

# Serenade: An Approach for Differentially Private Greedy Redescription Mining

Maiju Karjalainen, Esther Galbrun and Pauli Miettinen

*University of Eastern Finland*

## Abstract

We present an approach for mining differentially private redescrptions. Applying differential privacy to pattern mining approaches such as redescription mining is less straightforward than with numerical methods, in part because local patterns are more susceptible to the kind of noise inherent to differential privacy, and in part because the typically combinatorial algorithms require many operations and become computationally unaffordable when overloaded with privacy preserving procedures. Our solution, and the two algorithms we provide, address both the computational bottlenecks as well as the issues with noise. They might also provide inspiration for other differentially private pattern mining algorithms. Extensive experiments on real-world data validate the usefulness of our algorithms.

## 1. Introduction

As awareness about the rights of data subjects rises, so does the demand for privacy-preserving data mining methods. This trend, combined with concerns about black-box models being used in many domains, has lead to a need for methods that couple interpretability and privacy. Here we present algorithms answering this need, bringing differential privacy to redescription mining.

Differential privacy [1] is the most promising approach for privacy preserving data analysis. It allows to provide guarantees on the level of privacy offered, even in the presence of side information. Differentially private results are also secure against post-processing: no amount of post-processing will reveal any further private information from the results of differentially-private algorithms.

It is widely recognized that pattern and rule-based approaches, such as redescription mining [2], are highly interpretable data analysis approaches. However, integrating differential privacy with such approaches is not simple, as they are based on local properties of the data, which are more prone to privacy leakage than global properties. Moreover, the combinatorial nature of these approaches makes them typically more susceptible to noise than continuous optimization methods.


---

*KDID 2022: 20th anniversary of KDID Workshop*

✉ maiju.karjalainen@uef.fi (M. Karjalainen); esther.galbrun@uef.fi (E. Galbrun); pauli.miettinen@uef.fi (P. Miettinen)

🆔 0000-0001-5885-3143 (M. Karjalainen); 0000-0003-2585-9252 (E. Galbrun); 0000-0003-2271-316X (P. Miettinen)

© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

We tackle these obstacles and present two algorithms based on existing approaches for differential privacy and redescription mining: a greedy algorithm for redescription mining, and the exponential and Laplace mechanisms for differential privacy. Their successful combination requires developing novel methods. In particular, we present a general technique to lower the sensitivity of the quality measure to magnitudes well below 1, allowing much more accurate results to be drawn from the exponential mechanism without compromising the privacy. In addition, we use weighted reservoir sampling for memoization, yielding an efficient implementation of the exponential mechanism.

## 2. Background and Related Work

**Redescription mining.** Redescription mining is a data analysis task that aims at finding two different ways to describe almost the same observations. In our setting, the data set  $\mathcal{D}$  consists of a pair of tables, respectively denoted  $\mathbf{D}_L$  and  $\mathbf{D}_R$ . Both tables are over the same entities (rows) but disjoint collections of attributes (columns).

Such a pair of data tables, along with user-defined parameters, constitute the input of redescription mining. A redescription consists of a pair of *queries* denoted by  $q_L$  and  $q_R$ , one over either data table. A query consists of literals over Boolean, categorical or numerical attributes, combined with logical conjunction and disjunction operators, possibly involving negations. For instance, the query

$$q = (\neg \text{Hypertension} \wedge [\text{State} = \text{CA}]) \vee [\text{BirthY} \leq 1950]$$

describes entities, here individuals, who either do not suffer from hypertension and reside in California, or are born before 1950. The collection of entities described by the queries is called the *support* of the query, denoted  $\text{supp}(q)$ .

The main quality of a redescription is that the queries have similar supports. Referred to as its accuracy, it is measured using the Jaccard coefficient:

$$J(q_L, q_R) = |\text{supp}(q_L) \cap \text{supp}(q_R)| / |\text{supp}(q_L) \cup \text{supp}(q_R)| . \quad (1)$$

In addition to having a high accuracy, to be of interest a redescription should cover enough entities. This can be ensured by setting a minimum support threshold.

Since the introduction of redescription mining [3], various algorithms have been proposed for the task, including algorithms based on decision trees [3, 4, 5], on mining and combining frequent itemsets [6], or considering regression and correlation models [7]. We build on the REREMi algorithm [8], that iteratively grows the queries in a greedy fashion. Further details can be found in [2].

**Differential privacy.** Differential privacy ensures that it cannot be inferred, within a fixed probability, whether or not a particular individual is included in a data set, regardless of what other information the adversary might have. This is achieved by withholding direct access to the data and only returning randomized outputs.

Formally, if  $\mathcal{D}$  and  $\mathcal{D}'$  are two data sets that differ in exactly one row,  $A$  is a randomized algorithm operating on these tables,  $S$  is any set of potential results  $A$  might return, and

---

**Algorithm 1** SERENADECS, client algorithm with client-side decisions

---

**Input:** Privacy budget  $\varepsilon_{\text{tot}} = \varepsilon_{\text{init}} + \varepsilon_{\text{ext}} + \varepsilon_{\text{qual}}$ , number of initial pairs  $I$ , number of extensions  $L$ , minimum support  $z_o$ , minimum accuracy  $J_{\text{min}}$ , steepness  $s$

**Output:** A set of redescription together with their accuracy  $\mathcal{R} = \{(q_L, q_R, J)\}$

```
1: function SERENADECS( $\varepsilon_{\text{init}}, \varepsilon_{\text{ext}}, \varepsilon_{\text{qual}}, I, L, z_o, J_{\text{min}}, s$ )
2:    $\mathcal{R} \leftarrow \emptyset$ 
3:    $\mathcal{Q} \leftarrow \text{INITIALPAIRS}(\varepsilon_{\text{init}}/2, I, z_o, s)$ 
4:   for  $r_0 \in \mathcal{Q}$  do
5:      $\ell \leftarrow 0$ 
6:      $J_0 \leftarrow \text{COMPUTEQUALITY}(r_\ell, \varepsilon_{\text{qual}}/(L|\mathcal{Q}|))$ 
7:     while  $\ell < L$  and  $J_\ell > J_{\text{min}}$  do
8:        $\ell \leftarrow \ell + 1$ 
9:        $r_\ell \leftarrow \text{EXTENDONE}(\{r_{\ell-1}\}, \varepsilon_{\text{ext}}/(2L|\mathcal{Q}|), z_o, s)$ 
10:       $J_\ell \leftarrow \text{COMPUTEQUALITY}(r_\ell, \varepsilon_{\text{qual}}/(L|\mathcal{Q}|))$ 
11:      Add to  $\mathcal{R}$  the last  $r_\ell$  with  $J_\ell > J_{\text{min}}$ , if any
12:   return  $\mathcal{R}$ 
```

---

$\varepsilon \in [0, 1]$  is the privacy parameter (or budget), then  $A$  is  $\varepsilon$ -*differentially private* (later also simply  $\varepsilon$ -private) if

$$\Pr[A(\mathcal{D}) \in S] \leq e^\varepsilon \Pr[A(\mathcal{D}') \in S], \quad (2)$$

where the probability is over the randomness of  $A$  (see [9]).

