

Utility Cost of Formal Privacy for Releasing National Employer-Employee Statistics

Samuel Haney
Duke University
shaney@cs.duke.edu

Ashwin Machanavajjhala
Duke University
ashwin@cs.duke.edu

John M. Abowd
U.S. Census Bureau, U.S.A.
john.maron.abowd@census.gov

Matthew Graham
U.S. Census Bureau, U.S.A.
matthew.graham@census.gov

Mark Kutzbach
U.S. Census Bureau, U.S.A.
mark.j.kutzbach@census.gov

Lars Vilhuber
Cornell University
lars.vilhuber@cornell.edu

ABSTRACT

National statistical agencies around the world publish tabular summaries based on combined employer-employee (ER-EE) data. The privacy of both individuals and business establishments that feature in these data are protected by law in most countries. These data are currently released using a variety of statistical disclosure limitation (SDL) techniques that do not reveal the exact characteristics of particular employers and employees, but lack provable privacy guarantees limiting inferential disclosures.

In this work, we present novel algorithms for releasing tabular summaries of linked ER-EE data with formal, provable guarantees of privacy. We show that state-of-the-art differentially private algorithms add too much noise for the output to be useful. Instead, we identify the privacy requirements mandated by current interpretations of the relevant laws, and formalize them using the Pufferfish framework. We then develop new privacy definitions that are customized to ER-EE data and satisfy the statutory privacy requirements. We implement the experiments in this paper on production data gathered by the U.S. Census Bureau. An empirical evaluation of utility for these data shows that for reasonable values of the privacy-loss parameter $\epsilon \geq 1$, the additive error introduced by our provably private algorithms is comparable, and in some cases better, than the error introduced by existing SDL techniques that have no provable privacy guarantees. For some complex queries currently published, however, our algorithms do not have utility comparable to the existing traditional SDL algorithms. Those queries are fodder for future research.

Keywords

differential privacy, Pufferfish privacy, U.S. Census Bureau

1. INTRODUCTION

In this paper we present a case study in applying provably private algorithms for publishing tabular summaries of linked ER-EE data; i.e., data about business establishments and characteristics of

ACM acknowledges that this contribution was authored or co-authored by an employee, or contractor of the national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only. Permission to make digital or hard copies for personal or classroom use is granted. Copies must bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. To copy otherwise, distribute, republish, or post, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGMOD'17, May 14-19, 2017, Chicago, Illinois, USA

© 2017 ACM. ISBN 978-1-4503-4197-4/17/05...\$15.00

DOI: <http://dx.doi.org/10.1145/3035918.3035940>

their workforces. Such publications are used to compute national and local economic indicators, including job creation and destruction statistics. Canada, the United Kingdom, and most Eurostat countries have regularly published data from establishment-based surveys that directly measure employee characteristics. Statistics Sweden explicitly publishes such tabulations based on its Business Database (a frame containing lists of workplaces) combined with establishment-based surveys that measure employee characteristics. In the U.S., the Census Bureau publishes County Business Patterns (CBP) and Quarterly Workforce Indicators (QWI) using establishment frames that include characteristics of employees.

We focus on one establishment-based data product published by the U.S. Census Bureau's Longitudinal Employer-Household Dynamics Program (LEHD) [6]. The LEHD infrastructure files combine confidential census, survey and administrative records to tabulate and release public-use data called the LEHD Origin-Destination Employment Statistics (LODES) [13]. A LODES tabulation in 2011, for example, provides a snapshot of approximately 130 million jobs by workplace and residence location, as well as by a set of employer (industry sector and ownership type), employee (age, gender, race, ethnicity and education) and job characteristics [2].

These linked ER-EE data cannot be released as unaltered establishment-level microdata. They are subject to long-standing legal confidentiality protections that apply to both establishments and employees. Census Bureau publications must protect whether or not a specific individual is employed at a workplace and the reported count of employees at a workplace; however, the existence of one or more employers in a location, industrial sector, or ownership type does not require privacy protection.

The combination of detailed geography (at the level of a census block) with characteristics of establishments and employees results in *sparse* tabulations for which many cells have only a few contributing establishments and/or individuals. In addition, the outcome being tabulated, employment, is highly *right skewed* (i.e., has many large outlying values) at the establishment level. The combined effect of sparsity and skewness is the potential for re-identification attacks. Violating the statutory confidentiality pledge, for example by publishing data that permit re-identifying employees or inferring exact employer characteristics, can result in fines of up to \$250,000 and potential imprisonment for up to five years for those Census Bureau employees who authorize these publications.¹ Thus, privacy must be ensured as a matter of criminal law

¹These fines and prison terms change often, but since each data item is a separate violation, the ones cited in the text are really lower bounds. All authors on this paper take security training each year that includes the current values of these fines and prison terms.

with each employee and contractor of the Census Bureau bound by the law. With the stakes so high, we ask:

- “Can we develop algorithms for releasing such data that provably uphold the privacy requirements mandated by law?”, and
- “What is the loss of data usefulness (utility cost) when releasing summaries that provably ensure privacy?”

Statutory requirements obligate the statistical agency not to publish the raw microdata; however, they do not prescribe exactly how to publish the data [22]. Current publications of employment counts for detailed industries and geographies released by the Census Bureau and other agencies (including LODES) use a variety of confidentiality protection methods known collectively as *statistical disclosure limitation* (SDL) to limit potential re-identification attacks. These SDL methods are the *de facto* interpretations of the legally required confidentiality protections [21, 8]. For instance, in current LODES publications, workplace counts are protected by an input noise infusion system² that provably avoids exact disclosures of establishment level employment counts but has no guarantee that users are prevented from inferring actual establishment job counts to within some stated level of accuracy [10, 11, 27](we provide examples of vulnerabilities in Sec. 5.2). Moreover, due to the establishment-level skewness and the sparsity of some cells in the tables, these SDL techniques also limit the suitability of the published data for analyzing geographical and industrial characteristics of employers and employees in some cases. These limitations are the explicit utility loss from the input noise infusion SDL.

To overcome the aforementioned limitations of the current state of the art, we develop novel algorithms with provable privacy guarantees and measure the utility cost of such methods on the LODES dataset. Differential privacy [24] is the gold standard for provably private algorithms. It requires that the output of a private algorithm be insensitive to the presence of any single entity. A direct application of differentially private [25, 23] techniques to our data results in either insufficient privacy or poor utility when the entity is an establishment. Recent privacy frameworks generalizing differential privacy provide another avenue to address our problem. The Pufferfish framework [35] is a formal privacy framework that defines the privacy of an algorithm as the change in an adversary’s belief about a given set of secrets after seeing the output to the algorithm. The framework is provided with both the set of secrets (e.g. whether or not any individual is in the dataset in the case of differential privacy), and the assumptions about the adversary. Though we would ideally like a privacy guarantee that is assumption-free, seminal work of Kifer and Machanavajjhala [33] shows that such assumptions are required: one can not ensure privacy and utility simultaneously without making assumptions about an adversary’s beliefs about the data generation mechanism. For instance, it is known [33, 40, 41] that attacker-independent privacy notions like ϵ -differential privacy limit the ability of an attacker to learn properties of individual records in the input dataset if and only if the adversary believes records in the data are independent of one another. We use the Pufferfish framework to mathematically formulate privacy requirements based on current interpretations of relevant laws governing the release of LODES data, and develop novel privacy definitions and algorithms that meet these requirements. Our main deviation from differential privacy is a change to the information we want to keep private, not a change to the assumptions about

²In prior work [37], we developed an algorithm to protect employee residence locations that provably ensures probabilistic differential privacy. Residence locations are not a focus of this paper. In addition, other ER-EE products mentioned above do not contain information on employee residence locations.

the adversary under which disclosure is bounded. Our algorithms provably ensure that a privacy-loss limit is respected while achieving utility comparable to the current publication methodology.

Our contributions are as follows:

- (i) We formalize the privacy requirements for releasing ER-EE data using the Pufferfish framework [34, 35] (Sec. 4) to ensure that a strongly informed attacker cannot (a) infer whether an employee held a job, and (b) learn establishment sizes to within a pre-specified multiplicative factor from the output of a data release. The requirements are developed by reviewing the current interpretations of the legal regulations that the national statistical agencies are mandated to follow when releasing such data. For answering queries on employer attributes, we make exactly the same assumptions about the adversary that must be made under differential privacy to ensure that inferential disclosure is bounded. To answer queries over both employer and employee attributes, we must make additional assumptions about the adversary.
- (ii) Our data can be represented as a bipartite graph between employers and employees, each edge representing a job. Thus, edge-[31] and node-differentially private techniques [18, 20, 32] can be applied to protect the data (Sec. 6). We show that edge-differential privacy, which “hides” the presence of a single job, does not satisfy our privacy requirements for establishments. Node-differential privacy satisfies the privacy requirements but results in poor utility.
- (iii) We formulate novel attacker-independent privacy definitions that satisfy both the Pufferfish framework and the legal framework. These provably limit an informed attacker’s ability to make inferences about employees and employers (Sec. 7). Thus, our algorithms conform to current interpretations of the confidentiality laws governing these data.
- (iv) We develop algorithms for releasing counts that satisfy our new privacy definitions. We prove analytical bounds on the errors for each algorithm (Sec. 8 and 9). Our algorithms use an extension of the smooth sensitivity framework [38].
- (v) We empirically evaluate the utility cost of provable privacy using the production data from an ER-EE dataset maintained by the U.S. Census Bureau (Sec. 10). For releasing tabular summaries and rankings, we compare the error introduced by our algorithms to the error produced by the existing protection scheme. We show that we can release tabular summaries of establishment characteristics with additive error that is comparable (within a factor of 3) and in some cases smaller than the error introduced by the current SDL techniques for reasonable values of the privacy-loss parameters. On the other hand, node-differentially private techniques incur an additive error that is at least 10 times larger than that of SDL techniques (at $\epsilon = 4$), and the error does not decrease significantly when ϵ is increased. For counts and rankings, the relative error of our new algorithms is comparable to that of the existing SDL methods. Since tabular summaries of employment by establishment characteristics alone, such as in the Business Database statistics (Sweden) or the CBP (U.S.A.), are important publications, it is significant that we are able to achieve provable privacy protection with relatively little sacrifice in data utility as compared to existing methods.
- (vi) For tabulations involving both establishment and employee characteristics, under a weaker adversary model, we are able to show competitive error for releasing single queries and rankings. Our algorithms experience larger errors (within a factor of 100 compared to existing methods) for releasing tabular summaries involving multiple employee characteristics. Understanding whether the magnitude of this error is fundamental to provably private algorithms or whether better algorithm design could lower such errors is an avenue for future work.
- (vii) Our techniques are applicable to virtually all establishment-

based products released by statistical agencies for national production and employment statistics, including those produced by the Census Bureau, Bureau of Labor Statistics and Bureau of Economic Analysis in the U.S.A. as well as the inputs to the national income and product accounts, which are published for more than 200 countries. More broadly, our techniques are applicable to data derived from multiple entities connected by a bipartite graph such as user-purchases, reviewers-ratings, etc.

2. PRELIMINARIES

Database and Queries Let D be a table of records with schema (A_1, \dots, A_k) . The domain of each attribute A_i is denoted $\text{dom}(A_i)$. For a set of attributes $V = \{A_{i_1}, \dots, A_{i_m}\}$, let $\text{dom}(V)$ represent the multidimensional domain $\times_{A \in V} \text{dom}(A)$. For each record t in the table, we let $t[A_i] \in \text{dom}(A_i)$ be the value of attribute A_i . Let $n = |D|$ denote the size of the table; i.e., D has n records. A database with schema (S_1, \dots, S_m) is a collection of tables (D_1, \dots, D_m) , where D_i has schema S_i .

We will consider *marginal queries* over tables in this paper.

DEFINITION 2.1 (MARGINAL QUERY). *The marginal query $q_V(D)$ is defined as a vector of $|\text{dom}(V)|$ counts, one for each cell $\mathbf{v} = (v_1, \dots, v_m) \in \text{dom}(V)$. The count corresponding to cell \mathbf{v} , denoted by $q_V(D, \mathbf{v})$ is*

$$|\{t \in D \mid t[A_{i_1}] = v_1 \wedge \dots \wedge t[A_{i_m}] = v_m\}| \quad (1)$$

$q_{\emptyset}(D)$ returns a single cell whose count is the size of the table.

The marginal query can be succinctly expressed in SQL as:

```
Select Count(*) From D Group By  $A_{i_1}, \dots, A_{i_m}$ 
```

Differential Privacy A mechanism or algorithm is differentially private if its output is not significantly affected by the presence or absence of a single record from the input table. Let D and D' be tables that differ in the presence of a single record; i.e., $|(D \setminus D') \cup (D' \setminus D)| = 1$. We call such tables *neighbors*.

DEFINITION 2.2 ((ϵ, δ)-DIFFERENTIAL PRIVACY [26]). *Let \mathcal{M} be a randomized algorithm. Let the tables D and D' be neighbors with the same schema. Then \mathcal{M} satisfies (ϵ, δ) -differential privacy if for all D and D' and for all $S \subseteq \text{range}(\mathcal{M})$,*

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \cdot \Pr[\mathcal{M}(D') \in S] + \delta,$$

where δ allows for the ratio of probabilities to be unbounded with a small failure rate. Values of δ that are $\Omega(1/n)$ should be avoided, since algorithms that release more than a constant number of unperturbed records from the database satisfy such a definition. When $\delta = 0$, we refer to the condition as ϵ -differential privacy.

