
 295 in particular, for every  $\varepsilon$ -DP hypothesis test  $H$ , we have that  $H(\mathcal{M}_{\mathcal{C}_{\tilde{f}}}(P))$  and  $H(\mathcal{M}_{\mathcal{C}_{\tilde{f}}}(P'))$  are  
 296  $\varepsilon'$ -indistinguishable. Now, if we consider only hypothesis tests  $H : \mathbb{N} \rightarrow \{\tilde{f}(P), \tilde{f}(P')\}$  that simply  
 297 look at the size of the sample  $\mathcal{C}_{\tilde{f}}(\cdot)$ , then we can formalise this statement in the following way.

298 ► **Proposition 4.** *Suppose  $\tilde{f} : \mathcal{U}^* \rightarrow \mathbb{N}$  is such that for all  $\varepsilon$ -DP mechanisms  $\mathcal{M}$ , we have that*  
 299  *$\mathcal{M}_{\mathcal{C}_{\tilde{f}}}$  is  $\varepsilon'$ -DP. Then for all neighboring datasets  $P, P'$ , we have*

$$300 \quad \max \text{adv}(H; \tilde{f}(P), \tilde{f}(P')) \leq e^{\varepsilon'} - 1,$$

301 *where the optimisation is over all hypothesis tests  $H$  such that for all  $x \in \mathbb{N}$ , and  $b \in \{0, 1\}$ ,*  
 302  *$e^{-\varepsilon} \Pr(H(x) = b) \leq \Pr(H(x+1) = b) \leq e^\varepsilon \Pr(H(x) = b)$ .*

303 This result helps us build intuition for what type of survey designs could possibly amplify privacy.  
 304 If  $\tilde{f}$  results in privacy amplification then for any pair of neighbouring populations  $P$  and  $P'$ , the  
 305 distributions  $\tilde{f}(P)$  and  $\tilde{f}(P')$  must be close enough that they can not be distinguished between by  
 306 any hypothesis test  $H$  such that  $\log H$  is  $\varepsilon$ -Lipschitz. From this perspective the result in Section 4.1  
 307 follows from the fact that if  $\tilde{f}$  is deterministic with high sensitivity then we can define an appropriate  
 308 hypothesis test with large advantage based on  $\mathcal{M}^{\text{count}}$ . This is a useful perspective to keep in mind  
 309 throughout the remainder of the paper.

310 One consequence of this perspective is a lower bound on how well we can emulate a desired  
 311 deterministic function  $f$  while controlling or amplifying privacy. Suppose that absent privacy concerns,  
 312 an analyst has determined that they want to use a function  $f$  to determine the sample size. However,



313 to avoid privacy degradation they replace  $f$  with a randomised function  $\tilde{f}$ . How close can  $\tilde{f}$  get to  $f$   
 314 while maintaining or amplifying the original privacy level? We can obtain a lower bound on expected  
 315 closeness of  $f(P)$  and  $\tilde{f}(P)$  by relating it to the well studied problem of estimation lower bounds in  
 316 differential privacy.

317 ► **Proposition 5.** *Let  $f : \mathcal{U}^* \rightarrow \mathbb{R}$  and  $\varepsilon, \varepsilon' > 0$ . Suppose  $\tilde{f} : \mathcal{U}^* \rightarrow \mathbb{N}$  is a randomised  
 318 function such that for all  $\varepsilon$ -unbounded DP mechanisms  $\mathcal{M}$ , it holds that  $\mathcal{M}_{\tilde{f}}$  is  $\varepsilon'$ -bounded DP.  
 319 If  $\alpha \geq 0$  is such that for every  $\varepsilon'$ -unbounded DP mechanism  $\mathcal{A}$ , there exists a dataset  $P$  such that  
 320  $\mathbb{E}[|\mathcal{A}(P) - f(P)|^2] \geq \alpha$ , then there exists a dataset  $P$  such that*

$$321 \quad \mathbb{E}[|\tilde{f}(P) - f(P)|^2] \geq \alpha - \left(\frac{1}{\varepsilon}\right)^2.$$

322 The problem of lower bounding differentially private function estimation is well-studied [30, 4]  
 323 in the privacy literature. The lower bounds essentially arise from the fact that  $\mathcal{A}(P)$  and  $\mathcal{A}(P')$  must  
 324 be similar distributions for neighbouring databases, even if  $f(P)$  and  $f(P')$  are far apart. Since we  
 325 know from Proposition 4 that  $\tilde{f}(P)$  and  $\tilde{f}(P')$  must also be close, we obtain the related lower bound.  
 326 The slackness of  $(1/\varepsilon)^2$  is a result of the fact that while  $\mathcal{A}(P)$  and  $\mathcal{A}(P')$  must be indistinguishable  
 327 with respect to *any* hypothesis test,  $\tilde{f}(P)$  and  $\tilde{f}(P')$  need only be indistinguishable with respect to  
 328 any  $\varepsilon$ -DP hypothesis test.

### 329 4.3 Privacy amplification from randomised rounding

330 Many functions used to determine data-dependent sampling rates have high sensitivity, but at least  
 331 one common sampling method has low sensitivity: proportional sampling. In proportional sampling,  
 332 a constant, data-independent fraction of the population is sampled independently from each stratum.  
 333 This method is similar to simple random sampling, but a small amount of data dependence is  
 334 introduced by the fact that the total number of samples across all strata must be an integer. In this  
 335 section, we will show that while naïve implementations of proportional sampling can result in privacy  
 336 degradation, a minor change in the sampling size function results in privacy amplification comparable  
 337 to that afforded by simple random sampling.

338 Let  $r \in [0, 1]$  and  $f(P) = r|P|$  for some constant  $r \in (0, 1)$ . Since the output space of  $f$  is not  
 339  $\mathbb{N}$ , in practice, this is typically replaced with the deterministic function  $\tilde{f}_{\text{det},r}(P) = \text{round}(r|P|)$ ,  
 340 where  $\text{round}(\cdot)$  rounds its input to the nearest integer. Unfortunately, deterministic rounding can be  
 341 problematic for privacy, as we can see through a simple example. Suppose  $P$  and  $P'$  are neighbouring  
 342 populations such that  $|P| = 14$ ,  $|P'| = 15$ , and  $r = 1/10$ . Then, deterministic rounding always  
 343 results in one case being sampled from  $P$  and two cases being sampled from  $P'$ . As discussed in  
 344 Section 4.1, such a data-dependent deterministic function can never result in privacy amplification.

345 We propose a simple and practical change to the rounding process that *does* guarantee roughly  
 346 the expected level of privacy amplification. We replace the ideal function  $f$  with a randomised  
 347 rounding function  $\tilde{f}_{\text{rand},r}$ . That is, let  $p = r|P| - \lfloor r|P| \rfloor$  so  $\tilde{f}_{\text{rand},r}(P) = \lceil r|P| \rceil$  with probability  
 348  $p$ , and  $\tilde{f}_{\text{rand},r}(P) = \lfloor r|P| \rfloor$  with probability  $1 - p$ . The following proposition shows that, up to a  
 349 constant factor, randomised rounding recovers the expected factor of  $r$  in privacy amplification.

350 ► **Theorem 6 (Privacy Amplification from Randomised Rounding).** *Let  $r \in (0, 1)$ . Then for  
 351 every  $\varepsilon$ -unbounded DP mechanism  $\mathcal{M}$ , the mechanism  $\mathcal{M}_{\tilde{f}_{\text{rand},r}}$  is  $\varepsilon'$ -unbounded DP when restricted  
 352 to datasets of size at least  $1/r$ , where  $\varepsilon' = \log(1 + 2r(e^{2\varepsilon} - 1)) + \log(1 + r(e^{2\varepsilon} - 1)) \approx 6r\varepsilon$ .*

353 The approximation at the end of the proposition follows from applying (1) and (2), which give that  
 354  $\log(1 + 2r(\exp(2\varepsilon) - 1)) \approx 2r \cdot 2\varepsilon$  and  $\log(1 + r(\exp(2\varepsilon) - 1)) \approx r \cdot 2\varepsilon$ . The constant 6 can perhaps  
 355 be optimized through a more careful analysis. Randomised rounding is a practical modification since

## 23:10 Controlling Privacy Loss in Sampling Schemes

356 it does not change the size of the sample very much; if traditional proportional allocation would  
357 typically assign  $m$  samples, then the modified algorithm allocates at most  $m + 1$ .

### 358 **5 Cluster sampling**

359 In cluster sampling, the population is partitioned into disjoint subsets, called clusters. A subset of the  
360 clusters is sampled and data subjects are selected from within the chosen clusters. If the sampling  
361 scheme uses a single stage design, all data subjects contained in the selected clusters will be included  
362 in the sample. Otherwise, a random sample of data subjects might be selected from each of the  
363 selected clusters (multi-stage design). Cluster sampling produces accurate results when the clusters  
364 are mutually homogeneous; that is, when the distributions within each cluster are similar to the  
365 distribution over the entire population.

366 In the survey context, cluster sampling is often performed due to time or budgetary constraints  
367 which make sampling many units from a few clusters cheaper and/or faster than sampling a few  
368 units from each cluster. A typical example in the survey context is when clusters are chosen to be  
369 geographic regions. Sampling a few geographic clusters and interviewing everybody in those clusters  
370 saves traveling costs compared to interviewing the same number of people based on a simple random  
371 sample from the population. In algorithm design, cluster sampling is often performed to improve the  
372 performance and accuracy of classifiers. In this setting, sampling often involves a two-step approach  
373 where the data is first clustered, using some clustering classifier, and then a subset of the clusters  
374 is selected. Forms of cluster samplings have been applied in several learning areas, for example in  
375 federated learning [16] and active learning [27].

#### 376 **5.1 Privacy implications of single-stage cluster sampling with** 377 **simple random sampling**

378 We focus here on a simple cluster sampling design that is commonly used in survey sampling and  
379 which naïvely appears to be a good candidate for privacy amplification: simple random sampling  
380 without replacement of clusters. That is, suppose the dataset  $P$  is divided into  $k$  clusters,

$$381 \quad P = C_1 \sqcup \dots \sqcup C_k$$

382 and the sampling mechanism  $\mathcal{C}_\ell : \mathcal{U}^* \rightarrow \mathcal{U}^*$  chooses a random subset  $I \subset [k]$  of size  $\ell < k$ , then  
383 maps  $P$  to  $\sqcup_{i \in I} C_i$ .

384 Since simple random sampling at the individual level provides good privacy amplification, one  
385 might expect the same to happen when the clusters are sampled in a similar way. In fact, this is true  
386 when the size of each cluster is small. However, if the clusters are large this sampling design achieves  
387 less amplification than might be expected. This is characterized by the following theorem showing a  
388 lower bound in this setting.

389 **► Theorem 7 (Lower Bound on Privacy Amplification for Cluster Sampling).** *For any sequence*  
390  *$n_i > 0$  and privacy parameter  $\varepsilon > 0$ , there exist neighboring populations  $P = C_1 \sqcup \dots \sqcup C_i \sqcup \dots \sqcup C_k$*   
391 *and  $P' = C_1 \sqcup \dots \sqcup C'_i \sqcup \dots \sqcup C_k$  (with  $|C_i| = n_i$  and  $C'_i = C_i \cup \{x\}$  for some  $x \in \mathcal{U}$ ) and an*  
392  *$\varepsilon$ -unbounded DP mechanism  $\mathcal{M}$  such that if  $\mathcal{M}_{\mathcal{C}_\ell}(P)$  and  $\mathcal{M}_{\mathcal{C}_\ell}(P')$  are  $\varepsilon'$ -indistinguishable then*

$$393 \quad \varepsilon' \geq \ln \left( 1 + \frac{\frac{\ell}{k}}{\left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_i + n_{\min})\varepsilon}\right)} (e^\varepsilon - 1) \right),$$

394 where  $n_i = |C_i|$  and  $n_{\min} = \min_{j \in \{1, \dots, i-1\} \cup \{i+1, \dots, k\}} n_j$ .

