

DPXPlain: Privately Explaining Aggregate Query Answers

Yuchao Tao
Duke University
yctao@cs.duke.edu

Amir Gilad
Duke University
agilad@cs.duke.edu

Ashwin
Machanavajjhala
Duke University
ashwin@cs.duke.edu

Sudeepa Roy
Duke University
sudeepa@cs.duke.edu

ABSTRACT

Differential privacy (DP) is the state-of-the-art and rigorous notion of privacy for answering aggregate database queries while preserving the privacy of sensitive information in the data. In today’s era of data analysis, however, it poses new challenges for users to understand the trends and anomalies observed in the query results: Is the unexpected answer due to the data itself, or is it due to the extra noise that must be added to preserve DP? In the second case, even the observation made by the users on query results may be wrong. In the first case, can we still mine interesting explanations from the sensitive data while protecting its privacy? To address these challenges, we present a three-phase framework DPXPLAIN, which is the first system to the best of our knowledge for explaining group-by aggregate query answers with DP. In its three phases, DPXPLAIN (a) answers a group-by aggregate query with DP, (b) allows users to compare aggregate values of two groups and with high probability assesses whether this comparison holds or is flipped by the DP noise, and (c) eventually provides an explanation table containing the approximately ‘top-k’ explanation predicates along with their relative influences and ranks in the form of confidence intervals, while guaranteeing DP in all steps. We perform an extensive experimental analysis of DPXPLAIN with multiple use-cases on real and synthetic data showing that DPXPLAIN efficiently provides insightful explanations with good accuracy and utility.

PVLDB Reference Format:

Yuchao Tao, Amir Gilad, Ashwin Machanavajjhala, and Sudeepa Roy. DPXPlain: Privately Explaining Aggregate Query Answers. PVLDB, 16(1): 113 - 126, 2022.
doi:10.14778/3561261.3561271

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/yuchaotao/Private-Explanation-System>.

1 INTRODUCTION

Differential privacy (DP) [14, 40–42] is the gold standard for protecting privacy in query processing and is critically important for sensitive data analysis. It has been widely adopted by organizations like the U.S. Census Bureau [3, 38, 58, 84] and companies like Google [44, 97], Microsoft [29], and Apple [89]. The core idea behind DP is that a query answer on the original database cannot be distinguished from the same query answer on a slightly different

database. This is usually achieved by adding random noise to the query answer to create a small distortion in the answer. Recent works have made significant advances in the usability of DP, allowing for complex query support [32, 56, 59, 60, 69, 90, 97], and employing DP in different settings [32, 45, 48, 77, 90, 99]. These works assist in bridging the gaps between the functionality of non-DP databases and databases that employ DP.

Automatically generating meaningful *explanations* for query answers in response to questions asked by users is an important step in data analysis that can significantly reduce human efforts and assist users. Explanations help users validate query results, understand trends and anomalies, and make decisions about next steps regarding data processing and analysis, thereby facilitating data-driven decision making. Several approaches for explaining aggregate and non-aggregate query answers have been proposed in database research, including intervention [81, 82, 98], Shapley values [67], counterbalance [75], (augmented) provenance [5, 65], responsibility [73, 74], and entropy [43] (discussed in Section 6).

One major gap that remains wide open is to provide explanations for analyzing query answers from sensitive data under DP. Several new challenges arise from this need. First, in DP, the (aggregate) query answers shown to users are distorted due to the noise that must be added for preserving privacy, so the explanations need to separate the contributions of the noise from the data. Second, even after removing the effect of noise, new techniques have to be developed to provide explanations based on the sensitive data and measure their effects. For instance, standard explanation methods in non-DP settings are typically deterministic, while it is known that DP methods must be randomized. Therefore, no deterministic explanations can be provided, and even no deterministic scores or ranks of explanations can be displayed in response to user questions if we want to guarantee DP in the explanation system. Third, the system needs to ensure that the returned explanations, scores, and ranks still have high accuracy while being private.

In this paper, we propose DPXPLAIN, a novel three-phase framework that generates explanations¹ under DP for aggregate queries based on the notion of *intervention* [82, 98]². DPXPLAIN surmounts the aforementioned challenges and is the first system combining DP and explanations to the best of our knowledge. We illustrate DPXPLAIN through an example.

Example 1.1. Consider the Adult (a subset of Census) dataset [35] with 48,842 tuples. We consider the following attributes: age, workclass, education, marital-status, occupation, relationship, race, sex, native-country,

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 16, No. 1 ISSN 2150-8097.
doi:10.14778/3561261.3561271

¹The explanations we provided should not be considered causal explanations.

²A graphical user interface for DPXPLAIN is an ongoing work.

marital-status	occupation	...	education	high-income
Never-married	Machine-op-inspct	...	11th	0
Married-civ-spouse	Farming-fishing	...	HS-grad	0
Married-civ-spouse	Machine-op-inspct	...	Some-college	1
...

(a) Example of the Adult dataset.

Question-Phase-1:

SELECT marital-status, AVG(high-income) as avg-high-income
FROM Adult GROUP BY marital-status;

group marital-status	Priv-answer avg-high-income	True-answer (hidden)
Separated	0.064712	0.064706
Widowed	0.082854	0.084321
Married-spouse-absent	0.089988	0.092357
Divorced	0.101578	0.101161
Married-AF-spouse	0.463193	0.378378
Married-civ-spouse	0.446021	0.446133

Answer-Phase-1:

(b) Phase-1 of DPXPLAIN: Run a query and receive noisy answers by DP. True-answers are not visible to the user and for illustration only.

Question-Phase-2: Why avg-high-income of group "Married-civ-spouse" > that of group "Never-married"?

Answer-Phase-2: The 95% confidence interval of group difference is ¹0.399, 0.402°, hence the noise in the query is possibly not the reason.

(c) Phase-2 of DPXPLAIN: Ask a comparison question and receive a confidence interval of the comparison.

Answer-Phase-3:

explanation predicate	Rel Influ 95%-CI		Rank 95%-CI	
	L	U	L	U
occupation = "Exec-managerial"	3.25%	10.12%	1	9
education = "Bachelors"	2.93%	9.80%	1	8
age = "(40, 50]"	2.76%	9.63%	1	8
occupation = "Prof-specialty"	0.94%	7.81%	1	18
relationship = "Own-child"	-0.49%	6.38%	1	96

(d) Phase-3 of DPXPLAIN: Receive an explanation table from data for the previous question that passed Phase-2.

Figure 1: Database instance and the three phases of the DPXPLAIN framework.

and high-income, where high-income is a binary attribute indicating whether the income of a person is above 50K or not; some relevant columns are illustrated in Figure 1a.

In the **first phase (Phase-1)** of DPXPLAIN, the user submits a query and gets the results as shown in Figure 1b. This query is asking the fraction of people with high income in each marital-status group. As Figure 1b shows, the framework returns the answer with two columns: group and Priv-answer. Here group corresponds to the group-by attribute marital-status. However, since the data is private, instead of seeing the actual aggregate values avg-high-income, the user sees a perturbed answer Priv-answer for each group as output by some differentially private mechanism with a given privacy budget (here computed by the Gaussian mechanism with privacy budget $\rho = 0.1$ [14]). The third column True-answer shown in grey (hidden for users) in Figure 1b shows the **true aggregated output** for each group.

In the **second phase (Phase-2)** of DPXPLAIN, the user selects two groups to compare their aggregate values and asks for explanations. However, unlike standard explanation frameworks [43, 65, 75, 82, 98] where the answers to a query are correct and hence the question asked by the user is also correct, in the DP setting, the answers that the users see are perturbed. Therefore, the user question and the direction of comparison may not be valid. Hence our system first tests the validity of the question. If the question is valid, our system provides a data-dependent explanation of the user question. We explain this below with the running example.

First, consider the question in Figure 2 comparing the last two groups in Figure 1b (spouse in armed forces vs. a civilian). In this example, even though the noisy avg-high-income for "Married-AF-spouse" is larger than the noisy value for "Married-civ-spouse", this might not be true in the real data (as is the case in the True-answer column). Hence, our system tests whether the user question could potentially be explained just using the noise introduced by DP rather than from the data itself. To do

Question-Phase-2: Why avg-high-income of group "Married-AF-spouse" > that of group "Married-civ-spouse"?

Answer-Phase-2: The 95% confidence interval of group difference is ¹ 0.259, 0.460°, hence the noise in the query is possibly the reason.

Figure 2: A user question explained by high noise.

this, our system tests the validity of the user question by computing a confidence interval around the difference between these two outputs. In this case, the confidence interval is ¹ 0.259, 0.460°. Since it includes 0 and negative values, we cannot conclude with high probability that "Married-AF-support" > "Married-civ-spouse" is true in the original data. **Since the validity of the user question is uncertain, we know that any further explanation might not be meaningful and the user may choose to stop here.** In other words, the explanation for the comparison in the user question is primarily attributed to the added noise by the DP mechanism. If the user chooses to proceed to the next phase for further explanations from the data, they might not be meaningful.

Now consider the comparison between two other groups "Never-married" and "Married-civ-spouse", in Figure 1c. In this case, the confidence interval about the difference does not include zero and is tight around a positive number of 0.4, which indicates that the user question is correct with high probability. Notice that it is still possible for a valid question to have a confidence interval that includes zero given sufficiently large noise. Since the question is valid, the user may continue to the next phase.

In the **third phase (Phase-3)** of DPXPLAIN, for the questions that are likely to be valid, DPXPLAIN can provide a further detailed data-dependent explanation for the question. To achieve this again with DP, our framework reports an "Explanation Table"³ to the

³We note that our notion of explanation table is unrelated to that described by Gebaly et al. [43] for summarizing dimension attributes to explain a binary outcome attribute.

user as Figure 1d shows, which includes the top-5 *explanation predicates*. The explanation predicates explain the user question using the notion of *intervention* as done in previous work [82, 98] for explaining aggregate queries in the non-DP setting. Intuitively, if we intervene in the database by (hypothetically) removing tuples that satisfy the predicate, and re-evaluate the query, then the difference in the aggregate values of the two groups mentioned in the question will reduce. In the simplest form, explanation predicates are singleton predicates of the form “attribute = <value>”, while in general, our framework supports more complex predicates involving conjunction, disjunction, and comparison (>, etc.). In Figure 1d, the top-5 simple explanation predicates, as computed by DPXPLAIN, are shown out of 103 singleton predicates, according to their influences on the question but perturbed by noises to satisfy DP. The amount of noise is proportional to the sensitivity of the influence function, the maximum possible change of the influence of any explanation predicate when adding or removing a single tuple from the database. Once the top-5 predicates are selected, the explanation table also shows both their *relative influence* (intuitively, how much they affect the difference of the group aggregates in the question) and their *ranks* (that might be far away from the true top-5) in the form of confidence interval under DP.

From this table, `occupation = "Exec-managerial"` is returned as the top explanation predicate, indicating that the people with this job contribute more to the average high income of the married group compared to the never-married group. In other words, managers tend to earn more if they are married than those who are single, which probably can be attributed to the intuition that married people might be older and have more seniority, which is consistent with the third explanation `age = "(40, 50]"` in Figure 1d as well. Although these explanations are chosen at random, we observe that the first three explanations are almost constantly included. This is consistent with the narrow confidence interval of rank for the first three explanation predicates, which are all around $\approx 1, 8\%$. Looking at the confidence intervals of the relative influence and ranks in the explanation table, the user also knows that the first three explanations are likely to have some effect on the difference between the married and unmarried groups. However, for the last two explanations, the confidence intervals of influences are closer to 0 and the confidence intervals of ranks are wider, especially for the fifth one which includes negative influences in the interval and has a wide range of possible ranks (96 out of 103 simple explanation predicates in total).

Our contributions.

We develop DPXPLAIN, the first framework, to our knowledge, that generates explanations for query answers under DP adapting the notion of intervention [82, 98]. It explains user questions comparing two group-by aggregate query answers (COUNT, SUM, or AVG) with DP in three phases: private query answering, private user question validation, and private explanation table. We develop multiple novel techniques that allow DPXPLAIN to provide explanations under DP including (a) computing confidence intervals to check the validity of user questions, (b) choosing explanation predicates, and (c) computing confidence intervals around the influence and rank of the predicates.