Queries over tables can be answered while satisfying differential privacy by adding noise that is related to the *sensitivity* of the query.

DEFINITION 2.3 (SENSITIVITY). *Let \mathcal{I} denote the set of all tables with a given schema. Let $q : \mathcal{I} \rightarrow \mathbb{R}^d$ be a query function on tables from that set that outputs a vector of d real numbers. The sensitivity of q , denoted Δ_q , is*

$$\Delta_q = \max_{D, D' \text{ neighbors} \in \mathcal{I}} \|q(D) - q(D')\|_1.$$

The Laplace mechanism is a commonly used ϵ -differentially private technique.

DEFINITION 2.4 (LAPLACE MECHANISM [26]). *Let $q : \mathcal{I} \rightarrow \mathbb{R}^d$ be a query on a table. Let $\eta \sim \text{Lap}(\lambda)$ denote a random variable drawn from the Laplace distribution with pdf $\Pr[\eta = x] \propto$*

$e^{-|x|/\lambda}$. *The algorithm which returns $\tilde{q}(D) = q(D) + \eta^d$ satisfies ϵ -differential privacy, where η^d is a vector of d independently drawn Laplace random variables.*

DEFINITION 2.5 (EXPECTED L_p ERROR). *Let $q : \mathcal{I} \rightarrow \mathbb{R}^d$ be a query on a table, and $\tilde{q}(D)$ be the noisy answer returned by an algorithm. The expected L_p error of the algorithm is:*

$$\mathbb{E}(\|q(D) - \tilde{q}(D)\|_p) \quad (2)$$

where $\|x\|_p$ is the L_p norm, and expectation is over the randomness of the algorithm.

We use expected L_1 error to quantify the utility of algorithms.

The privacy loss increases when multiple queries are answered on the database, and we reason about this loss for differential privacy using the following composition rule:

THEOREM 2.1 (SEQUENTIAL COMPOSITION). *Let \mathcal{M}_1 and \mathcal{M}_2 be ϵ_1 - and ϵ_2 -differentially private algorithms. Releasing the outputs of $\mathcal{M}_1(D)$ and $\mathcal{M}_2(D)$ on the same input D results in $(\epsilon_1 + \epsilon_2)$ -differential privacy.*

Thus, the privacy-loss parameter is often called the *privacy budget* – the analyst is allowed to pose multiple queries as long as the total privacy loss from answering all queries is no greater than ϵ . In other words, a privacy-loss equal to ϵ exhausts the privacy budget.

3. THE DATASET

The LODES data are produced from the LEHD infrastructure files, which are composed of administrative records, census and survey data focused on the labor market, worker, and firm statistics. Confidential Census Bureau data, state Unemployment Insurance sources, Internal Revenue Service records, Social Security Administration records, and federal Office of Personnel Management records provide information on employment location and industry for jobs and residential location and demographic characteristics for workers. These data form the basis of the LODES tabulations. LODES are published as an annual cross-section of jobs held on April 1 of each year from 2002 onwards.

3.1 Table and Database Structure

The LODES data are organized as a relation with three database tables – **Job**, **Worker** and **Workplace** (see [2] for a complete description). The **Workplace** table contains one record per establishment and describes the following attributes – NAICS code (denoting industrial sector in which the establishment operates), ownership (public/private), geography (the Census block where the establishment is located). The **Worker** table contains one record for each individual working in any establishment at that point of time. Worker attributes include age, sex, race, ethnicity, and education. Finally, the **Job** table contains pairs (w, i) of worker and workplace IDs denoting that worker i works at establishment w . We assume each worker has exactly one job (although LODES allows for queries that include secondary jobs). We have not documented some of the attributes that do not feature in our queries.

We support *marginal queries* that output counts of employment of the current year’s cross section of jobs over workplaces where those establishments have been stratified by subsets of the available characteristics of employers and workers. Let V_I denote a subset of the worker attributes, and V_W denote a subset of workplace attributes. Let D denote the universal relation constructed by joining the **Job** table with the **Worker** and **Workplace** tables using the worker and workplace IDs, respectively. We call this the **WorkerFull** table as it contains one record for every record in **Worker**

(each worker has exactly one job). The records in `WorkerFull` contain all the worker and workplace attributes. Records that share the same workplace ID represent the workforce of that establishment.

Every cell $(\mathbf{v}_I, \mathbf{v}_W)$ in the marginal query $q_{V_I \cup V_W}(D)$ (as in Definition 2.1) represents the number of workers matching the criteria in \mathbf{v}_I who work in establishments matching the criteria in \mathbf{v}_W . The output is employment counts stratified by employee attributes in V_I and workplace attributes V_W .

3.2 Usage Scenarios

Researchers and analysts working with these data request answers to one or more marginal queries and use these summary tabulations as inputs to further data analysis. In this paper, we focus on analyzing the privacy and utility tradeoffs from releasing the output to one marginal query. Analyzing the privacy of multiple queries is a straightforward application of Theorem 2.1, since our privacy definitions also satisfy sequential composition like differential privacy. The examples of queries and protection parameters provided here are meant to illustrate our methods. We have made no recommendation to the Census Bureau regarding the set of queries, methods, and privacy parameters to use for a full implementation.

We highlight two application scenarios that use summary tables computed from LODES.

Resource Allocation: LODES data are used for a number of planning and development purposes. We consider an example in disaster assistance that motivates our use of L_1 error, as a measure for data accuracy. Experiments in Section 10 will explore the effect of privacy parameters on L_1 error.

The Federal Emergency Management Agency (FEMA) uses Census population count data to evaluate requests for a disaster declaration, which permit federal cost sharing of 75%, or even up to 90% in the most severe cases. Under the Robert T. Stafford Disaster Relief and Emergency Assistance Act, FEMA uses the cost estimate of \$3.50 per capita as an indicator that a disaster is of such size that it might warrant Federal assistance.³ For determinations, FEMA divides the Preliminary Damage Assessment by the 2010 Census population count for one or more counties. FEMA participates in a Census Bureau program, OnTheMap for Emergency Management [3], to release queries of LODES and population data concurrently with the designation of emergencies and disasters.⁴

If the threshold were applied to jobs rather than population, errors in the count would affect the hypothetical threshold for a disaster declaration for that area. Positive count errors would result in a higher damage threshold and would require a larger magnitude disaster for assistance, while negative errors would imply the opposite. Both positive and negative errors could result in the misallocation of funds relative to the intent of such a program, with each job in error having a net social cost of \$3.50.

Rankings: Users are often interested in comparing counts across a list of units. As a well-known example, Forbes magazine regularly publishes rank-ordered national lists of cities by various attributes within city-size classes. The OnTheMap web tool [4] invites users to rank the LODES job counts of a series of areas from largest to smallest. A data user may specify a domain of comparison by selecting a standard geography (e.g. state, Congressional District, metropolitan area) or hand drawing a polygon. Within the selection area, the user can perform an Area Comparison Analysis for a work area and further specify the types of areas to compare from

³For the present analysis, we use the county level threshold of \$3.50 per capita, which applies to major disasters declared on or after October 1, 2013. See [9].

⁴Note that we are not making any statement about the appropriate factors to consider for disaster declarations.

another list of standard geographies. For example, a business might be interested in the ranked order of Places (e.g., cities, towns, and Census Designated Places) by job count, within a state, for deciding where to open a new establishment. Given these parameters, OnTheMap returns the ranked list of areas, or Places in the example, in descending order by work area job count. Exercises in Section 10 will use Spearman’s rank-order correlation to explore the accuracy of ranked lists produced from LODES data.

4. PRIVACY REQUIREMENTS

We derive our privacy requirements from relevant U.S. laws. We first discuss the general privacy requirements, then formalize them using the Pufferfish framework.

4.1 Generic Legal Environment

Information on workers and firms is protected by several sections of the U.S. Code. Title 13 Section 9 [1, 5] mandates confidentiality protections for individual and business information collected by the Census Bureau. Under Title 13, the Census Bureau may not “make any publication whereby the data furnished by any particular establishment or individual under this title can be identified.” The Disclosure Review Board (DRB) at the Census Bureau approves the release of data that, in its view, satisfy these statutory confidentiality protection requirements.

The LEHD Infrastructure Files are derived from both employer and job (ER-EE) level data, and, thus, Title 13 Section 9 protections apply. Over the years, the DRB has interpreted the statute to require the following set of protections.

- The existence of a job held by a particular individual is confidential and must not be disclosed.
- The existence of an employer business as well as its type (or sector) and location is not confidential.⁵
- The data on the operations of a particular business must be protected. In our context, that means that characteristics of an establishment’s workforce (e.g., total employment and all disaggregations like number of female employees of age 20-25) must be protected.

Appendix A discusses how these statutory confidentiality requirements have been interpreted over the past half century, and the development of SDL techniques that uphold them. However, none of the prior interpretations or techniques have attempted to prove formal statements about the privacy that is guaranteed to individuals or businesses. We show (Section 5) that the SDL techniques currently used to publish ER-EE data may allow an informed attacker to infer confidential properties like the exact total employment of an establishment, or whether a certain employee is employed by a specific employer. We first present our formalization of the privacy requirements.

4.2 Formal Privacy Desiderata

We propose to design algorithms for releasing counts from ER-EE data that can provide provable guarantees of the following privacy desiderata. Our discussion focuses on formulating privacy requirements for a *one-time release* by the statistical agency using an algorithm denoted by A . Our privacy notions handle multiple releases through composition rules. We model the requirements based on the Pufferfish privacy framework [34, 35].

Informed Attacker: National statistical agencies are concerned about two kinds of attackers – uninformed and informed. Unin-

⁵This principle is also the basis of confidentiality protections for the Economic Census and the County Business Patterns [7].

formed attackers can access the output of the algorithm A , but may not possess detailed background knowledge about specific individuals and establishments in the data. Informed attackers are more powerful. They possess specific knowledge about individual employees or employers, or statistics about those in the dataset. Examples of such attackers include a group of employees who would like to determine a private attribute of their co-worker, or one (or more) employer(s) attempting to learn detailed statistics about a competing employer. Our goal is to ensure the confidentiality of employer and employee characteristics from such attackers.

We assume the adversary knows the set of all establishments (say \mathcal{E}), and their public attributes (location, industry code and ownership). The attacker also knows the universe of all workers U . Each worker $w \in U$ has a set of private attributes $A_1 \dots A_k$ (like age and sex). We add another attribute with domain $\mathcal{E} \cup \perp$ that represents whether w works in one of the establishments in \mathcal{E} , or not.

For each employee w , the attacker’s belief is defined as π_w , a probability distribution over all the values in $\mathcal{T} = (\mathcal{E} \cup \perp) \times A_1 \times A_2 \times \dots \times A_k$. $\theta = \prod_{w \in U} \pi_w$ represents the adversary’s belief about all employees in the universe U . That is, the adversary possesses no knowledge correlating employees. We denote by $\Theta = \{\theta\}$, the set of all possible adversarial beliefs that assume no correlations between employees and between employers. Nevertheless, Θ includes informed attackers who may know exact information about all but one employee, and those who know exact information about all but one employer. We note that Θ contains very strong attackers. Algorithms that can provably protect against such attackers while ensuring error comparable to current SDL techniques would underscore the possibility that provable privacy could be achieved at low utility cost.

We distinguish a subset of attackers $\Theta_{weak} \subset \Theta$ as *weak* attackers. Weak attackers have no prior knowledge over worker attributes – i.e., all workers are the same in their eyes. The weak attacker may still have the same detailed knowledge about establishments as our general attacker. We capture a weak adversary by requiring that the prior for each worker π_w be a product of $\pi_{(1,e)}$ (worker independent prior over establishments), and $\pi_{(2,w)}$ (a uniform prior over all worker attributes). We use these definitions to define a weaker privacy notion.

What should we protect? We now specify which properties of the data we need to protect against such adversaries.

1. No re-identification of individuals: We would like to ensure that adversaries do not learn too much additional information about any single employee in the dataset when an algorithm A operates on the dataset D . In particular, they should not be able to determine (i) whether or not an employee is in or out of the dataset (\perp versus not), (ii) whether or not an employee works at a specific (type of) employer (E versus $\mathcal{E} - E$, where $E \subseteq \mathcal{E}$), and (iii) whether or not the employee has certain characteristics (e.g., Hispanic with age greater than 35).

We formalize this as follows. For any pair of values $a, b \in \mathcal{T}$, we require that the ratio of the adversary’s posterior odds (after seeing the output $A(D)$) that a worker record takes the value $w = a$ vs $w = b$ to the adversary’s prior odds that $w = a$ vs $w = b$ be bounded at a known level. That is, we want to bound the Bayes factor: the ratio of the posterior odds to the prior odds, and this bound is the privacy-loss budget.

DEFINITION 4.1 (EMPLOYEE PRIVACY REQUIREMENT). For randomized algorithm A , if for some $\epsilon \in (0, \infty)$, and for every employee $w \in U$, for every adversary $\theta \in \Theta$, for every $a, b \in \mathcal{T}$ such that $Pr_\theta[w = a] > 0$ and $Pr_\theta[w = b] > 0$, and for every

output $\omega \in range(A)$:

$$\log \left(\frac{Pr_{\theta,A}[w = a | A(D) = \omega]}{Pr_{\theta,A}[w = b | A(D) = \omega]} \bigg/ \frac{Pr_\theta[w = a]}{Pr_\theta[w = b]} \right) \leq \epsilon \quad (3)$$

Then the algorithm A protects employees against informed attackers at privacy-loss level ϵ .