The core mechanism for ensuring differential privacy with numerical outputs is the *Laplace mechanism*, that perturbs numerical outputs with  $Laplace(0, \Delta A/\varepsilon)$ -distributed noise, where  $\Delta A$  is the *sensitivity* of  $A$ , i.e. the maximum change in the output of  $A$  given two inputs that differ in exactly one row. For other than numerical outputs, the *exponential mechanism* [10] is used. The space of all possible outputs of  $A$  must be endowed with a data-specific *quality function*  $f(A(\mathcal{D}))$ , and the answer is sampled from this space with probability  $\propto \exp(\varepsilon \cdot f(A(\mathcal{D})))$ . This ensures  $2\varepsilon\Delta f$ -privacy [10]. Further developments in differential privacy include concentrated differential privacy [11], which trades better behaviour under composition to weaker definition of privacy. Our method is trivially  $(\varepsilon \cdot (e^\varepsilon - 1)/2, \varepsilon)$ -concentrated differentially private [11], but further analysis of concentrated differential privacy is beyond the scope of this paper.

After the initial formulation [1] and the introduction of the exponential mechanism [10], differential privacy has seen significant interest in academia, and recent years have brought uptake by the industry as well (e.g. [12]). Real-world use of differential privacy has indicated that the typical assumption that  $\varepsilon < 1$  is often unachievable (e.g. the US Census Bureau used  $\varepsilon = 19.61$  for releasing the 2020 US Census results [13]). Differentially private algorithms have been proposed for mining frequent itemsets [14, 15], association rules [16, 17], as well as subgraphs [18] or sequences [19].

---

**Algorithm 2** SERENADEES, client algorithm with private extensions

---

**Input:** Privacy budget  $\varepsilon_{\text{tot}} = \varepsilon_{\text{init}} + \varepsilon_{\text{ext}} + \varepsilon_{\text{qual}}$ , number of initial pairs  $I$ , number of extensions  $K$ , minimum support  $z_o$ , steepness  $s$

**Output:** A set of redescrptions together with their accuracy  $\mathcal{R} = \{(q_L, q_R, J)\}$

```
1: function SERENADEES( $\varepsilon_{\text{init}}, \varepsilon_{\text{ext}}, \varepsilon_{\text{qual}}, I, K, z_o, s$ )
2:    $\mathcal{Q} \leftarrow \text{INTIALPAIRS}(\varepsilon_{\text{init}}/2, I, z_o, s)$ 
3:   for  $k = 1, \dots, K$  do
4:      $r \leftarrow \text{EXTENDONE}(\mathcal{Q}, \varepsilon_{\text{ext}}/2K, z_o, s)$ 
5:      $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{r\}$ 
6:     Remove extended redescription from  $\mathcal{Q}$ 
7:    $\mathcal{R} \leftarrow \emptyset$ 
8:   for all  $r \in \mathcal{Q}$  do
9:      $\mathcal{R} \leftarrow \mathcal{R} \cup \{(r, \text{COMPUTEQUALITY}(r, \varepsilon_{\text{qual}}/|\mathcal{Q}|))\}$ 
10:  return  $\mathcal{R}$ 
```

---

### 3. The Algorithms

In this section we present two algorithms for differentially private redescription mining: SERENADECS and SERENADEES. Both algorithms are based on the greedy redescription mining approach of REREMi and assume that one side of the data has only Boolean and categorical attributes, a common assumption in redescription mining (see, e.g. [8]).

Both algorithms consist of a client with no direct access to the data and a differentially private engine, with full access to the data, that the client queries. The two algorithms differ mainly on their client-side implementations, i.e. how they use the differentially private engine's methods. With SERENADECS, more of the decisions are made on the client side (CS), resulting in the ability to prune low-quality results at the cost of an increased budget expenditure. SERENADEES provides a thriftier alternative by letting the differentially private engine make the decisions (engine side, ES).

We will next explain how the clients work before explaining how the differentially private part has been implemented.

#### 3.1. Client side

The pseudocode for the algorithms is given in Algorithms 1 and 2. Both algorithms mine redescrptions through three main steps. The INTIALPAIRS method (Algorithm 4) returns redescrptions with a single literal in both queries. These pairs are then iteratively extended into longer redescrptions. The EXTENDONE method (Algorithm 3) takes a set of redescrptions as input and returns a new redescription, obtained by extending one of the input redescrptions by one literal. The COMPUTEQUALITY method computes the noisy Jaccard and support size of a redescription. The two algorithms differ on when the quality of redescrptions is computed and how to select the next redescription to extend.

SERENADECS closely follows the approach of REREMi: it considers one initial pair at a time (line 4 in Algorithm 1) for at most a fixed number of extensions (line 7). In each iteration, SERENADECS calls EXTENDONE with the current candidate as input (line 9), queries COMPUTEQUALITY to get the noisy Jaccard of the extension (line 10), and

---

**Algorithm 3** EXTENDONE, extending a redescription

---

**Input:** Data  $\mathcal{D}$ , collection of redescriptions  $\mathcal{Q}$ , privacy budget  $\varepsilon$ , min. support  $z_o$ , steepness  $s$

**Output:** Extended redescription  $(q'_L, q'_R)$

**Private information:** Data  $\mathcal{D} = (\mathbf{D}_L, \mathbf{D}_R)$ , sensitivity  $\Delta f$  given  $\mathcal{D}$ ,  $z_o$ , and  $s$