We design a low sensitivity influence function inspired by previous work on non-private explanations [98], which is the key to the accurate selection of the top-k explanation predicates.

We design an algorithm that uses a noisy binary search technique to find the confidence intervals of the explanation ranks, which overcomes the high sensitivity challenge of the rank function.

We have implemented a prototype of DPXPLAIN [2] to evaluate our approach. We include two case studies on a real and a synthetic dataset showing the entire process and the obtained explanations. We have further performed a comprehensive accuracy and performance evaluation, showing that DPXPLAIN correctly indicates the validity of the question with 100% accuracy for 8 out of 10 questions, selects at least 80% of the true top-5 explanation predicates correctly for 8 out of 10 questions, and generates descriptions about their influences and ranks with high accuracy.

2 PRELIMINARIES

We now give the necessary background for our model. The DPXPLAIN framework supports single-block SELECT - FROM - WHERE - GROUP BY queries with aggregates (Figure 3) on single tables⁴. Hence the database schema $A = {}^1A_1, \dots, A_m^\circ$ is a vector of attributes of a single relational table. Each attribute A_i is associated with a domain $\text{dom}^1A_i^\circ$, which can be continuous or categorical. A database (instance) D over a schema A is a bag of tuples (duplicate tuples are allowed) $t_i = {}^1a_1, \dots, a_m^\circ$, where $a_i \in \text{dom}^1A_i^\circ$ for all i . The domain of a tuple is denoted as $\text{dom}^1A^\circ = \text{dom}^1A_1^\circ \times \text{dom}^1A_2^\circ \times \dots \times \text{dom}^1A_m^\circ$. We denote $A_i^{\max} = \max\{|a_j| \mid a \in \text{dom}^1A_i^\circ\}$ as the maximum absolute value of A_i . The value of the attribute A_i of tuple t is denoted by $t.A_i$.

$q = \text{SELECT } A_{gb}, \text{agg}(A_{agg}) \text{ FROM } D \text{ WHERE } \phi \text{ GROUP BY } A_{gb};$

Figure 3: Group-by query with aggregates supported by DPXPLAIN. The true results are denoted by ${}^1\alpha_i, o_i^\circ$ and the noisy results released by a DP mechanism are denoted by ${}^1\alpha_i, \hat{o}_i^\circ$ where α_i is the value of A_{gb} and o_i, \hat{o}_i are aggregate values.

We consider group-by aggregate queries q of the form shown in Figure 3. Here A_{gb} is the group-by attribute and A_{agg} is the aggregate attribute, ϕ is a predicate without subqueries, and $\text{agg} \in \{\text{COUNT}, \text{SUM}, \text{AVG}\}$ is the aggregate function. When query q is evaluated on database D , its result is a set of tuples ${}^1\alpha_i, o_i^\circ$, where $\alpha_i \in \text{dom}^1A_{gb}^\circ$ and $o_i = \text{agg}^1\{t.A_{agg} \mid t \in D, \phi^1t^\circ = \text{true}, t.A_{gb} = \alpha_i\}^\circ$. For brevity, we will use ϕ^0 to denote $\{t \mid \phi^1t^\circ = \text{true}\}$ for any predicate ϕ^0 , and agg^1A_{agg}, D^0 , or simply agg^1D^0 when it is clear from context, to denote $\text{agg}^1\{t.A_{agg} \mid t \in D^0\}^\circ$ for any $D^0 \subseteq D$. Hence, $o_i = \text{agg}^1A_{agg}, g_i^1D^0$, where $g_i = \phi^1 \wedge {}^1A_{gb} = \alpha_i^\circ$.

Example 2.1. Consider Example 1.1. The schema is $A = (\text{marital-status}, \text{occupation}, \text{age}, \text{relationship}, \text{race}, \text{workclass}, \text{sex}, \text{native-country}, \text{education}, \text{high-income})$. All the attributes are categorical attributes and the domain of

⁴Unlike some standard explanation framework [98], in DP, we cannot consider materialization of join-result for multiple tables, since the privacy guarantee depends on *sensitivity*, and removing one tuple from a table may change the join and query result significantly. We leave it as an interesting future work.

high-income is f0, 1g. The query is shown in Figure 1b and the true result for each group is shown in the True-answer column. Here $A_{gb} = \text{marital status}$, $A_{agg} = \text{high income}$, and $agg = \text{AVG}$.

Differential Privacy. In this work, we consider query-answering and providing explanations using *differential privacy (DP)* [41] to protect private information in the data. In standard databases, a query result can give an adversary the option to find the presence or absence of an individual in the database, compromising their privacy. DP allows users to query the database without compromising the privacy by guaranteeing that the query result will not change too much (defined in the sequel) even if it is evaluated on any two different but *neighboring* databases defined below.

Definition 2.2 (Neighboring Database). Two databases D and D^0 are neighboring (denoted by $D \sim D^0$) if D^0 can be transformed from D by adding or removing⁵ a tuple in D .

In this paper, we consider a relaxation of DP called ρ -**zero-concentrated differential privacy (zCDP)** [14, 42] for several reasons, and refer to it simply as DP if not otherwise stated. First, we use Gaussian noise to perturb query answers and derive confidence intervals, which does not satisfy pure ϵ -DP [41] but satisfies approximate (ϵ, δ) -DP [41] and ρ -zCDP. Second, ρ -zCDP only has one parameter ρ , compared to (ϵ, δ) -DP which has two parameters, so it is easier to understand and control. Third, ρ -zCDP allows for tighter analyses for tracking the *privacy budget* (controlled by ρ) over multiple private releases, which is the case for this framework. A lower ρ value implies a lower privacy loss.

Definition 2.3 (Zero-Concentrated Differential Privacy (zCDP) [14]). A mechanism M is said to be ρ -zero-concentrated differential private, or ρ -zCDP for short, if for any neighboring datasets D and D^0 and all $\alpha \geq 1$, $1 \leq \alpha$ it holds that

$$D_{\alpha}^1 M^1 D^0 \ll M^1 D^{00} \quad \rho \alpha$$

where $D_{\alpha}^1 M^1 D^0 \ll M^1 D^{00}$ denotes the Rényi divergence of the distribution $M^1 D^0$ from the distribution $M^1 D^{00}$ at order α [76].

A popular approach for providing zCDP to a query result is to add Gaussian noise to the result before releasing it to a user. This approach is called *Gaussian mechanism* [14, 41].

Definition 2.4 (Gaussian Mechanism). Given a query q and a noise scale σ , Gaussian mechanism M^G is given as:

$$M^G(q; D, \sigma) = q(D) + N(0, \sigma^2)$$

where $N(0, \sigma^2)$ is a random variable from a normal distribution⁶ with mean zero and variance σ^2 .

Example 2.5. Suppose there is a database D with 100 tuples. Consider a query $q = \text{"SELECT COUNT(*) FROM D"}$, which counts the total number of tuples in a database D . Here $q(D) = 100$. Now we use Gaussian mechanism to release $q(D)$, which is to randomly sample a noise z from distribution $N(0, \sigma^2)$. Here we assume $\sigma = 1$.

⁵There are two variants of neighboring databases. The definition by addition/deletion of tuples is called "unbounded DP", and by updating tuples is called "bounded DP", since the size of data is fixed. In this work, we assume the unbounded version, while DPXPLAIN can be adapted also for the bounded version by adapting the noise scale.

⁶The probability density function of a normal distribution $N(\mu, \sigma^2)$ is given as $\exp(-\frac{1}{2\sigma^2}(x-\mu)^2)$.

Finally, we got a noisy result $\hat{q}(D) = 102.32$, which we may round to an integer in postprocessing without sacrificing the privacy guarantee (Proposition 2.9 below).

The privacy guarantee from the Gaussian mechanism depends on both the noise scale it uses and the sensitivity of the query. Query sensitivity reflects how sensitive the query is to the change of the input. More noise is needed for a more sensitive query to achieve the same level of privacy protection.

Definition 2.6 (Sensitivity). Given a scalar query q that outputs a single number, its sensitivity is defined as:

$$\Delta_q = \sup_{D, D^0} |q(D) - q(D^0)|$$

Example 2.7. Continuing Example 2.5, since the query q returns the database size, for any two neighboring databases, their sizes always differ by 1, so the sensitivity of q is 1.

THEOREM 2.8 (GAUSSIAN MECHANISM [14]). Given a query q with sensitivity Δ_q and a noise scale σ , its Gaussian mechanism M^G satisfies $(\Delta_q^2 \cdot 2\sigma^2, \rho)$ -zCDP. Equivalently, given a privacy budget ρ , choosing $\sigma = \Delta_q \cdot \sqrt{\frac{\rho}{2}}$ in Gaussian mechanism satisfies ρ -zCDP.

Composition Rules. In our analysis, we will use the following standard composition rules and other known results from the literature of DP [72] (in particular, zCDP [14]) frequently:

PROPOSITION 2.9. The following holds for zCDP [14, 72]:

Parallel composition: if mechanisms take disjoint data as input, the total privacy loss is the maximum privacy loss from each.

Sequential composition: if mechanisms take overlapping data as input, the total privacy loss is the sum of each privacy loss.

Post-processing: if we run a mechanism and post-process the result without accessing the data, the total privacy loss is only the privacy loss from the mechanism.

Private Query Answering. Recall that we have group-by aggregation query of the form $q = \text{SELECT } A_{gb}, \text{ agg}(A_{agg}) \text{ FROM } D \text{ WHERE } \phi \text{ GROUP BY } A_{gb}$, and it returns a list of tuples $\langle \alpha_i, o_i \rangle$ where $\alpha_i \in \text{dom}(A_{gb})$ and o_i is the corresponding aggregate value. Since no single tuple can exist in more than one group, adding or removing a single tuple can at most change the result of a single group. As mentioned earlier, Phase-1 returns noisy aggregate values \hat{o}_i for each α_i instead of o_i . The following holds:

Observation 2.1. According to the parallel composition rule (Proposition 2.9), if for each α_i , its (noisy) aggregate value \hat{o}_i is released under ρ_q -zCDP, the entire release of results including all groups $\{ \langle \alpha_i, \hat{o}_i \rangle : \alpha_i \in \text{dom}(A_{gb}) \}$ satisfies ρ_q -zCDP.

For a *COUNT* or *SUM* query, we use the Gaussian mechanism for each group α_i : $\hat{o}_i = o_i + N(0, \sigma^2)$, where the noise scale $\sigma = \Delta_q \cdot \sqrt{2\rho_q}$ to satisfy ρ_q -zCDP by Theorem 2.8. The sensitivity term Δ_q is 1 for *COUNT* and A_{agg}^{max} for *SUM*, the maximum absolute value of the aggregation attribute in its domain. For an *AVG* query, since $AVG = SUM \cdot COUNT$, we decompose it into a *SUM* and a *COUNT* query, privately answer each of them by half of the privacy budget $\rho_q \cdot 2$ to get \hat{o}_i^S and \hat{o}_i^C for each group α_i , and release $\hat{o}_i = \hat{o}_i^S \cdot \hat{o}_i^C$ as a post-processing step. The noisy query answers of the group-by query with *AVG* satisfy ρ_q -zCDP by the sequential composition rule (Proposition 2.9).

Confidence Level and Interval. Confidence intervals are commonly used to determine the error margin in uncertain computations and are used in various fields including machine learning [55] and DP [46]. In our context, we use confidence intervals to measure the uncertainty in the user question and our explanations.

Definition 2.10 (Confidence Level and Interval [96]). Given a confidence level γ and an unknown but fixed parameter θ , a random interval $I = [L, U]$ is said to be its confidence interval, or CI, with confidence level γ if the following holds:

$$\Pr_{\theta} [L \leq \theta \leq U] = \gamma$$

Example 2.11. Let $\theta = 0$. Suppose with probability 50% we have $I^L = 1$ and $I^U = 1$, and with another probability 50% we have $I^L = 1$ and $I^U = 2$. Therefore, $\Pr_{\theta} [L \leq \theta \leq U] = 50\%$, and we can conclude that the random interval $I = [L, U]$ is a 50% level confidence interval for θ .

3 PRIVATE EXPLANATIONS IN DPXPLAIN

In this section, we provide the model for private explanations of query results at the center of DPXPLAIN.