Definition 4.1 bounds the logarithm of the maximum Bayes factor an informed attacker can achieve. This implies, as a consequence of the general bound on privacy loss, that an informed attacker can’t learn any property of a worker record with probability 1 after seeing the output of the algorithms unless the attacker already knew that fact, as reflected in his prior odds.

2. No precise inference of establishment size: An informed attacker should not *infer* the total employment of a single establishment to within a multiplicative factor of α . We do not require stronger privacy of the form “presence of an establishment must not be inferred,” since (a) the existence of an employer establishment is considered public knowledge, (b) the data are an enumeration of all employer establishments, and (c) whether or not an establishment is big or small is well known. This requirement balances the legal need for protecting the operations of a business with widespread knowledge of approximate employment sizes of establishments.

We can formalize the employer-size privacy requirement as follows. For any establishment e , let $|e|$ denote the random variable representing the number of workers employed at e . We define the requirement for both informed and weak adversaries.

DEFINITION 4.2 (EMPLOYER SIZE REQUIREMENT). Let e be any establishment in \mathcal{E} . A randomized algorithm A protects establishment size against an informed attacker at privacy level (ϵ, α) if, for every informed attacker $\theta \in \Theta$, for every pair of numbers x, y , and for every output of the algorithm $\omega \in range(A)$,

$$\left| \log \left(\frac{Pr_{\theta,A}[|e| = x | A(D) = \omega]}{Pr_{\theta,A}[|e| = y | A(D) = \omega]} \bigg/ \frac{Pr_\theta[|e| = x]}{Pr_\theta[|e| = y]} \right) \right| \leq \epsilon \quad (4)$$

whenever $x \leq y \leq \lceil (1 + \alpha)x \rceil$ and $Pr_\theta[w = x], Pr_\theta[w = y] > 0$. We say that an algorithm weakly protects establishments against an informed attacker if the condition above holds for all $\theta \in \Theta_{weak}$.

As in Definition 4.1, this definition bounds the maximum Bayes factor the informed attacker can learn within the universe of allowable data tables. Unlike the case of individuals, Definition 4.2 does allow an adversary to learn about the gross size of an employer establishment.

3. No precise inference of establishment shape: An informed attacker cannot precisely infer the composition of a single establishment’s workforce (e.g., the fraction of males who have a bachelor’s degree or the fraction with Hispanic ethnicity). We call the distribution of an establishment’s workforce based on worker characteristics its *shape*. One can think of this requirement as protecting the distribution of characteristics of the workforce, whereas the previous requirement protected the magnitude of each characteristic. We believe this shape requirement implements the legally mandated confidentiality of an establishment’s operating characteristics. The definition bounds the maximum Bayes factor that the informed adversary can learn within the set of allowable inputs.

DEFINITION 4.3 (EMPLOYER SHAPE REQUIREMENT). Let e be any establishment in \mathcal{E} . Let $e_{\mathcal{X}}$ denote the subset of employees working at e who have values in $\mathcal{X} \subset A_1 \times \dots \times A_k$. A randomized algorithm A protects establishment shape against an informed attacker at a privacy level of (ϵ, α) , if for every informed attacker $\theta \in \Theta$, for every property of a worker record \mathcal{X} , for every pair of

Name	Does this satisfy requirement on		
	Individuals	Emp. Size	Emp. Shape
Input Noise Infusion (Sec. 5)	No	No	No
Differential Privacy (individuals, Sec. 6)	Yes	No	No
Differential Privacy (establishments, Sec. 6)	Yes	Yes	Yes
ER-EE-privacy (Sec. 7)	Yes	Yes	Yes
Weak ER-EE privacy (Sec. 7)	Yes	Yes*	Yes

Table 1: Privacy definitions and requirements they satisfy.
* Privacy requirement is satisfied under weak adversaries.

numbers $0 < p \leq q \leq \min(1, (1 + \alpha)p)$, for every output of the algorithm $\omega \in \text{range}(A)$, and for every natural number z ,

$$\left| \log \left(\frac{\Pr_{\theta, A}[|e_{\mathcal{X}}|/|e| = p, |e| = z | A(D) = \omega]}{\Pr_{\theta, A}[|e_{\mathcal{X}}|/|e| = q, |e| = z | A(D) = \omega]} \right) \right| \leq \epsilon \quad (5)$$

$$\frac{\Pr_{\theta}[|e_{\mathcal{X}}|/|e| = p, |e| = z]}{\Pr_{\theta}[|e_{\mathcal{X}}|/|e| = q, |e| = z]} \Bigg| \leq \epsilon$$

whenever $\Pr_{\theta} \left[\frac{|e_{\mathcal{X}}|}{|e|} = p, |e| = z \right], \Pr_{\theta} \left[\frac{|e_{\mathcal{X}}|}{|e|} = q, |e| = z \right] > 0$.

Table 1 summarizes whether certain SDL and formal privacy methods satisfy our privacy desiderata. The following sections provide definitions and analysis required to interpret this table.

5. PROTECTION BY CURRENT SDL

We discuss how LODS data are protected using the current SDL techniques, and how they avoid exact disclosures in Sec. 5.1. Examples of potential inference attacks are outlined in Section 5.2.

5.1 Input Noise Infusion

A popular technique for protecting ER-EE data is *input noise infusion* [10, 12], where the database is perturbed before answering queries. Every establishment w is assigned a unique distortion factor f_w , bounded away from 1. The unique, time-invariant, confidential distortion factor f_w , is within the union of the ranges $[1 - t, 1 - s] \cup [1 + s, 1 + t]$. The parameters $0 < s < t$ are kept confidential in order to limit inference attacks. Zero counts are left unmodified.

More formally, consider a table `WorkplaceFull` that has one row per workplace w , as well as a histogram $\mathbf{h}(w)$ of counts of workers employed at w cross-tabulated over all combinations of worker attributes. Let c denote one of the cells (combinations of worker attributes like males, age 16-18, Hispanic, etc.), and let $\mathbf{h}(w, c)$ denote the count for workplace w in cell c . Counts in $\mathbf{h}(w)$ are perturbed to get $\mathbf{h}^*(w)$ as follows.

$$\mathbf{h}^*(w, c) = f_w \cdot \mathbf{h}(w, c) \quad (6)$$

To limit re-identification of individual workers, additional output perturbation is employed for small counts. Specifically, when a marginal query q_V is posed to the system (say employments counts tabulated by age, sex, and place), both the true answer $q_V(D)$ as well as a noise infused answer $q_V^*(D)$ are computed. The latter is constructed by adding up appropriate counts from $\mathbf{h}^*(w)$ for the establishments that satisfy the workplace criteria. If for a cell \mathbf{v} in the output, the true count $q_V(D, \mathbf{v})$ lies within $(0, S)$, then the noise infused answer $q_V^*(D, \mathbf{v})$ is replaced by a sample drawn from a posterior predictive distribution that always outputs integers

$1, \dots, \lfloor S \rfloor$. The small cell limit S is set to 2.5 for our dataset. Note that zero counts are unperturbed.

Privacy Properties

- *No Exact Disclosures about Establishments*: As a direct consequence of the gap around 1 in the distortion factor $f_w \in [1 - t, 1 - s] \cup [1 + s, 1 + t]$, an establishment’s actual employment count is never used in any computations that produce tabular summaries. Hence, even if some cell count in a marginal query contains only one establishment, its employment count is not exactly revealed. The statutory requirement not to publish exact data about establishments is fulfilled by using the distortion factors s and t .
- *No Re-identification of Employees without Background knowledge*: Given the output of a single marginal query, an adversary can not re-identify the presence of a specific worker without additional background knowledge. This is ensured by replacing small counts using draws from a different distribution.

Nevertheless, the aforementioned scheme is vulnerable to inference attacks, especially in the presence of background knowledge, as discussed next.

5.2 SDL Vulnerabilities

The following two properties of the input noise infusion scheme allow inference about individuals, establishment sizes and establishment shapes:

- The same distortion factor f_w is used to perturb all the cells counts $\mathbf{h}(w)$ for an establishment w .
- If $\mathbf{h}(w) = 0$, then $\mathbf{h}^*(w) = 0$.

We first note that the privacy requirement on establishment shape (Definition 4.3) is not satisfied. Consider a marginal query $q_{V_I \cup V_W}$, where V_I are attributes of the employee and V_W are attributes of the establishment. Suppose there is one combination $\mathbf{v}_W \in \text{dom}(V_W)$ such that exactly one workplace w fits that criterion.⁶ Thus, counts output by the marginal query for all cells (\mathbf{v}_W, c) , for all $c \in \text{dom}(V_I)$ would represent employment counts for a single workplace. Whenever all these cell counts are greater than the small cell limit S , they are precisely the true count multiplied by the same (but unknown) noise factor f_w . This reveals the shape exactly.

Next, let us consider the privacy requirement on establishment size (Def. 4.2). Consider a combination $\mathbf{v}_W \in \text{dom}(V_W)$ such that exactly one workplace w fits that criterion. Additionally, suppose the attacker knows one of the cell counts (\mathbf{v}_W, c) truthfully (say, the attacker knows there are a 100 males, age 20-25). If this count is greater than the small cell limit S , then the adversary can reconstruct the noise factor f_w . Knowledge of f_w and the exact shape allows the attacker to reconstruct the counts in all other cells as well as the total size of the establishment’s workforce.

The agencies that use input noise distortion as a method of SDL understand this attack. It is not currently considered a violation of the relevant data protection statutes because the exact disclosure occurs only when one of the counts $\mathbf{h}(w, c)$ is known exactly for some cell c . The agencies assume that the only users who possess exact information on c are employees of a business that reported the data. If those employees are obligated to keep such information confidential, then it is the employer’s duty to prevent the attack or prosecute the attacker. If they are not, then the data item itself is no

⁶The number of establishments in a cell are not published for the dataset we consider in this paper. However, there are combinations of V_W that contain only one workplace, and an adversary could know this.

longer confidential, and statutory protection doesn't apply because the employer released the value.

Next, we show that individual employees could be re-identified by informed attackers (thus violating Definition 4.1). Suppose a marginal query $q_{V_I \cup V_W}$ with one combination $\mathbf{v}_W \in \text{dom}(V_W)$ fits exactly one workplace w . Additionally, suppose an adversary knows that there is only one employee in w with a college degree. If V_W contains the education attribute, then the only cells that correspond to having a college degree in $\mathbf{h}(w, c)$ with positive counts are those that correspond to the true values of the other attributes for that employee. Since zero counts are preserved in the current SDL, the attacker can infer the other true attributes for this employee by looking at the published counts. Current publications of the ER-EE data we consider are vulnerable to this attack, but can be thwarted by the algorithms we propose.

6. APPLYING DIFFERENTIAL PRIVACY

In this section we directly apply differential privacy to our problem by considering the entities in our problem (employers and employees) to be nodes in a *bipartite* graph connected by edges that represent jobs. Thus, work applying differential privacy to graphs can now be brought to bear on our problem [18, 20, 32, 31].

Two standard variants considered in the context of graphs are *edge* and *node* differential privacy. Edge differential privacy considers neighboring graphs that differ in the presence of a single edge. In our context, that corresponds to adding or removing a single job (hence, worker) from our database. We can show that this definition is sufficient to satisfy the employee privacy requirement (Definition 4.1). However, edge differential privacy does not ensure privacy of establishments (Definitions 4.2 and 4.3; see Claim B.1 in Appendix B). For instance, under this definition an adversary is allowed to compute the total employment count at a single establishment by adding to the true count noise drawn from $\text{Lap}(1/\epsilon)$. We can show that the noise added is at most $\frac{\log(1/p)}{\epsilon}$ with probability $1 - p$ (i.e., at most 5 for $\epsilon = 1$ and $p = 0.01$). Knowing that the total employment in an establishment is $10,000 \pm 5$ is almost as good as knowing the true count, and this inference continues to improve as the establishment size increases; hence, it cannot respect a fixed privacy-loss budget for establishments.

Node differential privacy considers neighboring graphs that differ in the presence of a single node and all the edges incident to it. In our context that corresponds to removing or *adding* a single employer along with all the workers who are employed at this employer. Node differential privacy is much stronger, and in the context of our problem will satisfy the employee and employer privacy requirements (Definitions 4.1, 4.2 and 4.3). However, this comes at a huge cost to utility.

Since there is no *a priori* bound on the number of edges incident on each node (other than the number of nodes in the graph), the Laplace mechanism is inapplicable for edge counting queries under node-differential privacy. Hence, an alternate method to perturb counts is *projection*. Projection techniques [18, 20, 32] modify the graph by adding or deleting edges and nodes until the maximum degree of a node is bounded by a small number θ . Edge counting queries on this bounded-degree graph have bounded sensitivity (of θ) and hence can be answered by adding noise from $\text{Laplace}(\theta/\epsilon)$. For instance, the truncation method [32] projects the graph by removing nodes until all nodes have degree less than θ . When θ is small, a significant fraction of the nodes with degree larger than θ will be excluded from the count query. In our context, this would severely distort the characteristics of large employers. Preserving properties of these establishments is important for economic stud-

ies. When θ is large, the noise added would also be very large, adding too much noise to cells with small businesses to be useful.

