

Butterfly: Protecting Output Privacy in Stream Mining

Ting Wang Ling Liu

College of Computing, Georgia Institute of Technology
801 Atlantic Drive, Atlanta, GA 30332
{twang, lingliu}@cc.gatech.edu

Abstract—Privacy preservation in data mining demands protecting both input and output privacy. The former refers to sanitizing the raw data itself before performing mining. The latter refers to preventing the mining output (model/pattern) from malicious pattern-based inference attacks. The preservation of input privacy does not necessarily lead to that of output privacy.

This work studies the problem of protecting output privacy in the context of frequent pattern mining over data streams. After exposing the privacy breaches existing in current stream mining systems, we propose **Butterfly**, a light-weighted countermeasure that can effectively eliminate these breaches without explicitly detecting them, meanwhile minimizing the loss of the output accuracy. We further optimize the basic scheme by taking account of two types of semantic constraints, aiming at maximally preserving utility-related semantics while maintaining the hard privacy and accuracy guarantee. We conduct extensive experiments over real-life datasets to show the effectiveness and efficiency of our approach.

I. INTRODUCTION

Recent years have witnessed increasing concerns about individual privacy in numerous data mining and management applications. Individuals were usually unwilling to provide their personal information if they knew that the privacy of the data could be compromised. A plethora of work has been done on preserving the *input privacy* for static data [1], [2], [3], [4], [5], which assumes untrusted mining service providers and enforces privacy regulations by sanitizing the raw data before sending it to the service providers. The mining algorithms are performed over the sanitized data. This scenario is shown as the first four steps of Fig. 1.

However, surprisingly limited attention has been given to preserving *output privacy* in data mining: the published mining output can be leveraged to infer properties possessed only by a unique or a small number of individuals, even though the models/patterns may be built over the sanitized data. This can be explained by the fact that input-privacy preserving techniques are designed to make the constructed models/patterns as close as possible to, if not identical to that built over the raw data, in order to guarantee the utility of the result. This no-outcome-change property is considered as a pillar of privacy preserving data mining [6]. As long as the significant statistical information of the raw data is preserved, there exists the risk of disclosure of private information. Therefore, the preservation of input privacy does not necessarily lead to that of output privacy.

Example 1: Consider a nursing-care records database that records the observed symptoms of the patients in a hospital. By mining such a database, one can discover valuable information

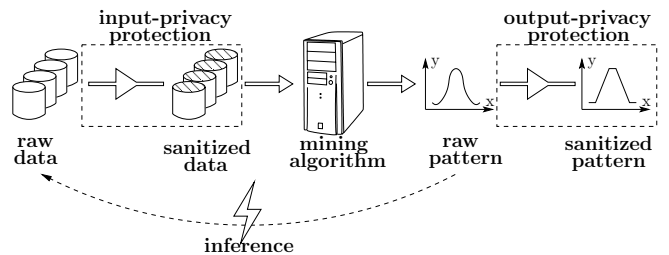


Fig. 1. Illustration of privacy protection in data mining applications.

about syndromes characterizing particular diseases. However, the released mining output can also be leveraged to uncover some combinations of symptoms that are so special that rare people match them (we will show how to achieve this in the following sections), which qualifies as severe threats to individuals' privacy.

Assume that Alice knows that Bob has certain symptoms a, b but not c (\bar{c}), and by analyzing the mining output, she finds that only one person in the hospital matches the specific combination of $\{a, b, \bar{c}\}$, and only one has all $\{a, b, \bar{c}, d\}$. She can then conclude that the one is Bob, who also has the symptom d . Further more, by studying other medical database, she may learn that the combination of $\{a, b, d\}$ is linked to a rare disease with high chance.

This output privacy issue is even severer in stream mining: The mining results need to be published in a continuous and in-time manner. Not only a single-time release may contain privacy breaches, but also multiple releases can potentially be exploited in combination, given the overlap of the corresponding input data. Taking the sliding window model as an example, in addition to the leakage in the output of a single window (*intra-window breach*), the output of multiple overlapping window can be combined to infer sensitive information (*inter-window breach*), even though each window itself contains no privacy breach per se. Therefore, one needs to consider protecting output privacy in stream mining as a unique problem.

A straightforward solution to preserving output privacy is to detect and eliminate all potential breaches, i.e., the *detecting-then-removing* strategy adopted in inference control of statistical databases and census data from 1970's. However, the results are usually negative in tone [7] for on-line stream mining systems: First, the detection of breaches usually requires complicated computations over entire dataset and the bookkeeping of voluminous history output; Second, even at such high cost, the operations of removing the found breaches,

e.g., suppression, addition [8], usually result in significant decrease of the utility of the output.

