

Appendix to “Serenade: An Approach for Differentially Private Greedy Redescription Mining”

Maiju Karjalainen, Esther Galbrun, and Pauli Miettinen^[0000–0003–2271–316X]

University of Eastern Finland `firstname.lastname@uef.fi`

1 Reproducibility

The full source code, together with XML-formatted settings files and scripts to run every experiment, can be obtained from <https://github.com/maijuka/Serenade>.

The `Mammals` data set can be obtained upon request from the authors. The `MIMIC` data set can be obtained from <https://physionet.org/content/mimiciii/1.4/>. The `VAA` data set can be downloaded from <http://urn.fi/urn:nbn:fi:fsd:T-FSD2701>. The `Psycho`, `Pref`, `SSC`, `SR`, and `RightW` data sets can be downloaded from https://openpsychometrics.org/_rawdata.

The steepness parameter for computing the quality function (s in Equation 4) are presented in Table A1. The optimal steepness was calculated for each data set by searching for the value resulting in the smallest sensitivity. The corresponding sensitivities are also listed in Table A1.

data set	s	Δf
<code>Mammals</code>	0.02	0.0059
<code>MIMIC</code>	0.01	0.0025
<code>VAA</code>	0.03	0.0075
<code>SSC</code>	0.01	0.0025
<code>Pref</code>	0.01	0.0025
<code>Psycho</code>	0.01	0.0025
<code>SR</code>	0.02	0.0050
<code>RightW</code>	0.02	0.0050

Table A1. Steepness parameters s for the quality function and corresponding sensitivities Δf .

2 Proof of Lemma 4

INITIALPAIRS (Algorithm 4) receives a budget of ε to choose I initial redescriptions. The budget is divided by the number of initial pairs I and by the sensitivity

Δf so that each redescription uses $\varepsilon/(I\Delta f)$. Then, we use the exponential mechanism to choose the redescription with a probability proportional to the budget (line 5), as shown in Equation 3. Since the exponential mechanism gives $2\varepsilon\Delta f$ -private results with a budget of ε , choosing one redescription is then $(\varepsilon/(I\Delta f)) \cdot (2\Delta f) = (2\varepsilon)/I$ -private. Repeating this for I redescriptions, the INITIALPAIRS algorithm is $((2\varepsilon)/I) \cdot I = 2\varepsilon$ -private.

EXTENDONE (Algorithm 3) receives a budget of ε to choose one extension to a given redescription. The budget is divided by the given sensitivity Δf , so that the extension is chosen using exponential mechanism with a probability proportional to $\varepsilon/\Delta f$, similarly as above. This makes extending a redescription once $(\varepsilon/\Delta f) \cdot (2\Delta f) = 2\varepsilon$ -private.

To compute the private quality for a redescription, we apply Laplace noise to the size of the intersection and the size of the union of the input queries' supports. With a budget ε , both operations receive a budget of $\varepsilon/2$, making the quality evaluation for one redescription ε -private.

Both algorithms SERENADEES and SERENADECS use a total budget ε_{tot} , and the user decides how to divide it into $\varepsilon_{\text{init}}$, ε_{ext} , and $\varepsilon_{\text{qual}}$.

SERENADEES (Algorithm 2) first uses budget $\varepsilon_{\text{init}}/2$ by calling INITIALPAIRS (line 2). Then, the loop starting on line 3 is executed K times, each time calling EXTENDONE with a budget $\varepsilon_{\text{ext}}/(2K)$, making the total budget spent in this loop $\varepsilon_{\text{ext}}/2$. Finally, the quality is computed for all redescriptions in \mathcal{Q} with a budget of $\varepsilon_{\text{qual}}/|\mathcal{Q}|$ to add Laplace noise to each redescription. As shown above for a budget ε both INITIALPAIRS and EXTENDONE were 2ε -private, and the quality computation is ε -private. With the budgets used in SERENADEES the algorithm is $(\varepsilon_{\text{init}}/2) \cdot 2 + (\varepsilon_{\text{ext}}/2) \cdot 2 + (\varepsilon_{\text{qual}}/|\mathcal{Q}|) \cdot |\mathcal{Q}| = \varepsilon_{\text{tot}}$ -private.

SERENADECS (Algorithm 1) computes the initial pairs in the same way as SERENADEES. The loop on line 4 is executed $|\mathcal{Q}|$ times, and the loop on line 9 is executed at most L times, which means EXTENDONE on line 9 is called at most $L|\mathcal{Q}|$ times with a budget of $\varepsilon_{\text{ext}}/(2L|\mathcal{Q}|)$ each time. Then the total budget spent for calling EXTENDONE is $(\varepsilon_{\text{ext}}/(2L|\mathcal{Q}|)) \cdot L|\mathcal{Q}| = \varepsilon_{\text{ext}}/2$. The quality of the redescription is evaluated after each EXTENDONE call (line 10) spending $\varepsilon_{\text{qual}}/(L|\mathcal{Q}|)$ budget each time, for at most $L|\mathcal{Q}|$ times, for total budget spent $\varepsilon_{\text{qual}}$. In total SERENADECS is $(\varepsilon_{\text{init}}/2) \cdot 2 + (\varepsilon_{\text{ext}}/2) \cdot 2 + (\varepsilon_{\text{qual}}/L|\mathcal{Q}|) \cdot L|\mathcal{Q}| = \varepsilon_{\text{tot}}$ -private.

