

Umбра: A Visual Analysis Approach for Defense Construction Against Inference Attacks on Sensitive Information

Xumeng Wang, Chris Bryan, Yiran Li, Rusheng Pan, Yanling Liu, Wei Chen, and Kwan-Liu Ma

Abstract—Collecting and analyzing anonymous personal information is required as a part of data analysis processes, such as medical diagnosis and restaurant recommendation. Such data should ostensibly be stored so that specific individual information cannot be disclosed. Unfortunately, inference attacks—integrating background knowledge and intelligent models—hinder classic sanitization techniques like syntactic anonymity and differential privacy from exhaustively protecting sensitive information. As a solution, we introduce a three-stage approach empowered within a visual interface, which depicts underlying inference behaviors via a Bayesian Network and supports a customized defense against inference attacks from unknown adversaries. In particular, our approach visually explains the process details of the underlying privacy preserving models, allowing users to verify if the results sufficiently satisfy the requirements of privacy preservation. We demonstrate the effectiveness of our approach through two case studies and expert reviews.

Index Terms—Privacy; inference attack; bayesian network; visual analytics.

1 INTRODUCTION

MORE than a century has passed since the first publication on the right to privacy [1]. But in today's digital age, the notion of a right to privacy has perhaps never been more prominent. There are myriad examples where private information has been disclosed (either inadvertently or maliciously), leading to personal repercussions including damaged reputation, fiduciary loss, and threats to personal safety. In May 2018, the European Union released the General Data Protection Regulation (GDPR) [2], [3]. This law requires **data owners**—persons who collect, hold, or process data—to put in place effective measures that guarantee privacy preservation for **data subjects**—persons whose data is collected and used.

Privacy preservation is necessary across many domains, including medicine [4] and behavioral research [5]. Unfortunately, achieving effective privacy protection remains a non-trivial task. A significant challenge is the evolution of adversarial techniques, whereby protection is not achievable simply by removing sensitive information or attributes from a data corpus. **Inference attacks** deduce sensitive information with high confidence by pairing a base dataset with outside information: news, third-party reports, additional datasets, etc. [6]. As an example, if an insurer charges a higher fee to people who are more prone to illness, an attacker could use monthly premiums to infer a person's state of health—high premiums corresponds to an unhealthy insured.

Inference attacks result in privacy violations when non-sensitive information is taken as a premise and sensitive information is derived. To complicate the problem, infer-

ences can change by varying background knowledge, making tracking and planning for such attacks akin to hitting a moving target. A key takeaway is that privacy cannot be effectively preserved without considering an attacker's potential background knowledge.

When mounting a defense against inference attacks, data owners must not only ask, “*how do I anonymize sensitive information in my dataset?*” but also, “*how do I simulate various inference attacks to verify data subjects are protected?*” To provide trust and control, data owners need to be able to understand the privacy preservation process—“*how can the process be intuitively explained?*”—and customize how it is applied: “*how can the user help direct the application of privacy preservation models to the dataset?*”

To address these issues, we develop a novel visual analytics approach based around interactively creating and manipulating Bayesian networks. Bayesian networks [7] are commonly used to model dependencies among variables, and have been shown effective for inference analysis across diverse application domains (e.g. [8], [9]).

Our approach is organized as a three-stage workflow: (1) inference initialization, (2) data sanitization, and (3) result verification. In the first stage, inferences that lead to privacy leaks are identified and visualized. In the second stage, targeted sanitization operations are customized to obfuscate data records and attributes, thereby reducing the confidence of potential adversarial inference attacks. To verify a sufficient level of protection, the third stage supports the interactive building and testing of models to simulate adversarial attacks on the sanitized data, including ones that make use of outside datasets. By providing feedback at each stage, we provide transparency into the privacy preservation process by showing *what* data is causing privacy issues, *why* sanitization operations are necessary, and *how* recommended solutions ensure dataset records are protected. The result is a powerful, highly customizable, and verifiable defense

- X. Wang, R. Pan, Y. Liu and W. Chen are with State Key Lab of CAD&CG, Zhejiang University. Wei Chen is the corresponding author. E-mail: {wangxumeng, panrusheng, 3150104312}@zju.edu.cn and chenwei@cad.zju.edu.cn.
- Y. Li and K.-L. Ma are with the University of California, Davis. E-mail: ranli@ucdavis.edu and ma@cs.ucdavis.edu.
- C. Bryan is with Arizona State University. E-mail: cbryan16@asu.edu.

solution against inference attacks.

To evaluate our workflow, we implement it as an interactive software system named *Umbra*. We present a set of case studies demonstrating how privacy risks can be defended. We also interviewed three domain experts in computer security to whom we demonstrated the system, and they provided feedback on its strengths and shortcomings.

2 RELATED WORK

This work is premised on the use of **inference**—drawing conclusions based on background knowledge and reasoning [6]—as a way to facilitate adversarial attacks to extract private information. To defend against such attacks, protection schemes modify the dataset in targeted ways. Such schemes must be tactfully chosen, as **sanitization**—obfuscation of data records and attributes—reduces the dataset’s **utility**—the ability of data to enable tasks such as analysis and decision-making [10]. Interactive visualization can effectively address this tactful scheme application by enabling human-in-the-loop sensemaking.

2.1 Inference Attacks

With the help of background knowledge such as association rules [11], adversaries can potentially access and/or reasonably assume sensitive information via inference attacks, even when the dataset is encrypted or anonymized. Although some databases employ precise query protocols to individually encrypt records, information can be recovered according to patterns given by range queries [12], [13].

Classic models for privacy preservation, such as syntactic anonymity and differential privacy, are vulnerable to inference attacks [14]. To construct adequate protection via anonymization approaches, a data owner must know both the background knowledge that adversaries will have access to and what attack behaviors they will employ [15]. Taking adversaries’ background knowledge into consideration, Sun et al. [16] constructed a privacy inference graph to describe potential privacy disclosures for k -anonymity. However, this is challenging, especially when predicting attacks from unknown adversaries, as they can access desired information simply by tracing paths between nodes in the graph [16]. As for differential privacy models, their independence assumption also provides a vulnerability for inference attacks [17], [18], since correlations can be observed in a majority of datasets. To break such limit, a feasible approach is applying Bayesian models [19].

In the past decade, Bayesian networks—which extract probabilistic relationships among variables as a graph model [7]—have emerged as a technique for inference identification. To adapt this idea to application scenarios, researchers tailor the graph model for particular perspectives. For example, Zhang and Song [20] proposed a graph-based posterior inference model in accordance with likelihood weighting. In considering various distinct possibilities in event sequences, their model includes nodes, corresponding to the resource information occupied by attackers, and edges, indicating underlying attacks. To reduce the learning cost of Bayesian Networks, Timmer et al. [21] designed an intuitive representation by extracting inference rules and related strengths. Based on inferential strength, the algorithms

presented by Timmer et al. can detect the undercutters against the inferences [21]. To our knowledge, our approach is one of the first to adopt Bayesian networks in visual analytics workflow for privacy preservation.

2.2 Defense against Inference Attacks

As one approach to personal identity privacy, Li et al. [22] proposed a flexible algorithm that combines two syntactic anonymity models, k -anonymity and l -diversity, by employing k as an individual number and l as diversity level. Unfortunately, both models provide insufficient protection against attacks that speculate on attribute distributions, e.g. skewness attacks [23]. Minimax filter [24] is a learning-based approach that uses independent assumptions between training data and test data. This method provides task-dependent protection by dimensionality reduction of raw features and can be extended to include noise by leveraging differential privacy approaches [24].

Utility is an important consideration in the sanitization process, because sharing sanitized data becomes pointless if all useful information is removed. Salamatian et al. [25] applied probabilistic privacy mapping to randomly perturb data values; they consider utility as one constraint in a convex optimization problem. Similarly, Chen et al. [26] balance privacy and utility as a knapsack problem using data values as weighted items. The item weight is quantified as the risk of privacy disclosure—calculated by a Naïve Bayes model—and item value is based on uniqueness and commonness. Given privacy-dependent attributes and utility-dependent attributes, Cai et al. [27] modify shared attributes by constructing a generic attribute hierarchy and removing the remaining privacy-dependent attributes.

However, the aforementioned strategies ignore the question of *how to simulate attacks*. Without specific simulation capabilities that account for the background knowledge of potential adversaries, they have limited adaptability. There are also other considerations, including the ability to enable task-oriented utility settings and designing for non-professional users (i.e., the question about *intuitively explaining the process*). We explicitly consider these issues in our design requirements and workflow (see Secs. 3 and 4).

2.3 Visualization for Privacy Preservation

Privacy preservation must be considered when using sensitive data in any manner, including visual analysis. Unfortunately, most prior work in this area only leverage classical models (see Sec. 2.1), meaning that it cannot provide guaranteed defend against inference attacks.

For example, k -anonymity is used to create visual clustering for parallel coordinates [28]. The identities of individual data points (data records) are unrecognizable due to the uncertainty caused by visual perception. Enabling custom parameters for visual clustering, Dasgupta et al. [29] extended their privacy preserving design to other charts types such as scatter plots. Subsequently, specific visual designs that can preserve privacy for various data types have been proposed, including event sequences [30], social networks [31], and trajectory data [32]. To quantify the ability of visualization techniques at obfuscating privacy leaks, a probabilistic model [33] is proposed to explain the

effects of attackers' background knowledge and the uncertainty in cluster-based charts. Empirical studies are also implemented to assess other visual expressions, including for alluvial diagrams [30], node-link graphs [34], and geo-based trajectories [35].

Interactive visualization has also been used to augment the privacy preservation process. ODD Visualizer [36] uses matrices to explain the re-identification risks that are caused by grouping attributes. The PER-Tree [37] measures potential privacy exposure risks via multiple syntactic anonymity models, allowing users to evaluate and remove the source of the risk. For social networks, GraphProtector [38] provides a transparent and intuitive privacy preserving pipeline by recommending sanitization solutions and providing a historical view of previous dataset updates, while allowing custom configurations that account for privacy and utility trade-offs. While these systems are effective in resisting subsets of attacks, as compared to the workflow introduced in this paper, none of them provides sufficient protection if inference attacks are employed.

3 DOMAIN AND TASK ANALYSIS

Our goal is to provide both guaranteed protection against inference attacks while simultaneously giving visibility and human-in-the-loop control into the backend data sanitization methods being used. As a domain abstraction, our target users are data owners who want to ensure that collected information is protected. As our approach is based around Bayesian modeling of inferences to identify and resolve privacy leaks, such users must also have an understanding of privacy protection and Bayesian network to reason the quantified risks [39], [40].