Using this technique on our data with $\theta = 1000$ in the context of ER-EE data results in removing all establishments with more than 1,000 employees; between 740 and 815 establishments would be removed.⁷ Moreover, in a tabular summary of counts by place, industry and ownership, over 93% of the counts have a total count less than a 1000. Adding Laplace noise with sensitivity of a 1000 (even with $\epsilon = 1$) would result in expected noise greater than the cell counts. We present further empirical evidence of the error incurred by this technique in Section 10.

7. FORMAL PRIVACY DEFINITION

We present two refinements of our privacy definitions that ensure protection against informed and weak adversaries, respectively, as we have defined them. Our new definitions are in Section 7.1; their privacy semantics are in Section 7.2. Our formal privacy definitions ensure that the requirements in Section 4 are satisfied.

7.1 Privacy for Employer-Employee Data

We denote the universe of establishments as \mathcal{E} and the universe of workers as \mathcal{U} .

DEFINITION 7.1 (STRONG α -NEIGHBORS). *Let D and D' be two ER-EE tables that differ in the employment attribute of exactly one record (say corresponding to establishment e). Let E denote the set of workers employed at e in D , and E' denote the set of workers employed at e in D' . Then D and D' are strong α -neighbors if $E \subseteq E'$, and $|E| \leq |E'| \leq \max((1 + \alpha)|E|, |E| + 1)$*

We refine our privacy definition using definition 7.1:

DEFINITION 7.2 ((α, ϵ) -ER-EE PRIVACY). *A randomized algorithm \mathcal{M} is said to satisfy (α, ϵ) -ER-EE Privacy, if for every set of outputs $S \subseteq \text{range}(\mathcal{M})$, and every pair of strong α -Neighbors D and D' , we have*

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S]$$

First, note that every pair of neighboring ER-EE tables must differ in the presence or absence of at least one worker. Next, note that neighboring tables do not differ in either the number of establishments or the values of their public attributes. Definition 7.1 bounds changes in a subset of the workforce of an establishment by α times the total workforce.

THEOREM 7.1. *Let \mathcal{M} be an algorithm satisfying (α, ϵ) -ER-EE privacy. Then, \mathcal{M} satisfies the individual privacy requirement at privacy level ϵ , and the establishment size and shape requirements at privacy level (ϵ, α) .*

When releasing counts over both establishment and worker attributes, this definition 7.2 is too strong to provide useful results. We further refine our definitions.

DEFINITION 7.3 (WEAK α -NEIGHBORS). *Let D and D' be two ER-EE tables such that they differ in the employment attribute of exactly one record (say corresponding to establishment e). Let $\phi : \mathcal{U} \rightarrow \{0, 1\}$ be any property of a worker record, and for any $S \subset \mathcal{U}$, let $\phi(S) = \sum_{r \in S} \phi(r)$. Let E denote the set of workers employed at e in D , and E' denote the set of workers employed at e in D' . D and D' are called weak α -neighbors if for every ϕ*

$$\phi(E) \leq \phi(E') \leq \max((1 + \alpha)\phi(E), \phi(E) + 1) \quad (7)$$

⁷The number of establishments with size > 1000 is a sensitive count. This count was computed using the Laplace mechanism with $\epsilon = 0.1$. The reported range is the 95% confidence interval.

Note that any property of the workforce can be represented using the function ϕ . The total employment count can be represented by the constant function that always outputs 1 for any record. The property “females with a college degree” can be represented by a ϕ that returns 1 for the records satisfying that property. Definition 7.3 bounds changes in every subset of the workforce corresponding to some of attribute values proportionally, by a factor of $(1 + \alpha)$. This neighboring definition can be used to give a weaker privacy notion.

DEFINITION 7.4 (WEAK (α, ϵ) -ER-EE PRIVACY). *A randomized algorithm \mathcal{M} is said to satisfy weak (α, ϵ) -ER-EE Privacy, if for every set of outputs $S \subseteq \text{range}(\mathcal{M})$, and every pair of weak α -Neighbors D and D' , we have*

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S]$$

THEOREM 7.2. *Let A be an algorithm satisfying weak (α, ϵ) -ER-EE privacy. Then, A satisfies the individual privacy requirement at privacy level ϵ and the establishment shape requirement at level (ϵ, α) . A satisfies the establishment size requirement at level (ϵ, α) for weak adversaries.*

The difference between the strong and weak variants of the requirement is the following. Suppose the attacker knows there are at least Δ 19 year-old employees in an establishment, and knows the exact counts of employees by age for all other ages (totaling to $x - \Delta$). Thus, the attacker’s only uncertainty is about the number of 19 year-olds. Then, Definition 7.2 requires that the attacker should not be able to distinguish between whether the number of 19 year-old employees is Δ or Δ' for all $\Delta \leq \Delta' \leq \Delta + \alpha \cdot x$. While one might expect few 19 year-olds in an establishment, algorithms satisfying Definition 7.2 will not be able to release that fact unless α is very small.

Under Definition 7.3 the attacker should not be able to distinguish between whether the number of 19 year-old employees is Δ or Δ' for all $\Delta \leq \Delta' \leq (1 + \alpha)\Delta$.

Keeping α fixed, larger ϵ values result in more privacy loss, since adversaries can better distinguish neighboring databases. On the other hand, when keeping ϵ fixed, larger α values result in less privacy loss.

7.2 Privacy Semantics

Setting α to 0 and ∞ results in neighboring tables that differ in the presence or absence of one worker and one establishment, respectively. These choices correspond to edge- and node-differential privacy discussed in Section 6

We can also bound the extent to which an adversary can infer properties of employees and establishments beyond just neighboring databases. Both neighboring definitions, 7.1 and 7.3, induce a metric $d(\cdot)$ over possible databases.

We can use this metric to reason about which inferences an adversary can make. For mechanism \mathcal{M} satisfying one of our privacy definitions,

$$\Pr[\mathcal{M}(D) \in S] \leq e^{\epsilon \cdot d(D, D')} \cdot \Pr[\mathcal{M}(D') \in S]. \quad (8)$$

In particular, there are two things to note. First, the distance between databases that differ in workplace attributes is infinite (we can disclose this information because the workplace attributes are public). Next, suppose D and D' are neighbors with different employment sets E_w and E'_w for establishment w , and $|E_w| = x$ and $|E'_w| = (1 + \alpha)^k \cdot x$. Then, an adversary can’t distinguish between D and D' based on an output with log-odds greater than $\epsilon \cdot k$; that is, $\epsilon \cdot k$ bounds the adversary’s Bayes factor.

7.3 Composition

THEOREM 7.3. *Let \mathcal{M}_1 and \mathcal{M}_2 be (α, ϵ_1) - and (α, ϵ_2) -ER-EE private algorithms. Releasing the outputs of $\mathcal{M}_1(D)$ and $\mathcal{M}_2(D)$ results in $(\alpha, \epsilon_1 + \epsilon_2)$ -ER-EE privacy. The same holds for weak (α, ϵ) -ER-EE privacy.*

Differentially private algorithms also satisfy parallel composition, which means that releasing the result of ϵ -differentially private algorithms on disjoint sets of records D_1 and D_2 also results in ϵ -differential privacy. The same is true for both Definitions 7.2 and 7.4 if the records in D_1 and D_2 pertain to distinct establishments.

THEOREM 7.4. *Let D_1 and D_2 represent subsets of records from the ER-EE dataset that pertain to distinct sets of establishments. Let \mathcal{M}_1 and \mathcal{M}_2 be (α, ϵ) - and (α, ϵ) -ER-EE private algorithms. Releasing the outputs of $\mathcal{M}_1(D_1)$ and $\mathcal{M}_2(D_2)$ results in (α, ϵ) -ER-EE privacy. The same holds for weak (α, ϵ) -ER-EE privacy.*

Parallel composition is nuanced when D_1 and D_2 could pertain to distinct sets of workers from the same sets of establishments (e.g., males in New York and females in New York).

THEOREM 7.5. *Let D_1 and D_2 represent subsets of records from the ER-EE dataset that pertain to distinct workers, but have records that arise from the same establishment. Let \mathcal{M}_1 and \mathcal{M}_2 be (α, ϵ) - and (α, ϵ) -ER-EE private algorithms. Releasing the outputs of $\mathcal{M}_1(D_1)$ and $\mathcal{M}_2(D_2)$ results in (α, ϵ) -ER-EE privacy. The same does not hold for weak (α, ϵ) -ER-EE privacy.*

8. ALGORITHMS

We next present algorithms for answering queries under both (α, ϵ) -ER-EE privacy and weak (α, ϵ) -ER-EE privacy. Our algorithms describe how to answer a single count query (e.g., number of workers in the age 25-35 with a college degree who are employed in publicly-owned establishments in New York). This would correspond to a single cell in a marginal query (e.g., {place, ownership, age, education}). We denote single counting queries as q_v , where $v \in \text{dom}(V)$ and q_v is the marginal query. These algorithms can be used to release all the counts in the marginal using the composition properties described in Section 7.3. algorithms for releasing single counts can be parallel-composed under (α, ϵ) -ER-EE privacy for all marginals. However, algorithms for releasing single counts can be parallel-composed under weak (α, ϵ) -ER-EE privacy for marginals containing only establishment attributes, since such cells aggregate over distinct subsets of employers. Using algorithms for releasing a single count under weak (α, ϵ) -ER-EE privacy to release a marginal containing worker attributes would result in an effective privacy-loss parameter of $d \cdot \epsilon$, where d is the domain size of the worker attributes in the marginal query.

8.1 Log-Laplace Algorithm

The global sensitivity of a count query under Definitions 7.2 and 7.4 is unbounded; if the count is x , the sensitivity can be as large as αx . However, the logarithm of the count has a low global sensitivity of $\ln(1 + \alpha)$. Thus the Log-Laplace mechanism (see Algorithm 1) adds Laplace noise to the log of the count.

THEOREM 8.1. *Suppose q_v is a query over only establishment attributes. Then, releasing q_v using Algorithm 1 satisfies (α, ϵ) -ER-EE privacy.*

Suppose q_v is a query over both establishment attributes and employee attributes. Then, releasing q_v using Algorithm 1 satisfies weak (α, ϵ) -ER-EE privacy.

Algorithm 1 Log-Laplace Mechanism

Input: n : the sum of employment counts for a set of cells, α, ϵ : privacy parameters

Output: \tilde{n} : the noisy employment count

Set $\gamma \leftarrow 1/\alpha$

$\ell \leftarrow \ln(n + \gamma)$

Sample $\eta \sim \text{Laplace}(2\ln(1 + \alpha)/\epsilon)$

$\tilde{n} \leftarrow e^{\ell + \eta} - \gamma$

All proofs are deferred to the Appendix. While the original Laplace mechanism is unbiased (the expectation of the noisy sum equals the true sum), the Log-Laplace mechanism is not. In particular we can show:

LEMMA 8.2. *Let x denote a real number, and \tilde{x} the random variable denoting the output of the Log-Laplace mechanism. Let $\lambda = 2\ln(\alpha + 1)/\epsilon$. Then, when $\lambda < 1$, $E[\tilde{x}] + \gamma = (x + \gamma)/(1 - \lambda^2)$. When $\lambda \geq 1$, $E[\tilde{x}]$ is unbounded.*

THEOREM 8.3. *For q_v , let \tilde{x} denote the output of the Log-Laplace mechanism and x denote the true answer of the query. The expected squared relative error of the Log-Laplace mechanism for q_v is bounded when $\lambda = 2\ln(\alpha + 1)/\epsilon$ is less than $1/2$, and is given by:*

$$\mathcal{E}_{rel}(q_v) = \max_D \left(\frac{(x - \tilde{x})^2}{x^2} \right) \leq \frac{(2\lambda^2 + 4\lambda^4)(1 + \gamma)^2}{(1 - 4\lambda^2)(1 - \lambda^2)} \quad (9)$$

8.2 Smooth Sensitivity-based Algorithms

We next derive a mechanism using an extension of the smooth sensitivity framework [38]. The smooth sensitivity framework adds noise based on *local sensitivity* of the input database rather than global sensitivity across all databases. While local sensitivity can be much smaller than global sensitivity, adding noise proportional to local sensitivity does not ensure differential privacy, and hence, local sensitivity must be “smoothed.”

DEFINITION 8.1 (LOCAL SENSITIVITY). *Let q be a query, \mathcal{I} be a domain of datasets, and $\text{nbrs}(x)$ denote the set of neighbors of x according to the appropriate definition; e.g. Definition 7.1 or 7.3. The local sensitivity of query q for a dataset $x \in \mathcal{I}$ is*

$$LS_q(x) = \max_{y \in \text{nbrs}(x)} \|q(x) - q(y)\|_1$$

Global sensitivity is the maximum local sensitivity over all databases. Nissim et al. [38] show that it is sufficient to add noise proportional to any smooth function that upper bounds local sensitivity. The smallest such upper bound is called the *smooth sensitivity*.