User Question and Standard Explanation Framework. In Phase-2 of DPXPLAIN, given the noisy results of a group-by aggregation query from Phase-1, users can ask questions comparing the aggregate values of two groups⁷:

Definition 3.1 (User Question). Given a database D , a group-by aggregate query q as shown in Figure 3, a DP mechanism \mathcal{M} , and two noisy answer tuples $\langle \alpha_i, \hat{\delta}_i \rangle, \langle \alpha_j, \hat{\delta}_j \rangle \in \mathcal{M}^1(D; q)$ where $\hat{\delta}_i > \hat{\delta}_j$, a user question has the form “why is the (noisy) aggregate value $\hat{\delta}_i$ of group α_i larger than the aggregate value $\hat{\delta}_j$ of group α_j ?”, which is denoted by “why $\langle \alpha_i, \alpha_j \rangle$?”.

Example 3.2. The question from Figure 1c is denoted as “why (‘Married-civ-spouse’, ‘Never-married’, >)?”.

To explain a user question, several previous approaches return top- k predicates that have the highest influences over the group difference in the question [43, 65, 82, 98]. We follow this paradigm and define explanation predicates.

Definition 3.3 (Explanation Predicate). Given a database D with a set of attributes A , a group-by aggregation query q (Figure 3) with group-by attribute A_{gb} and aggregate attribute A_{agg} and a predicate size l , an explanation predicate p is a Boolean expression of the form $p = \varphi_1 \wedge \dots \wedge \varphi_l$, where each φ_i has the form $A_i = a_i$ such that $A_i \in A \cap \{A_{gb}, A_{agg}\}$ is an attribute, and $a_i \in \text{dom}^1(A_i)$ is its value.

We assume $\text{dom}^1(A_i)$ is discrete, finite, and data-independent. We focus here on the conjunction of equality predicates. However, our framework can also handle predicates that contain disjunctions and inequalities of the form $A_i \neq a_i$ where $\neq \in \{>, <, =, \neq\}$ when the constant a_i is from a finite and data-independent set.

New challenges for explanations with DP. Unlike standard explanation framework on aggregate queries [65, 82, 98], the existing frameworks are not sufficient to support DP and need to be adapted: (i) the question itself might not be valid due to the noise injected

into the queries, (ii) the selection of top- k explanation predicates needs to satisfy DP, which further requires the influence function to have low sensitivity so that the selection is less perturbed, and (iii) since the selected explanation predicates are not guaranteed to be the true top- k , it is also necessary to output extra descriptions under DP for each selected explanation predicate about their actual influences and ranks. We detail the adjustments as follows.

Question Validation with DP (Phase-2). While the user is asking “why is $\hat{\delta}_i > \hat{\delta}_j$?”, in reality, it may be the case that the true results satisfy $\delta_i < \delta_j$, i.e., they have opposite relationship than the one observed by the user. This indicates that $\hat{\delta}_i > \hat{\delta}_j$ is the result of the noise being added to the results. In this scenario, one option to explain the user’s observation of $\hat{\delta}_i > \hat{\delta}_j$ will be releasing the true values (equivalently, the added exact noise values), which will violate DP. Instead, to provide an explanation in such scenarios, we generate a confidence interval for the difference of two (hidden) aggregate values $\delta_i - \delta_j$, which can include negative values (discussed in detail in Section 4.1). This leads to the first problem we need to solve in the DPXPLAIN framework:

Problem 1 (Private Confidence Interval of Question). Given a dataset D , a query q , a DP mechanism \mathcal{M} , a privacy budget ρ_q , a confidence level γ , and a user question $\langle \alpha_i, \alpha_j, > \rangle$ on the noisy query answers output by \mathcal{M} satisfying ρ_q -zCDP, find a confidence interval (see Definition 2.10) for the user question $\langle \alpha_i, \alpha_j, > \rangle$ for $\delta_i - \delta_j$ at confidence level γ without extra privacy cost.

In Phase-2, the framework returns a confidence interval of $\delta_i - \delta_j$ to the user. If it includes zero or negative numbers, it is possible that $\delta_i < \delta_j$, and the user’s observation of $\hat{\delta}_i > \hat{\delta}_j$ is the result of the noise added by the DP mechanism. In such cases, the user may stop at Phase-2. If the user is satisfied with the confidence interval for the validity of the question, she can proceed to Phase-3.

Influence Function (Phase-3). When considering DP, the order of the explanation predicates is perturbed by the noise we add to the influences according to the sensitivity of the influence function (discussed in detail in Section 4.3.1). To provide useful explanations, this sensitivity needs to be low, which means the influence does not change too much by adding or removing a tuple from the database. For example, a counting query that outputs the database size n has sensitivity 1, since its result can only change by 1 for any neighboring databases. Following this concept, we propose the second and a core problem for the DPXPLAIN framework, which is also critical to the subsequent problems defined below.

Problem 2 (Influence Function with Low Sensitivity). Find an influence function $INF : \mathcal{P} \rightarrow \mathbb{R}$ that maps an explanation predicate to a real number and has low sensitivity.

Private Top- k Explanations (Phase-3). In DPXPLAIN, to satisfy DP, in Phase-3 we output the top- k explanation predicates ordered by the noisy influences, and release the influences and ranks of these predicates in the form of confidence intervals to describe the uncertainty. To achieve this goal, we tackle the following three sub-problems.

Problem 3 (Private Top- k Explanation Predicates). Given a set of explanation predicates \mathcal{P} , an integer k , and a privacy parameter ρ_{Topk} , find the top- k highest influencing predicates p_1, p_2, \dots, p_k from \mathcal{P} while satisfying ρ_{Topk} -zCDP.

⁷Our framework can handle more general user questions involving single group or more than two groups; details are deferred to the full version [2].

Problem 4 (Private Confidence Interval of Influence). Given a confidence level γ , k explanation predicates p_1, p_2, \dots, p_k , and a privacy parameter ρ_{Influ} , find a confidence interval $I_{Influ} = [L_{Influ}^L, U_{Influ}^U]$ for influence $INF^1 p_u^0$ at confidence level γ for each $u \in \{1, \dots, k\}$ satisfying ρ_{Influ} -zCDP (overall privacy budget).

Problem 5 (Private Confidence Interval of Rank). Given a confidence level γ , k explanation predicates p_1, p_2, \dots, p_k , and a privacy parameter ρ_{Rank} , find a confidence interval $I_{Rank} = [L_{Rank}^L, U_{Rank}^U]$ for rank of p_u at confidence level γ for each $u \in \{1, \dots, k\}$ satisfying ρ_{Rank} -zCDP (overall privacy budget).

4 COMPUTING EXPLANATIONS UNDER DP

Next we provide solutions to problems 1, 2, 3, 4, and 5 in Sections 4.1, 4.2, 4.3.1, 4.3.2, and 4.3.3 respectively, and analyze their properties. We summarize the entire DPXPLAIN framework in Section 4.4.

4.1 Confidence Interval for a User Question

For **Problem 1**, the goal is to find a confidence interval of $o_i - o_j$ for the user question at the confidence level γ without extra privacy cost in Phase-2. We divide the solution into two cases. (1) When the aggregation is COUNT or SUM, the noisy difference $\hat{o}_i - \hat{o}_j$ follows Gaussian distribution, which leads to a natural confidence interval. (2) When the aggregation is AVG, the noisy difference does not follow Gaussian distribution, but we show that the confidence interval in this case can be derived through multiple partial confidence intervals. The solutions below only take the noisy query result as input, which does not incur extra privacy loss according to the post-processing property of DP (Proposition 2.9). The pseudo codes can be found in the full version [2].

Confidence interval for COUNT and SUM. For a COUNT or SUM query, recall from Section 2 that \hat{o}_i and \hat{o}_j are produced by adding Gaussian noises to o_i and o_j with some noise scale σ . Therefore, the difference between \hat{o}_i and \hat{o}_j also follows Gaussian distribution with mean $o_i - o_j$ and scale $\sqrt{2}\sigma$ (since the variance is $2\sigma^2$). Following the standard properties of Gaussian distribution, the interval with center c as $\hat{o}_i - \hat{o}_j$ and margin m as $\frac{1}{\sqrt{2}} \sqrt{2}\sigma \text{erf}^{-1}(\gamma)^8$, or $(c-m, c+m)$, is a γ level confidence interval of $o_i - o_j$ [96].

Confidence interval for AVG. For an AVG query, even the single noisy answer \hat{o}_i does not follow Gaussian distribution, because it is a division between two Gaussian variables as described in Section 2: $\hat{o}_i = \hat{o}_i^S \cdot \hat{o}_i^C$. However, we can still infer a range for \hat{o}_i based on the confidence intervals of \hat{o}_i^S and \hat{o}_i^C . More specifically, we first derive partial confidence intervals for o_i^S and o_i^C as discussed above, denoted by I^S and I^C , individually at some confidence level β . Let $I^A = I^S \cdot I^C := \{x \cdot y \mid x \in I^S, y \in I^C\}$ to be the set that includes all possible divisions between any numbers from I^S and I^C . If I^C contains zero, we return a trivial confidence interval $[-1, 1]$ that is always valid. Otherwise, I^A is a $2\beta - 1$ level confidence interval for the division, as stated in the following proposition.

LEMMA 4.1. Given I^S and I^C as two β level confidence intervals of o_i^S and o_i^C separately, the derived interval $I^A = \{x \cdot y \mid x \in I^S, y \in I^C\}$ is a $2\beta - 1$ level confidence interval of $o_i^S \cdot o_i^C$.

PROOF. The following holds:
 $Pr[o_i^S \cdot o_i^C \in I^A] = Pr[o_i^S \in I^S \wedge o_i^C \in I^C] = Pr[o_i^S \in I^S] \cdot Pr[o_i^C \in I^C] = \beta \cdot \beta = 2\beta - 1$
The first inequality above is due to fact that the second event is sufficient for the first event: if two numbers are from I^S and I^C , their division belongs to the set I^A by definition. The next inequality holds by applying the union bound. The third inequality is by definition. \square

Furthermore, the difference $\hat{o}_i - \hat{o}_j$ is a subtraction between two ratios of two Gaussian variables, which can be expressed as an arithmetic combination of multiple Gaussian variables: $\hat{o}_i - \hat{o}_j = X_i \cdot Y_i - X_j \cdot Y_j$, where $X_t = N(o_t^S, \sigma_S^2)$ and $Y_t = N(o_t^C, \sigma_C^2)$ for $t \in \{i, j\}$. Similar to Lemma 4.1, we can derive the confidence interval for $\hat{o}_i - \hat{o}_j$ based on 4 partial confidence intervals of $o_i^S, o_i^C, o_j^S,$ and o_j^C instead of 2. The confidence level we set for each partial confidence interval is $\beta = \frac{1 - \gamma}{4}$ by applying union bound on the failure probability $1 - \gamma$ that one of the four variables is outside its interval. After we have 4 partial confidence intervals $I_i^S, I_i^C, I_j^S,$ and I_j^C for $o_i^S, o_i^C, o_j^S,$ and o_j^C separately, similar to Lemma 4.1, we combine them together as $I^A = I_i^S \cdot I_i^C \cdot I_j^S \cdot I_j^C$ and derive the confidence interval for $o_i - o_j$ as $[\inf I^A, \sup I^A]$, which is guaranteed to be at confidence level γ . If 0 is included in either I_i^C or I_j^C , we set the confidence interval to be $[-1, 1]$ instead. Although there is no theoretical guarantee of the interval width, from two case studies in Section 5.2, we demonstrate narrow confidence intervals of AVG queries in practice, and observe no extreme case $[-1, 1]$ in the experiments.

4.2 Influence Function with Low Sensitivity

For **Problem 2**, the goal is to design an influence function that has low sensitivity. Inspired by PrivBayes [100], we start by adapting a known influence function to our framework.

Our influence function of an explanation predicate with respect to a comparison user question is inspired by the Scorpion framework [98], where the user questions seek explanations for outliers in the results of a group-by aggregate query. Scorpion identifies predicates on the input that cause the outliers to disappear from the output. Given the group-by aggregation query shown in Figure 3 and a group $\alpha_i \in \text{dom } A_{gb}$, recall from Section 2 that the true aggregate value for α_i is $o_i = \text{agg}^1 A_{agg}, g_i^1 D^{00}$, where $g_i = \phi^1 A_{gb} = \alpha_i^0$, i.e., $g_i^1 D^0$ denotes the set of tuples that contribute to the group α_i .