To motivate the specific tasks that our system should support, we conducted a task abstraction by discussing the topic of *interactively ensuring privacy* with an established researcher in computer security (a professor with 20+ years experience). Interactive dataset sanitization is not common in privacy preservation, so this discussion focused on identifying areas where standard automated approaches can be improved by adding interactivity and user control. For example, automated approaches are primarily black boxes, providing little transparency or explainability both to data owners and data subjects. Automated approaches provide little customization of defenses or details on demand about dataset specifics, which limits the ability of data owners to defend against distinct or novel attack behaviors. Based on this discussion, we derived seven design requirements (DRs) to support interactive and customized dataset protection against inference attacks.

DR1: Extract underlying inferences. Because future attack behaviors are unknown, any combination of underlying inferences can be considered as background knowledge that an adversary might leverage. To defend attacks, dependency relationships between all attribute value pairs must be understood [21]. *Therefore, underlying inferences should be extracted to facilitate subsequent summary and analysis.*

DR2: Allow custom inferences. Removing all potential inferences that relate to sensitive information will ostensibly provide comprehensive protection, but at the cost of severely reducing the utility of the dataset. Depending on

the context, a data owner might allow some inferences to remain, even if they technically lead to exposure risk (e.g., the inference is beyond an adversary's knowledge). New inferences may also be created by introducing additional datasets [41] or outside information [42]. *To ensure a full understanding of adversarial background knowledge, data owners need to specify new inferences and modify existing ones.*

DR3: Seek defenses against identified inference attacks. Privacy leaks should be identifiable and resolvable according to the inferences extracted from the dataset (DR1) and those customized by users (DR2). *To ensure comprehensive defense against attacks (which can be quite complex), computational-based approaches are necessary to guarantee the results of privacy preservation processes.*

DR4: Recommend sanitization solutions. Individual data records may be sanitized in multiple ways. For example, "a person whose birthday is July 1st" might be applied to one person (one data record) in a dataset, while vaguer descriptions ("a person born in July," "a person born on the 1st") might be applied to dozens. Removing either the month or day information preserves that person's anonymity, but which option is best? Automatic evaluation approaches, like using entropy-based indicators, can mislead in a majority of cases. For example, processed datasets with high indicator values may benefit little from practical applications such as machine learning. *To assist decision-making and enable human-in-the-loop sanitization, potential sanitization actions should be automatically evaluated, compared, and ranked as recommendations, while still allowing users to opt for a desired solution.*

DR5: Simulate outside attacks. To comprehensively verify privacy preserving levels, a common practice in related research [43] is to simulate attacks. Privacy exposure risks can be identified from the attack results. *For the sanitized dataset, new attack models can be built and applied; accuracy metrics can judge if the privacy preservation levels remain acceptable.*

DR6: Show data details. Due to a lack of data descriptions and result expressions, automatic methods face difficulties in parameter settings, model selection, result verification, etc. Besides, a data owner needs to verify a dataset's overall status before publishing it, potentially including low-level analysis of specific data subjects. *Therefore, details about individual records and the entire datasets should be reviewable in the sanitization process to see how they have been modified and sanitized.*

DR7: Explain the sanitization process. Mechanism explanation can help data owners judge if the sanitization model is reliable and trustworthy [44]. *Therefore, the entire sanitization process supported by automatic models (DR3) should be explained in an intuitive, user-friendly way [37], [38].*

4 MODELING AND DEFENDING INFERENCE

To support design requirements DR1–DR7, we adopt a three-stage workflow. Before describing the workflow in detail in Sec. 5, we first introduce several necessary background concepts and algorithms.

4.1 An Illustrative Dataset

To provide context and real-world applicability when describing our approach throughout Sections 4–6, we use a

dataset on post-student development published by McVicar and Anyadike-Dane [45]. This dataset is about the social statuses of persons in Northern Ireland. McVicar and Anyadike-Dane concluded that specific background characteristics—coming from disadvantaged areas, having an unemployed father, etc.—correlated to experiencing failure after graduation from educational institutions [45]. An adversary could therefore infer that an individual experienced failure if some or all background conditions were met, and odds are high this speculation would be correct.

Table 1 shows the attribute properties of this dataset. It contains 8 attributes about 712 individuals. For our purposes, the employment status (the `employ` attribute) is considered as sensitive.

TABLE 1

The post-student development dataset contains information about the social and economic status of 712 persons in Northern Ireland. The highlighted attribute `employ` is considered sensitive.

Attribute	Data Type	Description
gender	Categorical	The individual's gender.
residence	Categorical	Where the individual lives.
employ	Numerical	Total number of months employed over the prior 6 years.
grade	Categorical	Did the individual have five or more academic qualifications at grades A-C? (binary)
school	Categorical	The type of the individual's school.
cat	Categorical	Is the individual a Catholic?
fue	Categorical	Is the individual's father unemployed?
fmp	Categorical	Is the individual's father employed in a managerial position?

4.2 Key Terms for Bayesian-based Inference

There are a variety of approaches to inferring information (e.g., [25], [46]), among which, Bayesian-based approaches are representative and widely-applied [21], [26], [47]. Before describing the details, we define several key terms.

In a dataset, **inferences** are conditional probabilities among **states**, which are colloquial characteristics defined by single attribute values. As an example, a person (a single data point or **record**) in the post-student dataset could belong to a state called “success,” based on meeting some value threshold within the `employ` attribute, and “failure” if not. This threshold, called a **split point**, divides the numerical range of the `employ` attribute into categorical bins, one for each state. If the split point was set to 12, then the states $S_{success}$ (i.e., `employ`: [12~72]) and $S_{failure}$ (i.e. `employ`: [0~12]) would respectively denote success and failure. Alternatively, categorical attributes provide intrinsic state definitions: e.g., “lives in the north” can be defined by specific values for the `residence` attribute.

In the Umbra system (described in Sec. 6), we automatically identify split points for numerical attributes by parsing attribute distributions and finding sharp declines. If no such a decline is present, the medians are selected to equally divide the records. This prevents excessively unbalanced sample sizes, which hinders accurate inference generation.

States can be either sensitive or non-sensitive. **Sensitive states** are those that a data owner wishes to keep private, such as $S_{success}$ and $S_{failure}$. Other states, including `residence` ($S_{Northern}$, $S_{Southern}$, etc.) and `gender` (S_{Female} , S_{Male}), are considered **non-sensitive**, and their information

may be publicly known or published. We refer to a **group** as a set of records having the same set of non-sensitive states. To preserve privacy of the records in a group that has been identified as “at-risk” for exposure, we identify a set of candidate states for removal from the inference graph (see Sec. 4.4). This solution (the set of states to be removed to ensure the group’s privacy) is regarded as a **scheme**. When a scheme is applied, a percentage of state **occurrences**—records which satisfy a certain collection of states defined by a combination of one or more attributes—are removed, meaning that the state-related attribute values of the records are “blacked out” or marked as “unknown.”

4.3 Extracting Underlying Inferences as a Graph

The first design requirement (**DR1**) is to extract the underlying inferences in a dataset. The first step towards inference extraction is quantifying the probabilities of state occurrences. We first count state occurrences. Then, we calculate conditional probabilities as:

$$\Pr(S_0|S_1, \dots, S_n) = \Pr(S_0, S_1, \dots, S_n) / \Pr(S_1, \dots, S_n) \quad (1)$$

where S_i ($i = 0, \dots, n$) refers to different states.

An adversary may be aware of correlations between sensitive and non-sensitive states, thereby inferring sensitive states based on public information about non-sensitive states. Instead of Naïve Bayesian model [26], we use Bayesian network to learn state-to-state correlations and detect the underlying inferences thoroughly. Because individual records can only be binned into one state defined by an attribute (e.g., a student is either “pass” or “fail,” but cannot be both), we do not consider correlations between states defined by the same attributes.

To describe these correlations, inferences can be visualized as a network graph (see Fig. 2(c) and Fig. 10 for examples of inference graphs). States are represented as nodes and directed edges weight the “effect amount” from the source state to the target. We quantify these effects as $|\Pr(S_{target}|S_{source}) - \Pr(S_{target})|$, where S_{target} refers to the target state and S_{source} refers to the source state. If two nodes are connected through a path, there is an indirect correlation between the two states.

The difference between the conditional probability given a source state and the probability of a target state indicates the ability of source states to infer a target state. To allow for custom inferences (**DR2**), new edges can be added to the inference graph. Note that inaccurate inference simulation may lead to invalid privacy preservation. Thus, the dataset for inference graph construction and subsequent editing operations must be reliable.

4.4 Constructing Defenses Against Inference Attacks

To enable comprehensive defenses (**DR3**), we first learn the risks posed by underlying inferences to potential (i.e. simulated) attacks. The conditional probabilities of a sensitive state, given any set of non-sensitive states, are assumed known by adversaries. That is, an adversary may know the exact values of $\Pr(S_{success}|S_{Female}, S_{Southern})$, $\Pr(S_{success}|S_{Male}, S_{Northern})$, etc. Such an inference is used to assert if $S_{success}$ is true, violating privacy. Assume that the conditional probability using `residence` is much

higher than gender. If the published version of the dataset removes the residence information (so it is not known where the students come from), an adversary can only employ the conditional probabilities $\Pr(S_{success}|S_{Female})$ and $\Pr(S_{success}|S_{Male})$, which have much lower confidence.

An extreme solution for privacy preservation is to remove all non-sensitive states from the dataset (this is mentioned in (DR2)). By removing conditional information, an adversary cannot use conditional probabilities as a basis for inference attacks, and can only rely on less confident probabilities (i.e. $\Pr(S_{success})$ for all groups). If $\Pr(S_{success}|S_{Female})$ is close to $\Pr(S_{success})$, the gender information cannot be used for inference attacks and can therefore be included in the published dataset, which will indicate a higher overall utility. We describe a group of records as “privacy preserved” and ready to be shared when the conditional probability of each sensitive state lies within a “no-risk” range. The no-risk range is $\Pr(S_{sensitive}) \pm \delta$, where $\Pr(S_{sensitive})$ is calculated without any hint from public information, and δ is a user-defined parameter representing a privacy exposure risk limit. Such non-sensitive states may be safely published as they provide negligible help in inference attacks.

To assess the impact of a potential inference attack, we first verify if each sensitive group lies within the no-risk range. To resolve groups that are at-risk, one or more nodes in the inference graph must be deleted along paths to the sensitive state. In this way, inference attacks are defended by removing correlated states.