In this work, we study the problem of protecting output privacy in the context of frequent pattern mining over streams. Concretely, analogous to sanitizing raw data from leaking sensitive information, we propose the concept of “sanitized pattern”, and argue that by intelligently modifying the “raw pattern” outputted by mining algorithms, one can significantly reduce the risk of malicious inferences, while maximally preserving the utility of the raw patterns. This scenario is shown as the last step in Fig. 1.

Specifically, we present **Butterfly**, a light-weighted countermeasure against malicious inferences based on value perturbation. It possesses the following desirable features: (i) No need of explicit detection of privacy breaches; (ii) No need of bookkeeping of history output; (iii) Flexible control over the balance of multiple utility metrics and privacy guarantee.

Our Contributions (i) We articulate the problem of protecting output privacy in stream mining, and expose the privacy breaches existing in current stream mining systems; (ii) We propose a generic framework of protecting output privacy: On the first tier, it counters malicious inferences by amplifying the uncertainty of sensitive information; On the second tier, for the given privacy requirement, it maximally preserves output utility; (iii) We provide both theoretical analysis and experimental evaluation to validate our approach in terms of privacy guarantee, output utility and algorithm efficiency.

II. RELATED WORK

In this section, we outline the related work along three most relevant areas.

Disclosure Control in Statistical Database Extensive research has been done in statistical databases to provide statistical information without compromising sensitive information regarding individuals [9], [10], using the techniques of query restriction or data perturbation. Compared with the simple statistical information, e.g., min, max, avg, etc, the mining output (model/pattern) usually has more complex structures, leading to more complicated requirement for output utility, which makes it hard to directly apply these techniques in our scenario.

Input Privacy Preservation The work of [11], [1] paved the way for the rapidly expanding field of privacy preserving data mining. While a plethora of techniques have been developed, including data perturbation [11], [1], [2], [3], k -anonymity [4], [5] and secure multi-party computation [12], these techniques focus on protecting input privacy for static datasets, where the design goal is to provide sufficient privacy guarantee while minimizing the information loss in the mining output. A recent work [13] also addresses the problem of preserving input privacy for streaming data, by on-line analysis of correlation structure of multivariate streams.

Output Privacy Preservation Compared with the wealth of techniques developed for preserving input privacy, protecting mining output privacy has not received the attention it deserves. The work [14] proposes an empirical testing scheme to evaluate if the constructed classifier violates the privacy constraint. It is shown in [8] that the association rules can be

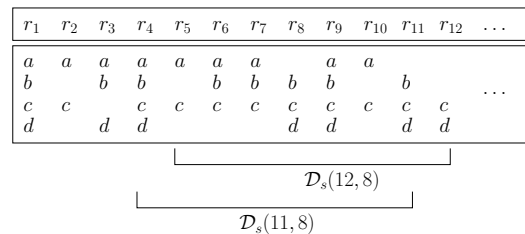


Fig. 2. Illustration of data stream and sliding window model.

exploited to infer information about individual transactions. The work [15] proposes a scheme to block the inference of sensitive patterns satisfying user-specified templates by suppressing certain raw transactions. A recent work [6] tries to consider input and output privacy in a unified framework, however it is not clear that it can prevent malicious inference over the mining output. To the best of our knowledge, none has addressed the problem of protecting output privacy in stream mining applications.

III. PROBLEM DEFINITION

In this section, we introduce the output privacy issue arising in this context of frequent pattern mining over data streams.

A. Frequent Pattern Mining

Consider a finite set of items $\mathcal{I} = \{i_1, i_2, \dots, i_M\}$. An itemset I is a subset of \mathcal{I} , i.e., $I \subseteq \mathcal{I}$. A database \mathcal{D} consists of a set of records, each corresponding to a non-empty itemset. The support of an itemset I w.r.t. \mathcal{D} , denoted by $T_{\mathcal{D}}(I)$, is the number of records in \mathcal{D} , which contain I as a subset.

A data stream \mathcal{D}_s is modeled as a sequence of records, (r_1, \dots, r_N) , where N is the current size of the stream. The sliding window model is introduced to deal with the potential of N going to infinity. Concretely, at each N , one considers only the window of most recent H records, r_{N-H+1}, \dots, r_N , denoted by $\mathcal{D}_s(N, H)$, where H is the size of the window. The problem of frequent stream pattern mining is to find in each $\mathcal{D}_s(N, H)$, all itemsets with support exceeding a user-defined threshold C , called the minimum support.