**Static variables:** Associative array  $M$ , priority queue  $P$

```
1: function EXTENDONE( $\mathcal{Q}, \varepsilon, z_o, s$ )
2:   for all  $(q_L, q_R) \in \mathcal{Q} \setminus M$  do
3:     for all one-literal extensions  $(q'_L, q'_R)$  do
4:       select the extensions with the highest key  $key(f(q_L, q_R, z_o, s), \varepsilon/\Delta f)$ 
5:        $M[(q_L, q_R)] \leftarrow (q'_L, q'_R)$  with age 1
6:        $P.push(k, (q_L, q_R))$ 
7:    $k, (q_L, q_R) \leftarrow P.pop()$ 
8:    $(q'_L, q'_R) \leftarrow M[(q_L, q_R)]$ ; delete  $M[(q_L, q_R)]$ 
9:   increment the age of every  $(q_L, q_R) \in M$  and remove expired from  $M$  and  $P$ 
10:  return  $(q'_L, q'_R)$ 
```

---

decides whether to continue with the extension process. This ensures we only attempt to extend redescriptions that are good enough, but it has two drawbacks: *i*) it needs to call COMPUTEQUALITY for every extension, using plenty of budget; *ii*) it might attempt to extend a redescription which has no good extensions, effectively wasting the associated budget. The number of redescriptions returned by SERENADECS is not predetermined, though it is never larger than the number of initial pairs  $I$ .

SERENADEES addresses these drawbacks of SERENADECS. At each of a fixed number of extensions, it sends *all* current candidate redescriptions that can be extended to EXTENDONE and receives one extension back (line 4 in Algorithm 2). The decision of which redescription to extend is made in the differentially private engine. Furthermore, the Jaccard is only computed for the redescriptions that are returned (line 9). This approach saves on the privacy budget compared to SERENADECS. It also allows some redescriptions to be extended multiple times while not wasting budget on redescriptions that cannot be extended. On the flip side, EXTENDONE might return a bad redescription, and this will only become apparent when it is too late to replace the candidate. SERENADEES always returns as many redescriptions as there were initial pairs.

### 3.2. Differentially Private Engine

Next, we present the three methods of the differentially private engine. We start with a general overview and description of the main ingredients of our algorithms.

**Reservoir sampling.** Greedy extension is a standard strategy for building redescriptions [2]. Extending redescription  $(q_L, q_R)$  consists in appending a literal to either of the queries using a disjunction or conjunction. The extension with the highest accuracy that satisfies constraints (e.g. minimum support threshold) is reported. For numerical attributes, all relevant intervals are evaluated in turn. A pruning trick [8] allows to identify relevant intervals, i.e. those that impact the accuracy, and avoid testing those

---

**Algorithm 4** INTIALPAIRS, generating initial candidates

---

**Input:** Privacy budget  $\varepsilon$ , number of initial pairs  $I$ , minimum support threshold  $z_o$ , steepness  $s$

**Output:** A set of initial redescrptions  $\mathcal{R} = \{(q_L, q_R)\}$

**Private information:** Data  $\mathcal{D} = (\mathbf{D}_L, \mathbf{D}_R)$ , sensitivity  $\Delta f$  given  $\mathcal{D}$ ,  $z_o$  and  $s$

```
1: function INTIALPAIRS( $\varepsilon, I, z_o, s$ )
2:    $\mathcal{R} \leftarrow \emptyset$ 
3:   Set  $S_i$  as size-1 reservoir sampler for  $i = 1, \dots, I$ 
4:   for all literals  $x$  over the non-numerical attributes of  $\mathbf{D}_L$ , numerical attributes  $y$  in  $\mathbf{D}_R$ ,
   and meaningful intervals  $[\lambda, \rho]$  of  $y$  do
5:     Add initial pair  $(x, [\lambda \leq y \leq \rho])$  to all reservoir samplers  $S_i$  where  $key(f(x, [\lambda \leq y \leq \rho]), z_o, s), \varepsilon/(I\Delta f))$  is greater than the current key
6:   for all reservoirs  $S_i$  do
7:     Add redescription from  $S_i$  to  $\mathcal{R}$ 
8:   return  $\mathcal{R}$ 
```

---

that do not (see Section 3.3).

A differentially private variant of such an extension strategy is best implemented with the exponential mechanism. Instead of reporting the best extension, the algorithm will randomly pick one of the candidates with probabilities proportional to their accuracies. The set of candidates will also include the current, unextended, redescription associated to a sampling probability reflecting its quality. That way, the structure of the queries will not reveal anything sensitive about the data (see Section 3.3 about the case of continuous-valued attributes).

A naïve approach would be to store all candidates, but this would use too much memory. Instead, we use *weighted reservoir sampling* [20, 21]. The idea is to keep a reservoir of size  $k$  and evaluate candidate extensions in turn. For each one, we sample a value  $u \sim Unif_{[0,1]}$ . If the probability of the candidate being selected by the exponential mechanism is  $\alpha$ , we compute its key as  $u^{1/\alpha}$  and store the candidate in the reservoir if this key is larger than the current smallest key in the reservoir (see [20, 21] for proofs of correctness). That is, for each candidate  $(q_L, q_R)$  we compute

$$key(f(q_L, q_R, \Theta), \varepsilon) = \log(Unif_{[0,1]}) \cdot \exp(\varepsilon \cdot f(q_L, q_R, \Theta)) , \quad (3)$$

where  $f(q_L, q_R, \Theta)$  is a function measuring the quality of  $(q_L, q_R)$  given parameters in  $\Theta$ , and retain the candidates with highest keys.

**Extending a redescription.** When called from SERENADEES, EXTENDONE (Algorithm 3) receives a collection of redescrptions  $\mathcal{Q}$ . In each call, the method samples from the distribution of all extensions of all redescrptions in  $\mathcal{Q}$ , and one selected candidate replaces the extended redescription. This proceeds for a fixed number of iterations, progressively extending and replacing redescrptions in  $\mathcal{Q}$ . Between two successive calls to EXTENDONE, only one redescription changes.

Our sampling approach also allows us to first draw a large sample and then take a smaller sample from it, by simply retaining the entries associated with the largest keys. As the collection  $\mathcal{Q}$  can be large, we use this strategy to save computational time.

Specifically, EXTENDONE stores for each redescription encountered so far the sampled extension for that redescription, together with its key. When a new redescription is added to  $\mathcal{Q}$ , EXTENDONE samples an extension for it and stores it. The redescription whose sampled extension has the largest key is selected and the corresponding extended redescription is returned.