To find a candidate set of states for removal, we first look at adjacent neighbor nodes to the sensitive states, starting with the node with the highest confidence in the inference (has the largest edge weight in the inference graph). We iteratively delete these states until the at-risk group’s privacy is preserved, record the deleted states as a scheme, and then restart and attempt other state combinations for deletion. If removing all states in a state set C_S fixes a privacy leak, we will not test any set that includes C_S , because there is no need to remove extra information (causing unnecessary utility loss). The set of identified schemes are candidate solutions for transforming the exposed group into a no-risk range. By exhaustively traversing all state combinations, all applicable schemes can be identified. Unfortunately, the time complexity of this operation is $O(n_S^{n_A})$, where n_A and n_S are, respectively, the number of attributes and the number of states defined by an attribute. (Keep in mind that multiple states from the same attribute cannot be in the same set, since a record can only have 1 state for each attribute.) We discuss system scalability in Sec. 8.2.

4.5 Recommending Solution Schemes

Multiple schemes might be applicable for a single identified privacy leak, therefore recommending a solution supports decision-making for users (DR4). For this process, it is necessary to know the impact that each scheme has on the resultant dataset. The hope is that sanitizing a dataset not only protects it from inference attacks, but it does so in a way that maintains high overall utility.

One consideration is that certain states may be more important to post hoc tasks, and their removal as part

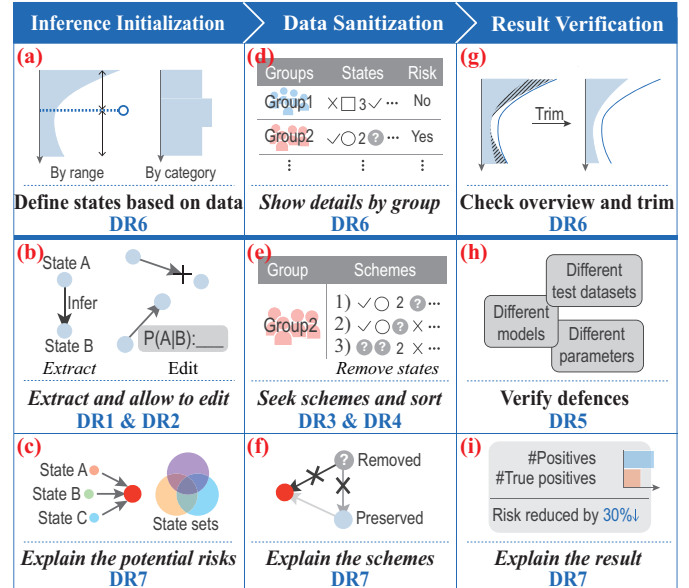


Fig. 1. The three-stage workflow for constructing defenses against inference attacks. Items in the top row allow users to seek appropriate sanitization based on dataset characteristics; the second and third rows allow users to review and refine the automatic models. Italicized text indicates actions that are automated by Umbra.

of a scheme might affect resultant attribute distributions. These states should be preserved as much as possible. As an example, assume a state S_1 has 2 occurrences and state S_2 has 1000. Removing one occurrence has a negligible impact on S_2 , but removes half the samples of S_1 .

We quantify the specific utility value of a non-sensitive state in two ways. Suppose that S is a non-sensitive state that happens with a probability of $\Pr(S)$. By calculating entropy based on information theory, S ’s utility value is:

$$u_E = -v_{Attr} \times \log(\Pr(S)) \quad (2)$$

where v_{Attr} ($0 < v_{Attr} \leq 1$) is a user-defined utility value of $Attr$, the attribute used to define S . However, many data owners are non-professional users (who may not be familiar with logarithms). Thus we provide a second utility metric (which is the default choice in Umbra):

$$u_E = v_{Attr} \times (1 - \Pr(S)) \quad (3)$$

Equations (2) and (3) compute states with fewer samples as having higher utility levels. We use this utility value to sort applicable schemes, ranking solutions by decreasing order of total utility loss (the sum of utility values of the states to be removed). The scheme with the lowest total utility loss for a group is the recommended choice.

5 WORKFLOW

Using the techniques described in Sec. 4, privacy preservation is achieved through a three-stage workflow. Fig. 1 shows this workflow, and labels the specific operations that support the design requirements outlined in Sec. 3. In this section, we briefly describe the workflow at a high level, according to operations shown in Figs. 1(a)–(i). Sec. 6 describes how each stage is implemented within Umbra.

5.1 Inference Initialization (Stage 1)

The Inference Initialization Stage is for defining states and constructing inferences. A dataset is loaded and its attributes are placed into two disjoint sets: those that can be shared and those considered sensitive. If additional datasets are desired (to simulate adversarial background knowledge), they can also be loaded.

Fig. 1(a): Based on the loaded data, users create states and learn about their correlations. State definitions may be customized according to user preferences (e.g., `age > 18` may signify a record is an “adult,” or alternatively “allowed to rent a vehicle”). Users can create, edit, and delete states by binning numerical ranges, adjusting split points, and merging categorical values, and can review attribute distributions and adjust the utility values (Equations 2 and 3) for non-sensitive attributes.

Fig. 1(b): Based on the defined set of states, underlying inferences are extracted (DR1) and the base inference graph is constructed. Custom inferences (DR2) are supported by manually editing the conditional probabilities.

Fig. 1(c): To understand exposure risks (DR4), inferences can be further analyzed. Because privacy exposure can be caused by a combination of multiple states (see Sec. 4.4), the conditional probabilities given in a state set can be compared to the probability of sensitive states.

5.2 Data Sanitization (Stage 2)

Based on the constructed inference graph, the Data Sanitization Stage recommends sanitization operations (schemes) on groups identified as being at-risk for exposure.

Fig. 1(e): For groups that are at-risk, applicable schemes are identified (DR3) and sorted (DR5) based on overall utility loss. When multiple schemes are available, the recommended solution can be compared to alternate schemes by reviewing the differences in effect on dataset utility.

Fig. 1(f): To understand how a scheme sanitizes an exposed group (DR4), the application of the scheme on the inference graph is previewed. Blocked by the removal of states (nodes), related inferences will either be invalidated or have much lower confidence.

Fig. 1(d): To understand how a group is affected by a scheme, specific states that will be deleted can be inspected in detail (DR6). If a scheme’s outcome is undesirable (for example, deleting the group will remove incidental states with high semantic values), an alternate solution may be chosen. Custom schemes are also allowed, wherein users specify a subset of the records within a group as the ones to be modified, because the selected records may have characteristics that are important for post hoc analysis needs.

5.3 Result Verification (Stage 3)

After the dataset has been sanitized by the application of schemes, the Result Verification Stage allows users to review, tweak, and validate the privacy preservation process.

Fig. 1(g): Because the application of a scheme removes one or more states, the distribution of attribute values in the resultant dataset may become skewed. Such a condition may provide hints to compute and complete missing values. Although we use distribution characteristics as the

basis for sorting the recommended solutions, “damaged” distributions can still occur due to strong correlations between states. To mitigate such situations, post-scheme attribute distributions can be checked (DR7) and “trimmed” to proportionally align with those in the original dataset. When trim is performed, we identify “excess” parts of the distribution, and randomly remove occurrences from these parts until the distribution matches the trim.

Fig. 1(h): To verify that the processed dataset has sufficient defenses (DR8), inference attacks can be simulated by interactively building and running classification models.

Fig. 1(i): The results of classification models are reported to provide insight into both the original and the sanitized dataset’s *sensitivity* (true positive rate) and *specificity* (true negative rate). Users may compare the two dataset reports to assess the success of the privacy preservation process (DR4).

6 THE UMBRA SYSTEM

Based on the workflow described in Sec. 5, we have built a visual analytics system, called Umbra. Each workflow stage is implemented as a dual-column interface, shown in Fig. 2. Each stage’s interface is similarly structured: relevant information about the dataset is shown in a left pane (Figs. 2(b), (d), and (f)), while model-related views, which explain a model’s state or show its results, are shown in a right pane (Figs. 2(c), (e), and (g)).

6.1 Interface for Inference Initialization Stage

The interface for the Inference Initialization Stage contains the State Initialization View in the left pane (Fig. 2(b)) and the Inference Simulation View in the right pane (Fig. 2(c)).

State Initialization View (Fig. 2(b)). A set of distribution charts visualizes the value histograms for each attribute. Numerical attributes and categorical attributes are shown as area charts and bar charts, respectively. Sensitive attributes (i.e., the `employ` attribute) are colored red, while non-sensitive attributes are colored blue.

To create and edit states for numerical attributes, clicking and dragging along the horizontal axis adjusts split points. In the figure, split points for `employ` have been defined as 12 and 48, signifying a record will fall into 1 of 3 states.

For categorical attributes, individual states are shown as vertical bars. For each chart of a non-sensitive attribute, a corresponding input box (labeled “Utility”) shows the parameter of the attribute’s utility, with a default value of 1. Because sensitive attributes must be removed for anonymization, they do not require utility values. For each non-sensitive state, opacity is mapped to the utility value (see Equations 2 and 3).

Inference Simulation View (Fig. 2(c)). As states are defined in the State Initialization View, the corresponding inference graph is updated. Each node represents one state and has the same hue and opacity encodings as in the distribution charts. Edges between nodes show inferences based on the source state to the target state. Its opacity encodes its magnitude and its line style encodes the effect sign (solid for positive and dashed for negative). A slider (labeled “Correlation Filter”) filters correlations associated with the magnitude.