Scorpion measures the influence of an explanation predicate p to some group α_i as the ratio between the change of output aggregate value and the change of group size:

$$\frac{\text{agg}^1 g_i^1 D^{00} - \text{agg}^1 g_i^1 p^1 D^{00}}{|g_i^1 p^1 D^{00}|} \quad (1)$$

Here $p^1 D^0$ denotes $D - p^1 D^0$, i.e., the set of tuples in D that do not satisfy the predicate p . To adapt this influence function to DPXPLAIN, we make the following two changes.

First, it should measure the influence w.r.t. the comparison from the user question $\alpha_i, \alpha_j, >$ instead of a single group. A natural extension is to change the target aggregate on g_i in the numerator in (1) to the difference between the aggregate values of two groups

⁸ erf^{-1} is the inverse function of the error function $\text{erf } z = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$.

⁹In the algorithm, we only need the maximum and the minimum of the set to construct the interval, which can be solved by a numerical optimizer.

g_i, g_j before and after applying the explanation predicate p , and change the denominator as the maximum change in g_i or g_j when p is applied, which gives the following influence function:

$$\frac{(agg^1 g_i^1 D^{00} \quad agg^1 g_j^1 D^{00}) \quad (agg^1 g_i^1: p^1 D^{000} \quad agg^1 g_j^1: p^1 D^{000})}{\max^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}|} \quad (2)$$

Second and more importantly, in DPXPLAIN, we need to preserve DP when we use influence function to sort and rank multiple explanation predicates, or to release the influence and rank of an explanation predicate. Therefore, **we need to account for the sensitivity of the influence function**, which is determined by the worst-case change of influence when a tuple is added or removed from the database. If the predicate only selects a small number of tuples, the denominator in (2) is small and thus changing the denominator in (2) by one (when a tuple is added or removed) can result in a big change in the influence as illustrated in the following example, making (2) unsuitable for DPXPLAIN.

Example 4.2 (The Issue of the Influence Sensitivity). Suppose there are two groups α_i and α_j in D with 1000 tuples in each, aggregate function $agg = SUM$ on attribute A_{agg} with domain $\gg 0, 100\%$, and the explanation predicate p matches only 1 tuple from the group α_i with $A_{agg} = 100$ and no tuple from α_j . Suppose $agg^1 g_i^1 D^{00} = 20,000$, $agg^1 g_j^1 D^{00} = 10,000$, then $agg^1 g_i^1: p^1 D^{000} = 19,900$ and $agg^1 g_j^1: p^1 D^{000} = 10,000$. Therefore, from Equation (2), the influence of p is $\frac{19,900 - 10,000}{19,900 + 10,000} = \frac{9,900}{29,900} \approx 0.33$ on the original database D . However, suppose a new tuple that satisfies p and belongs to group α_i is added with $A_{agg} = 2$. Now the influence in Equation (2) becomes $\frac{19,902 - 10,000}{19,902 + 10,000} = \frac{9,902}{29,902} \approx 0.33$. Note that while we added a tuple contributing only 2 to the sum, it led to a change of $100 - 51 = 49$ to the influence function because of the small denominator.

Therefore, we propose a new influence function that is inspired by Equation (2) but has lower sensitivity. Note that the denominator in Scorpion's influence function in Equation (2) acts as a normalizing factor, whose purpose is to penalize the explanation predicate that selects too many tuples, e.g., to prohibit the removal of the entire database by a dummy predicate. To have a similar normalizing factor with low sensitivity, we multiply the numerator in Equation (2) by $\frac{\min^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}|}{\max^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}| + 1}$. From this new normalizing factor, the numerator captures the minimum of the number of tuples that are not removed from each group, and the denominator keeps the normalizing factor in the interval $\gg 0, 1\%$ and does not change for different explanation predicates. Similar to Scorpion, if $p^1 D^0$ constitutes a large fraction of D (e.g., if $p^1 D^0 = D$), then the normalizing factor is small, reducing the value of the influence. Also note that, unlike standard SQL query answering where only non-empty groups are shown in the results, in DP, all groups from the actual domain have to be considered, hence unlike Equation (1), $g_i^1 D^0, g_j^1 D^0$ could be zero, hence 1 is added in the denominator to avoid division by zero. When $agg = AVG$, we remove the constant denominator to boost the signal of the influence and keep the sensitivity low, which will be discussed in the sensitivity analysis after Proposition 4.4 and in Example 4.5.

Definition 4.3 (Influence of Explanation Predicates). Given a database D , a query q as shown in Figure 3, and a user question $\langle \alpha_i, \alpha_j, \succ \rangle$, the influence of an explanation predicate p is defined as

$INF^1 p; \langle \alpha_i, \alpha_j, \succ \rangle, D^0$, or simply $INF^1 p^0$ when clear from context:

$$INF^1 p^0 = \frac{(agg^1 g_i^1 D^{00} \quad agg^1 g_j^1 D^{00}) \quad (agg^1 g_i^1: p^1 D^{000} \quad agg^1 g_j^1: p^1 D^{000})}{\begin{cases} \frac{\min^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}|}{\max^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}| + 1} & \text{for } agg \in \{COUNT, SUM\} \\ \min^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}| & \text{for } agg = AVG \end{cases}} \quad (3)$$

The next proposition summarizes the sensitivity of eq. (3).

PROPOSITION 4.4. [Influence Function Sensitivity] *Given an explanation predicate p and a user question with respect to a group-by query with aggregation agg , the following holds:*

- (1) If $agg = COUNT$, the sensitivity of $INF^1 p^0$ is 4.
- (2) If $agg = SUM$, the sensitivity of $INF^1 p^0$ is $4 A_{agg}^{max}$.
- (3) If $agg = AVG$, the sensitivity of $INF^1 p^0$ is $16 A_{agg}^{max}$.

We give an intuitive proof as follows, where the formal proofs are deferred to the full version [2] due to space restrictions. When $agg = COUNT$, we combine two group differences $(agg^1 g_i^1 D^{00} \quad agg^1 g_j^1 D^{00}) \quad (agg^1 g_i^1: p^1 D^{000} \quad agg^1 g_j^1: p^1 D^{000})$ into a single group difference as $agg^1 g_i^1: p^1 D^{000} \quad agg^1 g_j^1: p^1 D^{000}$, which is considered as a subtraction between two counting queries. We prove that the sensitivity of a counting query after a multiplication with the normalizing factor will multiply its original sensitivity by 2. Since we have two counting queries, the final sensitivity is 4. When $agg = SUM$, the proof is similar except we need to multiply the final sensitivity by A_{agg}^{max} , the maximum absolute domain value of A_{agg} . For AVG, we view it as a summation of 4 AVG queries that times with $\min^1 |g_i^1: p^1 D^{000}|, |g_j^1: p^1 D^{000}|$. Intuitively, we change AVG to SUM and bound the sensitivity. This sensitivity now becomes relatively small since we have amplified the influence.

Intuitively, the sensitivity of $INF^1 p^0$ is low compared to its value. When $agg = COUNT$, $INF^1 p^0$ is $O^1 n^0$ and Δ_{INF} is $O^1 1^0$, where n is the size of database. When $agg \in \{SUM, AVG\}$, $INF^1 p^0$ is $O^1 n A_{agg}^{max}$ and Δ_{INF} is $O^1 A_{agg}^{max}$. Therefore, the sensitivity of influence Δ_{INF} is low compared to the influence itself. However, as the example below shows, if we define the influence function for AVG the same way as COUNT or SUM, both $INF^1 p^0$ and Δ_{INF} will become $O^1 A_{agg}^{max}$, which makes the sensitivity (relatively) large.

Example 4.5 (The Issue with AVG Influence). Consider an AVG group-by query where the domain of the aggregate attribute is $\gg 0, 100\%$, and an explanation predicate p such that for group α_i we have 2 tuples with $AVG^1 g_i^1 D^{00} = 100 \cdot 2 = 50$, $AVG^1 g_i^1: p^1 D^{000} = 0 \cdot 1 = 0$, and for group α_j we have two tuples with $AVG^1 g_j^1 D^{00} = 100 \cdot 2 = 50$ and $AVG^1 g_j^1: p^1 D^{000} = 100 \cdot 2 = 50$. Suppose we define the influence function for AVG the same way as COUNT or SUM, therefore the influence of p in Equation (3) is $INF^1 p^0 = \frac{50 - 50}{50 + 50} = 0$. However, suppose we remove the single tuple from g_i , so $|g_i^1: p^1 D^{000}|$ becomes 0, now the influence in Equation (3) (for COUNT/SUM) becomes 0. Note that a single removal of a tuple completely changes the influence to 0, and this change is equal to the influence itself, which is relatively large and therefore is not a good choice for AVG.

Note that the user question " $\langle \alpha_i, \alpha_j, \succ \rangle$ " is asked based on the noisy results $\hat{o}_i > \hat{o}_j$, while the influence function uses the true results, i.e., even if $o_i < o_j$, we still consider $agg^1 g_i^1 D^{00} > agg^1 g_j^1 D^{00}$ in $INF^1 p^0$. Hence $INF^1 p^0$ can be positive or negative and

removing tuples satisfying p can make the gap smaller or larger. In the full version [2], we show that $\text{INF}^1 p^\circ$ is not monotone with p -s.

4.3 Private Top-k Explanations

In this section, we discuss the computation of the top-k explanation predicates and the confidence intervals of influences and ranks.

4.3.1 Problem 3: Private Top-k Explanation Predicates. The goal is to find with DP the top-k explanation predicates from a set of explanation predicates \mathcal{P} in terms of their (true) influences $\text{INF}^1 p^\circ$, which is the first step in Phase-3 of DPXPLAIN (Figure 1). Note that simply choosing the *true* top-k explanation predicates in terms of their $\text{INF}^1 p^\circ$ is not differentially private.

In DPXPLAIN, we adopt the **One-shot Top-k mechanism** [36, 37] to privately select the top-k. It works as follows. For each explanation predicate $p \in \mathcal{P}$, it adds a Gumbel noise¹⁰ to its influence with scale $\sigma = 2\Delta_{\text{INF}} \sqrt{k \cdot 8\rho_{\text{Top}k}^\circ}$, where Δ_{INF} is the sensitivity of the influence function (discussed in Proposition 4.4), reorders all the explanation predicates in descending order by their noisy influences, and outputs the first k explanation predicates. It satisfies $\rho_{\text{Top}k}$ -zCDP [16, 31, 36, 37, 79], since it is equivalent to iteratively applying k exponential mechanisms [41], where each satisfies $\epsilon^2 \cdot 8$ -zCDP [16, 31, 37, 79] and $\epsilon = \sqrt{8\rho_{\text{Top}k}^\circ \cdot k}$ [36, 37]. Therefore, in total it satisfies $1k\epsilon^2 \cdot 8^\circ$ -zCDP by the sequential composition property (Proposition 2.9) which is also $\rho_{\text{Top}k}$ -zCDP. The returned list of top-k predicates is close to that of the true top-k in terms of their influences; the proof is based on the utility proposition of the exponential mechanism in Theorem 3.11 of [41]. Since this algorithm iterates over each explanation predicate, the time complexity is proportional to the size of the explanation predicate set \mathcal{P} . By Definition 3.3, this number is $O\left(\binom{m}{l} N^l\right)$, where N is the maximum domain size of an attribute, l is the number of conjuncts in the explanation predicate and m is the number of attributes. In our experiments (Section 5), we fix $l = 1$ and use all the singleton predicates as the set \mathcal{P} , so its size is linear in the number of attributes. We summarize the properties of this approach in the following proposition and defer the pseudo codes and proofs to [2].