DEFINITION 8.2. *Let q be a query and b a smoothing parameter. Let \mathcal{I} denote the universe of all datasets. The b -smooth sensitivity of query q with respect to database x is defined as*

$$S_{q,b}^*(x) = \max_j e^{-jb} A^{(j)}(x),$$

$$\text{where } A^{(j)}(x) = \max_{y \in \mathcal{I}: d(x,y) \leq j} LS_q(y),$$

and $d(x, y)$ is the smaller integer ℓ such that there exist databases $x = x_0, x_1, \dots, x_\ell = y$, such that for all i , x_{i-1} and x_i are neighbors according to either Definition 7.1 or 7.3.

We next define *admissible* distributions, and show that adding noise proportional to an admissible distribution preserves privacy. Our definition is a slight generalization of [38] that is more flexible

and will result in better error for our application. In particular, the definition of admissible distributions in [38] splits the total ϵ -budget equally between the sliding and dilation properties. We extend this definition to allow a flexible split of the budget.

DEFINITION 8.3. *Let ϵ_1 and ϵ_2 be a division of the budget such that $\epsilon_1 + \epsilon_2 \leq \epsilon$. A probability distribution h is (a, b) -admissible with respect to ϵ , where a and b are functions δ , and of ϵ_1 and ϵ_2 respectively, if $\forall \lambda \in \mathbb{R}, \Delta \in \mathbb{R}^d$ with $|\lambda| \leq b$ and $\|\Delta\|_1 \leq a$, and $\forall S \subseteq \mathbb{R}^d$,*

$$\Pr_{Z \sim h} [Z \in S] \leq e^{\epsilon_1} \Pr_{Z \sim h} [Z \in S + \Delta] + \frac{\delta}{2}, \text{ and} \quad (10)$$

$$\Pr_{Z \sim h} [Z \in S] \leq e^{\epsilon_2} \Pr_{Z \sim h} [Z \in S \cdot e^\lambda] + \frac{\delta}{2}. \quad (11)$$

We can now adapt Lemma 2.6 from [38] to show that adding noise from admissible distributions, scaled by the smooth sensitivity, generates provably private algorithms.

THEOREM 8.4. *Suppose h is an (a, b) -admissible probability distribution with $\delta = 0$, and $Z \sim h$. For query q , let $S(x)$ be a b -smooth upper bound on the local sensitivity of q . Then, the algorithm $\mathcal{M}(x) = q(x) + \frac{S(x)}{a} \cdot Z$ satisfies (α, ϵ) -ER-EE privacy.*

We now compute the b -smooth sensitivity of our queries and describe an admissible distribution. For our problem, the local sensitivity itself is the smooth sensitivity.

LEMMA 8.5. *Let q_v be a query on x . Let x_v be the maximum number of workers belonging to a single workplace and matching the conditions in v . Then, the b -smooth sensitivity of x , $S_{v,b}^*(x)$, is*

$$S_{v,b}^*(x) = \begin{cases} \max(x_v \cdot \alpha, 1) & \text{if } e^b \geq (1 + \alpha), \\ \text{unbounded} & \text{otherwise.} \end{cases} \quad (12)$$

LEMMA 8.6. *Let $\epsilon_1 + \epsilon_2 \leq \epsilon$. $h(z) \propto \frac{1}{(1+|z|^\gamma)}$ is $(\frac{\epsilon_1}{1+\gamma}, \frac{\epsilon_2}{1+\gamma})$ -admissible for $\gamma > 0$ ($\delta = 0$).*

Combining Theorem 8.4 and Lemmas 8.5 and 8.6 gives us the following algorithm. We use $\gamma = 4$ to ensure that the mean and variance of the noise distribution are bounded. Note that privacy is guaranteed only when $\alpha + 1 < e^{\epsilon/5}$. Values of α and ϵ not satisfying this inequality are not allowed by the algorithm. When this condition is met, we set ϵ_2 such that b is as low as possible without violating the inequality in Equation 12. Since only a contributes directly to the error of the algorithm, this ensures that the ϵ budget is used efficiently to minimize error.

Algorithm 2 Smooth Gamma

Input: n : true count, α, ϵ : privacy parameters, $\alpha + 1 < e^{\epsilon/5}$

Output: \tilde{n} : noisy count

Sample $\eta \sim \frac{1}{(1+|z|^4)}$

$\epsilon_2 \leftarrow 5 \cdot \ln(\alpha + 1)$

$\epsilon_1 \leftarrow \epsilon - \epsilon_2$ ($\epsilon_1 > 0$ by the condition in the **Input**.)

$\tilde{n} \leftarrow n + \frac{S_{v,\epsilon_2/5}^*(x)}{\epsilon_1/5} \eta$,

LEMMA 8.7. *Suppose q_v is a query over only establishment attributes. Then releasing q_v using Algorithm 2 satisfies (α, ϵ) -ER-EE privacy.*

Suppose q_v is a query over both establishment and individual attributes. Then releasing q_v using Algorithm 2 satisfies weak (α, ϵ) -ER-EE privacy.

LEMMA 8.8. *Algorithm 2 is unbiased and has expected L_1 error of $O(\frac{x_v \cdot \alpha}{\epsilon} + \frac{1}{\epsilon})$.*

9. APPROXIMATING PRIVACY

A standard relaxation of differential privacy is to allow for a small failure probability of δ that the attacker can distinguish neighboring datasets based on an output. We can similarly define $(\alpha, \epsilon, \delta)$ -ER-EE privacy.

DEFINITION 9.1 ($(\alpha, \epsilon, \delta)$ -ER-EE PRIVACY). *A randomized algorithm \mathcal{M} is said to satisfy $(\alpha, \epsilon, \delta)$ -ER-EE Privacy, if for every set of outputs $S \subseteq \text{range}(\mathcal{M})$, and every pair of strong α -Neighbors D and D' , we have*

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S] + \delta$$

Weak $(\alpha, \epsilon, \delta)$ -employer employee privacy is defined analogously.

Exact records may be released when $\delta = \Omega(1/n)$ (where n is the number of records). δ is the probability that a mechanism gives no privacy guarantee. Therefore if $\delta > 1/n$, a mechanism which releases the exact values for a δ fraction of the records (and 0 for all other records) satisfies this privacy definition.

δ increases rapidly with database distance. We would like to show an analogue to Equation 8 to show how δ decays with distance. That is, suppose \mathcal{M} is a mechanism which satisfies $(\alpha, \epsilon, \delta)$ -ER-EE privacy, and suppose for two databases D and D' , $d(D, D') = d$. Then, $\forall S \subseteq \text{range}(\mathcal{M})$,

$$\Pr[\mathcal{M}(D) \in S] \leq e^{\epsilon d} \cdot \Pr[\mathcal{M}(D') \in S] + \Omega(\delta e^{(d-1)\epsilon}). \quad (13)$$

Note that for a high enough distance, the term $\Omega(\delta e^{(d-1)\epsilon})$ may become greater than one. This means that for database D' at least this far from D , we may allow an adversary to rule out without a doubt D' . This never happens under non-approximate privacy. An adversary must always have some amount of uncertainty between a pair of databases no matter how far apart.

Despite these drawbacks, allowing a small probability of failure can greatly increase the utility of an algorithm while still providing a provable guarantee. We use the framework developed in Section 8.2 to give an approximate algorithm.

LEMMA 9.1 ([38]). *The Laplace distribution, $h(z) \propto \frac{1}{2} \cdot e^{-|z|}$, is $(\epsilon/2, \frac{\epsilon}{2 \ln(1/\delta)})$ -admissible.*

Using Theorem 8.4 and Lemma 8.5, we obtain Algorithm 3.

Algorithm 3 Smooth Laplace

Input: n : true count, α, ϵ : privacy parameters, $\alpha + 1 \leq e^{\frac{\epsilon}{2 \ln(1/\delta)}}$.

Output: \tilde{n} : noisy count

$$\text{Sample } \eta \sim \text{Laplace}(1)$$

$$\tilde{n} \leftarrow n + \frac{S_{v^*}^* \frac{\epsilon}{2 \ln(1/\delta)}(x)}{\epsilon/2} - \eta,$$

LEMMA 9.2. *Suppose q_v is a query over only establishment attributes. Then releasing q_v using Algorithm 3 satisfies $(\alpha, \epsilon, \delta)$ -employer employee privacy.*

Suppose q_v is a query over both establishment and individual attributes. Then releasing q_v using Algorithm 3 satisfies weak $(\alpha, \epsilon, \delta)$ -ER-EE privacy.

LEMMA 9.3. *Algorithm 3 is unbiased and expected L_1 error is $O(\frac{x_v \cdot \alpha}{\epsilon} + \frac{1}{\epsilon})$.*

Note that the error of Algorithm 3 does *not* depend on δ . Therefore, the optimal δ for a fixed α and ϵ is the one which solves the inequality in Algorithm 3 with equality. In Table 2 (in Appendix C), we show some values of δ, ϵ , and α that work.

10. EMPIRICAL EVALUATION

Dataset: The data used for these experiments were a 3-state sample from the LODES data infrastructure to which standard edits and imputations had already been applied. The sample was taken from the 2011 snapshot in which April 1, the first day of Quarter 2, is the reference date. To be included jobs had to qualify as “beginning-of-quarter” jobs, which means that the job (person-firm relationship) had positive earnings in the reference quarter (Q2) as well as the previous quarter (Q1). Then the assumption is that the person was employed in the job on the first day of Q2. The count of jobs in the sample was 10.9 million jobs in about 527,000 establishments.

Queries and Quality Measures:

- **Workload 1** A marginal over all establishment characteristics: industry sector, ownership, and location at the resolution of places (e.g., cities and towns).
- **Workload 2** Single queries over all establishment attributes, and over the worker attributes of sex and education.
- **Workload 3** The marginal over all establishment attributes, and sex and education.

We report the cost of provable privacy as a ratio between the average L_1 error (over 20 independent trials) of our provably private algorithms divided by the L_1 error of current SDL algorithm.

In addition to reporting the overall error ratio, we also compute the error ratio stratified by place-size ranges. The strata we consider are cells in the marginals with a population⁸ of 0-100, 100-10k, 10k-100k, and 100k+, respectively. Our results for Workload 3 are discussed here, but appear in the Appendix.

Next we evaluate the cost of formal privacy in ranking tasks.

- **Ranking 1** Rank all the cells in the marginal over industry sector, ownership, and location by total count in descending order.
- **Ranking 2** Rank all the cells in the marginal over industry sector, ownership, and location by number of employees who are female with a college degree in each cell in descending order. The results of this experiment are discussed here, and the plots appear in the appendix.

We measure error as the Spearman rank-order correlation between the ordering based on noisy counts returned by our algorithm to the ordering based on the counts output by the current SDL algorithm. As in the L_1 error case, we also report error stratified by place size.

Algorithms: We compare the Log-Laplace, Smooth Gamma and Smooth Laplace algorithms. We present results for $\epsilon \in \{0.25, 0.5, 0.67, 1.0, 2.0, 4.0\}$ and for $\alpha \in \{0.01, 0.05, 0.1, 0.15, 0.2\}$. Recall that for $\alpha = 0.1$, for example, an adversary should not be able to tell the difference in employments within 10% of each other. We do not plot errors for the Log-Laplace mechanism when the expectation of the noisy count is unbounded (see Lemma 8.2).

We do not vary δ for the Smooth Laplace algorithm as part of our evaluation, since δ does not affect the amount of noise added (and consequently does not impact the accuracy). Nevertheless, as discussed in Section 9 and Table 2, for a given α, δ imposes a lower bound on ϵ and eliminates possible choices for (α, ϵ) pairs. In our figures, we report results for pairs of (α, ϵ) that are possible for a high failure probability of $\delta = 0.05$. The performance for smaller δ values can be read off the plots by checking whether the (α, ϵ) values are allowed for that δ .

We also implemented a node-differentially private algorithm (from Section 6), and call it “Truncated Laplace”. This suppresses establishments with counts $\geq \theta$, for $\theta \in \{2, 20, 50, 100, 200, 500\}$.

⁸Census Place total population (P0010001) as published in the 2010 Decennial Census.

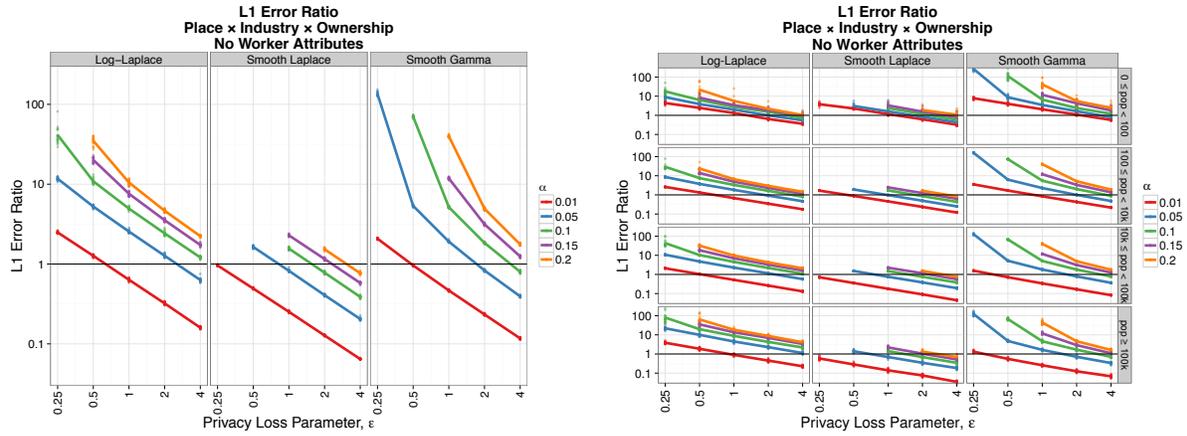


Figure 1: Average L_1 error ratio of releasing employment count for Census place by NAICS sector (industry) by ownership marginal compared to the current system.