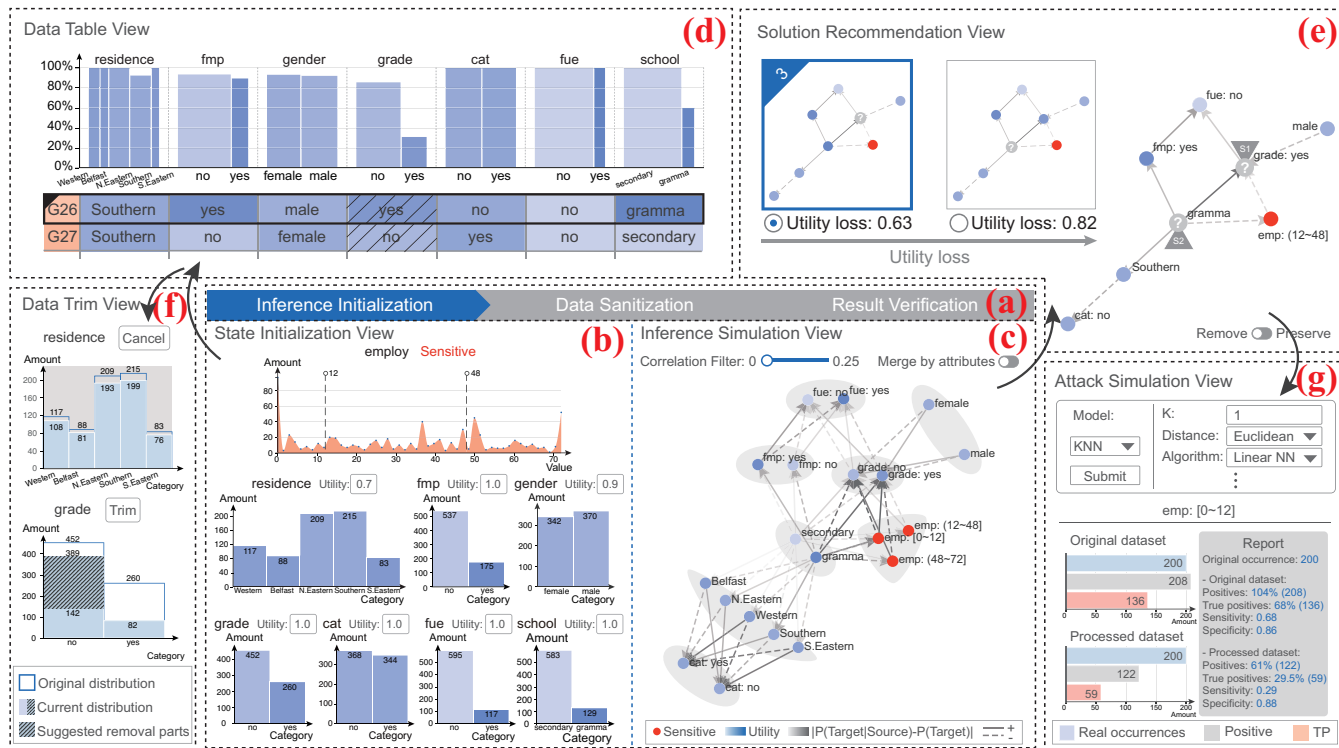


Fig. 2. An integrated layout of our system navigated by (a) a three-stage workflow. Each stage has a dual-column interface, with a data description view ((b), (d), and (f)) on the left, and a model exploration view ((c), (e), and (g)) on the right. The interface of the Inference Initialization Stage consists of a State Initialization View (b) that displays attribute distributions with utility-mapping opacity for event definition, and an Inference Initialization View (c) that depicts all underlying inferences among events by an inference graph. The second interface for the Data Sanitization Stage provides a selected group of records in a Data Table View (d), together with a list of candidate solutions, and an overview summarizing changes of all solutions in a Solution Recommendation View (e). The interface for the Result Verification Stage supports users in verifying the distributions in the Data Trim View (f), and evaluating the processed data's resistance to other attacks in Attack Simulation View (g).

States associated with the same attribute are bounded by grey convex hulls. State nodes can also be merged into aggregated nodes (the “Merge by attributes” toggle), meaning that graph nodes will now show attributes instead of states. Fig. 3(a) shows the merged inference graph of Fig. 2(c). In this case, the node opacity encodes the attribute’s utility value v_{Attr} . Clicking an edge in the merged graph unfolds an edge detail view (Fig. 3(b)), which shows the specific correlations between source and target states.

To see details about a specific inference, right-clicking a state node shows its conditional probabilities (see the tables in Fig. 10). Probabilities can be manually edited (or deleted) to enable custom inferences between states. In the main graph, new inferences can be created by dragging from one state node to another and then inputting both probabilities and conditional probabilities. To keep from overwhelming users, we restrict editing of a target state to only when it is given by one source state. If complex multi-state modifications are necessary, it is simpler to reconstruct the inference graph by loading additional datasets.

To analyze inferences based on multiple states, clicking a sensitive state node will pop up its state set chart (Fig. 4). The state set chart visualizes the conditional probabilities of all possible state combinations on the selected sensitive state. Each brown circle represents a set of non-sensitive states in the dataset. Note that a set can have only one state per attribute. A circle’s horizontal position represents the number of its contained states (Fig. 4. For instance, the

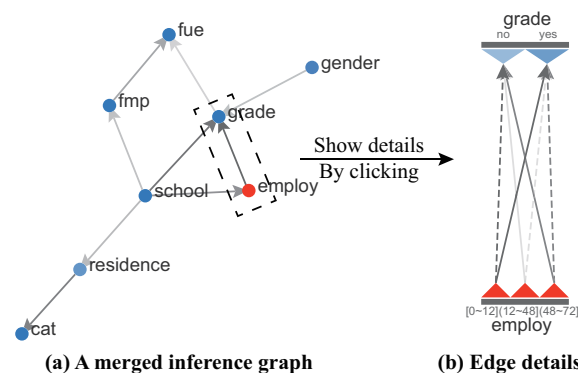


Fig. 3. (a) The inference graph from Fig. 2(c) shown with merged states. A strong correlation can be identified between grade and employment. (b) Edge details illustrating inferences between each state pair defined by the attributes; employment: [0~12] and employment: (48~72) can be inferred with high confidence based on grade.

state number in Fig. 4 varies from 1 to 7, because there are 7 non-sensitive attributes). Meanwhile, its vertical position indicates the conditional probability of the sensitive state given its state set. The horizontal green band shows the calculated no-risk range (see Sec. 4.4). Sets whose probabilities lie within the no-risk range are omitted from the chart, as users primarily care only about states that invoke risk; this means that the vertical space can be compressed. To

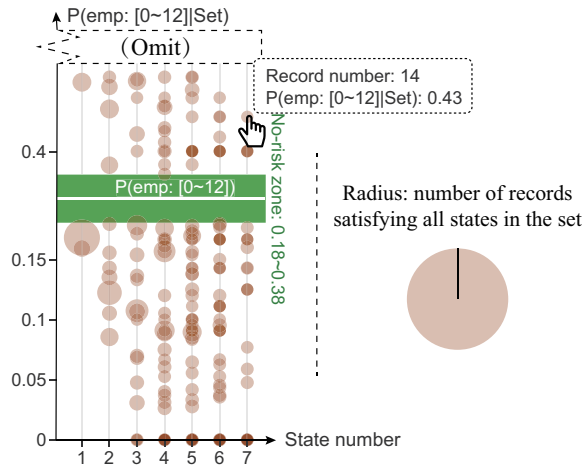


Fig. 4. The state set chart that shows the probabilities of the sensitive state employ: [0~12]. A state set out of the no-risk zone is checked. Related states are highlighted in the inference simulation view.

interpret this chart, the greater the distance of a circle from the no-risk range, the greater the risk caused by its state set. Hovering over a circle highlights the states that belong to it in the inference graph, and hovering on a state in the inference graph conversely highlights circles that include that state.

6.2 Interface for Data Sanitization Stage

Clicking on the navigator bar (Fig. 2(a)) trigger the transition from the Inference Initialization Stage to the Data Sanitization Stage. The second stage is composed of the left-pane Data Table View and right-pane Solution Recommendation View (Figs. 2(d) and (e), respectively).

Data Table View (Fig. 2(d)). Similar to the State Initialization View (Fig. 2(b)), a bar chart is shown for each non-sensitive attribute; individual bars represent states and color opacity indicates each state’s utility value. The height of a bar encodes the percentage of records in the state that remain after the application of currently selected schemes. A value of 100% indicates that no occurrences of that state will be removed. The width of a bar indicates the number of records belonging to that state—a wider bar indicates the state contains a higher percentage of records.

Below the bar charts, a table displays all groups (in Fig. 2(d), we only show two rows due to space limitations). If a group’s index cell (the first column) is colored red, the records within the group are considered at-risk for exposure, demanding the application of schemes to defend against attacks. Rows are ordered by the group exposure risk and size. Columns show the group’s state values.

By reviewing the bar charts, users can decide if a state will be overly modified by the application of schemes (i.e., too high a percentage of its occurrences will be removed). Clicking a state’s bar will sort the table to show related groups to the top. A striped pattern appeared in table cells indicates specific states that will be modified by currently selected schemes. To interact with the records in a group, clicking the group unfolds it to show its contained records (see Fig. 5(a), where group G26 contains 3 records). Users

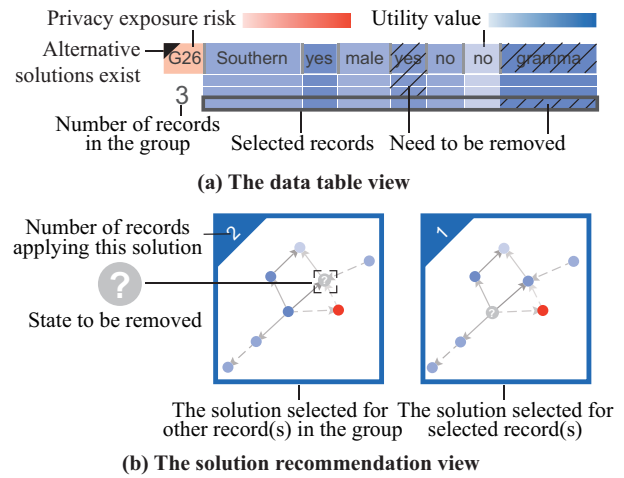


Fig. 5. Selected records and related solutions. (a) Group G26 is at risk, due to its red color. It contains three records, shown as rows beneath the main group row. Due to space, record values are blanked but can be shown by hovering on a cell. (b) The top-2 feasible solutions, shown as simplified inference graphs, demonstrate how inferences will be eliminated to protect the privacy of G26. The first solution will remove the “grade: yes” state and the second will remove “school: grammar” state. The sanitization results of the two solutions are summarized in the data table view in (a).

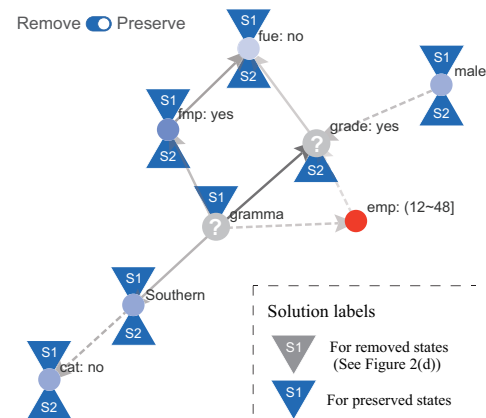


Fig. 6. The simplified inference graph for G26 in the Solution Recommendation View, with the toggle set to show states that will be “preserved.” Scheme S1 preserves all states except “grammar,” and S2 all states except “grade: yes.”

can brush and select records in the group to customize different schemes. Users can additionally review other schemes by clicking the group’s index cell—black triangles denote that alternate solutions exist.