395 We can compare the expression in the theorem above with the one we have for simple random  
396 sampling without replacement (cf. Theorem 14 from [6]):

$$397 \quad \varepsilon' = \ln \left( 1 + \frac{m}{n} (e^\varepsilon - 1) \right),$$

398 where  $m$  samples are drawn from a population of size  $n$ . We see that the two expressions coincide if  
399  $n_i + n_{\min} = 0$ , which is an unrealistic corner case. Let us instead consider the case in which all the  
400 clusters are small. In this case, the quantity  $n_i + n_{\min}$  will also be small, and if  $\varepsilon < 1$ , we can still  
401 expect some privacy amplification. However, as the clusters grow in size, the quantity  $n_i + n_{\min}$  will  
402 also increase, and the lower bound converges very quickly to  $\varepsilon$ , giving essentially no amplification.

403 Next, we present a corresponding upper bound.

404 ► **Theorem 8** (Upper Bound on Privacy Amplification for Cluster Sampling). *For any sequence*  
405  *$n_i > 0$ , privacy parameter  $\varepsilon > 0$ ,  $\varepsilon$ -unbounded DP mechanism  $\mathcal{M} : \mathcal{U}^* \rightarrow \mathcal{O}$ , and pair of*  
406 *neighboring populations  $P$  and  $P'$  such that  $P = C_1 \sqcup \dots \sqcup C_i \sqcup \dots \sqcup C_k$  and  $P' = C_1 \sqcup \dots \sqcup C'_i \sqcup \dots \sqcup C_k$*   
407 *(with  $|C_i| = n_i$  and  $C'_i = C_i \cup \{x\}$  for some  $x \in \mathcal{U}$ ), the mechanisms  $\mathcal{M}_{C_\ell}(P)$  and  $\mathcal{M}_{C_\ell}(P')$  are*  
408  *$\varepsilon'$ -indistinguishable where*

$$409 \quad \varepsilon' \leq \ln \left( 1 + \frac{\frac{\ell}{k}}{\left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_i + n_{\max})\varepsilon}\right)} (e^\varepsilon - 1) \right),$$

410 and  $n_{\max} = \max_{j \in \{1, \dots, i-1\} \cup \{i+1, \dots, k\}} n_j$ ,

411 Once again it is worth comparing the expression in the theorem above with the one we have for simple  
412 random sampling without replacement:

$$413 \quad \varepsilon' = \ln \left( 1 + \frac{m}{n} (e^\varepsilon - 1) \right).$$

414 Similar to the lower bound, the upper bound will quickly approach  $\varepsilon$  if the quantity  $n_i + n_{\max}$  is  
415 large. If each cluster contains a single data point, the two bounds are close. This is not surprising  
416 since in this case the type of cluster sampling we considered is just simple random sampling without  
417 replacement. Note that while  $\ell/k$  is the fraction of clusters included in the final sample and  $m/n$  is  
418 the fraction of data points, these are approximately the same when the clusters are small. If all the  
419 clusters are the same size, then  $n_{\max} = n_{\min}$  and the upper and lower bounds we gave above match.  
420 The proofs of these results are contained in the Appendix.

## 421 5.2 Discussion and hypothesis testing

422 Privacy amplification by subsampling is often referred to as *secrecy of the sample* due to the intuition  
423 that the additional privacy arises from the fact that there is uncertainty regarding which user's data is  
424 in the sample. The key intuition then for Theorem 7 is that the larger the clusters are, the easier it is  
425 for a differentially private algorithm  $\mathcal{M}$  to reverse engineer which clusters were sampled, breaking  
426 secrecy of the sample. Intuitively, if the clusters are different enough that a private algorithm can  
427 guess which clusters were chosen as part of the sample, then any amplification due to *secrecy of the*  
428 *sample* is negligible. We can formalize this intuition using once again using the lens of hypothesis  
429 testing. Note the framing in this section differs slightly from the framing in Section 4, although the  
430 underlying idea in both settings is that if a particular hypothesis test is effective, then there is a lower  
431 bound on the privacy parameter. In addition, note that privacy is also conserved in this setting, as  
432  $\mathcal{M}_{C_\ell}$  is at least as private as  $\mathcal{M}$ . The question is: when is  $\mathcal{M}_{C_\ell}$  more private than  $\mathcal{M}$ ?

## 23:12 Controlling Privacy Loss in Sampling Schemes

433 ► **Theorem 9.** Let  $\varepsilon > 0$ ,  $\ell \in [0, k]$ ,  $\mathcal{M} : \mathcal{U}^* \rightarrow \mathcal{O}$  be  $\varepsilon$ -DP and the sampling mechanism  $\mathcal{C}_\ell$  be as  
 434 defined in Section 5.1. Suppose there exists a hypothesis test  $\mathcal{H} : \mathcal{O} \rightarrow \{0, 1\}$  such that

$$435 \Pr(\mathcal{H}(\mathcal{M}_{\mathcal{C}_\ell}(P)) = 0 \mid C_i \in \mathcal{C}_\ell(P)) \geq e^{\varepsilon'} \Pr(\mathcal{H}(\mathcal{M}_{\mathcal{C}_\ell}(P)) = 0 \mid C_i \notin \mathcal{C}_\ell(P)).$$

436 Then there exists an event  $E$  in the output space of  $\mathcal{M}$  such that for any neighboring population  $P'$   
 437 that differs from  $P$  in  $C_i$ , if

$$438 \varepsilon'' = \log \frac{\Pr(\mathcal{M}_{\mathcal{C}_\ell}(P) \in E \mid C_i \in \mathcal{C}_\ell(P))}{\Pr(\mathcal{M}_{\mathcal{C}_\ell}(P') \in E \mid C_i \in \mathcal{C}_\ell(P'))} \in [0, \varepsilon],$$

439 and  $\mathcal{M}_{\mathcal{C}_\ell}(P)$  and  $\mathcal{M}_{\mathcal{C}_\ell}(P')$  are  $\tilde{\varepsilon}$ -indistinguishable, then

$$440 \tilde{\varepsilon} \geq \log \left( 1 + (e^{\varepsilon''} - 1) \frac{\ell/k}{\ell/k + e^{-\varepsilon'}(1 - \ell/k)} \right).$$

441 The key take-away of this theorem is that for any  $\varepsilon$ -DP mechanism  $\mathcal{M}$ , if there exists a hypothesis  
 442 test that, when given the output of  $\mathcal{M}_{\mathcal{C}_\ell}(P)$ , can confidently decide whether cluster  $C_i$  was chosen as  
 443 part of the final sample, then the privacy guarantee of  $\mathcal{M}_{\mathcal{C}_\ell}$  is no better than the privacy guarantee  
 444 would be if we knew for certain that  $C_i$  was chosen as part of the sample. That is, in this setting, we  
 445 gain no additional privacy as a result of secrecy of the sample. The parameter  $\varepsilon'$  controls how well  
 446 the hypothesis test can determine whether  $C_i \in \mathcal{C}_\ell$ . As  $\varepsilon'$  increases,  $\tilde{\varepsilon}$  approaches  $\varepsilon''$ , the privacy  
 447 parameter if  $C_i$  is known to be part of the sample, so privacy amplification is negligible.

448 This view is consistent with Theorem 7. Consider a population where only data points in cluster  $i$   
 449 have a particular property and let  $\mathcal{M}$  is an  $\varepsilon$ -DP mechanism that attempts to count how many data  
 450 points with the property are in the final sample. If cluster  $i$  is large, then it is easy to determine from  
 451 the output of the mechanism whether  $C_i$  is in the final sample. This example required cluster  $i$  to be  
 452 distinguishable from the remaining clusters using a private algorithm. While examples as extreme as  
 453 the one above may be uncommon in practice, clusters being different enough for a private algorithm  
 454 to distinguish between them is not an unrealistic assumption.

455 In Section 5.1, we analysed a single stage design. All subjects contained in the selected clusters  
 456 were included in the sample. In practice, multi-stage designs are common, where a random sample  
 457 of subjects are selected from within each chosen cluster. If the sampling within each cluster is  
 458 sufficiently simple then the privacy amplification from this stage can be immediately incorporated  
 459 into the upper bound in Theorem 8. For example, if each subject within the chosen clusters is sampled  
 460 with probability  $r$  and  $\mathcal{M}$  is  $\varepsilon$ -DP, i.e., we perform Poisson sampling with probability  $r$ , then we  
 461 immediately obtain an upper bound that is approximately  $r\varepsilon$ . However, if the sampling within clusters  
 462 is more complex, then further analysis is required. One can also imagine more complicated schemes  
 463 for selecting the chosen clusters. If these designs depend on properties of the data, then they are likely  
 464 to result in privacy degradation. We leave this study for future work.

## 465 **6 Stratified sampling**

466 Finally, we turn our attention to another common sampling design: stratified sampling. In stratified  
 467 sampling, the data is partitioned into disjoint subsets, called strata. A subset of data points is then  
 468 sampled from each stratum to ensure the final sample contains data points from every stratum.  
 469 Stratified sampling is common in survey sampling where it is used to improve accuracy and to ensure  
 470 sufficient representation of sub-populations of interest. A classic use case of stratified sampling is  
 471 business surveys, where businesses are typically stratified by industry and number of employees, or  
 472 by similar measures of establishment size. Stratification by establishment size results in substantial  
 473 gains in accuracy compared to simple random sampling, while stratification by industry ensures

474 that reliable estimates can be obtained at the industry level. Stratified sampling has several other  
 475 applications; for example it is used in algorithm design to improve performance [2, 24], in private  
 476 query design and optimization to improve accuracy [8], and to improve search and optimizations [25].

477 We focus here on *one-stage stratified sampling* using simple random sampling without replacement  
 478 within each stratum to select samples. We also assume that the stratum boundaries have been fixed in  
 479 advance. Given a target sample size  $m$ , the only design choice in this model is the *allocation function*,  
 480 which determines how many samples to take from each stratum. Different allocation functions are  
 481 used in practice. Which method is selected depends on the goals to be achieved (for example, ensuring  
 482 constant sampling rates across strata or minimizing the variance for a statistic of interest).

483 Before we describe allocation functions in detail, let us establish some notation for stratified  
 484 sampling. Suppose there are  $k$  strata in the population, and that each data point is a pair  $(s, x)$   
 485 where  $s \in [k]$  denotes which stratum the data subject belongs to, and  $x \in \mathcal{U}$  denotes their data.  
 486 Let  $\mathbf{f} = (f_1, \dots, f_k) : ([k] \times \mathcal{U})^* \rightarrow \mathbb{N}^k$  denote the allocation rule, so  $f_i(P)$  samples are drawn  
 487 uniformly at random without replacement from the  $i$ th stratum,  $P_i = \{(s, x) \in P \mid s = i\}$ . The final  
 488 sample  $S$  is the union of the samples from all the strata.

489 An important feature of stratified sampling is that the sampling rates can vary between the strata.  
 490 This means that data subjects in strata with low sampling rates may expect a higher level of privacy  
 491 than data subjects in strata with high sampling rates. This leads us to define a variant of differential  
 492 privacy that allows the privacy guarantee to vary between the strata. This generalisation of differential  
 493 privacy is tailored to stratified datasets and allows us to state more refined privacy guarantees than the  
 494 standard definition is capable of.