This memoization can lead to suboptimal choices if an extension of a redescription gets a very low key value even though it is relatively good, and vice versa. To counteract that, EXTENDONE prunes stored extensions that are too old, i.e. extensions that have not been selected in the past extension rounds. The maximum age of an extension is a user-supplied parameter. When called from SERENADECS, EXTENDONE operates on a single redescription and returns its candidate extension having the highest key, without any memoization.

**Computing initial pairs and redescription quality.** INITIALPAIRS (Algorithm 4) builds initial pairs by enumerating all pairs of attributes  $(x, y)$ , where  $x$  is an attribute from  $\mathbf{D}_L$  and  $y$  is an attribute from  $\mathbf{D}_R$ . The selection of the initial pairs from all pairs is based on a reservoir sampling approach. Assuming w.l.o.g. that the non-numerical attributes are in  $\mathbf{D}_L$ , we consider as initial query  $q_L$  every literal consisting of a Boolean attribute or an attribute–category pair. For  $q_R$ , we consider every literal (including all relevant intervals) from  $\mathbf{D}_R$  to match  $q_L$ , in the same way as when extending an existing redescription. The sampling is done over the queries  $q_R$ . In order to generate  $I$  initial pairs, we populate  $I$  different size-1 reservoirs, corresponding to  $I$  independent samples with replacement from the exponential mechanism, ensuring privacy.

COMPUTEQUALITY computes the size of the intersection and the size of the union of the input queries’ supports and returns them after applying Laplace noise. The noisy Jaccard of the redescription is computed as the noisy intersection size divided by the noisy union size.

### 3.3. Handling Continuous Attributes

Continuous attributes require extra care in both computational efficiency and privacy. When appearing in a query, a continuous attribute  $y \in \mathbb{R}$  will be bounded within an interval, as  $[\lambda \leq y \leq \rho]$  (where either  $\lambda$  or  $\rho$  can be infinite). Finding these intervals efficiently is the first problem. We follow the cut-point-based approach of REREMi [8]. This approach is based on the fact that only some values of  $y$  will have an effect on the support of the extended query when used as threshold. If an entity is not in  $\text{supp}(q_R)$ , whether or not it belongs to  $\text{supp}([\lambda \leq y \leq \rho])$  is indifferent w.r.t. the support of the extended query. Hence we only need to consider as thresholds those values that determine whether entities in  $\text{supp}(q_R)$  are included in  $\text{supp}([\lambda \leq y \leq \rho])$ . These values, along with  $-\infty$  and  $\infty$ , are called *cut-points*.

This approach extends to disjunctions and negations. It also helps us with privacy. Returning the exact value of cut-points can leak information about the data, as it reveals



that  $y$  takes this exact value at least once.<sup>1</sup> Let us assume that the lower cut-points and upper cut-points are  $-\infty = \lambda_1 < \lambda_2 < \dots < \lambda_l$  and  $\rho_1 < \rho_2 < \dots < \rho_p = \infty$ , respectively, and that we have chosen to use interval  $[\lambda_i \leq y \leq \rho_j]$ . Instead of reporting  $\lambda_i$ , we sample the value  $\lambda'_i$  we report as lower threshold from  $Unif_{(\lambda_{i-1}, \lambda_i]}$  and instead of  $\rho_j$ , we sample a value  $\rho'_j$  from  $Unif_{[\rho_j, \rho_{j+1})}$ . Note that if  $i = 1$  or  $j = p$ , we actually have a half-interval, which we write  $[y \leq \lambda'_i]$  or  $[\rho'_j \leq y]$ , and no sampling is necessary for that threshold. This sampling ensures that the accuracy of the redescription remains the same while no information is leaked about the exact values in the data.

### 3.4. Lower-Sensitivity Quality Function

A key element of the algorithm is the distribution of the budget between the steps. Both INTIALPAIRS and EXTENDONE use the exponential mechanism, which provides  $2\epsilon\Delta f$ -privacy for quality  $f$  of sensitivity  $\Delta f$ . For our application, the standard quality function is the Jaccard coefficient, which has sensitivity 1, since it takes values between 0 and 1. However, for interesting redescriptions, the support size is well above 0 and, as a result, adding or removing a single row will not change the Jaccard coefficient much. This motivates us to modify the quality function to take into account the support size. Specifically, we multiply the Jaccard coefficient with a logistic function centered at the user-defined minimum support threshold  $z_o$ . The quality function  $f$  becomes

$$f(q_L, q_R, z_o, s) = J(q_L, q_R) / (1 + \exp(-s \cdot (|\text{supp}(q_L) \cap \text{supp}(q_R)| - z_o))) , \quad (4)$$

where  $s$  is a parameter adjusting the steepness of the logistic curve.

For large values of  $s$  the sensitivity of  $f$  will still be 1. For smaller values of  $s$ , however, it is possible to significantly reduce the sensitivity of  $f$ . Given  $z_o$ ,  $s$ , and data set size  $n$ , the sensitivity of  $f$  can be calculated in  $O(n^2)$  time by considering all possible values for numerator and denominator in (1), plugging them into (4), and identifying the largest change that can result from the addition of a single row. Furthermore, the steepness can be optimized with a simple line search over different sensitivity parameters. With the data sets and parameters used in our experiments, the sensitivity was always less than 0.001, leaving more budget to be spent on initial pairs and extensions.

Notice that knowing the sensitivity releases information about the data as it depends on the data set size. This leaves a few options. The first is that the differentially private engine computes the optimal steepness and uses the corresponding sensitivity when distributing the budget. This approach is completely private, but means that the user has no knowledge of the actual noise levels used. The other extreme is to use a weaker form of differential privacy, where the user knows the size  $n$  of the data set. The user can then find the optimal  $s$  minimizing the sensitivity of (4). As an intermediate option, the user might query for  $n$ , use the noisy estimate to compute near-optimal steepness, and give that as a parameter to the algorithm. This way, at the cost of some privacy budget, the user retains some knowledge of the quality function and knows its steepness, but no privacy is lost.

---

<sup>1</sup>The fact that categorical attributes take each category at least once is not considered private information.