Solution Recommendation View (Fig. 2(e)). For the records selected in Data Table View, the top-recommended schemes are shown as simplified inference graphs (only states that contain records are shown). Fig. 5(b) shows an example of two recommended schemes. States that will be removed are labeled with question marks. (We choose this punctuation to denote uncertainty invoked by the removal of information.) The number of records that will be modified by the solution (i.e., by having attribute values set to “unknown”) is shown in the top-left corner. In Fig. 5(b), 1 and 2 records would be modified by each respective scheme. Recommended schemes are listed in the order of utility loss.

In the simplified inference graphs, two contrasting encodings are provided to demonstrate the impact of schemes (via a “Remove/Preserve” toggle). Triangles appended to nodes that indicate states that will either be preserved (blue triangles, see Fig. 6) or removed (grey triangles, Fig. 2(e)). Triangles are labeled with their respective schemes. In Fig. 2(e), the two grey triangles indicate that scheme S1 will remove the “grade: yes” state, and scheme S2 will remove the “gramma” (graduated from grammar school) state. Fig. 6 shows the same information from the “preserve” perspective.

6.3 Interface for Result Verification Stage

Navigating to the Results Verification Stage applies the selected schemes to the dataset (thus removing state occurrences). This last stage contains the Data Trim View (Fig. 2(f)) and Attack Simulation View (Fig. 2(g)).

Attack Simulation View (Fig. 2(g)). To verify that the applied schemes are a sufficient defense against inference attacks, users may simulate such attacks by interactively training and running binary classification models. The success of an attack on the sanitized data can be compared to the success of the attack on the original data.

Umbr supports several models for simulating inference attacks, including k-nearest neighbor (KNN), Bayesian network, SVM, random forest, and decision tree, with default parameters based on WEKA [48]. We report results in two ways: (1) bar charts provide visual summaries of classification results, and (2) text descriptions show detailed statistical information, including sensitivity and specificity. As an example, in Fig. 2(g), a KNN attack has been simulated, and the number of identified true positives (i.e. records with exposed privacy) is reduced from 136 to 59.

Data Trim View (Fig. 2(f)). To adjust the post-scheme distribution of attribute values, users may also trim attribute distributions. For each non-sensitive state, the original record count is visualized as a dark blue outline, the current (processed) count is shown as a light blue solid area, and the part that is recommended to be trimmed is shaded. Trimming removes the shaded area from the occurrences, changing the attribute’s distributions to be similar to what it was in the original dataset. After trimming, the current distributions will change, and the distribution charts will be covered by a grey rectangle. In Fig. 2(f), the residence attribute has been trimmed and greyed; the grade attribute recommends removing occurrences from the “no” state to match the original distribution.

7 CASE STUDIES

To demonstrate how Umbr can detect and defend against inference attacks, we present a set of case studies on two separate datasets (see Tables 2 and 3). Due to space constraints, screenshots for certain system interactions are omitted. Please watch the supplemental videos for more details.

7.1 Students’ Academic Performance

The **academic performance dataset** [49] contains background and academic information about students extracted from a learning management system (LMS). Suppose that,

TABLE 2
The academic performance dataset contains information about the background and academic behavior of 480 students. The highlighted **class** attribute is considered sensitive.

Attribute	Data Type	Description
placeOfBirth	Categorical	Student’s birth country.
grade	Categorical	Current grade level.
discuss	Numerical	Number of participations in discussion.
raiseHand	Numerical	Number of times the student raises a hand.
visitRes	Numerical	Number of visits to online course content.
absence	Categorical	Number of absences.
class	Categorical	Academic performance rating.
satisfy	Categorical	Satisfaction level of the student’s parents with the school.
relation	Categorical	Parent responsible for student.

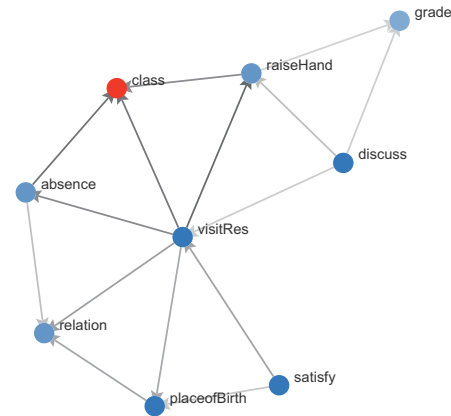


Fig. 7. The merged inference graph of the academic performance dataset shows that the absence, visitRes, and raiseHand attributes are neighbours of the class. Further analysis shows that two specific states (“absence: Above-7,” and “visitRes: [0~49]”) invoke the highest risk for exposure.

the LMS’s data owner wants to provide this data to a similar organization to analyze how certain patterns—e.g. encouraging students to take part in more discussions—helps to improve parental satisfaction. Before sharing, he needs to sanitize the dataset to protect students from having their academic performance exposed. This is represented by the **class** attribute, which shows academic performance and is considered sensitive.

To start, he loads the dataset in Umbr. Due to the low number of records in the dataset (only 80 students), inference confidences are already reduced. He therefore increases the privacy exposure limit from 0.1 (the default) to 0.3. In the inference initialization stage, he uses the State Initialization View to create split points for the numerical attributes based on their value distributions, and creates an “other” state for placeOfBirth data values that have less than 25 occurrences.

In the Inference Simulation View, the merged inference graph (shown in Fig. 7) shows that three attributes have a strong effect on the sensitive **class** attribute. After observing correlations between the attributes, he updates utility values for the attributes based on the length of the path from each attribute to satisfy. Longer paths indicate more indirect correlations, and thus these attributes contribute less to analyses on causal factors of parental satisfaction degrees.

He next loads the state set chart for the “class: L” state (L

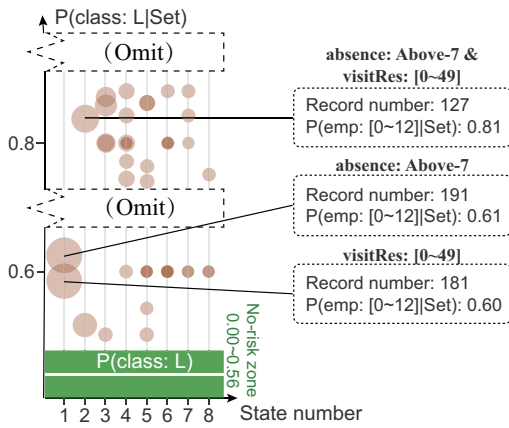


Fig. 8. The state set chart of sensitive state class: L in the academic performance dataset.

denoting low or failing academic performance), see Fig. 8. There are two single-state sets that are outside the green no-risk band: “absence: Above-7” (absent more than 7 times) and “visitRes: [0~49]” (visiting the course content less than 50 times). Surprisingly, when he hovers over these circles in the chart, most sets in the state set chart are highlighted, meaning that most sets that are at-risk contain at least one of these two states. Moreover, the two-state set with the highest exposure risk (0.81) is composed of these two states. From this, he deduces that the primary sources of privacy exposure come from states associated with the absence: Above-7 and visitRes: [0~49] attributes. Since they have direct, high correlations with the sensitive class attribute, an adversary could use them to infer values for that attribute.

Navigating to the Data Sanitization Stage, Umbra recommends removing about half the values of the two identified absence and visitRes states. States belonging to attributes that are “far away” in the inference graph (Fig. 7), such as satisfy and grade, have little correlation on the class attribute; hence, there is little need to modify them during defense construction.

In the Result Verification Stage, he applies a random forest model as a simulated attack—the results are shown in Fig. 9(a). In the original dataset, 119 records are exposed, but the processed data only shows that 58 records could be identified. To further decrease the amount of true positives, he uses data trimming to remove a small amount of occurrences (examples for the raiseHand and discuss attributes are shown in Fig. 9(b)). This slightly reduces the dataset’s sensitivity (original: 0.94, processed: 0.46, trimmed dataset: 0.40), and further reduces the number of true positives from the simulated attack model to 51 instances.

7.2 Home Insurance

The **home insurance dataset** [50] is an example of a Customers Relationship Management (CRM), which shows policy information for a home insurance company from 2007–2012. It stores information about customer homes and associated homeowner policies. To represent the dataset, we sampled 10,000 records from the full dataset (see Table 3 for a list of attributes used). As the CRM data owner, he is

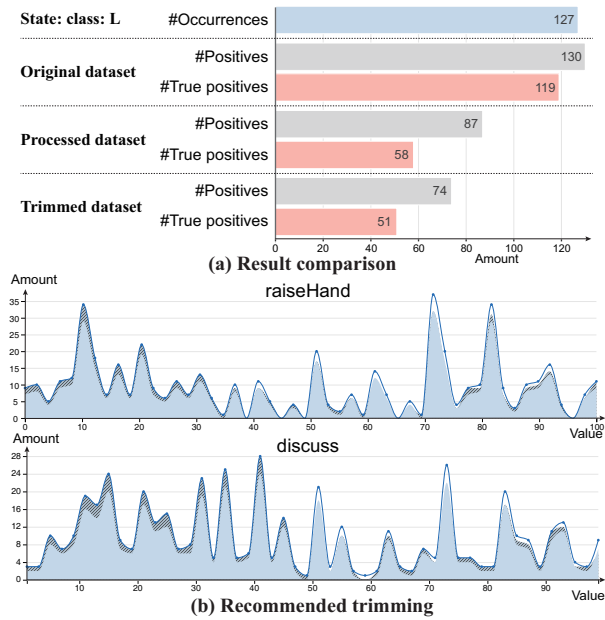


Fig. 9. Verifying the results of privacy preservation on the academic performance dataset. (a) The results of attack simulation implemented by *Random Forest* with default parameters. (b) The recommended trimming schemes for attributes raiseHand and discuss.

publishing the dataset for his company’s annual stakeholder meeting and must abide by the GDPR.

TABLE 3
Selected attributes from the home insurance dataset. The **unocc** and **claim3years** attributes are considered sensitive.