495 ► **Definition 10.** Let  $k \in \mathbb{N}$  and suppose there are  $k$  strata. A mechanism  $\mathcal{A}$  satisfies  $(\varepsilon_1, \dots, \varepsilon_k)$ -  
 496 stratified bounded differential privacy if for all datasets  $P$ , data points  $(s, x)$  and  $(s', x')$ ,  $\mathcal{A}(P \cup$   
 497  $\{(s, x)\})$  and  $\mathcal{A}(P \cup \{(s', x')\})$  are  $\max\{\varepsilon_s, \varepsilon_{s'}\}$ -indistinguishable. The mechanism  $\mathcal{A}$  satisfies  
 498  $(\varepsilon_1, \dots, \varepsilon_k)$ -stratified unbounded differential privacy if for all datasets  $P$ , data points  $(s, x)$ ,  $\mathcal{A}(P)$   
 499 and  $\mathcal{A}(P \cup \{(s, x)\})$  are  $\varepsilon_s$ -indistinguishable.

500 This definition is an adaptation of personalized differential privacy [22, 14, 3]. Note that it protects  
 501 not only the *value* of an individual's data point, but also which stratum they belong to.

## 502 6.1 Optimal allocation with privacy constraints

503 In this section, we will discuss how to think about choosing an allocation function when privacy is a  
 504 concern. A common goal when choosing an allocation  $\mathbf{f}$  is to minimise the variance of a particular  
 505 statistic. That is, suppose that  $\mathcal{C}_{\mathbf{f}}$  represents one-stage stratified sampling with allocation function  $\mathbf{f}$ .  
 506 Then, given a population  $P$  and desired sample size  $m$ , the *optimal allocation function*  $\mathbf{f}^*(P)$  with  
 507 respect to a statistic  $\theta$  is defined as

$$508 \quad \mathbf{f}^*(P) = \arg \min_{\mathbf{f}} \text{var}(\theta_{\mathcal{C}_{\mathbf{f}}}(P)), \quad (3)$$

509 where the randomness may come from both the allocation function and the sampling itself, and the  
 510 minimum is over all allocation functions such that  $\|\mathbf{f}(P)\|_1 \leq m$  for all  $P$ .<sup>3</sup>

511 A natural question then is: what is the optimal allocation when one wants to compute the statistic  
 512 of interest differentially privately? This is a simple yet subtle question. Our results in the previous  
 513 sections indicate that the landscapes of optimal allocations in the non-private and private settings

<sup>3</sup> As an aside, we note that the notion of optimal allocations implicitly assumes that the historic or auxiliary data,  $H$ ,  
 used to inform the sampling design and the population data  $P$  are the same, or at least similar enough that  $\mathbf{f}^*(H)$  is  
 a good proxy for  $\mathbf{f}^*(P)$ . This provides further justification for the assumption that  $H = P$  in our statements.

## 23:14 Controlling Privacy Loss in Sampling Schemes

514 may be very different. This is a result of the fact that allocation functions that do not amplify well  
 515 typically need to add more noise to achieve privacy (see discussion in Section 2.4). The additional  
 516 noise needed to achieve privacy may overwhelm any gains in accuracy for the non-private statistic.  
 517 Additionally, it is not immediately obvious how to define the optimal allocation in the private setting.

518 In this section, we formulate the notion of an optimal allocation under privacy constraints. Our  
 519 goal is to initiate the study of alternative allocation functions that may prove useful when privacy is a  
 520 concern. A full investigation of this question is outside the scope of this paper, but we provide some  
 521 intuition for why this may be an interesting and important question for future work.

522 Given a statistic  $\theta$ , we wish to define the optimal allocation for estimating  $\theta$  privately. Let  
 523  $\tilde{\theta}^\lambda$  be an  $\lambda$ -DP algorithm for estimating  $\theta$ , so  $\tilde{\theta}^\lambda(P)$  is an approximation of  $\theta(P)$ . The smaller  $\lambda$   
 524 is, the noisier  $\tilde{\theta}^\lambda$  is. The scale of  $\lambda$  needed to ensure that  $\tilde{\theta}_{\mathcal{C}_f}^\lambda$  is  $\varepsilon$ -DP depends on the allocation  
 525 function  $f$ . Allocation functions that are very sensitive to changes in the input dataset will require  
 526 more noise (smaller  $\lambda$ ) to mask changes in the allocation. For any allocation  $f$ , we will define the  
 527 optimal parameter  $\lambda$  as that which minimises the maximum variance of  $\tilde{\theta}_{\mathcal{C}_f}^\lambda$  over all datasets  $P$ , while  
 528 maintaining privacy:

$$529 \quad \lambda_f = \arg \min_{\lambda > 0} \sup_P \frac{\text{var}(\tilde{\theta}_{\mathcal{C}_f}^\lambda(P))}{\text{var}(\theta_{\mathcal{C}_f}(P))} \quad (4)$$

530 s.t.  $\tilde{\theta}_{\mathcal{C}_f}^\lambda$  is  $(\varepsilon_1, \dots, \varepsilon_k)$ -stratified DP.  
 531

532 Now, by definition,  $\tilde{\theta}_{\mathcal{C}_f}^\lambda$  is  $(\varepsilon_1, \dots, \varepsilon_k)$ -stratified DP for any allocation function  $f$ . We minimise  
 533 the multiplicative increase in variance so that the supremum is not dominated by populations  $P$  for  
 534 which  $\text{var}(\theta_{\mathcal{C}_f}(P))$  is large. Given privacy parameters  $\varepsilon_1, \dots, \varepsilon_k \geq 0$ , we now define the optimal  
 535 allocation as the allocation function that minimises the maximum variance over all populations  $P$ :

$$536 \quad f_\varepsilon^* = \arg \min_f \sup_P \text{var}(\tilde{\theta}_{\mathcal{C}_f}^{\lambda_f}(P)). \quad (5)$$

537

538 where the minimum again is over all allocations  $f$  such that  $\|f(P)\|_1 \leq m$  for all  $P$ , and the  
 539 supremum is over all populations of interest. This optimisation function has a different form to Eqn 3,  
 540 which performs the optimisation independently for each population  $P$ . This difference is necessary in  
 541 the private setting as we need to ensure that the choice of allocation function  $f_\varepsilon^*$  is not data dependent,  
 542 since this would introduce additional privacy concerns. We can view the optimal allocation as the  
 543 optimal balancing between the variance of the non-private statistic, and the scale of the noise needed  
 544 to maintain privacy.

545 We believe that examining the difference between the optimal allocation in the non-private setting  
 546 (Eqn (3)) and in the private setting (Eqn (5)) is an important question for future work. The main  
 547 challenge is computing the parameter  $\lambda_f$  for every allocation  $f$ . Analysing the privacy implications  
 548 of  $f$  in the style of the previous sections gives us an upper bound on  $\lambda_f$ , although this bound may be  
 549 loose for specific statistics  $\tilde{\theta}^\lambda$ . So, while the previous sections developed our intuition for  $\lambda_f$ , we  
 550 believe new techniques are required to understand this parameter enough to solve Eqn (5).

## 551 6.2 Challenges with optimal allocation

552 Optimal allocations are defined to perform well for a specific statistic of interest. However, in practice,  
 553 a wide variety of analyses will be performed on the final sample. The chosen allocation function may  
 554 be far from optimal for these other analyses. While this problem exists in the non-private setting, it  
 555 becomes more acute in the private setting. An allocation function that is optimal for one statistic may  
 556 result in privacy degradation (and hence low accuracy estimates) for another.

557 We illustrate this challenge using *Neyman allocation*, which is often employed for business  
 558 surveys. Neyman allocation is the optimal allocation method for the weighted mean [26]:

$$559 \quad \theta_\mu(S) = \frac{1}{|P|} \sum_{i=1}^k \frac{|P_i|}{|S_i|} \sum_{x \in S_i} x,$$

560 where  $|P_i|$  is the size of stratum  $i$ , and  $S_i = S \cap P_i$ . The estimator  $\theta_\mu(S)$  is an unbiased estimate of  
 561 the population mean for any stratified sampling design. Given a desired sample size  $m$ , let  $\mathbf{f}_{\text{Neyman}}$   
 562 be the allocation function corresponding to Neyman allocation. Provided each stratum is sufficiently  
 563 large,  $\mathbf{f}_{\text{Neyman}}(P) = (m_1, \dots, m_k)$ , where

$$564 \quad m_i = \frac{|P_i| \sigma(P_i)}{\sum_{j=1}^k |P_j| \sigma(P_j)} \cdot m,$$

565  $\sigma^2(P_i)$  is the empirical variance in stratum  $i$  and sufficiently large means that  $m_i \leq |P_i|$ . Neyman  
 566 allocation is deterministic and can be very sensitive to changes in the data due to its dependence  
 567 on the variance within each stratum. So, while it can provide accurate results for some statistics, it  
 568 provides very noisy results for other statistics of potential interest (e.g. privately computing strata  
 569 sizes).

570 To demonstrate the sensitivity of Neyman allocation, we analysed the sensitivity on a real data  
 571 set. The population is based on the County Business Patterns (CBP) data published by the U.S.  
 572 Census Bureau [15].<sup>4</sup> Each data point is an establishment and the establishments are stratified by  
 573 establishment size into  $k = 12$  strata. With a target final sample size of  $m = 10,000$ , and using the  
 574 weighted mean of the establishment size as the target statistic, the Neyman allocation for this popula-  
 575 tion is [1261, 621, 517, 1969, 833, 1947, 1058, 762, 257, 248, 306, 225]. We can find a neighbouring  
 576 population with Neyman allocation [1259, 620, 516, 1965, 831, 1943, 1056, 761, 257, 247, 306, 244].  
 577 While these allocations are not wildly different, they do differ by 19 samples in the top stratum, which  
 578 might not have a large impact on the weighted mean, but could lead to more substantial changes  
 579 for other statistics. As an illustrative example, we can consider the goal of privately estimating the  
 580 stratum sizes in the sample, for which this allocation would lead to significant privacy degradation.

### 581 6.3 Privacy amplification from proportional sampling

582 *Proportional sampling* is an alternative allocation function that is used to provide equitable repres-  
 583 entation of each sub-population, or stratum. Given a desired sample size  $m \in [n]$ , proportional  
 584 sampling samples an  $r = \frac{m}{n}$  fraction of the data points (rounded to an integer) from each stratum.  
 585 Proportional sampling is not an optimal allocation in the non-private setting but, when implemented  
 586 with randomised rounding, it has good privacy amplification. Now that we consider stratified sampling  
 587 with number of strata  $k \geq 1$ , we can state the following generalisation of Theorem 6.

588 ► **Theorem 11 (Privacy Amplification for Proportional Sampling).** *Let  $r \in [0, 1]$ ,  $\varepsilon > 0$ ,  $\mathcal{M}$*   
 589 *be an  $\varepsilon$ -DP mechanism, and  $P = S_1 \sqcup \dots \sqcup S_k$  and  $P' = S'_1 \sqcup \dots \sqcup S'_k$  be stratified neighboring*  
 590 *datasets that differ on stratum  $i$ . If for all  $j \in [k]$ ,  $r|S_j| \geq 1$  and  $r|S'_j| \geq 1$ , then  $\mathcal{M}_{\mathcal{C}_{\mathbf{f}_{r,\text{prop}}}}$  is  $\varepsilon'$ -DP*  
 591 *where*

$$592 \quad \varepsilon' \leq \log(1 + 2r(e^{2\varepsilon} - 1)) + \log(1 + r(e^{2\varepsilon} - 1)).$$

<sup>4</sup> The data released by the U.S. Census Bureau is a tabulated version of the true micro data from the Business Register (BR), a database of all known single and multi-establishment employer companies. The data set we use is micro data generated to be consistent with the tabulated version. Each data point in this population is the size of an establishment in the US. In order to compute the sensitivity, we need to top code the data, we top code the data at 10,000.