PROPOSITION 4.6. *Given an influence function INF with sensitivity Δ_{INF} , a set of explanation predicates \mathcal{P} , a privacy parameter $\rho_{\text{Top}k}$ and a size parameter k , the following holds:*

- (1) *One-shot Top-k mechanism finds k explanation predicates while satisfying $\rho_{\text{Top}k}$ -zCDP.*
- (2) *Denote by OPT^{1i° the i -th highest (true) influence, and by \mathcal{M}^{1i° the i -th explanation predicate selected by the One-shot Top-k mechanism. For $\forall t$ and $\forall i \in \{1, 2, \dots, k\}$, we have*

$$\Pr_{\mathcal{M}^{1i^\circ}} \left[\mathcal{M}^{1i^\circ} \in \text{OPT}^{1i^\circ} \pm \frac{2\Delta_{\text{INF}}}{\sqrt{8\rho_{\text{Top}k}^\circ \cdot k}} \ln \binom{k}{i} \right] \geq e^{-t} \quad (4)$$

Example 4.7. Reconsider the user question in Figure 1c. For this question, we have in total 103 explanation predicates as the set of explanation predicates. The privacy budget $\rho_{\text{Top}k} = 0.05$, the size parameter $k = 5$, and the sensitivity $\Delta_{\text{INF}} = 16$. For each of the explanation predicate, we add a Gumbel noise with scale $\sigma = 113$ to their influences. For example, for the predicates shown in Figure 1d,

¹⁰For a Gumbel noise $Z \sim \text{Gumbel}^1(\sigma)$, its CDF is $\Pr[Z \leq z] = \exp^{-\exp^{-z/\sigma}}$.

their noisy influences are 990, 670, 645, 475, 440, which are the highest 5 among all the noisy influences. The true influences for these five ones are 547, 501, 555, 434, 118. To see how close it is to the true top-5, we compare their true influences with the true highest five influences: 555, 547, 501, 434, 252, which shows the corresponding differences in terms of influence are 8, 46, 54, 0, 134. By Equation (4), the probability that such difference is beyond 864 is at most 5% for each explanation predicate. Finally, we sort explanation predicates by their noisy influences and report the top-k. These k predicates will be reordered as discussed in Section 4.4.

4.3.2 Problem 4: Private Confidence Interval of Influence.

The goal is to generate a confidence interval of influence $\text{INF}^1 p^\circ$ (Definition 4.3) of each explanation predicate $\text{INF}^1 p_1^\circ, \text{INF}^1 p_2^\circ, \dots, \text{INF}^1 p_k^\circ$ from the selected top-k (Section 4.3.1). For each $\text{INF}^1 p_i^\circ$, we apply the Gaussian mechanism (Theorem 2.8) with privacy budget $\rho_{\text{Inf}^1 p_i^\circ} \cdot k$ to release a noisy influence $\widehat{\text{INF}}_i$ with noise scale $\sigma = \Delta_{\text{INF}} \cdot \sqrt{2\rho_{\text{Inf}^1 p_i^\circ} \cdot k}$. The sensitivity term Δ_{INF} is determined by Proposition 4.4. Following the standard properties of Gaussian distribution, for each $\text{INF}^1 p_i^\circ$, we set the confidence interval by a center c as $\widehat{\text{INF}}_i$ and a margin m as $\frac{1}{2\sigma} \text{erf}^{-1}(\gamma)$, or $(c-m, c+m)$, as a γ level confidence interval of $\text{INF}^1 p_i^\circ$ [96]. Together, it satisfies $\rho_{\text{Inf}^1 p_i^\circ}$ -zCDP according to the composition property by Proposition 2.9. Pseudo codes can be found in the full version [2].

4.3.3 Problem 5: Private Confidence Interval of Rank.

The goal is to find the confidence interval of the rank of each explanation predicate from the selected top-k (Section 4.3.1). We denote $\text{rank}^1 p^\circ$ as the rank of $p \in \mathcal{P}$ by the natural ordering of the predicates imposed by their (true) influences according to the influence function INF , and denote rank^{1t° (for an integer $1 \leq t \leq |\mathcal{P}|$) as the predicate ranked in the t -th place according to INF . One trivial example of a confidence interval of rank is $\{1, |\mathcal{P}|\}$, which has no privacy loss and always includes the true rank.

Unlike the sensitivity of the influence function, the sensitivity of $\text{rank}^1 p^\circ$ is high, since adding one tuple could possibly change the highest influence to be the lowest and vice versa. Fortunately, we can employ a critical observation about rank and influence.

PROPOSITION 4.8. *Given a set of explanation predicates \mathcal{P} , an influence function INF with global sensitivity Δ_{INF} , and an integer $1 \leq t \leq |\mathcal{P}|$, $\text{INF}^1 \text{rank}^{1t^\circ}$ has sensitivity Δ_{INF} .*

The intuition behind this proof (details in [2]) is that, fixing an explanation predicate $p = \text{rank}^{1t^\circ}$, for a neighboring database, if its influence is increased, its rank will be moved to the top which pushes down other explanation predicates with lower influences, so the influence at the rank t in the neighboring database is still low. For a target explanation predicate p , since both $\text{INF}^1 p^\circ$ and $\text{INF}^1 \text{rank}^{1t^\circ}$ have low sensitivity as Δ_{INF} , intuitively we can check whether t is close to the rank of p by checking whether their influences $\text{INF}^1 p^\circ$ and $\text{INF}^1 \text{rank}^{1t^\circ}$ are close by adding a little noise to satisfy DP. Given this observation, we devise a binary-search-based strategy to find the confidence interval of rank.

Noisy binary search mechanism. We decompose the problem into finding two bounds of the confidence interval separately by a subroutine $\text{RANKBOUND}^1(p, \rho, \beta, \text{dir}^\circ)$ that guarantees that it will find a lower ($\text{dir} = _1$) or upper ($\text{dir} = _2$) bound of rank with

probability β for the explanation predicate p using privacy budget ρ . We divide the privacy budget ρ into two parts by a parameter η $2^{-1}, 1^0$ and return $\text{RANKBOUND}^1 p_u, \eta\rho, \beta, 1^0, \text{RANKBOUND}^1 p_u, 1^1 \eta^0 \rho, \beta, 1^0$ as the confidence interval of rank for each predicate p_u for $u \in \{1, \dots, k\}$, where $\rho = \rho_{\text{Rank}} \cdot k$ to divide the total privacy budget equally, and $\beta = \gamma \cdot 1^0 \cdot 2$ to ensure a confidence of γ .

The subroutine $\text{RANKBOUND}^1 p, \rho, \beta, \text{dir}^0$ works as follows. It is a noisy binary search with at most $N = \lceil \log_2 |P| \rceil$ loops. We initialize the search pointers $t_{\text{low}} = 1$ and $t_{\text{high}} = |P|$ as the two ends of possible ranks. Within each loop, we check the difference of influences at $t = \lfloor \frac{t_{\text{low}} + t_{\text{high}}}{2} \rfloor$ by adding a Gaussian noise:

$$\hat{s} = \text{INF}^1 p^0 - \text{INF}^1 \text{rank}^1 t^0 \sim \mathcal{N}(0, \sigma^2) \quad (5)$$

The noise scale is set as $\sigma = \sqrt{2} \Delta_{\text{INF}}^0 \cdot \sqrt{2^1 \rho \cdot N^0}$ to satisfy $\rho \cdot N$ -zCDP. Instead of comparing the noisy difference \hat{s} with 0 to check whether t is a close bound of $\text{rank}^1 p^0$, we compare it with the following slack constant ξ so that w.h.p. t is a true bound of $\text{rank}^1 p^0$.

$$\xi = \sigma \sqrt{2 \ln \frac{1}{\beta} \cdot 1} \quad \text{dir} \quad (6)$$

We update the binary search pointers by the comparison as follows: if $\hat{s} \geq \xi$, we set $t_{\text{high}} = \max\{t - 1, 1\}$, otherwise $t_{\text{low}} = \min\{t + 1, |P|\}$. The binary search stops when $t_{\text{high}} - t_{\text{low}} = 1$ and returns t_{high} as the rank bound. We defer the pseudo codes to [2].

Example 4.9. Figure 4 shows an example of RANKBOUND for finding the upper bound of the confidence interval for $\text{rank}^1 p^0$ for some explanation predicate p (with true rank 3 shown in red). The upper part of the figure shows the influences of all the explanation predicates in descending order, and the lower part shows the status of the binary search pointers in each loop. The search contains three loops starting from $t_{\text{low}} = 1$ and $t_{\text{high}} = 15$. Within each loop, to illustrate the idea, it is equivalent to adding a Gaussian noise to $\text{INF}^1 \text{rank}^1 t^0$, which is shown as a blue circle, compare it with $\text{INF}^1 p^0 - \xi$, which is shown as a dashed line, and update the pointers accordingly. For example, in loop 1, the blue circle 1 is in the green region, so the pointer t_{high} is moved from 15 to 7 (shown in the lower part). Finally, it breaks at $t_{\text{low}} = t_{\text{high}} = 5$.

We now show that **noisy binary search mechanism** satisfies the privacy requirement, and outputs valid confidence intervals. In Section 5, we show that the interval width is empirically small.

THEOREM 4.10. *Given a database D , a predicate space \mathcal{P} , an influence function INF with sensitivity Δ_{INF} , explanation predicates p_1, p_2, \dots, p_k , a confidence level γ , and a privacy parameter ρ_{Rank} , noisy binary search mechanism returns confidence intervals I_1, I_2, \dots, I_k such that*

- (1) Noisy binary search mechanism satisfies ρ_{Rank} -zCDP.
- (2) For $\forall u \in \{1, \dots, k\}$, I_u is a γ level confidence interval of $\text{rank}^1 p_u^0$.

The proof of item 1 follows from the composition theorem and the property of Gaussian mechanism [14]. The proof of item 2 is based on the property of the random binary search. We defer the formal proofs and a weak utility bound to the full version [2].

4.4 Putting it All Together

We now show how all the steps fit together into DPXPLAIN .

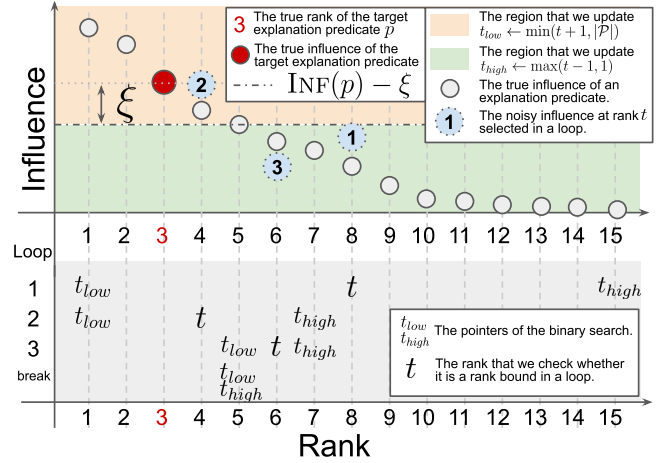


Figure 4: The execution of RANKBOUND for finding the upper bound of the confidence interval of rank for the predicate p (with true rank 3 shown in red) from a toy example.

Relative Influence. Recall that the influence defined by Definition 4.3 is the difference of $|o_i - o_j|$ before and after removing the tuples related to an explanation predicate (first term), and multiplies with a normalizer to penalize trivial predicates (second term). Since the absolute value of influence is hard to interpret, to help user better understand the confidence interval of influence, we show the *relative influence* compared to the original difference $|o_i - o_j|$ as a percentage. However, we cannot divide the influence by $|o_i - o_j|$ since using the actual data values will incur additional privacy loss, hence, for SUM and COUNT we divide the true influence by $|j \delta_i - j \delta_j|$ as an approximation since the normalizer in the second term is bounded in $[0, 1]$. However, when $\text{agg} = \text{AVG}$, the normalizer $\min\{j \delta_i^C, j \delta_j^C\}$ (second term) is not bounded in $[0, 1]$, so we further divide the influence by another constant, the minimum of the noisy counts/sizes of the groups, i.e., $\min\{\delta_i^C, \delta_j^C\}$ (approximating the upper bound $\min\{j \delta_i^C, j \delta_j^C\}$ of the normalizer to avoid additional privacy loss). In summary, we define the relative influence $\widetilde{\text{INF}}^1 p; \alpha_i, \alpha_j, >^0, D^0$, or simply $\widetilde{\text{INF}}^1 p^0$, as follows, which is only used for display purposes.

$$\widetilde{\text{INF}}^1 p^0 = \text{INF}^1 p^0 \cdot \begin{cases} \frac{|j \delta_i - j \delta_j|}{|j \delta_i^C - j \delta_j^C|} & \text{for } \text{agg} \in \{\text{COUNT}, \text{SUM}\} \\ \frac{|j \delta_i - j \delta_j|}{\min\{\delta_i^C, \delta_j^C\}} & \text{for } \text{agg} = \text{AVG} \end{cases}$$

Explanation Table. We define the explanation table as follows.