Attribute	Data Type	Description
lastPrem	Numerical	Total premium for the prior year.
yearBuild	Numerical	Year of house construction.
specPrem	Numerical	Premium for personal property.
unocc	Numerical	Number of days house is unoccupied.
alarm	categorical	Appropriate alarm.
lock	categorical	Appropriate lock.
riskRate	Numerical	Geographic risk for buildings.
claim3years	Categorical	Whether there was loss in last 3 years.
sex	Categorical	Customer sex.

Unfortunately, his adversary has a great deal of background knowledge, in the form of their own insurance dataset. We simulate this by sampling a second (non-repeating) subset of 10,000 records from [50]. His adversary will use this auxiliary dataset to infer private information from customers—specifically, unocc and claim3years values. Fortunately, the data owner has access to this second dataset, and will use it to construct his defenses.

He begins by loading both datasets into Umbra, accepting the default split points (for numerical attributes) and state categories (for categorical attributes). He then compares the inference graphs between the two datasets (see Fig. 10). Broadly, the state definitions and the conditional probabilities (i.e., nodes and edges) across the two graphs are very similar, indicating the auxiliary dataset will be good for inference attacks on his dataset. For example, if a premium was low in the past year (the lastPrem: [49.77~94.79] state), the customer’s house was probably in a non-occupied status for a long timespan (unocc: (4~181)). Alternatively,

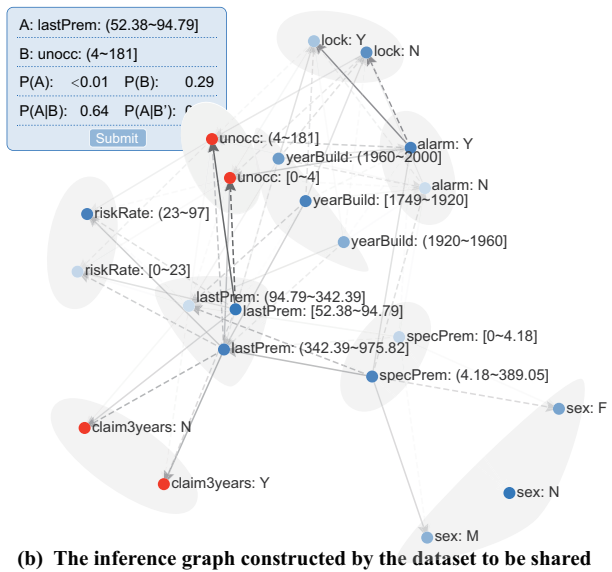
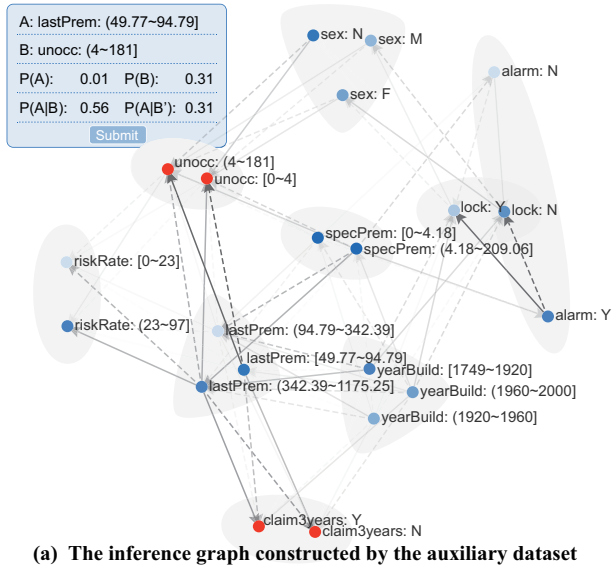


Fig. 10. The inference graphs constructed for the data owner on for the home insurance datasets.

some state definitions differ between the graphs. In his inference graph, the *specPrem* attribute has a range [0~389.05], while the auxiliary dataset’s maximum is only 209.06. To provide consistency, he has Umbra modify the state definitions in the auxiliary dataset’s to match the original.

Now, when he examines in detail the explicit inferences for *unocc* based on *lastPrem*, he notices the graphs have differing confidences: the conditional probability of *lastPrem*: (49.77~94.49] is 0.56 in the auxiliary dataset, lower than the 0.64 in his dataset. He pops up the conditional probability table in the auxiliary dataset’s inference graph (see Fig. 10(a)) and artificially increases the value of this state (the $P(A|B)$ value) to 0.64. In this way, if the true correlation between these states inadvertently becomes public, an adversary will not be able to take advantage of this higher correlation value.

Using inferences from the auxiliary dataset, he navigates to the Data Sanitization Stage and Umbra recommends a set of schemes. In particular, the schemes include the removal

of all occurrences in the state *lastPrem*: (49.77~94.79], and a majority of the occurrences in the state *lastPrem*: (342.39~1175.25] (see Fig. 11). Such an excessive amount of deletions is not desirable, so he checks for alternate solutions by clicking these states’ bars in the Data Table View.

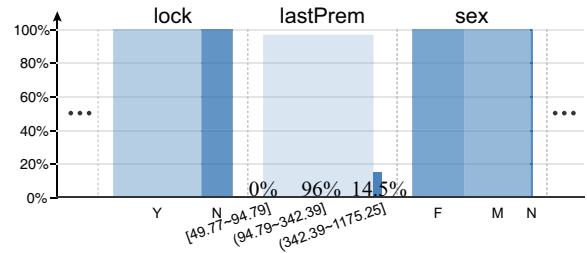


Fig. 11. The attribute distributions of the home insurance data processed by the recommended solutions.

For the state *lastPrem*: (49.77~94.79], no alternate solutions exist—the occurrences must be entirely deleted. However, for the state *lastPrem*: (342.39~1175.25], an alternate scheme does exist (see the triangles labeled S1 and S2 in Fig. 12). Unfortunately, neither scheme preserves the *lastPrem*: (342.39~1175.25], as states defined by the *lastPrem* attribute directly contribute to the inferences of both *unocc* and *claim3years* states. Hence, he has no choice but to accept the default recommendations. The takeaway is that the *lastPrem* attribute is highly linked to sensitive information; only the state with the lowest range of values (*lastPrem*: (52.38~94.79]) can be preserved.

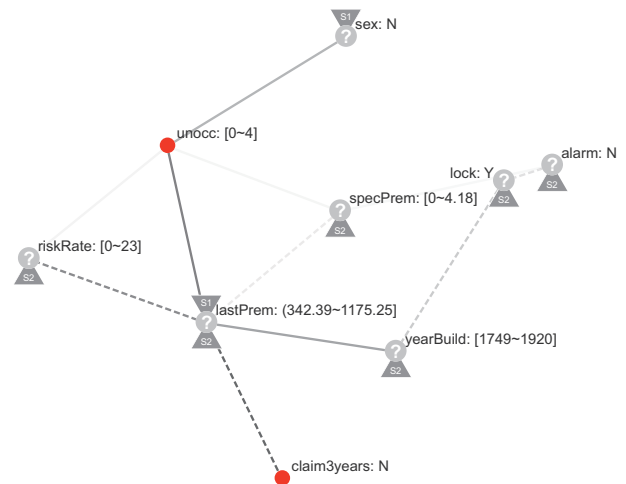


Fig. 12. The solution overview of a group with alternative solutions.

Finally, in the Results Verification Stage, he tests his constructed defenses using a Decision Tree model. The simulated attack—again using the auxiliary dataset—only returns a 0.1% true positive rate (3 total records) for *unocc*: (4~181] (a decrease from 1.0%, 30 records).

8 DISCUSSION

In this section, we summarize feedback from a trio of domain experts, and provide a discussion on the strengths and limitations of our approach.

8.1 Expert Reviews

To collect empirical feedback about Umbra, we conducted a set of expert review sessions with three professors (P_1 – P_3) with at least five years experience researching data privacy preservation (primarily via automated, non-interactive approaches). The intent in these sessions was to holistically assess our approach: not only Umbra’s interface designs, but also its three-stage workflow and use of Bayesian-based inference for defense construction. Each session was conducted as a three-part discussion: (1) discussion of automated approaches, (2) discussion of the validity of Bayesian-based defense construction, and (3) an interactive system demo with Umbra.

(1) Validating Bayesian inference for defense construction. We first explained how Umbra uses Bayesian network inferences to identify sensitive states which can be modified (or removed) via the application of schemes to preserve privacy. All experts confirmed it is reasonable and effective to fix privacy issues by removing combinations of states. One discussed concern of P_3 was scalability: enumerating all state combinations can be challenging when the dataset contains hundreds or thousands of states. P_3 also wondered a differential privacy approach would achieve superior results, compared to the data trimming we use (see Section 8.5 for a discussion of this).

(2) Interactive system demo. Next, we conducted a live system demo with the each expert. We began by walking through Umbra’s functionality, showing examples of how sensitive states can be identified, analyzed, and preserved. At any time, experts could ask questions about the system or provide feedback or critiques. While the demo was semi-structured in that we demonstrated both Umbra’s three-stage workflow and the specific visualizations and interactions in each of its stages, demonstration and interaction with Umbra’s specific features was allowed to freely flow based on the verbal discussion between ourselves and the expert.

All users were enthusiastic about the system. P_2 said, “The system is well-designed. We never thought about showing security and privacy in such a way!” He expressed a desire to use our system as a way explain privacy issues and privacy preserving processes to non-professionals, a comment we find particularly interesting as Umbra’s domain abstraction requires users with expertise in Bayesian modeling of inferences and privacy protection methods.

Notably, P_3 felt that our system has the ability to compensate for the shortcomings of automatic methods. For instance, he could easily identify which states needed to be deleted, thanks to the intuitive inference graph visualization of correlations between states. He also considered the data table and data trimming views as significantly helpful to the privacy preservation process. For the Attack Simulation View, he suggested using a percentage to summarize the results (as opposed to only the occurrence counts). Taking this advice, we updated Umbra and showed percentages and amount simultaneously in the reports of attack results (see Fig. 2(g)).

(3) System usage. P_1 was invited to use the system and give further feedbacks. After a three-hour usage, P_1 considered Umbra as a full-featured system for privacy

preserving. P_1 said, “The user interface of the system is simple to use and responsive in short time.” and “In multiple tests of the system, the data sanitized by the system can always reduce the accuracy of inferring sensitive attributes.”

TABLE 4
Computation time (in milliseconds) for the first two stages with different combinations of record amount and attribute amount.