593 Note that given a private statistic  $\tilde{\theta}^\lambda$  as defined as above, this allows us to set  $\lambda_{f_{r,\text{prop}}} \approx \frac{\varepsilon}{6r}$ , which  
 594 is considerably larger than  $\varepsilon$  for small sampling rates. Thus, while proportional sampling may not  
 595 minimise the variance of any single statistic, it may be a good choice since it performs reasonably  
 596 well for *all* statistics.

## 597 7 Conclusion

598 In this paper, we have considered the privacy guarantees of sampling schemes, extending previous  
 599 results to more complex and data-dependent sampling designs that are commonly used in practice. We  
 600 find that considering these sampling schemes requires developing more nuanced analytical tools. In  
 601 this work, we characterize the privacy impacts of randomized and data-dependent sampling schemes.  
 602 Then, we apply our insights to analyze cluster and stratified sampling and to consider the question  
 603 of optimal allocations under privacy. To the best of our knowledge, this work is the first to initiate  
 604 study into these designs. As such, we hope to see future work in three areas. First, future work  
 605 should tighten and optimize the constants in our theorems. Second, our results should be extended  
 606 from pure to approximate (and other variants) of differential privacy. Finally, we hope to see further  
 607 investigation into near-optimal allocations under privacy constraints.

## 608 — References —

- 609 1 Martín Abadi, Andy Chu, Ian J. Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and  
 610 Li Zhang. Deep learning with differential privacy. In Edgar R. Weippl, Stefan Katzenbeisser, Christopher  
 611 Kruegel, Andrew C. Myers, and Shai Halevi, editors, *Proceedings of the 2016 ACM SIGSAC Conference  
 612 on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*, pages 308–318. ACM,  
 613 2016. doi:10.1145/2976749.2978318.
- 614 2 Julaiti Alafate and Yoav S Freund. Faster boosting with smaller memory. In H. Wal-  
 615 lach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, edit-  
 616 ors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associ-  
 617 ates, Inc., 2019. URL: [https://proceedings.neurips.cc/paper/2019/file/  
 618 3fffebb08d23c609875d7177ee769a3e9-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/3fffebb08d23c609875d7177ee769a3e9-Paper.pdf).
- 619 3 Mohammad Alaggan, Sébastien Gambs, and Anne-Marie Kermarrec. Heterogeneous differential privacy.  
 620 *Journal of Privacy and Confidentiality*, 7, 04 2015.
- 621 4 Hilal Asi and John C. Duchi. Near instance-optimality in differential privacy, 2020. arXiv:2005.  
 622 10630.
- 623 5 Borja Balle, Gilles Barthe, and Marco Gaboardi. Privacy amplification by subsampling: Tight analyses via  
 624 couplings and divergences. In *Advances in Neural Information Processing Systems 31: Annual Conference  
 625 on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*,  
 626 pages 6280–6290, 2018.
- 627 6 Borja Balle, Gilles Barthe, and Marco Gaboardi. Privacy profiles and amplification by subsampling.  
 628 *Journal of Privacy and Confidentiality*, 10(1), Jan. 2020.
- 629 7 Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient  
 630 algorithms and tight error bounds. In *Proceedings of the 55th Annual IEEE Symposium on Foundations of  
 631 Computer Science, FOCS '14*, pages 464–473, Washington, DC, USA, 2014. IEEE Computer Society.
- 632 8 Johes Bater, Yongjoo Park, Xi He, Xiao Wang, and Jennie Rogers. SAQE: practical privacy-preserving  
 633 approximate query processing for data federations. *Proc. VLDB Endow.*, 13(11):2691–2705, 2020. URL:  
 634 <http://www.vldb.org/pvldb/vol13/p2691-bater.pdf>.
- 635 9 Amos Beimel, Shiva Prasad Kasiviswanathan, and Kobbi Nissim. Bounds on the sample complexity for  
 636 private learning and private data release. In Daniele Micciancio, editor, *Theory of Cryptography, 7th  
 637 Theory of Cryptography Conference, TCC 2010, Zurich, Switzerland, February 9-11, 2010. Proceedings*,  
 638 volume 5978 of *Lecture Notes in Computer Science*, pages 437–454. Springer, 2010. doi:10.1007/  
 639 978-3-642-11799-2\\_26.



- 640 **10** Amos Beimel, Kobbi Nissim, and Uri Stemmer. Characterizing the sample complexity of private learners.  
641 In Robert D. Kleinberg, editor, *Innovations in Theoretical Computer Science, ITCS '13, Berkeley, CA,*  
642 *USA, January 9-12, 2013*, pages 97–110. ACM, 2013. doi:10.1145/2422436.2422450.
- 643 **11** Mark Bun, Kobbi Nissim, Uri Stemmer, and Salil Vadhan. Differentially private release and learning of  
644 threshold functions. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer*  
645 *Science, FOCS '15*, pages 634–649, Washington, DC, USA, 2015. IEEE Computer Society.
- 646 **12** Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in  
647 private data analysis. In *Proceedings of the 3rd Conference on Theory of Cryptography, TCC '06*, pages  
648 265–284, Berlin, Heidelberg, 2006. Springer.
- 649 **13** Hamid Ebadi, Thibaud Antignac, and David Sands. Sampling and partitioning for differential privacy. In  
650 *14th Annual Conference on Privacy, Security and Trust, PST 2016, Auckland, New Zealand, December*  
651 *12-14, 2016*, pages 664–673. IEEE, 2016. doi:10.1109/PST.2016.7906954.
- 652 **14** Hamid Ebadi, David Sands, and Gerardo Schneider. Differential privacy: Now it's getting personal.  
653 *SIGPLAN Not.*, 50(1):69–81, January 2015.
- 654 **15** Fabian Eckert, Teresa C. Fort, Peter K. Schott, and Natalie J. Yang. Imputing missing values in the us  
655 census bureau's county business patterns. Technical report, National Bureau of Economic Research, 2021.
- 656 **16** Yann Fraboni, Richard Vidal, Laetitia Kameni, and Marco Lorenzi. Clustered sampling: Low-variance  
657 and improved representativity for clients selection in federated learning. In Marina Meila and Tong Zhang,  
658 editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24*  
659 *July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 3407–3416.  
660 PMLR, 2021. URL: <http://proceedings.mlr.press/v139/fraboni21a.html>.
- 661 **17** Marco Gaboardi, Michael Hay, and Salil Vadhan. A programming framework for opendp. *Manuscript*,  
662 2020.
- 663 **18** Marco Gaboardi, James Honaker, Gary King, Jack Murtagh, Kobbi Nissim, Jonathan Ullman, and Salil  
664 Vadhan. Psi ( $\psi$ ): A private data sharing interface. *arXiv preprint arXiv:1609.04340*, 2016.
- 665 **19** Antonious M. Girgis, Deepesh Data, Suhas N. Diggavi, Peter Kairouz, and Ananda Theertha Suresh.  
666 Shuffled model of differential privacy in federated learning. In Arindam Banerjee and Kenji Fukumizu,  
667 editors, *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April*  
668 *13-15, 2021, Virtual Event*, volume 130 of *Proceedings of Machine Learning Research*, pages 2521–2529.  
669 PMLR, 2021. URL: <http://proceedings.mlr.press/v130/girgis21a.html>.
- 670 **20** Jacob Imola and Kamalika Chaudhuri. Privacy amplification via bernoulli sampling. *arXiv preprint*  
671 *arXiv:2105.10594*, 2021.
- 672 **21** Joonas Jälkö, Antti Honkela, and Onur Dikmen. Differentially private variational inference for non-  
673 conjugate models. In Gal Elidan, Kristian Kersting, and Alexander T. Ihler, editors, *Proceedings of the*  
674 *Thirty-Third Conference on Uncertainty in Artificial Intelligence, UAI 2017, Sydney, Australia, August*  
675 *11-15, 2017*. AUAI Press, 2017. URL: [http://auai.org/uai2017/proceedings/papers/](http://auai.org/uai2017/proceedings/papers/152.pdf)  
676 [152.pdf](http://auai.org/uai2017/proceedings/papers/152.pdf).
- 677 **22** Z. Jorgensen, T. Yu, and G. Cormode. Conservative or liberal? personalized differential privacy. In *2015*  
678 *IEEE 31st International Conference on Data Engineering*, pages 1023–1034, 2015.
- 679 **23** Shiva Prasad Kasiviswanathan, Homin K. Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith.  
680 What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826, 2011.
- 681 **24** Wouter Kool, Herke van Hoof, and Max Welling. Estimating gradients for discrete random variables by  
682 sampling without replacement. In *International Conference on Learning Representations*, 2020. URL:  
683 <https://openreview.net/forum?id=rklEj2EFvB>.
- 684 **25** Levi H. S. Lelis, Roni Stern, Shahab Jabbari Arfaee, Sandra Zilles, Ariel Felner, and Robert C. Holte.  
685 Predicting optimal solution costs with bidirectional stratified sampling in regular search spaces. *Artif.*  
686 *Intell.*, 230:51–73, 2016. doi:10.1016/j.artint.2015.09.012.
- 687 **26** Jerzy Neyman. On the two different aspects of the representative method: The method of stratified  
688 sampling and the method of purposive selection. *Journal of the Royal Statistical Society*, 97(4):558–606,  
689 1934.
- 690 **27** Hieu Tat Nguyen and Arnold W. M. Smeulders. Active learning using pre-clustering. In Carla E. Brodley,  
691 editor, *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004), Banff*,

## 23:18 Controlling Privacy Loss in Sampling Schemes

- 692        *Alberta, Canada, July 4-8, 2004*, volume 69 of *ACM International Conference Proceeding Series*. ACM,  
693        2004. doi:10.1145/1015330.1015349.
- 694    **28**    Mijung Park, James R. Foulds, Kamalika Chaudhuri, and Max Welling. Variational bayes in private  
695        settings (VIPS). *J. Artif. Intell. Res.*, 68:109–157, 2020. doi:10.1613/jair.1.11763.
- 696    **29**    Adam Smith. Differential privacy and the secrecy of the sample, Feb 2010. URL: [https://](https://adamsmith.wordpress.com/2009/09/02/sample-secrecy/)  
697        [adamsmith.wordpress.com/2009/09/02/sample-secrecy/](https://adamsmith.wordpress.com/2009/09/02/sample-secrecy/).
- 698    **30**    Salil Vadhan. The complexity of differential privacy. In Yehuda Lindell, editor, *Tutorials on the Founda-*  
699        *tions of Cryptography: Dedicated to Oded Goldreich*, chapter 7, pages 347–450. Springer International  
700        Publishing AG, Cham, Switzerland, 2017.
- 701    **31**    Yu-Xiang Wang, Borja Balle, and Shiva Prasad Kasiviswanathan. Subsampled renyi differential privacy  
702        and analytical moments accountant. In *The 22nd International Conference on Artificial Intelligence*  
703        *and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of*  
704        *Machine Learning Research*, pages 1226–1235. PMLR, 2019. URL: [http://proceedings.mlr.](http://proceedings.mlr.press/v89/wang19b.html)  
705        [press/v89/wang19b.html](http://proceedings.mlr.press/v89/wang19b.html).
- 706    **32**    Yu-Xiang Wang, Stephen E. Fienberg, and Alexander J. Smola. Privacy for free: Posterior sampling  
707        and stochastic gradient monte carlo. In Francis R. Bach and David M. Blei, editors, *Proceedings of*  
708        *the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*,  
709        volume 37 of *JMLR Workshop and Conference Proceedings*, pages 2493–2502. JMLR.org, 2015. URL:  
710        <http://proceedings.mlr.press/v37/wang15.html>.
- 711    **33**    Yu-Xiang Wang, Jing Lei, and Stephen E. Fienberg. Learning with differential privacy: Stability,  
712        learnability and the sufficiency and necessity of ERM principle. *J. Mach. Learn. Res.*, 17:183:1–183:40,  
713        2016.
- 714    **34**    Yu-Xiang Wang, Jing Lei, and Stephen E. Fienberg. A minimax theory for adaptive data analysis. *arXiv*  
715        *preprint arXiv:1602.04287*, 2016.
- 716    **35**    Larry Wasserman and Shuheng Zhou. A statistical framework for differential privacy. *Journal of the*  
717        *American Statistical Association*, 105(489):375–389, 2010.