3 Proof of Lemma 2

Without memoization, the claim follows from Lemma 1. Indeed, the process is same, except that instead of selecting an extension of one redescription, all redescriptions in \mathcal{Q} are tested.

With memoization, we can think of the process as sampling from a stream into which we add new elements (the extensions of the new redescriptions) in successive iterations. The proofs of correct sampling [2, 1] still hold, as the random values for each extension, are independent from each other, though constant over different runs for a memoized redescription.

Table A2. Example redescriptions obtained from the data sets.

data set	algorithm	q_L	q_R	noisy J	true J
Psycho	SERENADEES	Do you think tobacco	$[27.0 \leq \text{age}]$	0.317	0.231
		has hurt you			
		Do you find it difficult	$[\text{age} \leq 33.0]$	0.343	0.408
SSC	SERENADEES	to pass urine in the presence of others			
		$[3.0 \leq \text{I have positive feelings about the way I approach my own sexual needs and desires}]$	$[3.0 \leq \text{I am very aware of the sexual aspects of myself eg habits}]$	0.428	0.375
		$\wedge [23.0 \leq \text{age}]$			
RightW	SERENADEES	$[5.0 \leq \text{You have to admire those who challenged the law and the majority's view by protesting for women's abortion rights, for animal rights, or to abolish school prayer}]$	$[\text{voted} = \text{No}]$	0.428	0.538
Pref	SERENADEES	$[\text{Q5} = \text{Write a movie script}]$	$([\text{Q9} = \text{Invent a story}] \wedge [\text{country} = \text{US}]) \vee [\text{Q2} = \text{Be a movie critic}]$	0.315	0.344
SR	SERENADECS	$[3.0 \leq \text{I am happiest when I am in my bed}]$	$[3.0 \leq \text{I give people handmade gifts}] \wedge [14.0 \leq \text{age} \leq 32.0]$	0.651	0.725
MIMIC	SERENADECS	Atrial fibrillation	$[\text{Uric Acid} = 1] \vee [\text{Creatinine Kinase (CK)} = 0]$	0.424	0.227
VAA	SERENADEES	$[\text{Gender} = \text{M}]$	$[\text{CountQ21:Fb} \leq 1.0] \vee [\text{W5:SeniorCare} \leq 0.0]$	0.637	0.445

4 Further experimental results

4.1 Results from ReReMi

Figure A1 shows the results obtained with REREMI on different data sets. These can be compared to results shown in Figure 1(c–e).

4.2 Results with the questionnaire datasets

We also evaluated the algorithms on three data sets consisting of answers to personality tests.¹ The Woodworth Psychoneurotic Inventory (Psycho data set) is a personality test originally designed to measure the emotional stability and

¹ Data sets available at https://openpsychometrics.org/_rawdata/.

data set	n , entities	\mathbf{D}_L attributes	\mathbf{D}_R attributes
Mammals	2 575	194 Boolean	48 numerical
MIMIC	46 065	107 Boolean	104 Boolean
VAA	1 656	9 categorical	107 categorical
SSC *	17 685	102 numerical	—
Pref *	13 502	13 categorical	—
Psycho *	6 019	119 Boolean	—
SR *	4 184	53 numerical	—
RightW *	2 863	61 numerical	—

* These data sets consist of a single table common to both sides.

Table A3. Data set properties.

the risk for shell shock of World War I soldiers. The Wagner Preference Inventory (**Pref** data set) is used to measure the hemispheric brain dominance with questions on preferred activities. The Multidimensional Sexual Self-Concept Questionnaire (**SSC** data set) consists of questions on sexual behaviours. These three 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. Figure A2 shows the results obtained with these datasets.

4.3 Scatter plots of true vs. noisy Jaccards

Figure A3 shows scatter plots of true vs. noisy Jaccards.

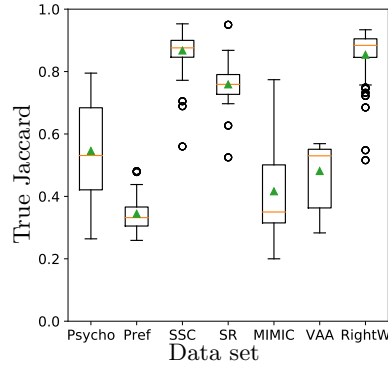


Fig. A1. Accuracy per data set with the REREMi algorithm.