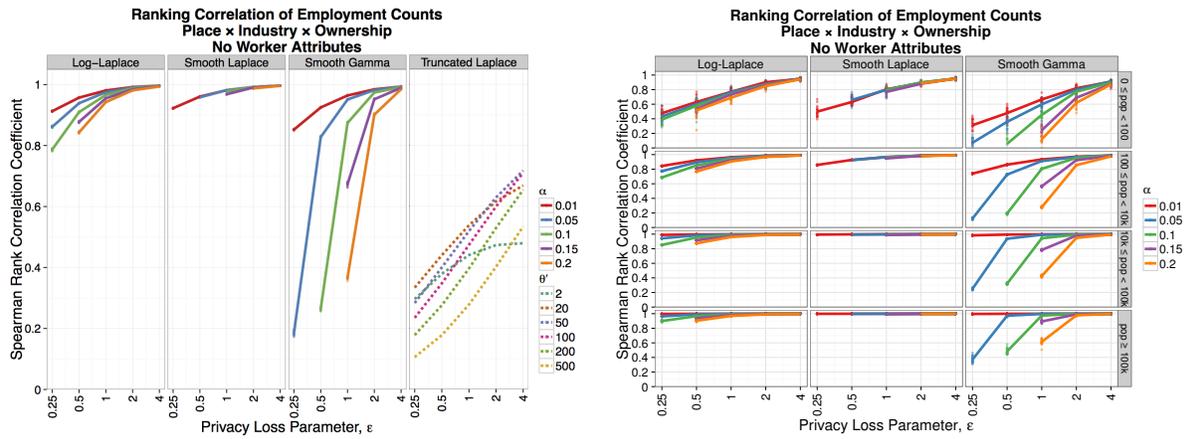


Figure 2: Spearman correlation between tested model and input noise infusion on the count of employment ranked for Census place by NAICS sector (industry) by ownership.

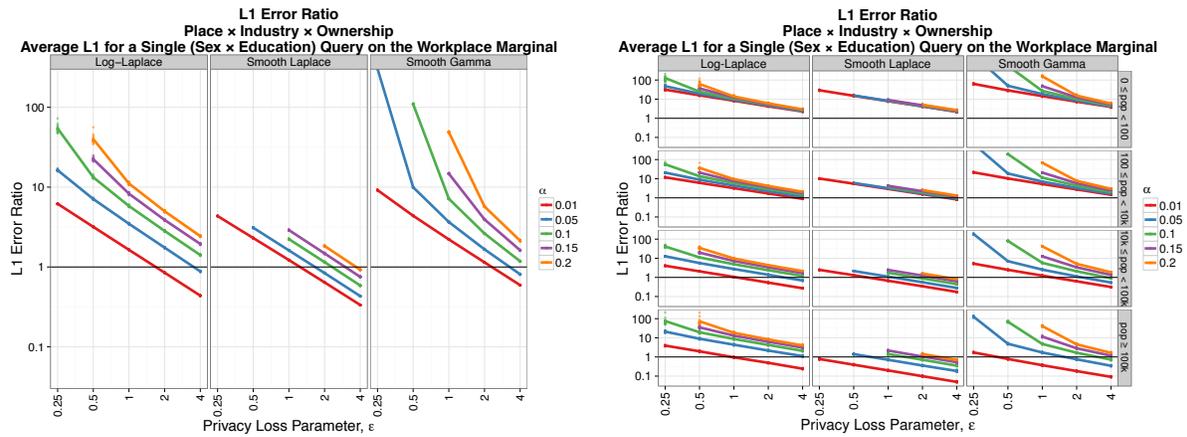


Figure 3: Average L_1 error ratio of releasing single queries of employment in the Census place by NAICS sector (industry) by ownership by sex by education marginal, compared to the current system.

Finding 1: For releasing marginals over only establishment attributes, our algorithms perform comparably or better than the current protection system while satisfying the strong privacy notion given in Def. 7.2. We use $\epsilon = 2$ and $\alpha = 0.1$ as a baseline. As can be seen in Figure 1, the average error of the Log-Laplace and Smooth Gamma algorithms is within a factor of 3 of the error of the current protection system. The Smooth Laplace algorithm performs better than the current protection system at $\epsilon = 2$ and $\alpha = 0.1$. Ranking results for establishment-attribute queries are shown in Figure 2. Overall, the Smooth Laplace algorithm has a correlation close to 1 when ϵ is at least 2, and the other two algorithms are close to 1 for $\epsilon \geq 4$. For rankings involving only larger population sizes, Smooth Laplace is close to 1 for all values of ϵ tested, and the Log-Laplace algorithm is close to 1 for ϵ at least 1.

We also observe that our algorithms have relative error comparable to that of the current protection system for a majority of the cells. For Log-Laplace, the relative L_1 is within 10 percentage points of the relative error of SDL for 65% of the counts at $\alpha = 0.1$ and $\epsilon = 2$. Smooth Laplace and Smooth Gamma are within 10 percentage points for 75% and 29% of the counts, respectively.

Finding 2: For releasing individual queries over establishment and worker attributes, our algorithms perform comparably or better than the current protection system while satisfying the weaker privacy notion given in Def. 7.4. As shown in Figure 3, the average error of the Log-Laplace algorithm is within a factor of 3 of the error of the current protection system. The Smooth Laplace algorithm has nearly the same average error as the current protection system. At $\epsilon = 4$, the Smooth Laplace algorithm outperforms the current protection system for all values of α that we tested. Ranking results are shown in Figure 5. For the overall ranking, only the Smooth Laplace algorithm approaches relative error of 1 for ϵ at least 4. Restricted to larger population sizes, both Log-Laplace and Smooth Laplace perform well for all tested values of ϵ .

Finding 3: For releasing marginal queries over establishment and worker attributes, our algorithms perform worse than the current protection system, but can have acceptable performance at high values of ϵ and low values of α . As shown in Figure 4, when $\alpha \leq 0.05$ and $\epsilon \geq 4$, the Log-Laplace algorithm has average L_1 error within a factor of 10 of the current protection system. The Smooth Laplace algorithm is within a factor 10 for all tested values of α at $\epsilon = 4$. At the lowest tested value of α ($\alpha = 0.01$), the Smooth Laplace algorithm is within a factor of 3 of the error of the current system.

Finding 4: All three of our algorithms perform better as population size grows. This can be seen in both results that measure L_1 error (Figures 1, 4, and 3) as well as those that measure ranking (Figures 5 and 2). In the results that measure L_1 error, the algorithms have lower error with respect to the current protection system as the population size grows. In the ranking results, the rank order of our algorithms is more similar to the ranking order of the current protection system as the population size grows. In all cases, the improvement from the small population range (0-100) to the next smallest (100-10k) results in the largest increase in performance. Further increases in population size have a smaller effect.

Finding 5: Our Smooth Laplace algorithm performs the best of the three, and the Log-Laplace and Smooth Gamma algorithms perform similarly. Smooth Laplace performs best in all experiments, and this is not surprising since it satisfies a weaker notion of privacy. However, the ordering of the performance of the other two algorithms is not consistent. Log-Laplace generally performs better at lower values of ϵ .

Finding 6: For all counting workloads, Truncated Laplace had a significantly higher cost than our techniques. E.g., for Workload 1

it incurs an additive error that is at least 10 times larger than that of SDL techniques (at $\epsilon = 4$), and the error does not decrease significantly when ϵ is increased. Truncated Laplace also performs poorly for ranking workloads. E.g., for Ranking 1 (Figure 2), it has a correlation coefficient of no better than 0.7 for any value of θ at every value of ϵ that we tested. For a fixed θ , increasing ϵ past some point only provides small utility gains (this can be seen in the plot for $\theta = 2$ and happens at larger values of ϵ for larger θ). This is because much of the error introduced is bias from removing large establishments, and increasing ϵ does not change this error.

Summary: Our empirical results suggest that there are a number of settings of (α, ϵ) that allow for publishing cell counts with little or no additional cost to accuracy in return for better and provable privacy protection. For marginals over only establishment attributes, all algorithms perform well at $\epsilon = 2$ and $\alpha = 0.1$. For individual queries of worker and establishment attributes, our algorithms perform well at $\epsilon = 2$ and $\alpha = 0.1$. For marginals over worker and establishment attributes, the Smooth Laplace algorithm performs well when $\epsilon \geq 4$. All algorithms perform better when the queries are over places with greater population size. Allowing for a small failure probability results in a significant reduction in error as seen with the Smooth Laplace algorithm. Counts output by our algorithms can also be used for ranking with high accuracy for $\epsilon \geq 1$.

Choosing an algorithm: The three algorithms, Smooth Gamma, Smooth Laplace, and Log-Laplace, have minor differences that can make one better than another in different scenarios. Smooth Laplace differs from the others in that it allows a small probability that the privacy guarantee is not met. The drawbacks of this relaxation are discussed in detail in Section 9, and the algorithm can be used to achieve significantly lower error than the others if these drawbacks are acceptable. If instead we want privacy with $\delta = 0$, Smooth Gamma and Log-Laplace also have slight differences. Values of α and ϵ disallowed by the Smooth Gamma algorithm are still allowed by the Log-Laplace algorithm, though the Log-Laplace algorithm typically results in large errors for these values. The Log-Laplace algorithm slightly outperforms the Smooth Gamma algorithm for small ϵ , but is outperformed by Smooth Gamma for large ϵ .

11. CONCLUSIONS

We considered the problem of releasing ER-EE data with provable privacy guarantees and measured the utility cost of this privacy protection. We identified privacy requirements based on current interpretation of laws pertaining to the release of these data and mathematically formalized them using the Pufferfish framework. We showed that current SDL techniques do not satisfy these strong privacy requirements. Direct adaptations of ϵ -differential privacy either do not satisfy the privacy requirements or output data with very limited utility. We develop novel privacy definitions that provably satisfy our privacy requirements. For the task of releasing marginals over establishment attributes, releasing single counts over marginals that include worker attributes, and ranking queries, our algorithms incur error that is comparable and in some cases less than the error incurred by current SDL algorithms for reasonable privacy parameters. Our results suggest these data can be released using provably private algorithms with a low utility cost.

Acknowledgements. Disclaimer: Any opinions and conclusions expressed herein are those of the authors and do not necessarily represent the views of the U.S. Census Bureau. All results have been reviewed to ensure that no confidential information is disclosed. Authors acknowledge support from NSF grants 1253327, 1408982, 1443014, BCS-0941226, TC-1012593 and SES-1131848, DARPA & SPAWAR N66001-15-C-4067, and Alfred P. Sloan Foundation.