## A Basic facts about indistinguishability

718

719 ► **Definition 12.** Let the LCS distance between two data sets  $P$  and  $P'$ , denoted  $d_{\text{LCS}}(P, P')$ , be  
 720 the minimal  $k$  such that if we let  $P = P_0$  and  $P' = P_k$ , there exist data sets  $P_1, P_2, \dots, P_{k-1}$  where  
 721 for all  $i = 0, \dots, k-1$ ,  $P_i$  and  $P_{i+1}$  are unbounded neighbors.

722 ► **Lemma 13.** [12] Let  $X, Y$  and  $Z$  be random variables. For any  $\varepsilon, \varepsilon' > 0$ , if  $X$  and  $Y$  are  $\varepsilon$ -  
 723 indistinguishable, and  $Y$  and  $Z$  are  $\varepsilon'$ -indistinguishable, then  $X$  and  $Z$  are  $\varepsilon + \varepsilon'$ -indistinguishable.

724 Many of our proofs use couplings so let us briefly describe on the main method we will use to  
 725 construct a coupling of two random variables. Let  $X$  be a random variable taking values in  $\Omega_X$   
 726 and  $Y$  be a random variable taking values in  $\Omega_Y$ . Suppose there exists a (possibly randomised)  
 727 transformation  $f : \Omega_X \rightarrow \Omega_Y$  such that  $Y = f(X)$ . That is, for all  $y \in \Omega_Y$ ,

$$728 \quad \Pr(Y = y) = \sum_{x \in \Omega_X} \Pr(X = x) \Pr(f(x) = y).$$

729 Then we can construct a coupling of  $X$  and  $Y$  by  $\mu(x, y) = \Pr(X = x) \Pr(f(x) = y)$ . A short  
 730 calculation confirms that this defines a coupling. Further, notice that  $\mu(x, y) \neq 0$  if and only  
 731  $\Pr(f(x) = y) \neq 0$ .

732 ► **Lemma 14.** Let  $X$  and  $Y$  be random variables taking values in  $\mathcal{U}^*$  such that there exists a  
 733 coupling  $\mu$  such that if  $\mu(x, y) \neq 0$  then the LCS distance between  $x$  and  $y$  is at most  $A$ . Then if  $\mathcal{M}$   
 734 is  $\varepsilon$ -unbounded DP then  $\mathcal{M}(X)$  and  $\mathcal{M}(Y)$  are  $A\varepsilon$ -indistinguishable.

**Proof.**

$$735 \quad \Pr(\mathcal{M}(X) \in E) = \sum_{x \in \mathcal{U}^*} \Pr(X = x) \Pr(\mathcal{M}(x) \in E)$$

$$736 \quad = \sum_{x \in \mathcal{U}^*} \sum_{y \in \mathcal{U}^*} \mu(x, y) \Pr(\mathcal{M}(x) \in E)$$

$$737 \quad \leq \sum_{y \in \mathcal{U}^*} \sum_{x \in \mathcal{U}^*} \mu(x, y) e^{A\varepsilon} \Pr(\mathcal{M}(y) \in E)$$

$$738 \quad = e^{A\varepsilon} \sum_{y \in \mathcal{U}^*} \Pr(Y = y) \Pr(\mathcal{M}(y) \in E)$$

$$739 \quad = e^{A\varepsilon} \Pr(\mathcal{M}(Y) \in E).$$

741

742 ► **Lemma 15** (Advanced joint convexity, [6]). Let  $X$  and  $Y$  be random variables satisfying  
 743  $X = (1 - q)X_0 + qX_1$  and  $Y = (1 - q)Y_0 + qY_1$  for some  $q \in [0, 1]$  and random variables  
 744  $X_0, X_1, Y_0$  and  $Y_1$ . If

745 ■  $X_0$  and  $Y_0$  are  $\varepsilon$ -indistinguishable,

746 ■  $X_1$  and  $Y_1$  are  $\varepsilon + \varepsilon'$ -indistinguishable and

747 ■  $X_1$  and  $Y_0$  are  $\varepsilon + \varepsilon'$ -indistinguishable

748 ■  $X_0$  and  $Y_1$  are  $\varepsilon + \varepsilon'$ -indistinguishable

749 then  $X$  and  $Y$  are  $\varepsilon + \log(1 + q(e^{\varepsilon'} - 1))$ -indistinguishable.

## 23:20 Controlling Privacy Loss in Sampling Schemes

**Proof.**

$$\begin{aligned}
 750 \quad \Pr(X \in E) &= (1 - q) \Pr(X_0 \in E) + q \Pr(X_1 \in E) \\
 751 \quad &= (1 - q) \Pr(X_0 \in E) + qe^{-\varepsilon'} \Pr(X_1 \in E) \\
 752 \quad &\quad + q(1 - e^{-\varepsilon'})(1 - q) \Pr(X_1 \in E) + q(1 - e^{-\varepsilon'})q \Pr(X_1 \in E) \\
 753 \quad &\leq (1 - q)e^\varepsilon \Pr(Y_0 \in E) + qe^\varepsilon \Pr(Y_1 \in E) \\
 754 \quad &\quad + q(1 - e^{-\varepsilon'})(1 - q)e^{\varepsilon+\varepsilon'} \Pr(Y_0 \in E) + q(1 - e^{-\varepsilon'})qe^{\varepsilon+\varepsilon'} \Pr(Y_1 \in E) \\
 755 \quad &= (e^\varepsilon + qe^\varepsilon(e^{\varepsilon'} - 1)) \Pr(Y \in E)
 \end{aligned}$$

757

### **B** Randomized data-independent sampling

759 **► Lemma 16.** Given  $m \in \mathbb{N}$ , define  $\mathcal{C}_m : \mathcal{U}^* \rightarrow \mathcal{U}^m$  be defined as follows: given a dataset  $P$ , form  
 760 a sample  $S$  by sampling  $m$  data points randomly without replacement from  $P$ , then  $\mathcal{C}_m(P) = S$ . Let  
 761  $P$  and  $P'$  be unbounded neighboring datasets and  $m, m' \in \mathbb{N}$ , then  $\mathcal{M}_{\mathcal{C}_m}(P)$  and  $\mathcal{M}_{\mathcal{C}_{m'}}(P')$  are

$$762 \quad \left( \log \left( 1 + \frac{m}{|P|+1} (e^{2\varepsilon} - 1) \right) + |m - m'| \varepsilon \right) \text{ -indistinguishable.}$$

763 **Proof.** Let  $P' = P \cup \{x\}$ . First, let us focus on the case where  $m' = m$ . Now,

$$\begin{aligned}
 764 \quad \mathcal{M}_{\mathcal{C}_m}(P') &= \frac{\binom{|P|}{m}}{\binom{|P|+1}{m}} \mathcal{M}_{\mathcal{C}_m}(P) + \left( 1 - \frac{\binom{|P|}{m}}{\binom{|P|+1}{m}} \right) \mathcal{M}(\mathcal{C}_m(P')|_{x \in S}) \\
 765 \quad &= \left( 1 - \frac{m}{|P|+1} \right) \mathcal{M}_{\mathcal{C}_m}(P) + \frac{m}{|P|+1} \mathcal{M}(\mathcal{C}_m(P')|_{x \in S}),
 \end{aligned}$$

767 where  $\mathcal{C}_m(P')|_{x \in S}$  denotes the random variable  $\mathcal{C}_m(P')$  conditioned on the event that  $x \in S$ . Now,  
 768 we can define a coupling of  $\mathcal{C}_m(P)$  and  $\mathcal{C}_m(P')|_{x \in S}$  by first sampling  $S$  from  $\mathcal{C}_m(P)$ , then replacing  
 769 a random element of  $S$  by  $x$ . This coupling has LCS distance at most 2, so by Lemma 14,  $\mathcal{M}_{\mathcal{C}_m}(P)$   
 770 and  $\mathcal{M}(\mathcal{C}_m(P')|_{x \in S})$  are  $2\varepsilon$ -indistinguishable. Thus, by Lemma 15,  $\mathcal{M}_{\mathcal{C}_m}(P)$  and  $\mathcal{M}_{\mathcal{C}_m}(P')$  are

$$771 \quad \log \left( 1 + \frac{m}{|P|+1} (e^{2\varepsilon} - 1) \right) \text{ -indistinguishable.}$$

772 Next, let us consider the case  $|m - m'| = 1$  and  $P = P'$ . We can define a coupling of  
 773  $\mathcal{C}_m(P)$  and  $\mathcal{C}_{m'}(P)$  as follows: first sample  $S$  from  $\mathcal{C}_m(P)$ , then add a random element of  $P \setminus S$   
 774 to  $S$ . This coupling has LCS distance at most 1, so by Lemma 14,  $\mathcal{M}_{\mathcal{C}_m}(P)$  and  $\mathcal{M}_{\mathcal{C}_{m'}}(P)$  are  
 775  $\varepsilon$ -indistinguishable.