#Attribute \ #Record	5		10		15	
	State1	Stage2	State1	Stage2	State1	Stage2
500	21	36	599	262	380	1,495
1,000	113	193	712	2,495	521	4,929
5,000	689	1,323	1,376	28,134	3,257	125,321
10,000	854	1,857	2,328	24,931	4,857	196,544

8.2 Scalability and Robustness

We tested the efficiency of our approach by automatically running the first two workflow stages with various-sized sampled subsets of the Home Insurance Dataset [50]—between 500 and 10,000 records, with 5, 10, and 15 attributes each. To limit the number of states, we extracted categorical attributes with two categories and defined numerical states with medians as split points. We tested our system on a desktop with 16G memory and 4 Intel Core i7 4770 at 3.4 GHz processors. For each subset, the average computation time over ten runs is listed in Table 4.

The results demonstrate that when excessive states exist, it is impossible to get timely recommendations. This confirms P_3 ’s observation that enumerating all combinations severely reduces efficiency. Fortunately, the process of recommending schemes—after generating an inference graph—can easily be streamlined by using state-of-the-art software engineering practices, like the incremental approach based on MapReduce proposed by Yue et al. [51]. In practice, the scheme with the lowest utility loss always removes the states that have high correlations with sensitive states. If we need to provide solutions for a dataset with excessive states, we can start construction of the state set from the states close to sensitive states in the inference graph. For instance, we can limit the search space to the third-order neighbors of sensitive states, meaning that only a minority of sets need to be considered.

Unfortunately, arbitrarily limiting the search space may cause the optimum solution to be missed in special cases. In Umbra, we do not limit the search space when processing datasets, which guarantees optimum recommendations. We additionally exclude sets whose subsets meet privacy preserving requirements with a lower number of states, which improves performance.

8.3 Defining States: Specific or General?

Our current workflow allows users to define states only by the values of a single attribute. In real-world scenarios, inferences may involve subtle and/or complex relationships across multiple attributes. For instance, an inference could be written as, “people who are cheerful and capable of working have more chances to get a satisfying job.” The state set of “being cheerful and capable of working” can actually be regarded as two states—personality: cheerful and

ability: good. In the inference extraction process, we consider complex states as a group of one-attribute states that occur simultaneously.

A comprehensive state definition method, supporting complex (and potentially even fuzzy) states, can guide our system to focus on certain inferences. However, such methods need more specific definitions and a more decentralized description to explain inferences, which will increase users' learning and interaction load. Moreover, users may miss important inferences, because rules-based definitions limit the search space of the backend algorithm. That is the reason why we take the current state definition method for Umbra.

8.4 Limitation on Inference Depiction

Although edges with opacity are capable of explaining the inferences based on one state, we found that paths are not efficient for depicting inferences. This is because conditional probabilities $\Pr(S_A|S_B, S_C)$ are obtained by counting the occurrences of S_A given S_B and S_C from the data sets, instead of operating $\Pr(S_A|S_B)$ and $\Pr(S_A|S_C)$. Inferences with multiple states can hardly be presented by means of opacity-modulated paths.

Distinct non-sensitive condition combinations can lead to different occurrence probabilities of sensitive states, which trigger different privacy exposure risks. Thus, there is a strong need to make up for this deficiency. Neither the state set chart or the data table view can intuitively explain the effect of increasing or decreasing a state on an inference. In the future, we plan to elaborate more details of inferences in an inference view.

8.5 Sanitization Alternatives

Umbra currently focuses on sanitizing datasets by removing certain sets of attribute values (data trimming), which is proved to be effective [26]. As mentioned in Sec. 2.2, there exist several sanitization alternatives, such as differential privacy approaches [19], [52] (which were brought up by P_3 in his expert review session). However, differential approaches contain their own complexities. For example, explaining data changes caused by differential privacy sanitization has its own interpretation difficulties and would require a significantly different visualization designs. In addition, much like this initial iteration of Umbra requires users with knowledge of Bayesian techniques, systems that leverage differential privacy require users with expert knowledge about those types of approaches.

Other sanitization approaches, such as adding noise to attribute values, can mean that no states are actually removed or merged. Unfortunately, randomization approaches mean that the values in each state have a certain probability of being fake, and as a result the dataset contains some amount of error (incorrect data records) within it. In contrast, the sanitized datasets produced by Umbra only contain true attribute values. Depending on a data owner's desires as well as the application context of the dataset that is currently being sanitized (a data owner might ask, "can a small amount of error be allowed?"), alternative sanitization schemes might be preferred [37]. That said, a full investigation of alternative sanitization schemes is beyond the scope of the current work, though it should certainly be investigated in the future.

8.6 What About Data Subjects and Novice Users?

Umbra is intended for data owners with understanding of privacy protection models and Bayesian statistics. While data subjects could potentially use systems like Umbra for personal data management, such as verifying the protection status of their own data records (DR6), it's likely such systems are overwhelming in terms of privacy models, Bayesian statistics, and visual encodings. Prior research suggests that novice users can have trouble reasoning about complex concepts such as privacy modeling [53], [54], Bayesian probabilities and statistics [39], [40], and even novel visualization encodings [55], [56], [57]. Communicating highly technical concepts to such users—not only the vast majority of data subjects, but likely a large number of data owners without sufficient technical expertise—in ways that are easily digestible and interpretable is a non-trivial task, and one that is outside the scope of this current work.

9 CONCLUSION

Inference attacks are intractable due to distinct inferences applied by adversaries with different background knowledge. We present a visual analytics approach as a three-stage workflow to simulate the inference behaviors of adversaries and seek appropriate processing to effectively block the underlying inferences about sensitive states. Intuitive visual designs allow users to learn the privacy preserving process without excessive learning costs.

We demonstrate in case studies and expert reviews that Umbra meets the needs of users to understand automatic models and follow the changes of datasets that occur during sanitization. In particular, the inference graph view is especially effective for visualizing the sources of privacy exposure risks from multiple perspectives and customizing utility criteria reasonably. Umbra thus shows promise to defend against inference attacks.

ACKNOWLEDGMENTS

The authors would like to thank the experts who we consulted/interviewed for this work. Xumeng Wang and Wei Chen are supported by National Natural Science Foundation of China (61772456, 61761136020).

REFERENCES

- [1] S. D. Warren and L. D. Brandeis, "The right to privacy," *Harvard Law Review*, pp. 193–220, 1890.
- [2] P. Voigt and A. Von dem Bussche, *The EU General Data Protection Regulation (GDPR)*. Springer, 2017, vol. 18.
- [3] "Eugdpr - information portal," <https://eugdpr.org/>.
- [4] H. Liu, J. Zhou, Q.-L. Feng, H.-T. Gu, G. Wan, H.-M. Zhang, Y.-J. Xie, and X.-S. Li, "Fetal echocardiography for congenital heart disease diagnosis: a meta-analysis, power analysis and missing data analysis," *European Journal of Preventive Cardiology*, vol. 22, no. 12, pp. 1531–1547, 2015.
- [5] M. J. Keith, S. C. Thompson, J. Hale, P. B. Lowry, and C. Greer, "Information disclosure on mobile devices: Re-examining privacy calculus with actual user behavior," *International Journal of Human-computer Studies*, vol. 71, no. 12, pp. 1163–1173, 2013.
- [6] C. Li, H. Shirani-Mehr, and X. Yang, "Protecting individual information against inference attacks in data publishing," in *Proceedings of International Conference on Database Systems for Advanced Applications*. Springer, 2007, pp. 422–433.
- [7] N. Friedman, D. Geiger, and M. Goldszmidt, "Bayesian network classifiers," *Machine Learning*, vol. 29, no. 2-3, pp. 131–163, 1997.