12. REFERENCES

- [1] Data protection and privacy, U.S. Census Bureau. <http://www.census.gov/privacy/>.
- [2] LEHD Origin-Destination Employment Statistics (LODES) Technical Document. <http://lehd.ces.census.gov/data/loDES/LODESTechDoc7.1.pdf>.
- [3] OnTheMap for Emergency Management. <http://onthemap.ces.census.gov/em/>.
- [4] OnTheMap Web Tool. <http://onthemap.ces.census.gov/>.
- [5] Title 13 - protection of confidential information, U.S. Census Bureau. http://www.census.gov/about/policies/privacy/data_protection/title_13_-_protection_of_confidential_information.html.
- [6] U.S. Census Bureau's Longitudinal Employer-Household Dynamics Program. <http://lehd.ces.census.gov/>.
- [7] 1996 Technical documentation, County Business Patterns. Technical report, U.S. Census Bureau, Administrative and Customer Services Division, Microdata Access and Use Branch, 1998.
- [8] Report on statistical disclosure limitation methodology. Research Report WP-22, Federal Committee on Statistical Methodology, Dec 2005.
- [9] Notices: Notice of adjustment of countywide per capita impact indicator. Federal Register, Vol. 78, No. 208, pp 64231-64232, Monday, October 28, 2013.
- [10] J. M. Abowd, R. K. Gittings, K. L. McKinney, B. E. Stephens, L. Vilhuber, and S. Woodcock. Dynamically consistent noise infusion and partially synthetic data as confidentiality protection measures for related time-series. In *Federal Committee on Statistical Methodology Research Conference*, 2012.
- [11] J. M. Abowd and I. Schmutte. Economic analysis and statistical disclosure limitation. *Brookings Papers on Economic Activity*, Spring 2015.
- [12] J. M. Abowd, B. E. Stephens, and L. Vilhuber. Confidentiality protection in the census bureaus quarterly workforce indicators. Technical Report TP-2006-02, U.S. Census Bureau, LEHD Program, December 2006.
- [13] J. M. Abowd, B. E. Stephens, L. Vilhuber, F. Andersson, K. L. McKinney, M. Roemer, and S. Woodcock. The LEHD infrastructure files and the creation of the Quarterly Workforce Indicators. In T. Dunne, J. Jensen, and M. Roberts, editors, *Producer Dynamics: New Evidence from Micro Data*, pages 149–230. Chicago: University of Chicago Press for the National Bureau of Economic Research, 2009.
- [14] M. Anderson and W. Seltzer. Challenges to the confidentiality of US federal statistics, 1910-1965. *Journal of Official Statistics*, 23(1):1–34, 2007.
- [15] M. Anderson and W. Seltzer. Federal statistical confidentiality and business data: Twentieth century challenges and continuing issues. *Journal of Privacy and Confidentiality*, 1(1):7–52, 2009.
- [16] M. E. Andrés, N. E. Bordenabe, K. Chatzikokolakis, and C. Palamidessi. Geo-indistinguishability: differential privacy for location-based systems. In *ACM CCS*, 2013.
- [17] K. Chatzikokolakis, M. E. Andrés, N. E. Bordenabe, and C. Palamidessi. Broadening the scope of differential privacy using metrics. In *PETS*, 2013.
- [18] S. Chen and S. Zhou. Recursive mechanism: Towards node differential privacy and unrestricted joins. In *ACM SIGMOD*, 2014.
- [19] T. Dalenius. Towards a methodology for statistical disclosure control. *Statistik Tidskrift*, 15:429–444, 1977.
- [20] W.-Y. Day, N. Li, and M. Lyu. Publishing graph degree distribution with node differential privacy. In *ACM SIGMOD*, 2016.
- [21] G. T. Duncan, M. Elliot, and J.-J. Salazar-González. *Statistical Confidentiality Principles and Practice*. Statistics for Social and Behavioral Sciences. Springer New York, 2011.
- [22] G. T. Duncan, T. B. Jabine, and V. A. de Wolf, editors. *Private Lives and Public Policies: Confidentiality and Accessibility of Government Statistics*. National Research Council: Committee on National Statistics, Washington, DC, 1993.
- [23] C. Dwork. Differential privacy. In *ICALP (2)*, pages 1–12, 2006.
- [24] C. Dwork. The differential privacy frontier (extended abstract). In *TCC*, pages 496–502, 2009.
- [25] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, pages 265–284, 2006.
- [26] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3,4):211–407, Aug. 2014.
- [27] T. Evans, L. Zayatz, and J. Slanta. Using noise for disclosure limitation of establishment tabular data. *Journal of Official Statistics*, 14(4):537–551, 1998.
- [28] I. P. Fellegi. On the question of statistical confidentiality. *Journal of the American Statistical Association*, 67(337):pp. 7–18, 1972.
- [29] S. Haney, A. Machanavajjhala, and B. Ding. Design of policy-aware differentially private algorithms. *Proceedings of the VLDB Endowment*, 9(4):264–275, 2015.
- [30] X. He, A. Machanavajjhala, and B. Ding. Blowfish privacy: Tuning privacy-utility trade-offs using policies. In *ACM SIGMOD*, pages 1447–1458, 2014.
- [31] V. Karwa, S. Raskhodnikova, A. Smith, and G. Yaroslavtsev. Private analysis of graph structure. *PVLDB*, 4(11):1146–1157, 2011.
- [32] S. Kasiviswanathan, K. Nissim, S. Raskhodnikova, and A. Smith. Analyzing graphs with node differential privacy. In *TCC*, 2013.
- [33] D. Kifer and A. Machanavajjhala. No free lunch in data privacy. In *SIGMOD*, pages 193–204, 2011.
- [34] D. Kifer and A. Machanavajjhala. A rigorous and customizable framework for privacy. In *PODS*, pages 77–88, 2012.
- [35] D. Kifer and A. Machanavajjhala. Pufferfish: A framework for mathematical privacy definitions. *ACM Transactions on Database Systems (TODS)*, 39(1):1–36, 2014.
- [36] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian. l-diversity: Privacy beyond k-anonymity. In *ICDE*, page 24, 2006.
- [37] A. Machanavajjhala, D. Kifer, J. M. Abowd, J. Gehrke, and L. Vilhuber. Privacy: Theory meets practice on the map. In *ICDE*, pages 277–286, 2008.
- [38] K. Nissim, S. Raskhodnikova, and A. Smith. Smooth sensitivity and sampling in private data analysis. In *STOC*, pages 75–84. ACM, 2007.
- [39] Y. Wang, S. Song, and K. Chaudhuri. Privacy-preserving analysis of correlated data. *CoRR*, abs/1603.03977, 2016.
- [40] Y. Xiao and L. Xiong. Protecting locations with differential privacy under temporal correlations. In *ACM CCS*, pages 1298–1309, 2015.
- [41] B. Yang, I. Sato, and H. Nakagawa. Bayesian differential privacy on correlated data. In *ACM SIGMOD*, pages 747–762, 2015.

APPENDIX

A. RELATED WORK

Prior work on interpreting privacy requirements: As reviews like [14] and the similar review for business data in [15] make clear, the privacy provisions in the statutes were designed to protect individuals and businesses from uses of the data that specifically identified them and, moreover, harmed them by subjecting them to punishment or competitive disadvantage in ways that could not be accomplished as easily without the confidential data. For protecting individuals, [28] formalized these requirements as ensuring that no exact disclosure of records in the underlying data. Fellegi derived the necessary and sufficient conditions for a set of published tables to be fully resistant to a subtraction attack that could expose one or more records. Thus, the primary goal of existing SDL techniques has been to prevent exact re-identification. For instance, a method called primary and complementary suppression [8] implements Fellegi's conditions [28], and has long been considered by statistical agencies around the world as compliant with confidentiality protection laws [21].

With regards to business data, protecting the characteristics of an establishment has been interpreted and implemented as described in Section 5 as ensuring that (a) the true counts of the workforce characteristics are never released or used to compute aggregates,

and (b) employment counts of a workplace are perturbed by a confidential multiplicative factor unique to that workplace.

Despite avoiding exact disclosures, data publications might violate individual or business privacy by allowing too precise an inference about the true values, given the published values. This idea was first formulated in the SDL literature [19] as well as rediscovered and popularized in the computer science literature (e.g., [36, 34, 23]). We show in Section 5.2 examples of such inferences that can be made by an adversary, especially in the presence of background knowledge. Our goal is to protect ER-EE data as per the aforementioned requirements while ensuring a formal privacy notion that can limit both inferential and exact disclosures.

Customizing differential privacy: Kifer and Machanavajjhala [33] prove a no-free-lunch theorem for privacy that states that one cannot achieve privacy and utility simultaneously without making assumptions on the attacker's prior knowledge. This means that no single privacy notion (including differential privacy) ensures sufficient privacy protection for all applications and data types. Hence, recent work has focused on generalizations, both strengthening and relaxing parts of the differential privacy framework. The Pufferfish framework [34, 35], which we use in this paper, generalizes differential privacy by specifying what information should be kept secret, and the attacker's prior knowledge. He et al. [30] propose the Blowfish framework, which also generalizes differential privacy and is inspired by Pufferfish, and Haney et al. [29] give a general method for creating algorithms for any Blowfish policy graph. Our privacy requirements are an instantiation of Pufferfish, and our privacy definitions can be thought of as instantiations of Blowfish. Chatzikokolakis et al [17] investigate notions of privacy that can be defined as metrics over the set of databases. These have led to the design of application specific privacy notions (e.g., [16, 39]).

Differential Privacy for complex entities: Section 6 shows that our problem can be considered as one of differential privacy on graphs, but the use of existing techniques either does not satisfy our privacy requirements or results in a complete loss of utility.

B. PROOFS

CLAIM B.1. *Differential privacy on establishments satisfies all three privacy requirements given in Section 4. Differential privacy on individuals does not satisfy all requirements.*

PROOF. The proof that differential privacy on establishments satisfies the three privacy requirements is essentially the same as the proof of Theorem 7.1.

Differential privacy on individuals, however, does not satisfy the requirements. For D and D' differing in d individuals, the following holds for all outputs S :

$$\Pr[A(D) \in S] \leq e^{d\epsilon} \Pr A(D') \in S. \quad (14)$$

When an establishment e with size $|e|$ changes in size by a factor of α , the change in the number of individuals is $\alpha|e|$. As $|e|$ goes to infinity, the change in the number of individuals becomes arbitrarily large, and therefore $d\epsilon$ is not a constant. \square

PROOF OF THEOREM 7.1. Requirement 1: We will denote $\Pr[\mathcal{M}(D) = w]$ as $\Pr[w]$. For some individual r ,

$$\log \frac{\Pr[w \mid r \in D]}{\Pr[w \mid r \notin D]} \leq \epsilon. \quad (15)$$

Then, the individual privacy requirement of Definition 4.1 follows from Bayes' theorem.

Requirement 2: Let x and y be two numbers such that $x \leq y \leq (1 + \alpha)x$. Then for all establishments e (where $|e|$ denotes the

number of employees at e), we want to show that

$$\log \frac{\Pr[w \mid |e| = x]}{\Pr[w \mid |e| = y]} \leq \epsilon. \quad (16)$$

Additionally,

$$\Pr[w \mid |e| = x] = \sum_{E \subset U} \Pr[w \mid e = E] \Pr[e = E \mid |e| = x], \quad (17)$$

and

$$\Pr[w \mid |e| = y] \quad (18)$$

$$= \sum_{E' \subset U} \Pr[w \mid e = E'] \Pr[e = E' \mid |e| = y] \quad (19)$$

$$= \sum_{E \subset U: |E|=x} \sum_{E' \supset E} \Pr[w \mid e = E'] \Pr[e = E' \mid |e| = y] \quad (20)$$

$$\geq \sum_{E \subset U: |E|=x} \left[\min_{E' \supset E} \Pr[w \mid e = E'] \sum_{E' \supset E} \Pr[e = E' \mid |e| = y] \right] \quad (21)$$

$$= \sum_{E \subset U} \left[\min_{E' \supset E} \Pr[w \mid e = E'] \Pr[e = E \mid |e| = x] \right]. \quad (22)$$

Then,

$$\frac{\sum_{E \subset U} \Pr[w \mid e = E] \Pr[e = E \mid |e| = x]}{\sum_{E \subset U} [\min_{E' \supset E} \Pr[w \mid e = E'] \Pr[e = E \mid |e| = x]]} \quad (23)$$

$$\leq \max_{E \subset U} \frac{\Pr[w \mid e = E] \Pr[e = E \mid |e| = x]}{\min_{E' \supset E} \Pr[w \mid e = E'] \Pr[e = E \mid |e| = x]} \quad (24)$$

$$= \max_{E \subset U} \frac{\Pr[w \mid e = E]}{\min_{E' \supset E} \Pr[w \mid e = E']} \leq e^\epsilon. \quad (25)$$

Requirement 3: Requirement 3 follows from requirement 2 for an adversary. For some p , q , and z , let y be $(q - p)|e|$.

$$\frac{\Pr[w \mid |e_\chi|/|e| = p, |e| = z]}{\Pr[w \mid |e_\chi|/|e| = p, |e| = z - y]} \leq e^\epsilon. \quad (26)$$

Additionally,

$$\frac{\Pr[w \mid |e_\chi|/|e| = q, |e| = z]}{\Pr[w \mid |e_\chi|/|e| = p, |e| = z - y]} \leq e^\epsilon. \quad (27)$$

Both of these hold because the difference in sizes is at most $\alpha|e|$. Therefore,

$$\frac{\Pr[w \mid |e_\chi|/|e| = q, |e| = z]}{\Pr[w \mid |e_\chi|/|e| = p, |e| = z]} \leq e^\epsilon. \quad \square \quad (28)$$

PROOF OF THEOREM 7.2. A proof of requirement 1 is the same as in the proof of Theorem 7.1. For weak adversaries, proof of requirement 2 is also the same as Theorem 7.1. Requirement 3 follows for weak adversaries.

We can also prove requirement 3, even for strong adversaries. We want to show

$$\frac{\Pr[w \mid |e_\chi|/|e| = q, |e| = z]}{\Pr[w \mid |e_\chi|/|e| = p, |e| = z]} \leq e^\epsilon. \quad (29)$$

for all χ , and for $0 < p \leq q \leq \min((1 + \alpha)p, 1)$. Let $x = |e|p$ and $y = |e|q$. Then we have

$$\frac{\Pr[w \mid |e_\chi|/|e| = q, |e| = z]}{\Pr[w \mid |e_\chi|/|e| = p, |e| = z]} = \frac{\Pr[w \mid |e_\chi| = y]}{\Pr[w \mid |e_\chi| = x]}. \quad (30)$$

But

$$\frac{\Pr[w \mid |e_\chi| = y]}{\Pr[w \mid |e_\chi| = x]} \leq e^\epsilon, \quad (31)$$

for all \mathcal{X} when $x \leq y \leq (1 + \alpha)x$ by a similar argument to the proof of Requirement 2 of Theorem 7.1. \square

PROOF OF THEOREMS 7.3, 7.4, AND 7.5. The proof of sequential composition follows from Pufferfish (Theorem 9.1 of [35]).