776 Finally, we'll use Lemma 13 to complete the proof. Note that  $\mathcal{M}_{\mathcal{C}_m}(P)$  and  $\mathcal{M}_{\mathcal{C}_m}(P')$  are  
 777  $\log \left( 1 + \frac{m}{|P|+1} (e^{2\varepsilon} - 1) \right)$ -indistinguishable. Then there exist  $m_1, \dots, m_{\ell-1}$  such that if we set  
 778  $m_0 = m$  and  $m_{|m-m'|} = m'$  then for all  $i$ ,  $|m_i - m_{i-1}| \leq 1$  and so  $\mathcal{M}(\mathcal{C}_{m_{i-1}}(P'))$  and  $\mathcal{M}(\mathcal{C}_{m_i}(P'))$   
 779 are  $\varepsilon$ -indistinguishable. Therefore, by Lemma 13,  $\mathcal{M}_{\mathcal{C}_m}(P)$  and  $\mathcal{M}_{\mathcal{C}_{m'}}(P')$  are

$$780 \quad \left( \log \left( 1 + \frac{m}{|P|+1} (e^{2\varepsilon} - 1) \right) + |m - m'| \varepsilon \right) \text{ -indistinguishable.}$$

781

782 ► **Definition 17** (log-Lipschitz functions). A function  $q : [n] \rightarrow \mathbb{R}_{\geq 0}$  is  $\varepsilon$ -log-Lipschitz if for all  
783  $m \in \{0, 1, \dots, n-1\}$ ,

$$784 \quad |\log q(m) - \log q(m+1)| \leq \varepsilon$$

785 ► **Lemma 18.** Let  $w : [n] \rightarrow \mathbb{R}_{\geq 0}$  be nondecreasing, and let  $p : [n] \rightarrow \mathbb{R}_{\geq 0}$  be any function. Then,

$$786 \quad \max_{q: [n] \rightarrow \mathbb{R}_{\geq 0} \text{ is } \varepsilon\text{-log-Lipschitz}} \frac{\sum_{m=0}^n q(m)w(m)p(m)}{\sum_{m=0}^n q(m)p(m)} \leq \frac{\sum_{m=0}^n e^{\varepsilon m} w(m)p(m)}{\sum_{m=0}^n e^{\varepsilon m} p(m)}$$

787 **Proof.** We will show by induction on  $k = 0, 1, \dots, n$  that we can assume w.l.o.g. that the maximizer  
788 has the form  $q(0) = 1, q(1) = e^\varepsilon, \dots, q(k) = e^{\varepsilon k}$ . This holds for  $k = 0$  by simply normalizing. Then,  
789 assuming it holds for some  $k > 0$ , and given any  $\varepsilon$ -log-Lipschitz  $q$  such that  $q(0) = 1, \dots, q(k) = e^{\varepsilon k}$ ,  
790 let us define  $q'$  as follows.

$$791 \quad q'(m) = \begin{cases} q(m) & \text{for } m = 0, 1, \dots, k \\ \frac{e^\varepsilon q(k)q(m)}{q(k+1)} & \text{for } m = k+1, \dots, n \end{cases}$$

793 By construction,  $q'$  is  $\varepsilon$ -log-Lipschitz. In particular,  $q'(k+1) = e^\varepsilon q(k) = e^{(k+1)\varepsilon}$ . In addition, since  
794  $q$  is  $\varepsilon$ -log-Lipschitz, we have that

$$795 \quad \frac{e^\varepsilon q(k)}{q(k+1)} \geq 1,$$

797 which means that  $q'(k+1) \geq q(k+1)$ .

798 Next, we use the inequality (a slight generalization of the mediant inequality) that for  $a, b, c, d > 0$   
799 and  $t \geq 1$  such that  $a/b \leq c/d$ ,

$$800 \quad \frac{a+c}{b+d} \leq \frac{a+tc}{b+td}$$

801 Let  $a = \sum_{m=0}^k q(m)w(m)p(m)$ ,  $b = \sum_{m=0}^k q(m)p(m)$ ,  $c = \sum_{m=k+1}^n q(m)w(m)p(m)$ , and  $d =$   
802  $\sum_{m=k+1}^n q(m)p(m)$ , and  $t = e^\varepsilon q(k)/q(k+1)$ . By the non-decreasing property of  $w$ , we have that

$$803 \quad a/b \leq w(k) \leq w(k+1) \leq c/d.$$

805 Therefore, the inequality above, and by definition of  $q'$ , we have that

$$806 \quad \frac{\sum_{m=0}^n q(m)w(m)p(m)}{\sum_{m=0}^n q(m)p(m)} \leq \frac{\sum_{m=0}^k q(m)w(m)p(m) + e^\varepsilon (q(k)/q(k+1)) \cdot \sum_{m=k+1}^n q(m)w(m)p(m)}{\sum_{m=0}^k q(m)p(m) + e^\varepsilon (q(k)/q(k+1)) \cdot \sum_{m=k+1}^n q(m)p(m)}$$

$$807 \quad = \frac{\sum_{m=0}^n q'(m)w(m)p(m)}{\sum_{m=0}^n q'(m)p(m)}.$$

809 So, by induction, we can assume that the maximizer has the form  $q(0) = 1, q(1) = e^\varepsilon, \dots, q(n) =$   
810  $e^{\varepsilon n}$ , which completes the proof. ◀

811 **Proof of Theorem 3.** Let  $\mathcal{C}_m : \mathcal{U}^* \rightarrow \mathcal{U}^m$  be the sampling scheme that given a dataset  $P$ , returns  
812  $S$  where  $S$  is a uniformly random subset of  $P$  of size  $m$  (drawn without replacement). Let  $y \in \mathcal{O}$  be

## 23:22 Controlling Privacy Loss in Sampling Schemes

813 any outcome, and let  $P \sim P'$  be neighboring datasets. Then, we have that

$$\begin{aligned}
 814 \quad \Pr[\mathcal{M}_{\mathcal{C}}(P) = y] &= \sum_{m=0}^n \Pr[\mathcal{M}_{\mathcal{C}_m}(P) = y] \cdot \Pr[|\mathcal{C}(P)| = m] \\
 815 &\leq \sum_{m=0}^n \left(1 + \frac{m}{n}(e^\varepsilon - 1)\right) \cdot \Pr[\mathcal{M}_{\mathcal{C}_m}(P') = y] \cdot t(m) \\
 816 &\leq \frac{\sum_{m=0}^n \left(1 + \frac{m}{n}(e^\varepsilon - 1)\right) \cdot e^{\varepsilon m} \cdot t(m)}{\sum_{m=0}^n e^{\varepsilon m} t(m)} \cdot \sum_{m=0}^n \Pr[\mathcal{M}_{\mathcal{C}_m}(P') = y] \cdot t(m) \\
 817 &= \left(1 + \frac{\mathbb{E}_{m \sim \tilde{t}}[m]}{n}(e^\varepsilon - 1)\right) \cdot \Pr[\mathcal{M}_{\mathcal{C}}(P') = y] \\
 818
 \end{aligned}$$

819 where the first inequality follows from Lemma 16. Then, note that  $(1 + (m/n)(e^\varepsilon - 1))$  is non-  
 820 decreasing, and that  $\Pr[\mathcal{M}_{\mathcal{C}_m}(P') = y]$  is  $\varepsilon$ -log-Lipschitz by definition, so the second inequality  
 821 follows by Lemma 18. After rearranging and simplifying, we obtain the desired result.

822 Finally, for the lower bound, suppose the data universe  $\mathcal{U} = [0, 1]$ . Let  $P = \{1, \dots, 1\}$  consist  
 823 of  $n$  1s and  $P'$  be the neighboring dataset  $P' = P \setminus \{1\} \cup \{0\}$ . Let  $\mathcal{M} : \mathcal{U}^* \rightarrow \mathbb{R}$  be defined by  
 824  $\mathcal{M}(S) = \sum_{x \in S} \mathbb{1}\{x = 1\} + \text{Lap}(1/\varepsilon)$  so  $\mathcal{M}$  is  $\varepsilon$ -unbounded DP. Then

$$\begin{aligned}
 825 \quad \frac{\Pr(\mathcal{M}_{\mathcal{C}}(P') = n)}{\Pr(\mathcal{M}_{\mathcal{C}}(P) = n)} &= \frac{\sum_{m=0}^n \Pr(t = m) \left(\frac{m}{n} e^{-(n-m+1)\varepsilon} + \left(1 - \frac{m}{n}\right) e^{-(n-m)\varepsilon}\right)}{\sum_{m=0}^n \Pr(t = m) e^{-(n-m)\varepsilon}} \\
 826 &= \frac{\sum_{m=0}^n \Pr(t = m) \left(\frac{m}{n} e^{(m-1)\varepsilon} + \left(1 - \frac{m}{n}\right) e^{m\varepsilon}\right)}{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon}} \\
 827 &= \frac{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon} \left(1 - \frac{m}{n}(1 - e^{-\varepsilon})\right)}{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon}} \\
 828 &= 1 - \frac{1}{n}(1 - e^{-\varepsilon}) \frac{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon} m}{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon}}. \\
 829
 \end{aligned}$$

830 Thus, taking the reciprocal,

$$831 \quad \log \frac{\Pr(\mathcal{M}_{\mathcal{C}}(P) = n)}{\Pr(\mathcal{M}_{\mathcal{C}}(P') = n)} = -\log \left(1 - \frac{1}{n}(1 - e^{-\varepsilon}) \frac{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon} m}{\sum_{m=0}^n \Pr(t = m) e^{m\varepsilon}}\right).$$

832

### 833 **C** Data-dependent sampling

834 **Proof of Proposition 4: hypothesis testing perspective.** Let  $H : \mathbb{N} \rightarrow \{0, 1\}$  be the hy-  
 835 pothesis test such that for all  $x \in \mathbb{N}$ , and  $b \in \{0, 1\}$ ,  $e^{-\varepsilon} \Pr(H(x) = b) \leq \Pr(H(x+1) =$   
 836  $b) \leq e^\varepsilon \Pr(H(x) = b)$ . Then  $H' : \mathcal{U}^* \rightarrow \{0, 1\}$  defined by  $H'(S) = H(|S|)$  is  $\varepsilon$ -unbounded DP.  
 837 By assumption,  $H'_{\tilde{f}}$  is  $\varepsilon'$ -DP. This implies that  $H(\tilde{f}(P))$  and  $H(\tilde{f}(P'))$  are  $\varepsilon'$ -indistinguishable.  
 838 Therefore,

$$839 \quad \text{adv}(H) = \Pr[H(\tilde{f}(P)) = 0] - \Pr[H(\tilde{f}(P')) = 0] \leq \Pr[H(\tilde{f}(P')) = 0](e^{\varepsilon'} - 1) \leq e^{\varepsilon'} - 1.$$

840 The result follows from taking the supremum over all  $\varepsilon$ -DP  $H$ . ◀

841 **Proof of Theorem 6: proportional allocation with randomized rounding.** Let  $P$  be a data-  
 842 set,  $x$  be a data point and  $P' = P \cup \{x\}$ . Let  $m = r|P|$ ,  $m' = r|P'|$ ,  $m^L = \lfloor m \rfloor$ ,  $m'^L = \lfloor m' \rfloor$ ,  
 843  $p = m - m^L$  and  $p' = m' - m'^L$ . Now,  $m' - m = r < 1$  so we have two cases,  $m^L = m'^L$  or  
 844  $m^L = m'^L - 1$ .

845 As in Lemma 16, let  $\mathcal{C}_m : \mathcal{U}^* \rightarrow \mathcal{U}^m$  be the sampling scheme that given a dataset  $P$ , returns  $S$   
 846 where  $S$  is a uniformly random subset of  $P$  of size  $m$  (drawn without replacement). Note that by  
 847 Theorem 2, for  $m, m' \in \mathbb{N}$ ,  $\mathcal{M}_m(P)$  and  $\mathcal{M}_{m'}(P)$  are  $|m - m'|\varepsilon$ -indistinguishable, and  $\mathcal{M}_{\mathcal{C}_m}(P)$   
 848 and  $\mathcal{M}_{\mathcal{C}_{m'}}(P')$  are  $\log\left(1 + \frac{m}{|P|+1}(e^{2\varepsilon} - 1)\right) + |m - m'|\varepsilon$ -indistinguishable.