### 3.5. Other Parameters

The two key parameters of the algorithms are the number of initial pairs,  $I$ , and the number of extensions,  $K$  or  $L$ . The higher  $I$  is set, the more redescrptions will be returned, but each of them will receive a smaller portion of the privacy budget, both when generating the initial pairs and when computing the privatized supports and accuracies. On the other hand, the number of extensions ( $K$  or  $L$ ) dictates how complex the redescrptions can become. SERENADECS will return no more than  $I$  redescrptions, and none of them will have more than  $L$  extensions (i.e. they have at most  $L + 2$  literals in total). SERENADEES will return exactly  $I$  redescrptions, but some of them might have just two literals, or some of them might have  $K + 2$  literals. Hence, if SERENADECS is set to perform  $L$  extensions, a comparable setting for SERENADEES is  $K = I \cdot L$  extensions.

### 3.6. Privacy

The privacy of the proposed approach is mainly based on the exponential mechanism. For this, we need to sample from every possible extension with the correct probability. Next we show that this is the case. We then continue by analysing the use of the budget, and conclude that the methods are indeed private.

**Lemma 1.** *Given a redescription  $r_{\ell-1} = (q_L, q_R)$ , the probability that EXTENDONE in line 9 of Algorithm 1 returns an extension  $r_\ell = (q'_L, q'_R)$ , where either i)  $q'_L$  is an arbitrary single-literal extension of  $q_L$  and  $q_R = q'_R$ , ii)  $q'_R$  is an arbitrary single-literal extension of  $q_R$  and  $q_L = q'_L$ , or iii)  $q_L = q'_L$  and  $q_R = q'_R$ , is  $\propto \exp(\varepsilon \cdot f(q'_L, q'_R, z_o, s))$ .*

*Proof.* We can assume w.l.o.g. that  $q_L = q'_L$ . Assume first that  $\mathbf{D}_R$  does not contain continuous attributes. To select the extension of  $q_R$ , EXTENDONE can be considered to go through every extension, computing their quality according to  $f$ , and placing the extension into a stream together with a weight parameter that is set to  $\exp(\varepsilon \cdot f)$ . Finally, the algorithm will also place the unextended  $q_R$  into the stream with weight corresponding to the quality of  $(q_L, q_R)$ . The weighted reservoir sampler will select one item from this stream, so that the probability of an item being selected is proportional to its weight [21, 20], thus sampling the extension with correct probability.

For continuous attributes, the algorithm further perturbs the boundaries as explained in Section 3.3. As this does not affect the quality, all extensions with the same quality have equal probability to be selected, concluding the proof.  $\square$

**Lemma 2.** *The extended redescription  $r = (q'_L, q'_R)$  returned by EXTENDONE when called in line 4 of Algorithm 2 is selected with probability proportional to  $\exp(\varepsilon \cdot f(q'_L, q'_R, z_o, s))$  from all possible at-most-single literal extension to all redescrptions in  $\mathcal{Q}$ .*

**Lemma 3.** *Given a pair of single-literal queries  $(q_L, q_R)$  over data sets  $\mathbf{D}_L$  and  $\mathbf{D}_R$ , the probability that INTIALPAIRS returns  $(q_L, q_R)$  is  $\propto I \cdot \exp(\varepsilon \cdot f(q_L, q_R, z_o, s))$ .*

**Lemma 4.** *The total budget used by SERENADECS and SERENADEES never exceeds  $\varepsilon_{tot}$ .*

The proof for Lemma 3 follows the same lines as for Lemma 1, except that the sampling is repeated  $I$  times, and is omitted. The proofs for Lemma 2 and Lemma 4 are postponed to the full version of the paper. <sup>2</sup>

**Theorem 5.** *SERENADECS and SERENADEES are  $\varepsilon_{tot}$ -differentially private.*

*Proof.* Per Lemma 3, both algorithms generate initial pairs privately. Each initial pair is an independent query to the data, spending altogether budget  $\varepsilon_{init}$ . Every extension in SERENADECS is an independent, differentially private query (Lemma 1). Due to the composability of differential privacy [9], they do not reveal more information than permitted by their allocated budget. The same holds for extensions in SERENADEES. Finally, quality calculations in COMPUTEQUALITY use the standard Laplace mechanism and are private, in line with their allocated budget. The total budgets do not exceed the allocated  $\varepsilon_{tot}$  (Lemma 4).  $\square$

## 4. Experimental Evaluation

To test our proposed methods, we conducted both quantitative and qualitative experiments. The purpose of the quantitative experiments is to investigate the behavior of the algorithms with different parameters and carry out an ablation study. The qualitative experiments are meant to assess the usefulness of our approach for real-world data analysis. For this reason, we did not do any parameter tuning as that would consume some of the privacy budget.

All algorithms are implemented in Python and the experiments are run with Python 3.6.8 on a machine with 2 AMD EPYC 7702 processors with 64 cores each and 1 TB of main memory.

### 4.1. Data and Other Methods

The `Mammals` data set contains information about which mammal species inhabit which areas of Europe on one side [22], and climate information on the other side [23]. This data set has often been used in redescription mining literature (e.g. [8, 24, 4, 25, 5, 26]). The data contains 194 species and 48 climate attributes as its columns, with 2575 localities as its rows.

The MIMIC data set [27]<sup>3</sup> is an example of health data where patient privacy is critical. It contains de-identified health information about hospital patients. Our dataset contains two views: the diagnoses view and the lab events view. The aim when mining this data set is to find out what kind of lab events are associated with various diagnoses of the examined patients. The data set contains 107 Boolean attributes in the diagnoses view (1 denotes a positive diagnosis) and 104 ternary attributes in the lab events view, with 46 065 patients as its rows.

---

<sup>2</sup>The source code and the appendix containing all details that are postponed to the full version of the paper are presented in <https://github.com/maijuka/Serenade>

<sup>3</sup>Data set available at <https://physionet.org/content/mimiciii/1.4/>.

The VAA data set [28] contains the background information and answers to questions regarding political opinions of candidates to the Finnish parliamentary elections of 2011. The data was collected by an online voting advice application (VAA). We removed candidates who did not provide answers to all questions, leaving 1656 candidates, 9 background attributes and 107 opinion attributes.

The Open Sex-Role Inventory [29] (SR data set) is used to study gender roles by measuring masculinity and femininity. The Right-Wing Authoritarianism Scale [30] (RightW data set) was developed to measure the psychology of the followers of fascist regimes.<sup>4</sup> Columns recording to the time spent answering the questions were removed from the RightW and the SR data sets. These data sets have the same variables on both sides; INTIALPAIRS and EXTENDONE are implemented to ensure the same attribute cannot be used in both queries of the same redescription.