To satisfy parallel composition, we need the following: Let D_1 and D_2 be disjoint parts of the domain. Then for all databases A there exists a B such that

$$d(A \cap D_1, B \cap D_1) + d(A \cap D_2, B \cap D_2) \leq d(A, B), \quad (32)$$

where $d(\cdot)$ is the distance metric induced by our neighboring definition. In particular, we should show that if $A \cap D_1, B \cap D_1$ are neighbors, and $A \cap D_2, B \cap D_2$ are neighbors, then A, B are not neighbors. This is clearly the case when D_1 and D_2 are over different sets of establishments, since neighbors can only differ in a single establishment regardless of whether we are considering strong or weak privacy. Next, consider some pair D_1 and D_2 over workers with shared establishments, and let e be the largest such shared establishment. Suppose in $A \cap D_1$, e has E workers, and in $B \cap D_1$, e has $(1 + \alpha)E$ workers. Similarly, in $A \cap D_2$, e has E workers, and in $B \cap D_2$, e has $(1 + \alpha)E$ workers. Then A and B differ in $2\alpha E$ workers, and are therefore not neighbors under our strong definition. The same does not hold under the weak definition. \square

PROOF OF THEOREM 8.1. Consider two neighboring datasets that differ in one establishment e . S' and S be the two sets of employees of e in the two databases. Let n_- denote the sum of all other counts. We just consider the case where $S \subseteq S'$. Let $y = |S|$ and $x = |S'|$. We need to ensure privacy for two cases: (i) $x = (1 + \alpha) \cdot y$, and (ii) $x = y + 1$. In case (i),

$$\begin{aligned} \frac{P(M(D_1) = o)}{P(M(D_2) = o)} &= \frac{P(M(x + n_-) = o)}{P(M(y + n_-) = o)} \\ &= \frac{P(\eta = \ln(o + \gamma) - \ln(x + n_- + \gamma))}{P(\eta = \ln(o + \gamma) - \ln(y + n_- + \gamma))} \\ &\leq \exp\left(\frac{\epsilon}{\ln(1 + \alpha)} \ln\left(\frac{(1 + \alpha) \cdot y + n_- + \gamma}{y + n_- + \gamma}\right)\right) \\ &\leq \exp\left(\frac{\epsilon}{\ln(1 + \alpha)} \cdot \ln(1 + \alpha)\right) = e^\epsilon \end{aligned}$$

In case (ii),

$$\begin{aligned} \frac{P(M(D_1) = o)}{P(M(D_2) = o)} &= \frac{P(M(1 + y + n_-) = o)}{P(M(y + n_-) = o)} \\ &\leq \exp\left(\frac{\epsilon}{\ln(1 + \alpha)} \ln\left(\frac{1 + y + n_- + \gamma}{y + n_- + \gamma}\right)\right) \\ &\leq \exp\left(\frac{\epsilon}{\ln(1 + \alpha)} \ln(1 + \alpha)\right) = e^\epsilon \end{aligned}$$

\square

PROOF OF LEMMA 8.2.

$$E[\tilde{x}] = -\gamma + (x + \gamma) \cdot E[e^\eta]$$

where $\eta \sim \text{Laplace}(\lambda)$. $E[e^\eta]$ corresponds to the value of the moment generating function $M_\eta(1)$.

$$M_\eta(1) = E[e^\eta] = 1 + \sum_{n=1}^{\infty} E[\eta^n]/n!$$

Since $\text{Laplace}(\lambda)$ is an even distribution, for all odd i , $E[\eta^i] = 0$. Moreover, $E[\eta^{2n}] = 2n!\lambda^{2n}$. Therefore, when $\lambda < 1$

$$\begin{aligned} E[\tilde{x}] &= -\gamma + (x + \gamma) \cdot \sum_{n=1}^{\infty} \lambda^{2n} \\ &= -\gamma + (x + \gamma)/(1 - \lambda^2) \end{aligned}$$

When λ is not bounded by 1, then the expected value is not bounded. Thus this mechanism is good only when $\lambda < 1$. \square

PROOF OF THEOREM 8.3. Let y denote $x + \gamma$, where $q_v(D) = x$ is the true sum. Similarly, let \tilde{y} denote $\tilde{x} + \gamma$, where \tilde{x} is the output of the log-laplace mechanism. We will show that

$$E\left(\left(\frac{y - \tilde{y}}{y}\right)^2\right) = \frac{2\lambda^2 + 4\lambda^4}{(1 - 4\lambda^2)(1 - \lambda^2)}$$

The result in the theorem directly follows.

$$\begin{aligned} E((y - \tilde{y})^2/y^2) &= E(\tilde{y}^2)/y^2 - 2E(\tilde{y})/y + 1 \\ &= E(\tilde{y}^2)/y^2 - 2/(1 - \lambda^2) + 1 \end{aligned}$$

$E(\tilde{y}^2)/y^2 = E[e^{2\eta}]$, where $\eta \sim \text{Laplace}(\lambda)$. $E[e^{2\eta}]$ corresponds to the value of the moment generating function $M_\eta(2)$.

$$\begin{aligned} E(\tilde{y}^2)/y^2 &= E[e^{2\eta}] = M_\eta(2) = 1 + \sum_{n=1}^{\infty} 2^n E[\eta^n]/n! \\ &= 1 + \sum_{n=1}^{\infty} (2\lambda)^{2n} \\ &= 1/(1 - 4\lambda^2) \quad \text{when } \lambda < 1/2 \end{aligned}$$

Therefore, we have:

$$\begin{aligned} E((y - \tilde{y})^2/y^2) &= 1/(1 - 4\lambda^2) - 2/(1 - \lambda^2) + 1 \\ &= \frac{2\lambda^2 + 4\lambda^4}{(1 - 4\lambda^2)(1 - \lambda^2)} \quad \square \end{aligned}$$

PROOF OF THEOREM 8.4. For $y \in \text{nbrs}(x)$, we show that

$$\Pr[\mathcal{M}(x) \in S] \leq e^\epsilon \cdot \Pr[\mathcal{M}(y) \in S]. \quad (33)$$

We have

$$\Pr[\mathcal{M}(x) \in S] = \Pr_{Z \sim h} \left[Z \in \frac{S - q(x)}{S(x)/\alpha} \right] \quad (34)$$

$$\leq \Pr_{Z \sim h} \left[Z \in \frac{S - q(y)}{S(x)/\alpha} \right] \cdot e^{\epsilon_1} + \frac{\delta}{2} \quad (35)$$

$$\leq \Pr_{Z \sim h} \left[Z \in \frac{S - q(y)}{S(y)/\alpha} \right] \cdot e^{\epsilon_1 + \epsilon_2} + \delta \quad (36)$$

$$\leq \Pr[\mathcal{M}(y) \in S] \cdot e^\epsilon + \delta. \quad (37)$$

(35) holds by the first property of Definition 8.3. (36) holds by the second property of Definition 8.3. \square

PROOF OF LEMMA 8.5. First, we show the following general claim: For all x , let $N_1(x)$ and $N_2(x)$ be sets of neighbors of x such that $N_1(x) \cup N_2(x) = N(x)$. Then, let $LS_1(x)$ be the local sensitivity of x over $N_1(x)$ and let $LS_2(x)$ be the local sensitivity over $N_2(x)$. Then, $LS = \max(LS_1, LS_2)$. Let S_1 and S_2 be smooth upper bounds on LS_1 and LS_2 . We claim that the function $S = \max(S_1, S_2)$ is a smooth upper bound on LS .

1. We know that S is an upper bound on both LS_1 and LS_2 .
2. S is smooth because for any neighboring pair x and y , the difference between $S(x)$ and $S(y)$ is at most the difference between $S_1(x)$ and $S_1(y)$ and likewise for S_2 .

We apply this to our problem by letting $N_1(x)$ be the set of neighbors of x that differ in size by exactly 1, and N_2 be the set of neighbors such that the size of an establishment's employment differs by a factor of at most ϵ . For N_1 , the global and local sensitivity is 1, which is smooth. For N_2 , we must give a bound on the smooth sensitivity to complete the proof.

The local sensitivity of q_v with respect to x and N_2 is the maximum amount by which any firm's (matching the criteria in v) count of employees (matching the criteria in v) can change. Note that x_v is the largest such count, and therefore $LS_{q_v}(x) = x_v \cdot (1 + \alpha) - x_v = x_v \cdot \alpha$. Then, we have

$$A^{(j)}(x) = \max_{y \in D: d(x,y) \leq j} y_v \cdot \alpha.$$

This value is maximized by maximizing y_v . The maximum value for y_v is $x_v(1 + \alpha)^j$. Therefore, $A^{(j)}(x) = x_v \cdot \alpha(1 + \alpha)^j$. Our smooth sensitivity is therefore

$$S_{v,b}^*(x) = \max_j \left(\frac{1 + \alpha}{e^b} \right)^j x_v \alpha.$$

Our databases do not have a fixed size, so j can be any positive integer, and therefore the smooth sensitivity is not necessarily bounded. When $e^b < (1 + \alpha)$,

$$S_{v,b}^*(x) = \max_j \left(\frac{1 + \alpha}{e^b} \right)^j x_v \alpha = \lim_{j \rightarrow \infty} \left(\frac{1 + \alpha}{e^b} \right)^j x_v \alpha,$$

which is unbounded. When $e^b \geq (1 + \alpha)$, this limit is bounded, and in this case $S_{v,b}^*(x) = x_v \cdot \alpha$. \square

PROOF OF LEMMA 8.6. Let $|\lambda| \leq \frac{\epsilon_1}{\gamma+1}$. We must show that $\ln \left(\frac{e^\lambda h(e^\lambda z)}{h(z)} \right) \leq \epsilon_1$. This follows because

$$\ln \left(\frac{e^\lambda h(e^\lambda z)}{h(z)} \right) = \ln \left(\frac{e^\lambda (1 + (e^\lambda |z|)^\gamma)}{1 + |z|^\gamma} \right) \quad (38)$$

$$\leq \ln \left(\frac{(e^\lambda |z|)^\gamma}{|z|^\gamma} \right) \leq \lambda(\gamma + 1). \quad (39)$$

For the sliding property, we must show $\ln \left(\frac{h(z+\Delta)}{h(z)} \right) \leq \epsilon_2$. This follows from [38]. \square

PROOF OF LEMMAS 8.7 AND 9.2. These follow directly from Theorem 8.4 and Lemma 8.5. \square

PROOF OF LEMMAS 8.8 AND 9.3. We show that the that the random variable η drawn in each algorithm has expectation 0, and constant expected L_1 . It is well known that the Laplace mechanism is unbiased with Laplace(λ) having expected error of λ . We prove here that $h(z) \propto \frac{1}{(1+|z|^\gamma)}$ is unbiased with bounded error for $\gamma = 4$. $h(z)$ is unbiased since

$$\mathbb{E}[h(z)] = \int \frac{z}{1 + |z|^4} dz = 0.$$

The expected L_1 error of $h(z)$ is given by

$$\int \frac{|z|}{1 + |z|^4} dz = \frac{\pi}{2} \approx 1.57. \quad \square$$

C. ADDITIONAL FIGURES

δ	α	ϵ	δ	α	ϵ
.05	.01	.105	5×10^{-4}	.01	.15
.05	.10	1.01	5×10^{-4}	.10	1.45
.05	.20	1.932	5×10^{-4}	.20	2.13

Table 2: Minimum values of ϵ given α and δ

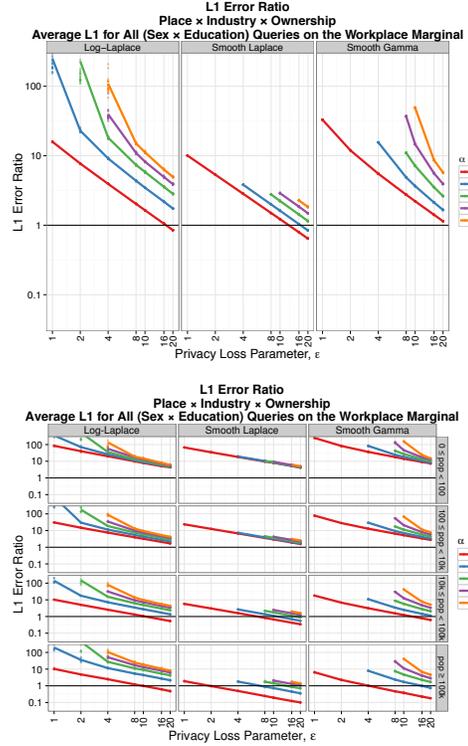


Figure 4: Average L_1 error ratio of worker and workplace attribute marginal compared to the current system.

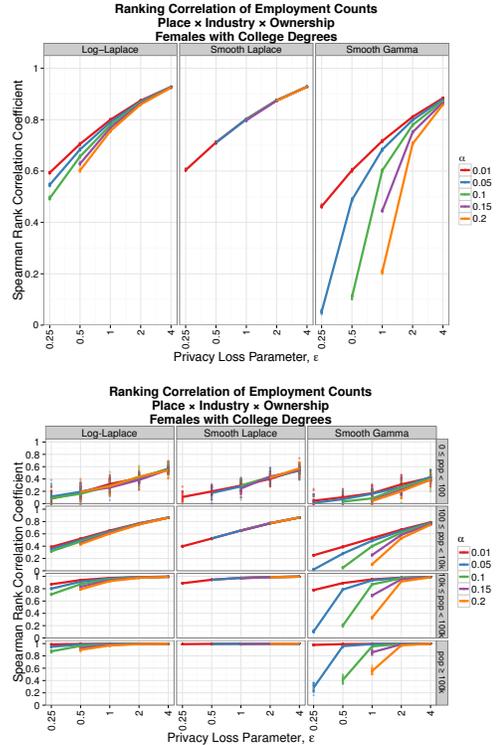


Figure 5: Spearman correlation between tested model and input noise infusion on count of female workers with bachelors degree or higher, ranked by Census place by NAICS sector (industry) by ownership.