849 Firstly, suppose  $m^L = m'^L$ . Let

$$850 \mu_0 = \frac{1}{1-r}((1-p-r)\mathcal{M}_{\mathcal{C}_{m^L}}(P) + p\mathcal{M}_{\mathcal{C}_{m^L+1}}(P)),$$

$$851 \mu'_0 = \frac{1}{1-r}((1-p-r)\mathcal{M}_{\mathcal{C}_{m^L}}(P') + p\mathcal{M}_{\mathcal{C}_{m^L+1}}(P')),$$

$$852 \mu_1 = \mathcal{M}_{\mathcal{C}_{m^L}}(P),$$

$$853 \mu'_1 = \mathcal{M}_{\mathcal{C}_{m^L+1}}(P').$$

855 Notice that

$$856 \mathcal{M}_{\mathcal{C}_r}(P) = (1-r)\mu_0 + r\mu_1 \quad \text{and} \quad \mathcal{M}_{\mathcal{C}_r}(P') = (1-r)\mu'_0 + r\mu'_1.$$

858 Now, by Lemma 15 and Lemma 14,  $\mu_0$  and  $\mu'_0$  are  $\log(1 + \frac{m^L+1}{|P|+1}(e^{2\varepsilon} - 1))$ -indistinguishable. Further,  
 859 all the pairs  $(\mu'_0, \mu_1)$ ,  $(\mu_1, \mu'_1)$  and  $(\mu_0, \mu'_1)$  are  $(\log(1 + \frac{m^L+1}{|P|+1}(e^{2\varepsilon} - 1)) + \varepsilon)$ -indistinguishable.  
 860 Therefore, by Lemma 15,  $\mathcal{M}_{\mathcal{C}_r}(P)$  and  $\mathcal{M}_{\mathcal{C}_r}(P')$  are  $\varepsilon'$ -indistinguishable where

$$861 \varepsilon' \leq \log\left(1 + \frac{m^L+1}{|P|+1}(e^{2\varepsilon} - 1)\right) + \log(1 + r(e^\varepsilon - 1))$$

$$862 \leq \log\left(1 + \left(r + \frac{1}{|P|+1}\right)(e^{2\varepsilon} - 1)\right) + \log(1 + r(e^\varepsilon - 1)).$$

864 Next, suppose  $m'^L = m^L + 1$ . Let  $1 - q = \min\{p, 1 - p'\}$  and

$$865 \mu_0 = \mathcal{M}_{\mathcal{C}_{m^L+1}}(P),$$

$$866 \mu'_0 = \mathcal{M}_{\mathcal{C}_{m^L+1}}(P'),$$

$$867 \mu_1 = \frac{1}{q}((p-1+q)\mathcal{M}_{\mathcal{C}_{m^L+1}}(P) + (1-p)\mathcal{M}_{\mathcal{C}_{m^L}}(P)),$$

$$868 \mu'_1 = \frac{1}{q}((1-p'-1+q)\mathcal{M}_{\mathcal{C}_{m^L+1}}(P') + p'\mathcal{M}_{\mathcal{C}_{m^L+2}}(P')).$$

870 Notice that

$$871 \mathcal{M}_{\mathcal{C}_r}(P) = (1-q)\mu_0 + q\mu_1 \quad \text{and} \quad \mathcal{M}_{\mathcal{C}_r}(P') = (1-q)\mu'_0 + q\mu'_1.$$

873 Now, by Lemma 2,  $\mu_0$  and  $\mu'_0$  are  $\log\left(1 + \frac{m^L+1}{|P|+1}\right)$ -indistinguishable. Further, all the pairs  $(\mu'_0, \mu_1)$ ,  
 874  $(\mu_1, \mu'_1)$  and  $(\mu_0, \mu'_1)$  are  $(\log(1 + \frac{m^L+1}{|P|+1}(e^{2\varepsilon} - 1)) + 2\varepsilon)$ -indistinguishable. Also, note that  $q \leq r$ .  
 875 Then by Lemma 15,  $\mathcal{M}_{\mathcal{C}_r}(P)$  and  $\mathcal{M}_{\mathcal{C}_r}(P')$  are  $\varepsilon'$ -indistinguishable where

$$876 \varepsilon' \leq \log\left(1 + \frac{m^L+1}{|P|+1}(e^{2\varepsilon} - 1)\right) + \log(1 + p(e^{2\varepsilon} - 1))$$

$$877 \leq \log\left(1 + \left(r + \frac{1}{|P|+1}\right)(e^{2\varepsilon} - 1)\right) + \log(1 + r(e^{2\varepsilon} - 1)).$$

879

**D Cluster sampling**

**Proof of Theorem 8.** Without loss of generality, let  $i = 1$ . Notice that conditioned on cluster  $1 \notin I$ , the distribution of outputs of  $\mathcal{M}_C(P)$  and  $\mathcal{M}_C(P')$  are identical. Let  $E$  be a set of outcomes. Then

$$\begin{aligned} \Pr(\mathcal{M}_C(P) \in E) &= \frac{\ell}{k} \Pr(\mathcal{M}_C(P) \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P) \in E \mid 1 \notin I) \\ &= \frac{\ell}{k} \Pr(\mathcal{M}_C(P) \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I). \end{aligned}$$

Now, we have that

$$\begin{aligned} \frac{\ell}{k} \Pr(\mathcal{M}_C(P) \in E \mid 1 \in I) &= \frac{\ell}{k} \sum_{|I|=\ell, 1 \in I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\ &\leq \frac{\ell}{k} \sum_{|I|=\ell, 1 \in I} \frac{1}{\binom{k}{\ell}} e^\varepsilon \Pr(\mathcal{M}(P'_I) \in E) \\ &= \frac{\ell}{k} e^\varepsilon \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I), \end{aligned}$$

where the inequality follows from the fact that the LCS distance between  $P_I$  and  $P'_I$  is 1. Thus,

$$\begin{aligned} \Pr(\mathcal{M}_C(P) \in E) &\leq \frac{\ell}{k} e^\varepsilon \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I) \\ &= \Pr(\mathcal{M}_C(P') \in E) + \frac{\ell}{k} (e^\varepsilon - 1) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I). \end{aligned}$$

Now, we need to relate  $\Pr(\mathcal{M}_C(P') \in E \mid 1 \in I)$  to  $\Pr(\mathcal{M}_C(P) \in E)$ . For a set  $I$  such that  $1 \notin I$  and index  $i \in I$ , let  $I \cup \{1\} \setminus \{i\}$  be the set where index  $i$  has been replaced with 1. Then,

$$\begin{aligned} \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I) &= \sum_{|I|=\ell, 1 \notin I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\ &= \sum_{|I|=\ell, 1 \notin I} \sum_{i \in I} \frac{1}{\ell} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\ &\geq \sum_{|I|=\ell, 1 \notin I} \sum_{i \in I} \frac{1}{\ell} \frac{1}{\binom{k}{\ell}} e^{-(n_1+n_i)\varepsilon} \Pr(\mathcal{M}(P_{I \cup \{1\} \setminus \{i\}}) \in E) \\ &\geq e^{-(n_1+n_{\max})\varepsilon} \frac{1}{\ell} \sum_{|I|=\ell, 1 \notin I} \sum_{i \in I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_{I \cup \{1\} \setminus \{i\}}) \in E), \end{aligned}$$

where the first inequality follows from the fact that the LCS distance between  $P_I$  and  $P_{I \cup \{1\} \setminus \{i\}}$  is at most  $n_1 + n_i$ . Now, notice that the sets  $I \cup \{1\} \setminus \{i\}$  in the above sum all contain 1, and each index  $I'$  such that  $|I'| = \ell$  and  $1 \in I'$  appears in the sum  $k - \ell$  times (corresponding to the  $k - \ell$  possible choices for the swapped index  $i$ ). Therefore, we can rewrite the sum as

$$\begin{aligned} \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I) &\geq e^{-(n_1+n_{\max})\varepsilon} \frac{k - \ell}{\ell} \sum_{|I|=\ell, 1 \in I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\ &= e^{-(n_1+n_{\max})\varepsilon} \frac{k - \ell}{\ell} \frac{\ell}{k} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) \\ &= e^{-(n_1+n_{\max})\varepsilon} \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I). \end{aligned}$$



911 Thus,

$$\begin{aligned}
 912 \quad \Pr(\mathcal{M}_C(P') \in E) &= \frac{\ell}{k} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I) \\
 913 \quad &\geq \frac{\ell}{k} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\max})\varepsilon} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) \\
 914 \quad &= \left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\max})\varepsilon}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I). \\
 915
 \end{aligned}$$

916 Finally,

$$\begin{aligned}
 917 \quad \Pr(\mathcal{M}_C(P) \in E) &\leq \Pr(\mathcal{M}_C(P') \in E) + \frac{\ell}{k} (e^\varepsilon - 1) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) \\
 918 \quad &\leq \Pr(\mathcal{M}_C(P') \in E) + \frac{\ell}{k} (e^\varepsilon - 1) \frac{1}{\left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\max})\varepsilon}\right)} \Pr(\mathcal{M}_C(P') \in E) \\
 919 \quad &\leq \left(1 + \frac{\ell}{k} (e^\varepsilon - 1) \frac{1}{\left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\max})\varepsilon}\right)}\right) \Pr(\mathcal{M}_C(P') \in E) \\
 920
 \end{aligned}$$

921 ◀

922 Now we turn our attention to the lower bound.

923 **Proof of Theorem 7.** Let  $C_1 = \{1, \dots, 1\}$  and  $C_j = \{-1, \dots, -1\}$  for all  $j \in \{2, \dots, k\}$ .  
 924 Let  $C'_1 = C_1 \setminus \{1\} \cup \{-1\}$  be the same as  $C_1$  except with one 1 switched to a -1. Let  $\mathcal{M}(S) =$   
 925  $\sum_{x \in S} x + \text{Lap}(1/\varepsilon)$ , so  $\mathcal{M}$  is  $\varepsilon$ -unbounded DP. Notice that  $\mathcal{M}$  has the property that if  $\sum_{x \in S'} x =$   
 926  $\sum_{x \in S} x + a$ , for some  $a \in \mathbb{R}$  then  $\Pr(\mathcal{M}(S) = \sum_{x \in S} x) = e^{|a|\varepsilon} \Pr(\mathcal{M}(S') = \sum_{x \in S} x)$ . This  
 927 equality allows us to tighten many of the inequalities that appeared in the proof of Theorem 8, and  
 928 give a lower bound.

$$\begin{aligned}
 929 \quad \Pr(\mathcal{M}_C(P) = n_1 + 1) &= \frac{\ell}{k} \Pr(\mathcal{M}_C(P) = n_1 + 1 \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P) = n_1 + 1 \mid 1 \notin I) \\
 930 \quad &= \frac{\ell}{k} e^\varepsilon \Pr(\mathcal{M}_C(P') = n_1 + 1 \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') = n_1 + 1 \mid 1 \notin I) \\
 931 \quad &= \Pr(\mathcal{M}_C(P') = n_1 + 1) + \frac{\ell}{k} (e^\varepsilon - 1) \Pr(\mathcal{M}_C(P') = n_1 + 1 \mid 1 \in I).
 \end{aligned}$$

933 Now,

$$\begin{aligned}
 934 \quad \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I) &= \sum_{|I|=\ell, 1 \notin I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\
 935 \quad &= \sum_{|I|=\ell, 1 \notin I} \sum_{i \in I} \frac{1}{\ell} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\
 936 \quad &= \sum_{|I|=\ell, 1 \notin I} \sum_{i \in I} \frac{1}{\ell} \frac{1}{\binom{k}{\ell}} e^{-(n_1+n_i)\varepsilon} \Pr(\mathcal{M}(P_{I \cup \{1\} \setminus \{i\}}) \in E) \\
 937 \quad &\leq e^{-(n_1+n_{\min})\varepsilon} \frac{1}{\ell} \sum_{|I|=\ell, 1 \notin I} \sum_{i \in I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_{I \cup \{1\} \setminus \{i\}}) \in E) \\
 938 \quad &= e^{-(n_1+n_{\max})\varepsilon} \frac{k-\ell}{\ell} \sum_{|I|=\ell, 1 \in I} \frac{1}{\binom{k}{\ell}} \Pr(\mathcal{M}(P_I) \in E) \\
 939 \quad &= e^{-(n_1+n_{\max})\varepsilon} \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I). \\
 940
 \end{aligned}$$