- [8] G. Li, J. Shi, and J. Zhou, "Bayesian adaptive combination of short-term wind speed forecasts from neural network models," *Renewable Energy*, vol. 36, no. 1, pp. 352–359, 2011.
- [9] B. Yet, A. Constantinou, N. Fenton, M. Neil, E. Luedeling, and K. Shepherd, "A bayesian network framework for project cost, benefit and risk analysis with an agricultural development case study," *Expert Systems with Applications*, vol. 60, pp. 141–155, 2016.
- [10] G. T. Duncan and S. L. Stokes, "Disclosure risk vs. data utility: The ru confidentiality map as applied to topcoding," *Chance*, vol. 17, no. 3, pp. 16–20, 2004.
- [11] K. Sathiyapriya and G. S. Sadasivam, "A survey on privacy preserving association rule mining," *International Journal of Data Mining & Knowledge Management Process*, vol. 3, no. 2, p. 119, 2013.
- [12] J. L. Dautrich Jr and C. V. Ravishankar, "Compromising privacy in precise query protocols," in *Proceedings of the 16th International Conference on Extending Database Technology*. ACM, 2013, pp. 155–166.
- [13] M. S. Islam, M. Kuzu, and M. Kantarcioglu, "Access pattern disclosure on searchable encryption: Ramification, attack and mitigation." in *Proceedings of the 19th Annual Network and Distributed System Security Symposium*, vol. 20, 2012, p. 12.
- [14] P. Zhao, H. Jiang, C. Wang, H. Huang, G. Liu, and Y. Yang, "On the performance of k-anonymity against inference attacks with background information," *IEEE Internet of Things Journal*, 2018.
- [15] T. Li, N. Li, and J. Zhang, "Modeling and integrating background knowledge in data anonymization," in *Proceedings of IEEE International Conference on Data Engineering*, 2009, pp. 6–17.
- [16] Y. Sun, L. Yin, L. Liu, and S. Xin, "Toward inference attacks for k-anonymity," *Personal and Ubiquitous Computing*, vol. 18, no. 8, pp. 1871–1880, 2014.
- [17] C. Liu, S. Chakraborty, and P. Mittal, "Dependence makes you vulnerable: Differential privacy under dependent tuples." in *Proceedings of the 23rd Annual Network and Distributed System Security Symposium*, vol. 16, 2016, pp. 21–24.
- [18] G. Cormode, "Personal privacy vs population privacy: learning to attack anonymization," in *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2011, pp. 1253–1261.
- [19] B. Yang, I. Sato, and H. Nakagawa, "Bayesian differential privacy on correlated data," in *Proceedings of the 2015 ACM SIGMOD international conference on Management of Data*, pp. 747–762.
- [20] S. Zhang and S. Song, "A novel attack graph posterior inference model based on bayesian network," *Journal of Information Security*, vol. 2, no. 01, p. 8, 2011.
- [21] S. T. Timmer, J.-J. C. Meyer, H. Prakken, S. Renooij, and B. Verheij, "Inference and attack in bayesian networks," in *Proceedings of the 25th Benelux Conference on Artificial Intelligence*. Citeseer, 2013, pp. 199–206.
- [22] Z. Li, G. Zhan, and X. Ye, "Towards an anti-inference (k, l)-anonymity model with value association rules," in *Proceedings of International Conference on Database and Expert Systems Applications*. Springer, 2006, pp. 883–893.
- [23] N. Li, T. Li, and S. Venkatasubramanian, "t-closeness: Privacy beyond k-anonymity and l-diversity," in *Proceedings of IEEE 23rd International Conference on Data Engineering*, 2007, pp. 106–115.
- [24] J. Hamm, "Minimax filter: learning to preserve privacy from inference attacks," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 4704–4734, 2017.
- [25] S. Salamatian, A. Zhang, F. du Pin Calmon, S. Bhamidipati, N. Fawaz, B. Kveton, P. Oliveira, and N. Taft, "Managing your private and public data: Bringing down inference attacks against your privacy." *IEEE Journal of Selected Topics in Signal Processing*, vol. 9, no. 7, pp. 1240–1255, 2015.
- [26] J. Chen, J. He, L. Cai, and J. Pan, "Disclose more and risk less: Privacy preserving online social network data sharing," *IEEE Transactions on Dependable and Secure Computing*, 2018.
- [27] Z. Cai, Z. He, X. Guan, and Y. Li, "Collective data-sanitization for preventing sensitive information inference attacks in social networks," *IEEE Transactions on Dependable and Secure Computing*, vol. 15, no. 4, pp. 577–590, 2018.
- [28] A. Dasgupta and R. Kosara, "Adaptive privacy-preserving visualization using parallel coordinates," *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 12, pp. 2241–2248, 2011.
- [29] A. Dasgupta, M. Chen, and R. Kosara, "Measuring privacy and utility in privacy-preserving visualization," *Computer Graphics Forum*, vol. 32, no. 8, pp. 35–47, 2013.
- [30] J.-K. Chou, Y. Wang, and K.-L. Ma, "Privacy preserving event sequence data visualization using a sankey diagram-like representation," in *Proceedings of SIGGRAPH ASIA 2016 Symposium on Visualization*, p. 1.
- [31] J.-K. Chou, C. Bryan, and K.-L. Ma, "Privacy preserving visualization for social network data with ontology information," in *Proceedings of IEEE Pacific Visualization Symposium*, 2017, pp. 11–20.
- [32] J. Oksanen, C. Bergman, J. Sainio, and J. Westerholm, "Methods for deriving and calibrating privacy-preserving heat maps from mobile sports tracking application data," *Journal of Transport Geography*, vol. 48, pp. 135–144, 2015.
- [33] A. Dasgupta, R. Kosara, and M. Chen, "Guess me if you can: A visual uncertainty model for transparent evaluation of disclosure risks in privacy-preserving data visualization."
- [34] J. Chou, C. Bryan, J. Li, and K. Ma, "An empirical study on perceptually masking privacy in graph visualizations," in *2018 IEEE Symposium on Visualization for Cyber Security (VizSec)*, Oct 2018, pp. 1–8.
- [35] X. Wang, T. Gu, X. Luo, X. Cai, T. Lao, W. Chen, Y. Wu, J. Yu, and W. Chen, "A user study on the capability of three geo-based features in analyzing and locating trajectories," *IEEE Transactions on Intelligent Transportation Systems*, 2018.
- [36] C.-H. Kao, C.-H. Hsieh, Y.-F. Chu, Y.-T. Kuang, and C.-K. Yang, "Using data visualization technique to detect sensitive information re-identification problem of real open dataset," *Journal of Systems Architecture*, vol. 80, pp. 85–91, 2017.
- [37] X. Wang, J.-K. Chou, W. Chen, H. Guan, W. Chen, T. Lao, and K.-L. Ma, "A utility-aware visual approach for anonymizing multi-attribute tabular data," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 1, pp. 351–360, 2018.
- [38] X. Wang, W. Chen, J.-K. Chou, C. Bryan, H. Guan, W. Chen, R. Pan, and K.-L. Ma, "GraphProtector: A visual interface for employing and assessing multiple privacy preserving graph algorithms," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 193–203, 2019.
- [39] A. Ottley, E. M. Peck, L. T. Harrison, D. Afergan, C. Ziemkiewicz, H. A. Taylor, P. K. Han, and R. Chang, "Improving Bayesian reasoning: The effects of phrasing, visualization, and spatial ability," *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 1, pp. 529–538, 2015.
- [40] L. Micallef, P. Dragicevic, and J.-D. Fekete, "Assessing the effect of visualizations on bayesian reasoning through crowdsourcing," *IEEE Transactions on Visualization and Computer Graphics*, vol. 18, no. 12, pp. 2536–2545, 2012.
- [41] T. Xiao, H. Li, W. Ouyang, and X. Wang, "Learning deep feature representations with domain guided dropout for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1249–1258.
- [42] J. P. Vandenbroucke, A. Broadbent, and N. Pearce, "Causality and causal inference in epidemiology: the need for a pluralistic approach," *International Journal of Epidemiology*, vol. 45, no. 6, pp. 1776–1786, 2016.
- [43] S. Ji, W. Li, P. Mittal, X. Hu, and R. A. Beyah, "Secgraph: A uniform and open-source evaluation system for graph data anonymization and de-anonymization." in *USENIX Security Symposium*, 2015, pp. 303–318.
- [44] B. Shneiderman, "Human-centered artificial intelligence: Reliable, safe & trustworthy," *International Journal of Human-Computer Interaction*, vol. 36, no. 6, pp. 495–504, 2020.
- [45] D. McVicar and M. Anyadike-Danes, "Predicting successful and unsuccessful transitions from school to work by using sequence methods," *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, vol. 165, no. 2, pp. 317–334, 2002.
- [46] Z. Cai, Z. He, X. Guan, and Y. Li, "Collective data-sanitization for preventing sensitive information inference attacks in social networks," *IEEE Transactions on Dependable and Secure Computing*, vol. 15, no. 4, pp. 577–590, 2016.
- [47] R. Heatherly, M. Kantarcioglu, and B. Thuraisingham, "Preventing private information inference attacks on social networks," *IEEE Transactions on Knowledge and Data Engineering*, vol. 25, no. 8, pp. 1849–1862, 2012.
- [48] F. Eibe, M. Hall, and I. Witten, "The weka workbench. online appendix for data mining: Practical machine learning tools and techniques," *Morgan Kaufmann*, 2016.
- [49] I. Aljarah, "Students' academic performance dataset," <https://www.kaggle.com/aljarah/xAPI-Edu-Data#xAPI-Edu-Data.csv>.

- [50] Y. Canario, "Home insurance," <https://www.kaggle.com/ycanario/home-insurance>.
- [51] K. Yue, Q. Fang, X. Wang, J. Li, and W. Liu, "A parallel and incremental approach for data-intensive learning of bayesian networks," *IEEE transactions on cybernetics*, vol. 45, no. 12, pp. 2890–2904, 2015.
- [52] J. He, L. Cai, and X. Guan, "Preserving data-privacy with added noises: Optimal estimation and privacy analysis," *IEEE Transactions on Information Theory*, vol. 64, no. 8, pp. 5677–5690, 2018.
- [53] R. Kang, L. Dabbish, N. Fruchter, and S. Kiesler, "my data just goes everywhere: user mental models of the internet and implications for privacy and security," in *Eleventh Symposium On Usable Privacy and Security ({SOUPS} 2015)*, 2015, pp. 39–52.
- [54] A. Rao, F. Schaub, N. Sadeh, A. Acquisti, and R. Kang, "Expecting the unexpected: Understanding mismatched privacy expectations online," in *Twelfth Symposium on Usable Privacy and Security ({SOUPS} 2016)*, 2016, pp. 77–96.
- [55] E. M. Peck, S. E. Ayuso, and O. El-Etr, "Data is personal: Attitudes and perceptions of data visualization in rural pennsylvania," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2019, pp. 1–12.
- [56] A. V. Maltese, J. A. Harsh, and D. Svetina, "Data visualization literacy: Investigating data interpretation along the noviceexpert continuum," *Journal of College Science Teaching*, vol. 45, no. 1, pp. 84–90, 2015.
- [57] J. Boy, R. A. Rensink, E. Bertini, and J.-D. Fekete, "A principled way of assessing visualization literacy," *IEEE transactions on visualization and computer graphics*, vol. 20, no. 12, pp. 1963–1972, 2014.



Rusheng Pan is a postgraduate student in the State Key Lab of CAD&CG at Zhejiang University, Hangzhou. He earned the B.S. degree in control science and engineering from Zhejiang University in 2018. His research interests are privacy preservation and graph visualization.



Yanling Liu received a bachelors degree in Software Engineering in 2019 from Zhejiang University. Her interests include information visualization and visual analytics.



Xumeng Wang is a Ph.D. candidate in the State Key Lab of CAD&CG at Zhejiang University, Hangzhou. She earned the B.S. degree in information and computing science from Zhejiang University in 2016. Her research interests are visual analytics and privacy preservation.



Wei Chen is a professor in the State Key Lab of CAD&CG, Zhejiang University. His research interests include visualization and visual analysis, and has published more than 70 IEEE/ACM Transactions and IEEE VIS papers. He actively served as guest or associate editors of the ACM Transactions on Intelligent System and Technology, the IEEE Computer Graphics and Applications, and Journal of Visualization.



Chris Bryan is an assistant professor of computer science in the School of Computing, Informatics, and Decision Systems Engineering at Arizona State University, where he directs the Sonoran Visualization Laboratory (SVL @ ASU). He received the PhD degree in computer science from the University of California, Davis, in 2018. His research areas include information visualization, human-computer interaction, and virtual reality. He is a member of the IEEE.



Kwan-Liu Ma is a distinguished professor of computer science at the University of California, Davis. He directs VIDL Labs and UC Davis Center of Excellence for Visualization. Professor Ma received his PhD degree in computer science from the University of Utah in 1993. His research interests include visualization, computer graphics, high-performance computing, and human-computer interaction. Professor Ma was a recipient of the NSF PECASE award in 2000 and the IEEE VGTC 2013 Visualization Technical Achievement Award. He is an IEEE Fellow. He presently serves as the AEIC of IEEE CG&A.



Yiran Li is a Ph.D. Student at the Department of Computer Science in UC Davis. She earned the B.S. degree in Mathematical Sciences from Zhejiang University in 2018. Her research interest is visual analytics.