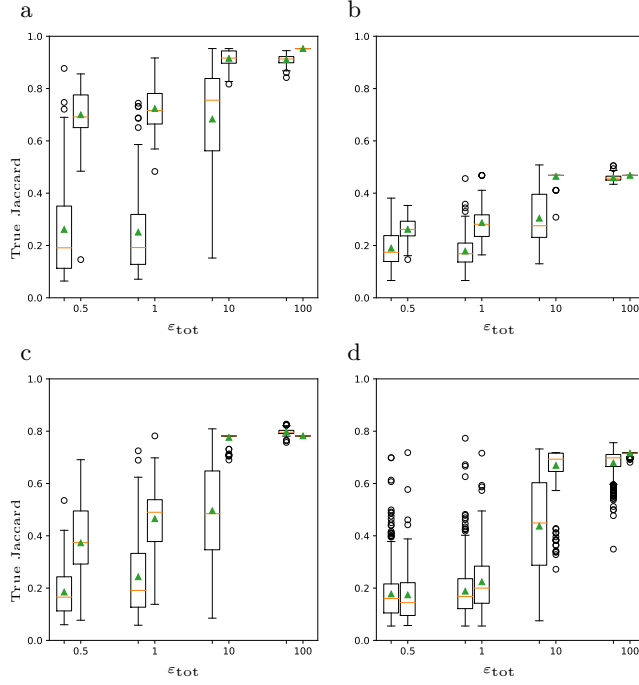


Fig. A2. Effects of the budget in terms of the true accuracy on different data sets. y -axis is always true Jaccard. x -axis: (a) Accuracy of the initial pairs, varying ϵ_{init} , with $\epsilon_{\text{tot}} = 1$. (b)–(h) varying ϵ_{tot} , budget distribution d_3 . This is the same as in Figure 1(a) for **Mammals**, but without **NoSENS**. The data sets are (b) **SR**, (c) **RightW**, (d) **VAA**, (e) **SSC**, (f) **Pref**, (g) **Psycho**, and (h) **MIMIC**. Boxes in a group represent (left to right) SERENADEES and SERENADECS. The horizontal bar and triangle within each box represent the median and the mean, respectively. Results are collected over five restarts.

4.4 Raw qualitative results

Table A2 shows some redescrptions mined from different datasets, as anecdotal examples of results one can expect by our algorithms with these kinds of data sets. A few of them are presented in Section 4.3.

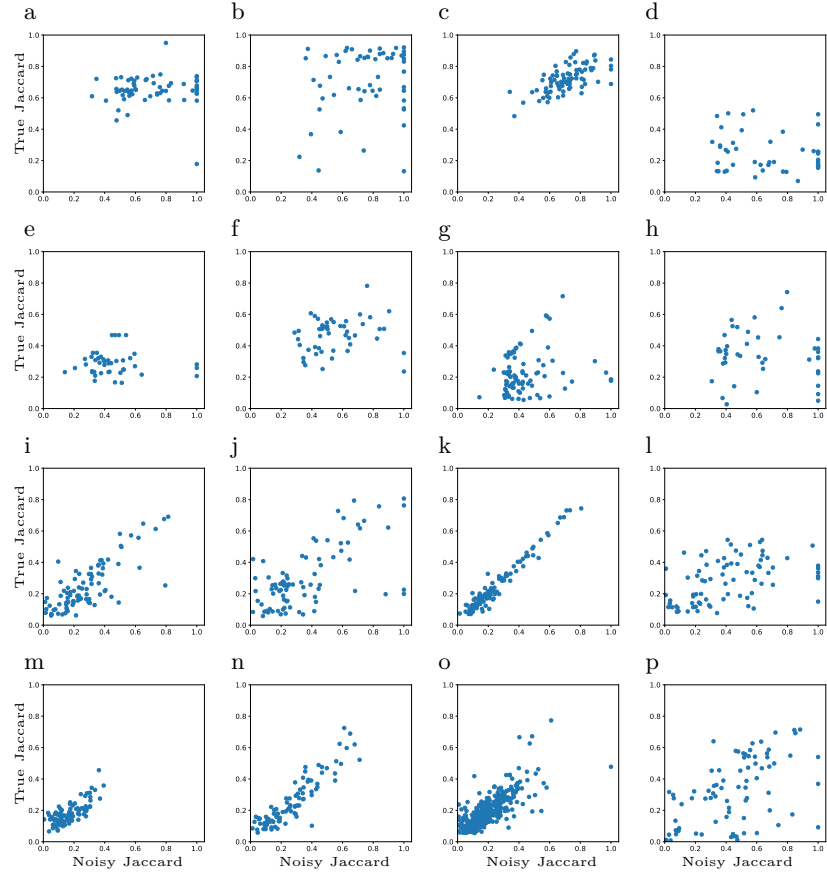


Fig. A3. Total budget 1, distribution d_3 . The vertical axis represents the true Jaccard whereas the horizontal axis represents the noisy Jaccard. Top two rows: SERENADECS. Bottom two rows: SERENADEES. (a) and (i) **SR**, (b) and (j) **RightW**, (c) and (k) **SSC**, (d) and (l) **VAA**, (e) and (m) **Pref**, (f) and (n) **Psycho**, (g) and (o) **MIMIC**, and (h) and (p) **Mammals**.