Definition 4.11 (Explanation Table containing top- k explanations). Given a database D , a group-by aggregate query q as shown in Figure 3, a user question $\langle \alpha_i, \alpha_j, >^0$, a predicate space \mathcal{P} , a confidence level γ , and an integer k , a table of top- k explanations is a list of k 5-element tuples $\langle p_u, \text{relinflu}_u^L, \text{relinflu}_u^U, \text{rank}_u^L, \text{rank}_u^U \rangle$ for $u = 1, 2, \dots, k$ such that p_u is an explanation predicate, $\text{relinflu}_u^L, \text{relinflu}_u^U$ is a confidence interval of relative influence $\widetilde{\text{INF}}^1 p_u^0$ with confidence level γ , and $\text{rank}_u^L, \text{rank}_u^U$ is a confidence interval of $\text{rank}^1 p_u^0$ with confidence level γ .

Sorting the explanations in the explanation table. Since this table contains the bounds of the influences and ranks it is natural to

present the table as a sorted list. Since the numbers in the table are generated by random processes, each column may imply a different sorting. In this paper, we sort the selected top-k explanations by the upper bound of the relative influence CI (the third column in Figure 1d) in descending order; if there is a tie, we break it using the upper bound of the rank confidence interval (the fifth column in Figure 1d). Finding a principled way for sorting the explanation predicates is an intriguing subject of future work.

Overall DP guarantee. We summarize the privacy guarantee of DPXPLAIN as follows: (i) the private noisy query answers returned by Gaussian mechanism in Phase-1 satisfy ρ_q -zCDP together (see Section 2); (ii) Phase-2 only returns the confidence intervals of the noisy answers in Phase-1 with zero additional privacy loss (discussed in Section 4.1); (iii) Phase-3 returns k explanation predicates and their upper and lower bounds on relative influence and ranks given a required confidence interval with three privacy parameters ρ_{Topk} , ρ_{Influ} , ρ_{Rank} (discussed in Section 4.3.1, 4.3.2 and 4.3.3). The following theorem summarizes the total privacy guarantee.

THEOREM 4.12. *Given a group-by query q and a user question comparing two aggregate values in the answers of q , the DPXPLAIN framework guarantees $^1\rho_q \succ \rho_{Topk} \succ \rho_{Influ} \succ \rho_{Rank}^0$ -zCDP.*

5 EXPERIMENTS

In this section, we evaluate the quality and efficiency of the explanations generated by DPXPLAIN. To our knowledge, there are no existing benchmarks for explanations for query answers (even without privacy consideration) in the database research literature. We have implemented DPXPLAIN [1] in Python 3.7.4 using the Pandas [92], NumPy [51], and SciPy [95] libraries. All experiments were run on Intel i7-7700 CPU @ 3.60GHz with 32 GB of RAM.

5.1 Experiment Setup

We first detail the data, queries, questions, and parameters.

Datasets. We consider two datasets in our experiments.

IPUMS-CPS (real data): A dataset of Current Population Survey from the U.S. Census Bureau [47] with 1,146,552 tuples from the year 2011 to 2019. The dataset contains 8 categorical attributes where domain sizes vary from 3 to 36 and one numerical attribute. The attribute AGE is discretized as 10 years per range, e.g., [0,10] is considered a single value. To set the domain of numerical attributes, we only include tuples with attribute INCTOT (the total income) smaller than 200k as a domain bound.

German-Credit (synthetic data): A corrected collection of credit data [50]. It includes 20 attributes where the domain sizes vary from 2 to 11 and a numerical attribute. Attributes duration, credit-amount, and age are discretized. The domain of attribute good-credit is zero or one. We synthesize the dataset to 1 million rows by combining a Bayesian network learner [7] and XGBoost [13] following the strategy of QUAIL [80].

Queries and Questions. The queries and questions used on the experiments are shown in Table 1.

Default setting of DPXPLAIN. Unless mentioned otherwise, the following default parameters are used (also for the motivating example): $\rho_q = 0.1$, $\rho_{Topk} = 0.5$, $\rho_{Influ} = 0.5$, $\rho_{Rank} = 1.0$, $\gamma = 0.95$, $k = 5$, $\eta = 0.1$, and the number of conjuncts in explanation predicates $l = 1$ (Definition 3.3). We choose $\eta = 0.1$ to allocate more

Table 1: Queries and questions for the experiments; Valid indicates if it is a valid question on the hidden true data.

Data	Query	Question	Valid
IPUMS-CPS	q_1 : AVG(INCTOT) by SEX	I1: Why Male > Female ?	Yes
	q_2 : INCTOT by RELATE	I2: Why Grandchild > Foster children ? I3: Why Head/householder > Spouse ?	Yes No
	q_3 : INCTOT by EDUC	I4: Why Bachelor > High school ? I5: Why Grade 9 > None or preschool ?	Yes No
German-Credit	q_4 : AVG(good-credit) by status	G1: Why no balance > no chk account ?	Yes
	q_5 : AVG(good-credit) by purpose	G2: Why car (new) > car (used) ? G3: Why business > vacation ?	Yes No
	q_6 : AVG(good-credit) by residence	G4: Why "< 1 yr" > ">= 7 yrs" ? G5: Why "[1, 4) yrs" > "[4, 7) yrs" ?	Yes No

explanation predicate	Rel Influ 95%-CI		Rank 95%-CI		Rel Influ (hidden)	Rank (hidden)
	L	U	L	U		
RELATE = "Head/householder"	12.18%	12.52%	1	1	12.41%	1
EDUC = "Bachelor's degree"	7.10%	7.45%	2	3	7.32%	2
RACE = "White"	6.41%	6.75%	2	5	6.54%	3
RELATE = "Spouse"	5.70%	6.04%	2	5	6.01%	4
CLASSWKR = "NIU"	3.83%	4.18%	2	6	4.22%	5

Figure 5: Phase-3 of DPXPLAIN for the case IPUMS-CPS.

privacy budget for the rank upper bound by our observation that the scores of explanation predicates have a long and flat tail, which intuitively means that a tight rank upper bound indicates a precise score and, thus, costs more privacy. For the total privacy budget, which is 2.1 by default, we provide experiments to show that reducing the budget of each component can still lead to a high utility for all questions except I2 and I5 in Table 1 (Figures 7, 8a, 9a, 9b).

5.2 Case Studies

Case-1, IPUMS-CPS. In **Phase-1**, the user submits a query q_1 from Table 1, and gets a noisy result: ("Female", 31135.25) and ("Male", 45778.46). The hidden true values are ("Female", 31135.78) and ("Male", 45778.39). Next, in **Phase-2**, since there is a gap of 14643.21 between two groups, the user asks a question I1 from Table 1. The framework then quantifies the noise in the question by reporting a confidence interval of the gap as (14636.63, 14649.79). Since the interval does not include zero, DPXPLAIN suggests that this is a valid question, which is correct. Finally, in **Phase-3**, the framework presents top-5 explanations to the user as Figure 5 shows. The last two columns are the true relative influences and ranks. We correctly find the top-5 explanation predicates, and the first and fourth explanations together suggests that a married man tends to earn more than a married woman, which is supported by the wage disparities in the labor market [94]. The second and third explanations also match the wage disparities within the educated group and white people. The total runtime for preparing the explanations in Phase-2 and Phase-3 is 67 seconds.

Case-2, German-Credit. In **Phase-1**, the user submits a query q_4 from Table 1, and gets a noisy result: ("no checking account", 0.526571) and ("no balance", 0.574447). The true hidden result is ("no checking account", 0.526574) and ("no balance",

explanation predicate	Rel Influ		95%-CI Rank		95%-CI	Rel Influ (hidden)	Rank (hidden)
	L	U	L	U			
existing-credits = "1"	77.90%	78.99%	1	1		78.16%	1
job = "skilled employee / official"	71.21%	72.29%	1	2		71.83%	2
sex-marst = "male : married/widowed"	54.34%	55.42%	2	4		55.10%	3
credit-amount = "(500, 2500)"	50.01%	51.10%	2	5		50.27%	4
credit-history = "no credits"	49.07%	50.16%	4	5		49.14%	5
taken/all credits paid back duly"							

Figure 6: Phase-3 of DPXPLAIN for the case German-Credit.

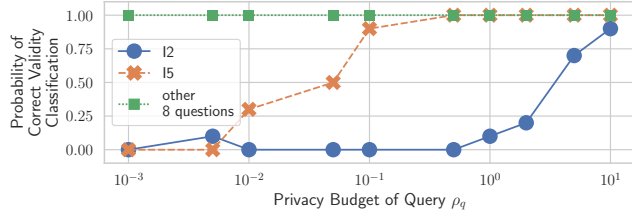


Figure 7: The probability of correctly validating user questions. All questions except I2 and I5 (Figure 7) are at 100%.

0.574466). Next, in **Phase-2**, since there is a gap of 0.047876 between two groups, the user asks a question G1 from Table 1. The framework then quantifies the noise in the question by reporting a confidence interval of the gap as (0.047786, 0.047967). Since the interval does not include zero, the framework suggests that this is a valid question, which is correct. Finally, in **Phase-3**, the framework presents top-5 explanations to the user as Figure 6 shows. The last two columns are the true relative influences and ranks. We correctly find the top-5 explanations, and the first explanation suggests that for a person who already has a credit in the bank, the bank tends to mark the credit as good with a higher probability than the case of no account if she has a checking account even with zero balance, which follows the intuition that a person having a credit account but no checking account is risky to the bank. The total runtime for preparing the explanations in Phase-2 and Phase-3 is 40 seconds.

5.3 Accuracy and Performance Analysis

We detail our experimental analysis for the different questions and configurations of DPXPLAIN. All results are averaged over 10 runs. **Correctness of noise interval.** In Phase-2 of DPXPLAIN, the validity of the question is suggested as follows: if the confidence interval contains non-positive numbers, the question is invalid, otherwise valid. From Figure 7, we find that 8 out of 10 questions (plotted together for clarity) from Table 1 are classified correctly with an accuracy of 100% given a wide range of privacy budget of query ρ_q . However, there are two questions, I2 and I5, only show high accuracy given a large privacy budget of $\rho_q = 10$. One reason is that the minimum group size involved in I2 and I5 is at least 600 and 60 times smaller compared to other questions, and, therefore, the partial confidence intervals in the denominators of the AVG query are low, which makes the final confidence intervals wider including negative numbers when it should not.

Accuracy of top-k explanation predicates. In Phase-3 of DPXPLAIN, we first select top-k explanation predicates. We measure the accuracy of the selection by Precision@k [52], the fraction of

the selected top-k explanation predicates that are actually ranked within top-k. Another experiment on the full ranking is included in the full version [2]. From Figure 8a, we find that the privacy budget of top-k selection ρ_{Topk} has a positive effect to Precision@k at $k = 5$ for various questions. When $\rho_{Topk} = 1.0$, all the questions except I2 and I5 have Precision@k = 0.8. The selection accuracy of question I2 and I5 are generally lower because of small group sizes, and, therefore, the influences of explanation predicates are small and the rankings are perturbed by the noise more significantly.

From Figure 8b, we find that the trend of Precision@k by k is different across questions and there is no clear trend that Precision@k increases as k increases. For example, for G3, it first decreases from $k=3$ to $k=5$, but increases from $k=5$ to $k=6$. When $k = 3$, most questions have high Precision@k; this is because the highest three influences are much higher than the others, which makes the probability high to include the true top three. With larger k, explanation predicates that have similar scores have an equal probability to be included in top-k and therefore the top-k selected by the algorithm are different from the true top-k selections. The relationship between Precision@k and k depends on the distribution of all the explanation predicate influences.