We also conducted experiments with other datasets. The results were not significant and are postponed to the full version of the paper.

The main studied algorithms are SERENADECS and SERENADEES. In addition, for ablation studies, we used a variation of SERENADEES called NOSENS. It works like SERENADEES, but uses the standard Jaccard (1) instead of the lower-sensitivity version (4) for computing keys (3). We also did a comparison to REREMI, the baseline non-private redescription mining algorithm [8].

For all experiments, we set  $J_{\min} = 0.3$  for SERENADECS and  $z_o$  to 10% of data set size  $n$ .

## 4.2. Quantitative Experiments

In this empirical evaluation, we investigate the effects of parameters and carry out an ablation study using NOSENS.

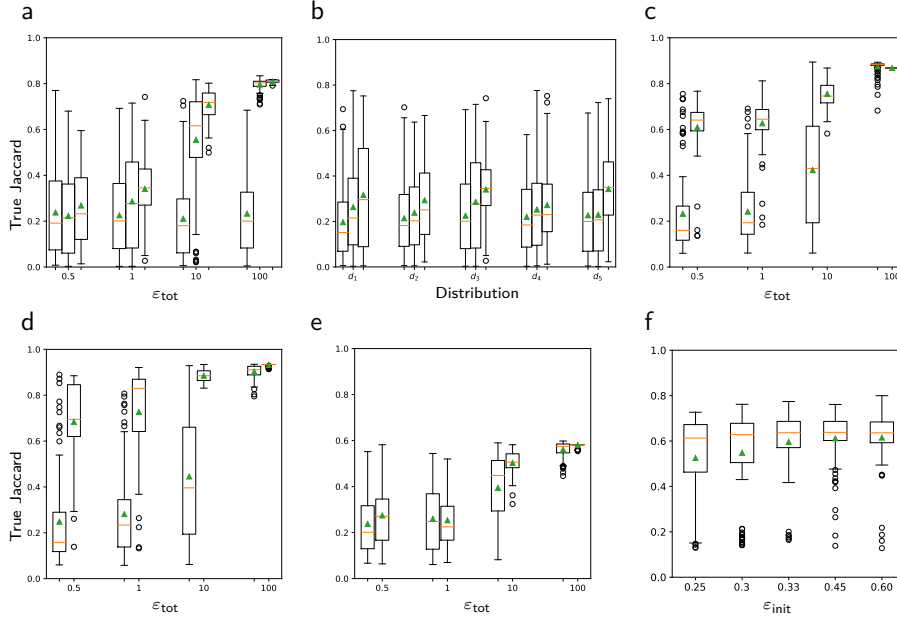
**Privacy budget.** A typical value for the privacy budget is  $\varepsilon_{\text{tot}} = 1$ , though much higher total budgets are sometimes used (e.g. [31, 32, 13]). Hence we first studied the effect of varying the total budget  $\varepsilon_{\text{tot}}$  along with different ways to distribute it between the different tasks. Total budgets of 0.5, 1, 10, and 100 were used. Of these,  $\varepsilon_{\text{tot}} = 0.5$  and  $\varepsilon_{\text{tot}} = 1$  are typical values found in the literature,  $\varepsilon_{\text{tot}} = 10$  is closer to values sometimes used in practice [13], and  $\varepsilon_{\text{tot}} = 100$  can be seen as a baseline of how the algorithms perform when privacy requirements are lifted.

Another parameter to consider is the distribution of the total budget between  $\varepsilon_{\text{init}}$ ,  $\varepsilon_{\text{ext}}$ , and  $\varepsilon_{\text{qual}}$ . We used five different ways of distributing  $\varepsilon_{\text{tot}}$  between  $\varepsilon_{\text{init}}$ ,  $\varepsilon_{\text{ext}}$ , and  $\varepsilon_{\text{qual}}$ , respectively in the following proportions:  $d_1 = (0.33, 0.33, 0.33)$ ,  $d_2 = (0.45, 0.45, 0.1)$ ,  $d_3 = (0.6, 0.3, 0.1)$ ,  $d_4 = (0.3, 0.6, 0.1)$ , and  $d_5 = (0.25, 0.25, 0.5)$ . The maximum age of an extension was set to 40,  $I = 20$ ,  $L = 4$ , and  $K = 80$ . Each experiment was repeated five times.

Results for these experiments on the Mammals data set are presented in Fig. 1(a-b) where we see that both SERENADECS and SERENADEES are capable of returning very

---

<sup>4</sup>Data sets available at [https://openpsychometrics.org/\\_rawdata/](https://openpsychometrics.org/_rawdata/).



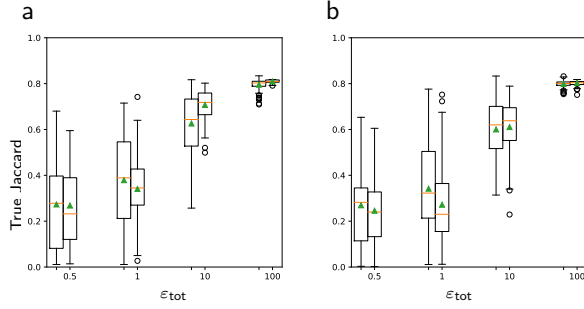
**Figure 1:** Effects of the budget on the accuracy.  $y$ -axis is true Jaccard.  $x$ -axis: (a) Varying  $\epsilon_{\text{tot}}$ , budget distribution  $d_3$ . (b) Varying budget distribution,  $\epsilon_{\text{tot}} = 1$ . Boxes represent (left to right) NOSENS, SERENADEES, and SERENADECS. (c–e) same as in (a) with data sets (c) SR, (d) RightW and (e) VAA, but without NOSENS. (f) Accuracy of the initial pairs, varying  $\epsilon_{\text{init}}$ , with  $\epsilon_{\text{tot}} = 1$ . The horizontal bar: median, triangle: mean.

good results when allocated sufficiently high budgets. On the other hand, NOSENS struggles even when allocated very high budgets. Given its inefficient use of the budget, this is not surprising. Results for similar experiments on three other data sets are presented in Fig. 1(c–e) where we see that SERENADECS returns clearly better results with almost all parameters for these datasets. We can also see that VAA (Fig 1(e)) yielded poor results regardless of budget and algorithm. This is explained by the fact that the non differentially private redescription mining algorithm RREMI<sup>5</sup> also fails to mine good quality redescriptions from this data set. These results are postponed to the full version of the paper. From Fig. 1(f) we see that we get slightly better initial pairs with a higher budget, which could be the reason why the distributions  $d_3$  and  $d_2$ , having the highest budgets for initial pairs, perform slightly better than the other budget distributions.