## 23:26 Controlling Privacy Loss in Sampling Schemes

941 Thus,

$$\begin{aligned}
 942 \quad \Pr(\mathcal{M}_C(P') \in E) &= \frac{\ell}{k} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \notin I) \\
 943 \quad &\leq \frac{\ell}{k} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\min})\varepsilon} \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) \\
 944 \quad &= \left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\min})\varepsilon}\right) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I). \\
 945
 \end{aligned}$$

946 Finally,

$$\begin{aligned}
 947 \quad \Pr(\mathcal{M}_C(P) \in E) &= \Pr(\mathcal{M}_C(P') \in E) + \frac{\ell}{k}(e^\varepsilon - 1) \Pr(\mathcal{M}_C(P') \in E \mid 1 \in I) \\
 948 \quad &\geq \Pr(\mathcal{M}_C(P') \in E) + \frac{\ell}{k}(e^\varepsilon - 1) \frac{1}{\left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\min})\varepsilon}\right)} \Pr(\mathcal{M}_C(P') \in E) \\
 949 \quad &= \left(1 + \frac{\ell}{k}(e^\varepsilon - 1) \frac{1}{\left(\frac{\ell}{k} + \left(1 - \frac{\ell}{k}\right) e^{-(n_1+n_{\min})\varepsilon}\right)}\right) \Pr(\mathcal{M}_C(P') \in E). \\
 950
 \end{aligned}$$

951

952 **Proof of Theorem 9.** Let  $D = \mathcal{C}_\ell(P)$  and  $D' = \mathcal{C}_\ell(P')$ . For an event  $E \in \mathcal{O}$ , define the  
 953 probabilities  $p, q, p'$  and  $q'$  as follows.

$$\begin{aligned}
 954 \quad p &= \Pr(\mathcal{M}(D) \in E \mid C_1 \in D) & q &= \Pr(\mathcal{M}(D) \in E \mid C_1 \notin D) \\
 955 \quad p' &= \Pr(\mathcal{M}(D') \in E \mid C_1 \in D') & q' &= \Pr(\mathcal{M}(D') \in E \mid C_1 \notin D')
 \end{aligned}$$

957 By the existence of  $\mathcal{H}$  described in the lemma statement, there must exist an event  $E$  such that  
 958  $q \leq e^{-\varepsilon'} p$ . Since  $P$  and  $P'$  only differ on  $C_1$ , the distributions of  $\mathcal{M}(D)|_{C_1 \notin D}$  and  $\mathcal{M}(D')|_{C_1 \notin D'}$  are  
 959 identical, which means that  $q = q'$ . Then, we can compute a lower bound on the indistinguishability  
 960 of  $\mathcal{M}(D)$  and  $\mathcal{M}(D')$  as follows. Without loss of generality, assume  $p' > p$

$$\begin{aligned}
 961 \quad \frac{\Pr(\mathcal{M}(D') \in E)}{\Pr(\mathcal{M}(D) \in E)} &= \frac{p' \cdot \Pr(C_1 \in D') + q \cdot \Pr(C_1 \notin D')}{p \cdot \Pr(C_1 \in D) + q \cdot \Pr(C_1 \notin D)} \\
 962 \quad &= \frac{p' \cdot \frac{\ell}{k} + q \cdot \left(1 - \frac{\ell}{k}\right)}{p \cdot \frac{\ell}{k} + q \cdot \left(1 - \frac{\ell}{k}\right)} \\
 963 \quad &= \frac{p \cdot \frac{\ell}{k} + q \cdot \left(1 - \frac{\ell}{k}\right) + (p' - p) \cdot \frac{\ell}{k}}{p \cdot \frac{\ell}{k} + q \cdot \left(1 - \frac{\ell}{k}\right)} \\
 964 \quad &= 1 + \frac{(p' - p) \frac{\ell}{k}}{p \cdot \frac{\ell}{k} + q \cdot \left(1 - \frac{\ell}{k}\right)} \\
 965 \quad &\geq 1 + \frac{(p' - p) \frac{\ell}{k}}{p \cdot \left(\frac{\ell}{k} + e^{-\varepsilon'} \left(1 - \frac{\ell}{k}\right)\right)} \\
 966 \quad &= 1 + \left(\frac{p'}{p} - 1\right) \frac{\frac{\ell}{k}}{\frac{\ell}{k} + e^{-\varepsilon'} \left(1 - \frac{\ell}{k}\right)} \\
 967
 \end{aligned}$$

968 where the final inequality follows from the fact that  $\mathcal{M}$  is  $\varepsilon$ -DP, so  $p'/p \geq e^{\varepsilon'}$  by definition. ◀

### 969 **E Stratified sampling**

970 **Proof of Theorem 11: proportional allocation for stratified sampling.** Given  $\mathcal{M} : ([k] \times$   
 971  $\mathcal{U})^* \rightarrow \mathcal{Y}$ , for all datasets  $T_2, \dots, T_k \in \mathcal{U}^*$ , define  $\mathcal{M}^{T_2, \dots, T_k} : \mathcal{U}^* \rightarrow \mathcal{Y}$  by  $\mathcal{M}^{T_2, \dots, T_k}(S) =$

972  $\mathcal{M}(S \sqcup T_2 \sqcup \dots \sqcup T_k)$ . Then since  $\mathcal{M}$  was  $(\varepsilon, \dots, \varepsilon)$ -stratified unbounded DP,  $\mathcal{M}^{T_2, \dots, T_k}$  is  $\varepsilon$ -  
 973 unbounded DP. Let  $\mathcal{C}_r$  be as in Lemma 6 so for all  $S, S'$  unbounded neighbours such that  $r|S| \geq 1$   
 974 and  $r|S'| \geq 1$ ,  $\mathcal{M}_{\mathcal{C}_r}^{T_2, \dots, T_k}(S)$  and  $\mathcal{M}_{\mathcal{C}_r}^{T_2, \dots, T_k}(S')$  are  $\varepsilon'$ -indistinguishable where

$$975 \quad \varepsilon' \leq \log(1 + 2r(e^{2\varepsilon} - 1)) + \log(1 + r(e^{2\varepsilon} - 1)).$$

976 Now, let  $P = S_1 \sqcup S_2 \sqcup \dots \sqcup S_k$  and  $P' = S'_1 \sqcup S_2 \sqcup \dots \sqcup S_k$  be unbounded stratified neighboring  
 977 datasets that differ in the first stratum. Since  $S_2 \sqcup \dots \sqcup S_k$  are shared between  $P$  and  $P'$ , and the  
 978 datasets  $T_i$  only dependent on strata  $S_i$ , the distribution of  $T_2, \dots, T_k$  are identical given inputs  $P$  and  
 979  $P'$ . Let  $q$  be the distribution of  $T_2, \dots, T_k$  so  $q(T_2, \dots, T_k) = \Pr(\mathcal{C}_r(S_2) = T_2, \dots, \mathcal{C}_r(S_k) = T_k)$ .  
 980 Then given an event  $E$ ,

$$\begin{aligned} 981 \quad \Pr(\mathcal{M}_{\mathcal{C}_{f_{\text{prop}}, r}}(P) \in E) &= \int_{T_2, \dots, T_k} q(T_2, \dots, T_k) \Pr(\mathcal{M}_{\mathcal{C}_r}^{T_2, \dots, T_k}(S_1) \in E) \\ 982 &\leq \int_{T_2, \dots, T_k} q(T_2, \dots, T_k) e^{\varepsilon'} \Pr(\mathcal{M}_{\mathcal{C}_r}^{T_2, \dots, T_k}(S'_1) \in E) \\ 983 &= e^{\varepsilon'} \Pr(\mathcal{M}_{\mathcal{C}_{f_{\text{prop}}, r}}(P') \in E). \end{aligned}$$

985

## 986 **F** Instantiating the estimation-based lower bound

987 For any function  $f$ , the accuracy to which one can privately estimate  $f$  depends on the *sensitivity* of  
 988  $f$ , i.e., how much changing the input dataset can change the output.

989 **Proof of Proposition 5.** Define  $\mathcal{M}_{SS} : \mathcal{U}^* \rightarrow \mathbb{N}$  as follows. For all  $P \in \mathcal{U}^*$ ,

$$990 \quad \mathcal{M}(P) = |P| + \text{Lap}(1/\varepsilon).$$

991 Then  $\mathcal{M}$  is  $\varepsilon$ -unbounded DP. Suppose that  $\tilde{f} : \mathcal{U}^* \rightarrow \mathbb{N}$  is such that for all  $\varepsilon$ -unbounded DP  
 992 mechanisms  $\mathcal{A}$ ,  $\mathcal{A}_{\tilde{f}}$  is  $\varepsilon'$ -bounded DP. This implies that

$$993 \quad \mathcal{M}_{\mathcal{C}_{\tilde{f}}}(P) = \tilde{f}(P) + \text{Lap}(1/\varepsilon)$$

994 is  $\varepsilon'$ -bounded DP. Therefore, by the definition of  $\alpha$ , there exists a population  $P$  such that

$$995 \quad \sup_{P \in \mathcal{U}^n} \mathbb{E}[|\mathcal{M}_{\mathcal{C}_{\tilde{f}}}(P) - f(P)|^2] \geq \alpha.$$

996 Also

$$\begin{aligned} 997 \quad \alpha &\leq \mathbb{E}[|\mathcal{M}_{\mathcal{C}_{\tilde{f}}}(P) - f(P)|^2] = \mathbb{E}[|\tilde{f}(P) + \text{Lap}(1/\varepsilon) - f(P)|^2] \\ 998 &= \mathbb{E}[(\tilde{f}(P) - f(P))^2 + \text{Lap}(1/\varepsilon)(\tilde{f}(P) - f(P)) + (\text{Lap}(1/\varepsilon))^2] \\ 999 &= \mathbb{E}[|\tilde{f}(P) - f(P)|^2] + (1/\varepsilon)^2. \end{aligned}$$

1001 After a small amount of rearranging we arrive at the result. ◀

## 1002 **G** Experimental Setup

### 1003 **G.1** Figure 2

1004 In this experiment, recall that we let  $\mathcal{M}$  be an  $\varepsilon$ -DP exponential mechanism that returns a noisy count  
 1005 of the number of ones in a dataset from  $\{0, 1\}^n$ , and  $\Pr[\mathcal{M}(P) = i] \propto \exp(\varepsilon \cdot |i - \#\text{ones in } P|/2)$ .

## 23:28 Controlling Privacy Loss in Sampling Schemes

1006 Then, we let  $t$  be the following randomised function which determines the size of the sample. For a  
1007 given  $\gamma \in (0, 1)$ , we define  $t$  as follows.

$$1008 \quad t = \begin{cases} 2 & \text{w.p. } 1 - \gamma \\ n & \text{w.p. } \gamma \end{cases}$$

1009

1010 We define neighboring datasets  $P, P'$ , where  $P$  has  $n/2$  zeroes and  $n/2$  ones, and  $P'$  has  $n/2$  zeroes  
1011 and  $n/2 + 1$  ones. Let  $\mathcal{C}_t(P)$  (resp.  $\mathcal{C}_t(P')$ ) sample  $m \sim t$  records randomly without replacement  
1012 from  $P$  (resp.  $P'$ ). In Figure 2, we vary  $\gamma$  from  $10^{-25}$  and  $10^{-15}$  and compute the empirical privacy  
1013 loss  $\varepsilon'$  of  $(\mathcal{M} \circ \mathcal{C}_t)(P)$  and  $(\mathcal{M} \circ \mathcal{C}_t)(P')$ .