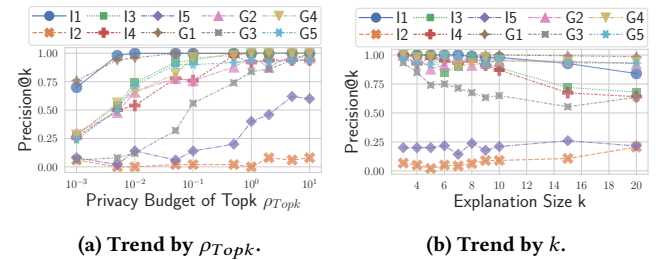


Figure 8: Precision@k of top-k selection by DPXPLAIN.

Precision of relative influence and rank confidence Interval (CI). In Phase-3, the last step is to describe the selected top-k explanation predicates by a CI of relative influence and rank for each. To measure the precision of the description, we adopt the measure of **interval width** [46]. Figure 9 illustrates the average width of k CIs of relative influence and rank. From Figure 9a and 9b, we find that the increase of privacy budget ρ_{Influ} and ρ_{Rank} shrinks the interval width of relative influence CI and rank CI separately. In particular, when $\rho_{Influ} = 0.5$, 6 out of 10 questions have the interval width of relative influence CI = 0.025; when $\rho_{Rank} = 1.0$, 2 questions have the interval width of rank CI = 2 and 6 questions have this number = 10. We also measure the **effect of confidence level** γ to the CI by changing γ from 0.1 to 0.9 by step size 0.1 and from 0.95 and 0.99. Figures can be found in the full version [2]. The results show that it has a non-significant effect to the interval width, as it changes < 0.03 for the influence CI of 6 questions, and changes < 5 for the rank CI of 8 questions.

Runtime analysis. We analyze the runtime of DPXPLAIN for generating Phase-2 and Phase-3. Figure 10a shows a runtime breakdown on average for all the questions from Table 1 with total runtime of 32 seconds on average. 88% of the time is used for the top-k explanation predicate selection procedure, especially on computing the influences for all the explanation predicates. The next highest runtime is for computing the confidence interval of influence, which

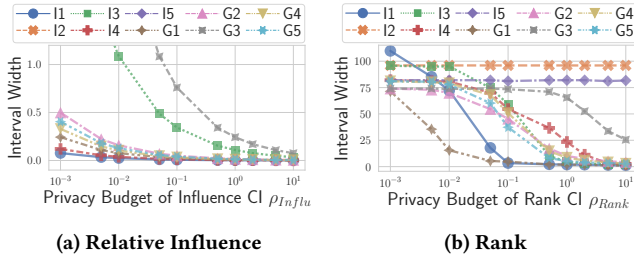


Figure 9: The width of confidence intervals by DPXPLAIN. The numbers are beyond 2 for the relative influence of I2 and I5.

needs to evaluate each sub queries. For the step noise quantification and confidence interval of rank, the time usage is not significant since the first only needs to find the image of two intervals and the second is a binary search. Figure 10b, shows that the runtime is linearly proportional to the size of explanations k , and the difference between questions is due to the difference of group sizes. We also find the runtime grows exponentially with the number of conjuncts l as the number of explanation predicates grows exponentially: for $l = 1, 2, 3$, the runtime about question I1 is 67, 3078 and 79634, and for question G1 it is 40, 1587 and 39922 seconds.

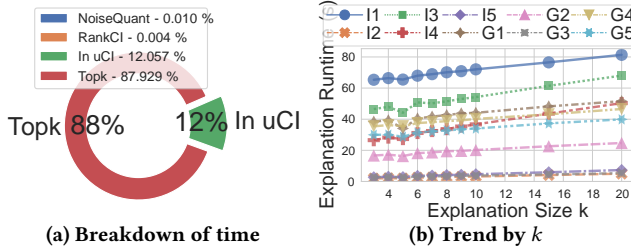


Figure 10: Runtime analysis of DPXPLAIN.

6 RELATED WORK

We next survey related work in the fields of DP and explanations for query results. *To the best of our knowledge, DPXPLAIN is the first work that explains aggregate query results while satisfying DP.*

Explanations for query results. The database community has proposed several approaches to explaining aggregate and non-aggregate queries in multiple previous works. Proposed approaches include provenance [17, 26, 27, 53, 54, 62, 63, 93], intervention [81, 82, 98], entropy [43], responsibility [73, 74], Shapley values [67, 78], counterbalance [75] and augmented provenance [65], and several of these approaches have used predicates on tuple values as explanations like DPXPLAIN, e.g., [43, 65, 82, 98]. We note that any approaches that consider individual tuples or explicit tuple sets in any form as explanations (e.g., [26, 63, 67, 73]) cannot be applied in the DP setting since they would violate privacy. Among the other summarization or predicate-based approaches, Scorpion [98] explains outliers in query results with the intervention of most influential predicates. Our influence function (Section 4.2) is inspired by the influence function of Scorpion, but has been modified to deliver accurate results while satisfying DP. Another intervention-based work [82] that also uses explanation predicates,

models inter-dependence among tuples from multiple relations with causal paths. DPXPLAIN does not support joins in the queries, which is a challenging future work (see Section 7).

Differential privacy. Private SQL query answering systems [32–34, 56, 59, 60, 69, 90, 97] consider a workload of aggregation queries with or without joins on a single or multi-relational database, but none supports explanation under differential privacy. The selection of private top-k candidates is well-studied by the community [8, 10, 11, 15, 18, 30, 37, 61, 66, 70, 71, 77, 91]. We adopt One-shot Top-k mechanism [77] since it is easy to understand. Private confidence interval is a new trend of estimating the uncertainty under differential privacy [12, 21, 45], however, the current bootstrap based methods measure the uncertainty from both the sampling process and the noise injection, while we only focus on the second part which is likely to give tighter intervals. The most relevant work to the private rank estimation is private quantile [4, 20, 39, 48, 57, 64, 86], which is to find the value given a position such as median, but the problem of rank estimation in our setting is reversed.

Privacy and provenance. As mentioned earlier, data provenance is often used for explaining query results, mainly for non-aggregate queries. Within the context of provenance privacy [6, 9, 19, 22, 23, 83, 85, 88], one line of work [22–24] studied the preservation of workflow privacy (privacy of data transferred in a workflow with multiple modules or functions), with a privacy criterion inspired by l -diversity [68]. A recent work [28] explored what can be inferred about the *query* from provenance-based explanations and found that the query can be reversed-engineered from the provenance in various semirings [49]. To account for this, a follow-up paper [25] proposed an approach for provenance obfuscation that is based on abstraction. This work uses k -anonymity [87] to measure how many ‘good’ queries can generate concrete provenance that can be mapped to the abstracted provenance, thus quantifying the privacy of the underlying query. Devising techniques for releasing provenance of non-aggregate and aggregate queries while satisfying DP is an interesting research direction.

7 FUTURE WORK

There are several interesting future directions. Extending DPXPLAIN to more general queries (like joins) and questions is an important future work. Unlike standard explanation frameworks like [98] where the join results can be materialized before running the explanation mechanism, a careful sensitivity analysis of adding/removing tuples from multiple tables is needed in the DP settings [90]. Second, the complexity of the top-k selection algorithm links to the number of explanation predicates that could be exponentially large, leaving room for future improvements. Additionally, other interesting notions of explanations for query answers (e.g., [65, 67, 75]) can be explored in the DP setting. Finally, evaluating our approach with a comprehensive user study and examining different metrics of understandability of the explanations generated by DPXPLAIN is also an important direction for future investigation.

ACKNOWLEDGMENTS

This work was supported by the NSF awards IIS-2147061, IIS-2016393, IIS-2008107, IIS-1703431, and IIS-1552538.