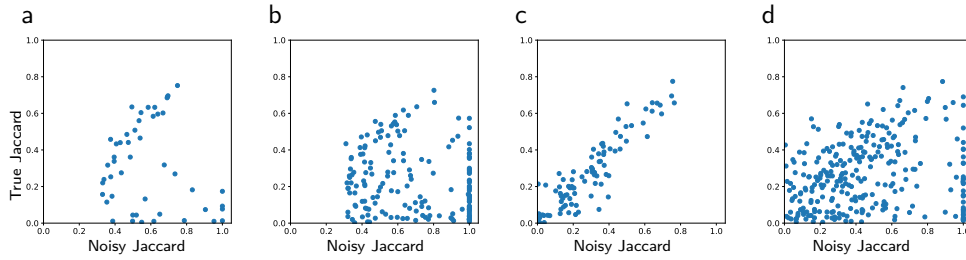
SERENADEES typically has higher variance on the accuracy of the results, giving accuracies both higher and lower than SERENADECS. This is in line with the expected behavior. Of the different budget distributions,  $d_3$  seems to be consistently giving the best results, although the differences are not significant.

An analyst is typically not interested in results with very low reported accuracy (noisy Jaccard). Figure 2 repeats Fig. 1(a–b) for SERENADEES and SERENADECS after removing redescriptions with noisy Jaccard below 0.3. The average and median for SERENADEES

<sup>5</sup>Source code available at <https://pypi.org/project/python-siren/> [8]



**Figure 2:** Effects of pruning results with low noisy Jaccard. (a) Varying  $\epsilon_{\text{tot}}$ , budget distribution  $d_3$ . (b) Varying  $\epsilon_{\text{tot}}$ , budget distribution  $d_1$ . Left box: SERENADEES, right box: SERENADECS. The horizontal bar: median, triangle: mean.



**Figure 3:** Effects of the number of initial pairs. The  $y$ -axis represents the true Jaccard whereas the  $x$ -axis represents the noisy Jaccard. (a) SERENADECS and  $I = 20$ . (b) SERENADECS and  $I = 80$ . (c) SERENADEES and  $I = 20$ . (d) SERENADEES and  $I = 80$ .

improve, and especially the lower end of the tail rises, as expected. We can also see that this pruning leaves some redescrptions with *true* accuracy below 0.3, but does not actually discard many of those reported to be good. The discrepancy between the reported and true accuracy is an important feature: if the system reports a redescription as having a high accuracy, the user should be able to trust that it truly is so.

**Number of initial pairs and reliability of noisy Jaccards.** Next we investigate the effects of the number of initial pairs. For these tests, the maximum age was always set to be twice the number of initial pairs, and the number of extension rounds was four times the number of initial pairs. Budget distribution  $d_1$  was used throughout. Recall that more initial pairs means the budget is divided between more redescrptions. This shows especially strongly in the differences between the true and reported accuracies (Figure 3).

Figure 3 also shows that SERENADECS returns fewer redescrptions than SERENADEES and that it occasionally reports overly optimistic Jaccard values even with  $I = 20$ . This is due to its mediocre use of the budget, leaving just slivers for relevant quality computations. Note that reported accuracies for SERENADECS are cut above  $J_{\text{min}} = 0.3$  through filtering. With 80 initial pairs, reported Jaccards become less correlated with true Jaccards for both methods, as budget per redescription dwindles.

**Maximum age.** The final quantitative study probes the effects of the maximum age memoization in EXTENDONE in SERENADEES. We investigated this by setting  $K = 80$  and trying maximum ages 0 (computing extensions anew in every iteration), 5, 10, 20, 40, 60, and 81 (stored extensions never expire). We found that the memoization has a clear beneficial effect on the running time. Using maximum age 20 brings the running time down to about 25 min, from over 70 min without memoization (maximum age 0). Most redescrptions are selected for extension within about 20 steps, so the running time does not improve noticeably beyond that. On the other hand, the maximum age did not seem to have a significant effect on the quality.

### 4.3. Qualitative Experiments

Next, we present a few results mined from different datasets, as anecdotal examples of results one can expect by our algorithms with these kinds of data sets. We did five runs of the algorithms for the MIMIC, VAA, RightW, SR, SSC, Psycho, and Pref data sets with total budget 1 and budget distribution  $d_3$ . The results were chosen over five runs for both algorithms, so the total budget per data set is 10. The results with VAA, SSC, Psycho and Pref datasets, and the raw redescrptions are postponed to the full version of the paper.

An example result  $(q_L, q_R)$  from MIMIC with SERENADECS is

$$(\text{Atrial Fibrillation}, \text{Uric Acid} = 1 \vee \text{Creatine Kinase} = 0),$$

meaning that the set of patients who had atrial fibrillation (a type of arrhythmia, abnormal heart rhythm) corresponds to the set of patients whose uric acid (waste product in blood) test results were abnormal or whose creatine kinase (enzyme related to energy production) test results returned normal values. The reported accuracy of this redescription is 0.424 and its true accuracy is 0.227. Such results can be used to understand the procedures at this hospital, and as they are privatized, the results can also be communicated more widely without putting the privacy of patients in jeopardy.

A result with SERENADEES from the RightW data set tells that the set of participants who moderately or strongly agreed that “You have to admire those who challenged the law and the majority’s view by protesting for women’s abortion rights, for animal rights, or to abolish school prayer” corresponds to the set of participants who did not vote. The reported accuracy of this redescription is 0.428 and its true accuracy is 0.538.

A result from the SR data set with SERENADECS the set of participants who agreed or slightly agreed that they are the happiest when they are in their bed corresponds to the set of participants who agreed or slightly agreed that they give people handmade gifts and are between 14 and 32 years of age. The reported accuracy of this redescription is 0.651 and its true accuracy is 0.725.

## 5. Conclusions

It is the nature of differential privacy and local pattern mining that the user can never be entirely certain whether a pattern reported by the algorithm actually represents a strong pattern in the data. But careful selection of the parameters allows to avoid most of the

pitfalls, and when a privacy budget higher than 1 is available, the problem essentially vanishes. Even with a privacy budget of at most 1, the results can be very useful. Of the two approaches presented, SERENADECS gives results with higher (true and reported) Jaccard, but with lower correlation between the two, meaning that a good-looking result from SERENADECS might turn out not to be so good in reality. SERENADEES, on the other hand, returns more truthful Jaccards, thanks to its better use of budget, but its extension method might lead to generally somewhat inferior results.

Overall, the reported strong patterns are reliably present in the data. This means that a user who obtains interesting strong patterns via these differentially private approaches can seek further means to confirm them, for example by applying for a research permit for the data set or by conducting a separate confirmatory study. The techniques presented in this paper might also be used to implement other differentially private pattern mining algorithms in the future.

## References

- [1] C. Dwork, F. McSherry, K. Nissim, A. Smith, Calibrating noise to sensitivity in private data analysis, in: TCC'06, 2006, pp. 265–284.
- [2] E. Galbrun, P. Miettinen, Redescription Mining, Springer, 2018.
- [3] N. Ramakrishnan, D. Kumar, B. Mishra, M. Potts, R. F. Helm, Turning CARTwheels: An alternating algorithm for mining redescription, in: KDD'04, 2004, pp. 266–275.
- [4] T. Zinchenko, E. Galbrun, P. Miettinen, Mining Predictive Redescriptions with Trees, in: ICDM'15 Workshops, 2015, pp. 1672–1675.
- [5] M. Mihelčić, S. Džeroski, N. Lavrač, T. Šmuc, A framework for redescription set construction, Expert Syst. Appl. 68 (2017) 196–215.
- [6] A. Gallo, P. Miettinen, H. Mannila, Finding subgroups having several descriptions: Algorithms for redescription mining, in: SDM'08, 2008, pp. 334–345.
- [7] F. I. Stamm, M. Becker, M. Strohmaier, F. Lemmerich, Redescription model mining, in: KDD'21, 2021, pp. 1521–1529.
- [8] E. Galbrun, P. Miettinen, From black and white to full color: Extending redescription mining outside the Boolean world, Stat. Anal. Data Min. 5 (2012) 284–303.
- [9] C. Dwork, A. Roth, The algorithmic foundations of differential privacy, Found. Trends Theor. Comput. Sci. 9 (2014) 211–407.
- [10] F. McSherry, K. Talwar, Mechanism design via differential privacy, in: FOCS'07, 2007, pp. 94–103.
- [11] C. Dwork, G. N. Rothblum, Concentrated differential privacy, arXiv preprint arXiv:1603.01887 (2016).
- [12] B. Ding, J. Kulkarni, S. Yekhanin, Collecting Telemetry Data Privately, in: NIPS'17, 2017, pp. 3571–3580.
- [13] US Census Bureau, Census bureau sets key parameters to protect privacy in 2020 census results, 2021. URL: <https://www.census.gov/newsroom/press-releases/2021/2020-census-key-parameters.html>.



- [14] C. Zeng, J. F. Naughton, J.-Y. Cai, On differentially private frequent itemset mining, *Proc. VLDB Endowment* 6 (2012) 25–36.
- [15] J. Lee, C. W. Clifton, Top-k frequent itemsets via differentially private FP-trees, in: *KDD '14*, 2014, pp. 931–940.
- [16] Y.-T. Tsou, H. Zhen, X. Jiang, Y. Huang, S.-Y. Kuo, DPARM: Differentially private association rules mining, *IEEE Access* 8 (2020) 142131–142147.
- [17] M. Maruseac, G. Ghinita, Precision-enhanced differentially-private mining of high-confidence association rules, *IEEE Trans. Dependable Secure Comput.* 17 (2018) 1297–1309.
- [18] S. Xu, S. Su, L. Xiong, X. Cheng, K. Xiao, Differentially private frequent subgraph mining, in: *ICDE'16*, 2016, pp. 229–240.
- [19] L. Bonomi, L. Xiong, A two-phase algorithm for mining sequential patterns with differential privacy, in: *CIKM'13*, 2013, pp. 269–278.
- [20] P. S. Efraimidis, P. G. Spirakis, Weighted random sampling with reservoir, *Inf. Proc. Lett.* 97 (2006) 181–185.
- [21] M. Kolonko, D. Wäsch, Sequential reservoir sampling with a nonuniform distribution, *ACM Trans. Math. Soft.* 32 (2006) 257–273.
- [22] A. J. Mitchell-Jones, G. Amori, W. Bogdanowicz, B. Krystufek, P. J. H. Reijnders, F. Spitzenberger, M. Stubbe, J. B. M. Thissen, V. Vohralik, J. Zima, *The Atlas of European Mammals*, Academic Press, 1999.
- [23] R. J. Hijmans, S. E. Cameron, L. J. Parra, P. G. Jones, A. Jarvis, Very High Resolution Interpolated Climate Surfaces for Global Land Areas, *Int. J. Climatol.* 25 (2005) 1965–1978. URL: [www.worldclim.org](http://www.worldclim.org).
- [24] E. Galbrun, P. Miettinen, Siren: An interactive tool for mining and visualizing geospatial redescrptions [demo], in: *KDD'12*, 2012, pp. 1544–1547.
- [25] J. Kalofolias, E. Galbrun, P. Miettinen, From sets of good redescrptions to good sets of redescrptions, in: *ICDM'16*, 2016, pp. 211–220.
- [26] E. Galbrun, P. Miettinen, Mining Redescrptions with Siren, *ACM Trans. Knowl. Discov. Data* 12 (2018).
- [27] A. E. Johnson, T. J. Pollard, L. Shen, L. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, R. G. Mark, MIMIC-III, a freely accessible critical care database, *Scientific data* 3 (2016) 160035.
- [28] Helsingin Sanomat, Parliamentary elections 2011: Candidate responses to helsingin sanomat candidate selector, 2016. URL: <http://urn.fi/urn:nbn:fi:fsd:T-FSD2701>, version 2.1.
- [29] T. E. Malloy, *Bem sex role inventory*, Wiley Online Library, 2010.
- [30] B. Altemeyer, *Right-wing authoritarianism*, Univ. of Manitoba Press, 1983.
- [31] R. Balu, T. Furon, Differentially private matrix factorization using sketching techniques, in: *IH&MMSec'16*, 2016, pp. 57–62.
- [32] S. Fletcher, M. Z. Islam, Decision tree classification with differential privacy: A survey, *ACM Comput. Surv.* 52 (2019) 83:1–83:33.