REFERENCES

- [1] Codebase of DPXPlain. <https://github.com/yuchaotao/Private-Explanation-System>.
- [2] DPXPlain: Explaining query results under differential privacy. <https://arxiv.org/abs/2209.01286>.
- [3] J. M. Abowd. The us census bureau adopts differential privacy. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2867–2867, 2018.
- [4] D. Alabi, O. Ben-Eliezer, and A. Chaturvedi. Bounded space differentially private quantiles. *CoRR*, abs/2201.03380, 2022.
- [5] Y. Amsterdamer, D. Deutch, and V. Tannen. Provenance for aggregate queries. In *PODS*, pages 153–164, 2011.
- [6] P. Anderson and J. Cheney. Toward provenance-based security for configuration languages. In U. A. Acar and T. J. Green, editors, *4th Workshop on the Theory and Practice of Provenance, TaPP*, 2017.
- [7] A. Ankan and A. Panda. pgmpy: Probabilistic graphical models using python. In *Proceedings of the 14th python in science conference (scipy 2015)*, pages 6–11. Citeseer, 2015.
- [8] M. Bafna and J. Ullman. The price of selection in differential privacy. In *Conference on Learning Theory*, pages 151–168. PMLR, 2017.
- [9] E. Bertino, G. Ghinita, M. Kantarcioglu, D. Nguyen, J. Park, R. S. Sandhu, S. Sultana, B. M. Thuraisingham, and S. Xu. A roadmap for privacy-enhanced secure data provenance. *J. Intell. Inf. Syst.*, 43(3):481–501, 2014.
- [10] R. Bhaskar, S. Laxman, A. Smith, and A. Thakurta. Discovering frequent patterns in sensitive data. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 503–512, 2010.
- [11] L. Bonomi and L. Xiong. Mining frequent patterns with differential privacy. *Proceedings of the VLDB Endowment*, 6(12):1422–1427, 2013.
- [12] T. Brawner and J. Honaker. Bootstrap inference and differential privacy: Standard errors for free. *Unpublished Manuscript*, 2018.
- [13] J. Brownlee. *XGBoost With python: Gradient boosted trees with XGBoost and scikit-learn*. Machine Learning Mastery, 2016.
- [14] M. Bun and T. Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pages 635–658. Springer, 2016.
- [15] R. S. Carvalho, K. Wang, L. Gondara, and C. Miao. Differentially private top-k selection via stability on unknown domain. In *Conference on Uncertainty in Artificial Intelligence*, pages 1109–1118. PMLR, 2020.
- [16] M. Cesar and R. Rogers. Bounding, concentrating, and truncating: Unifying privacy loss composition for data analytics. In *Algorithmic Learning Theory*, pages 421–457. PMLR, 2021.
- [17] A. Chapman and H. V. Jagadish. Why not? In *SIGMOD*, pages 523–534, 2009.
- [18] K. Chaudhuri, D. Hsu, and S. Song. The large margin mechanism for differentially private maximization. *arXiv preprint arXiv:1409.2177*, 2014.
- [19] J. Cheney. A formal framework for provenance security. In *CSF*, pages 281–293, 2011.
- [20] G. Cormode, T. Kulkarni, and D. Srivastava. Answering range queries under local differential privacy. *Proceedings of the VLDB Endowment*, 12(10):1126–1138, 2019.
- [21] C. Covington, X. He, J. Honaker, and G. Kamath. Unbiased statistical estimation and valid confidence intervals under differential privacy. *arXiv preprint arXiv:2110.14465*, 2021.
- [22] S. B. Davidson, S. Khanna, T. Milo, D. Panigrahi, and S. Roy. Provenance views for module privacy. In *PODS*, pages 175–186, 2011.
- [23] S. B. Davidson, S. Khanna, S. Roy, J. Stoyanovich, V. Tannen, and Y. Chen. On provenance and privacy. In *ICDT*, pages 3–10, 2011.
- [24] S. B. Davidson, S. Khanna, V. Tannen, S. Roy, Y. Chen, T. Milo, and J. Stoyanovich. Enabling privacy in provenance-aware workflow systems. In *CIDR*, pages 215–218, 2011.
- [25] D. Deutch, A. Frankenthal, A. Gilad, and Y. Moskovitch. On optimizing the trade-off between privacy and utility in data provenance. In *SIGMOD*, pages 379–391, 2021.
- [26] D. Deutch, N. Frost, and A. Gilad. Explaining natural language query results. *VLDB J.*, 29(1):485–508, 2020.
- [27] D. Deutch, N. Frost, A. Gilad, and T. Haimovich. Explaining missing query results in natural language. In *EDBT*, pages 427–430, 2020.
- [28] D. Deutch and A. Gilad. Reverse-engineering conjunctive queries from provenance examples. In *EDBT*, pages 277–288, 2019.
- [29] B. Ding, J. Kulkarni, and S. Yekhanin. Collecting telemetry data privately. *Advances in Neural Information Processing Systems*, 30, 2017.
- [30] Z. Ding, D. Kifer, T. Steinke, Y. Wang, Y. Xiao, D. Zhang, et al. The permute-and-flip mechanism is identical to report-noisy-max with exponential noise. *arXiv preprint arXiv:2105.07260*, 2021.
- [31] J. Dong, D. Durfee, and R. Rogers. Optimal differential privacy composition for exponential mechanisms. In *International Conference on Machine Learning*, pages 2597–2606. PMLR, 2020.
- [32] W. Dong, J. Fang, K. Yi, Y. Tao, and A. Machanavajjhala. R2t: Instance-optimal truncation for differentially private query evaluation with foreign keys. In *Proc. ACM SIGMOD International Conference on Management of Data*, 2022.
- [33] W. Dong and K. Yi. Residual sensitivity for differentially private multi-way joins. In *Proceedings of the 2021 International Conference on Management of Data*, SIGMOD '21, page 432–444, New York, NY, USA, 2021. Association for Computing Machinery.
- [34] W. Dong and K. Yi. A nearly instance-optimal differentially private mechanism for conjunctive queries. In *Proceedings of the 41st ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, PODS '22, page 213–225, New York, NY, USA, 2022. Association for Computing Machinery.
- [35] D. Dua and C. Graff. UCI machine learning repository, 2017.
- [36] D. Durfee and R. Rogers. One-shot dp top-k mechanisms. *Differential Privacy.org*, 08 2021. <https://differentialprivacy.org/one-shot-top-k/>.
- [37] D. Durfee and R. M. Rogers. Practical differentially private top-k selection with pay-what-you-get composition. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 3527–3537, 2019.
- [38] C. Dwork. Differential privacy and the us census. In *Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, pages 1–1, 2019.
- [39] C. Dwork and J. Lei. Differential privacy and robust statistics. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pages 371–380, 2009.
- [40] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [41] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- [42] C. Dwork and G. N. Rothblum. Concentrated differential privacy. *arXiv preprint arXiv:1603.01887*, 2016.
- [43] K. El Gebaly, P. Agrawal, L. Golab, F. Korn, and D. Srivastava. Interpretable and informative explanations of outcomes. *Proc. VLDB Endow.*, 8(1):61–72, sep 2014.
- [44] Ú. Erlingsson, V. Pihur, and A. Korolova. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security*, pages 1054–1067, 2014.
- [45] C. Ferrando, S. Wang, and D. Sheldon. General-purpose differentially-private confidence intervals. *arXiv preprint arXiv:2006.07749*, 2020.
- [46] C. Ferrando, S. Wang, and D. Sheldon. Parametric bootstrap for differentially private confidence intervals, 2021.
- [47] S. Flood, M. King, R. Rodgers, S. Ruggles, J. R. Warren, and M. Westberry. Integrated public use microdata series, current population survey: Version 9.0 [dataset]. *Minneapolis, MN: IPUMS*, 2021. <https://doi.org/10.18128/D030.V9.0>.
- [48] J. Gillenwater, M. Joseph, and A. Kulesza. Differentially private quantiles. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 3713–3722. PMLR, 2021.
- [49] T. J. Green, G. Karvounarakis, and V. Tannen. Provenance semirings. In *PODS*, pages 31–40, 2007.
- [50] U. Grömping. South german credit data: Correcting a widely used data set. *Rep. Math., Phys. Chem., Berlin, Germany, Tech. Rep.*, 4:2019, 2019.
- [51] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, Sept. 2020.
- [52] J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. T. Riedl. Evaluating collaborative filtering recommender systems. *ACM Transactions on Information Systems (TOIS)*, 22(1):5–53, 2004.
- [53] M. Herschel and M. A. Hernández. Explaining missing answers to SPJUA queries. *PVLDB*, 3(1):185–196, 2010.
- [54] J. Huang, T. Chen, A. Doan, and J. F. Naughton. On the provenance of non-answers to queries over extracted data. *PVLDB*, 1(1):736–747, 2008.
- [55] B. Jiang, X. Zhang, and T. Cai. Estimating the confidence interval for prediction errors of support vector machine classifiers. *J. Mach. Learn. Res.*, 9:521–540, 2008.
- [56] N. Johnson, J. P. Near, and D. Song. Towards practical differential privacy for sql queries. *Proceedings of the VLDB Endowment*, 11(5):526–539, 2018.
- [57] H. Kaplan, S. Schnapp, and U. Stemmer. Differentially private approximate quantiles. *CoRR*, abs/2110.05429, 2021.
- [58] C. T. Kenny, S. Kuriwaki, C. McCartan, E. T. Rosenman, T. Simko, and K. Imai. The use of differential privacy for census data and its impact on redistricting: The case of the 2020 us census. *Science advances*, 7(11):eabk3283, 2021.
- [59] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajjhala, M. Hay, and G. Miklau. Privatesql: a differentially private sql query engine. *Proceedings of the VLDB Endowment*, 12(11):1371–1384, 2019.

- [60] I. Kotsogiannis, Y. Tao, A. Machanavajjhala, G. Miklau, and M. Hay. Architecting a differentially private sql engine. In *CIDR*, 2019.
- [61] J. Lee and C. W. Clifton. Top-k frequent itemsets via differentially private fp-trees. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 931–940, 2014.
- [62] S. Lee, S. Köhler, B. Ludäscher, and B. Glavic. A sql-middleware unifying why and why-not provenance for first-order queries. In *ICDE*, pages 485–496, 2017.
- [63] S. Lee, B. Ludäscher, and B. Glavic. PUG: a framework and practical implementation for why and why-not provenance. *VLDB J.*, 28(1):47–71, 2019.
- [64] J. Lei. Differentially private m-estimators. *Advances in Neural Information Processing Systems*, 24, 2011.
- [65] C. Li, Z. Miao, Q. Zeng, B. Glavic, and S. Roy. Putting things into context: Rich explanations for query answers using join graphs. In *SIGMOD*, pages 1051–1063, 2021.
- [66] N. Li, W. H. Qardaji, D. Su, and J. Cao. Priv’basis: Frequent itemset mining with differential privacy. *Proc. VLDB Endow.*, 5(11):1340–1351, 2012.
- [67] E. Livshits, L. E. Bertossi, B. Kimelfeld, and M. Sebag. The shapley value of tuples in query answering. In *ICDT*, volume 155, pages 20:1–20:19, 2020.
- [68] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkitasubramaniam. L-diversity: Privacy beyond k -anonymity. *TKDD*, 1(1):3, 2007.
- [69] R. McKenna, G. Miklau, M. Hay, and A. Machanavajjhala. Optimizing error of high-dimensional statistical queries under differential privacy. *Proc. VLDB Endow.*, 11(10):1206–1219, 2018.
- [70] R. McKenna and D. R. Sheldon. Permute-and-flip: A new mechanism for differentially private selection. *Advances in Neural Information Processing Systems*, 33:193–203, 2020.
- [71] F. McSherry and K. Talwar. Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS’07)*, pages 94–103. IEEE, 2007.
- [72] F. D. McSherry. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, pages 19–30, 2009.
- [73] A. Meliou, W. Gatterbauer, K. F. Moore, and D. Suciu. The complexity of causality and responsibility for query answers and non-answers. *Proc. VLDB Endow.*, 4(1):34–45, 2010.
- [74] A. Meliou, W. Gatterbauer, S. Nath, and D. Suciu. Tracing data errors with view-conditioned causality. In T. K. Sellis, R. J. Miller, A. Kementsietsidis, and Y. Velegarakis, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2011, Athens, Greece, June 12-16, 2011*, pages 505–516. ACM, 2011.
- [75] Z. Miao, Q. Zeng, B. Glavic, and S. Roy. Going beyond provenance: Explaining query answers with pattern-based counterbalances. In *SIGMOD*, pages 485–502, 2019.
- [76] I. Mironov. Rényi differential privacy. In *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*, pages 263–275. IEEE, 2017.
- [77] G. Qiao, W. J. Su, and L. Zhang. Oneshot differentially private top-k selection. *arXiv preprint arXiv:2105.08233*, 2021.
- [78] A. Reshef, B. Kimelfeld, and E. Livshits. The impact of negation on the complexity of the shapley value in conjunctive queries. In D. Suciu, Y. Tao, and Z. Wei, editors, *PODS*, pages 285–297, 2020.
- [79] R. Rogers and T. Steinke. A better privacy analysis of the exponential mechanism. DifferentialPrivacy.org, 07 2021. <https://differentialprivacy.org/exponential-mechanism-bounded-range/>.
- [80] L. Rosenblatt, X. Liu, S. Pouyanfar, E. de Leon, A. Desai, and J. Allen. Differentially private synthetic data: Applied evaluations and enhancements. *arXiv preprint arXiv:2011.05537*, 2020.
- [81] S. Roy, L. J. Orr, and D. Suciu. Explaining query answers with explanation-ready databases. *Proc. VLDB Endow.*, 9(4):348–359, 2015.
- [82] S. Roy and D. Suciu. A formal approach to finding explanations for database queries. In C. E. Dyreson, F. Li, and M. T. Özsu, editors, *SIGMOD*, pages 1579–1590, 2014.
- [83] P. Ruan, G. Chen, A. Dinh, Q. Lin, B. C. Ooi, and M. Zhang. Fine-grained, secure and efficient data provenance for blockchain. *Proc. VLDB Endow.*, 12(9):975–988, 2019.
- [84] S. Ruggles, C. Fitch, D. Magnuson, and J. Schroeder. Differential privacy and census data: Implications for social and economic research. In *AEA papers and proceedings*, volume 109, pages 403–08, 2019.
- [85] J. L. C. Sanchez, J. B. Bernabé, and A. F. Skarmeta. Towards privacy preserving data provenance for the internet of things. In *WF-IoT*, pages 41–46, 2018.
- [86] A. Smith. Privacy-preserving statistical estimation with optimal convergence rates. In *Proceedings of the forty-third annual ACM symposium on Theory of computing*, pages 813–822, 2011.
- [87] L. Sweeney. K-anonymity: A model for protecting privacy. *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*, 10(5):557–570, 2002.
- [88] Y. S. Tan, R. K. L. Ko, and G. Holmes. Security and data accountability in distributed systems: A provenance survey. In *HPCC/EUC*, pages 1571–1578, 2013.
- [89] J. Tang, A. Korolova, X. Bai, X. Wang, and X. Wang. Privacy loss in apple’s implementation of differential privacy on macos 10.12. *arXiv preprint arXiv:1709.02753*, 2017.
- [90] Y. Tao, X. He, A. Machanavajjhala, and S. Roy. Computing local sensitivities of counting queries with joins. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pages 479–494, 2020.
- [91] A. G. Thakurta and A. Smith. Differentially private feature selection via stability arguments, and the robustness of the lasso. In *Conference on Learning Theory*, pages 819–850. PMLR, 2013.
- [92] The pandas development team. pandas-dev/pandas: Pandas, Feb. 2020.
- [93] Q. T. Tran and C.-Y. Chan. How to conquer why-not questions. In *SIGMOD*, pages 15–26, 2010.
- [94] G. Vandenbroucke. Married men sit atop the wage ladder. 24, 2018.
- [95] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020.
- [96] L. Wasserman. *All of statistics: a concise course in statistical inference*, volume 26. Springer, 2004.
- [97] R. J. Wilson, C. Y. Zhang, W. Lam, D. Desfontaines, D. Simmons-Marengo, and B. Gipson. Differentially private sql with bounded user contribution. *arXiv preprint arXiv:1909.01917*, 2019.
- [98] E. Wu and S. Madden. Scorpion: Explaining away outliers in aggregate queries. *Proc. VLDB Endow.*, 6(8):553–564, 2013.
- [99] Z. Yan, G. Li, and J. Liu. Private rank aggregation under local differential privacy. *International Journal of Intelligent Systems*, 35(10):1492–1519, 2020.
- [100] J. Zhang, G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao. Privbayes: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)*, 42(4):1–41, 2017.