One can further generalize the concept of itemset by introducing the negation \bar{i} of an item i : A record is said to contain \bar{i} if it does not contain i . Following, we will use the term *pattern* to denote a set of items or negation of items, e.g., $\bar{a}\bar{b}c$. We use the notation \bar{I} to represent the negation of an itemset I , i.e., $\bar{I} = \{\bar{i} | i \notin I\}$.

Analogously, one can define the support of a pattern p w.r.t. a database \mathcal{D} : We say a record satisfies p if it contains all the items and negations of items in p . The support of p w.r.t. \mathcal{D} is the number of records in \mathcal{D} that satisfy p .

Example 2: Consider a data stream with current size $N = 12$, window size $H = 8$, shown in Fig. 2, where $a \sim h$ and $r_1 \sim r_{12}$ represent the set of items and records respectively. The pattern $\bar{a}\bar{b}c$ has support 2 w.r.t. $\mathcal{D}_s(12, 8)$, because only records r_8 and r_{11} match it.

B. Mining Output Privacy

The output privacy refers to the requirement that the output (model/pattern) of a mining process does not disclose any sensitive information regarding an individual or a small number of records.

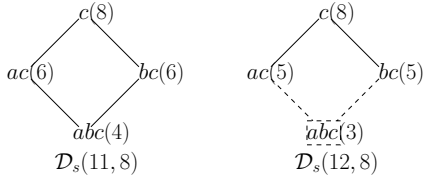


Fig. 3. Illustration of intra-window and inter-window breaches. The lattice structure \mathcal{X}_c^{abc} is shown, with the itemsets and their support in two windows $\mathcal{D}_s(11, 8)$ and $\mathcal{D}_s(12, 8)$ respectively.

Examples of such sensitive information in the context of frequent pattern mining are usually in the form of patterns with low support. Recall Example 1 of nursing-care records in Section I, clearly the disclosure of such patterns through output inference can lead to uncovering sensitive information regarding few individuals.

Therefore we introduce the concept of *vulnerable pattern*. Intuitively, vulnerable patterns are those with very low support w.r.t. the given database, i.e., only few individual records match them. To quantitatively measure this, we introduce a threshold K ($K \ll C$), called the vulnerable support. We have the following classification of patterns based on their support values.

Definition 1 (Pattern Classification): Given a database \mathcal{D} , let \mathcal{P} be the set of patterns appearing in \mathcal{D} , then all $p \in \mathcal{P}$ can be classified into three disjoint classes

$$\begin{cases} \text{Frequent Pattern :} & \mathcal{P}_f = \{p | T_{\mathcal{D}}(p) \geq C\} \\ \text{Hard Vulnerable Pattern :} & \mathcal{P}_{hv} = \{p | 0 < T_{\mathcal{D}}(p) \leq K\} \\ \text{Soft Vulnerable Pattern :} & \mathcal{P}_{sv} = \{p | K < T_{\mathcal{D}}(p) < C\} \end{cases}$$

for the given thresholds C and K .

The frequent patterns (\mathcal{P}_f) expose the significant statistics of the underlying data, and are often the candidate in the mining process. Actually the frequent itemsets found by the mining process are a subset of \mathcal{P}_f . The hard vulnerable patterns (\mathcal{P}_{hv}) represent the properties possessed by only a very small number of individuals, so it is unacceptable that they are disclosed or inferred from the mining output. The soft vulnerable patterns (\mathcal{P}_{sv}) neither demonstrate statistical significance, nor violate the privacy of individual records.

Therefore, the problem of protecting output privacy in stream frequent pattern mining can be stated as follows:

Definition 2 (Problem Formulation): For each sliding window $\mathcal{D}_s(N, H)$, output privacy protection prevents the disclosure or inference of any hard vulnerable patterns w.r.t. $\mathcal{D}_s(N, H)$ from the mining output.

Although the output of frequent pattern mining contains only those patterns with their support exceeding C ($C \gg K$), as we will show in the next section, an adversary may still be able to infer certain patterns with support below K from the released frequent itemsets and their associated support.

IV. OUTPUT PRIVACY BREACHES

For ease of presentation, we will use the following notations: for two itemsets I and J , $I \cup J$ denotes their union, $J \setminus I$ the difference of J and I , and $|I|$ the size of I .

A. Attack Techniques

Lattice Structure As a special case of multi-attribute aggregation, computing the support of $I \subset J$ can be considered

as generalization over all the attributes of $J \setminus I$. Therefore one can apply the standard work of computing multi-attribute aggregation, a lattice structure. Without ambiguity, we use the notation $\mathcal{X}_I^J = \{X | I \subseteq X \subseteq J\}$ to represent both the set of itemsets and their corresponding lattice structure. An example of lattice \mathcal{X}_c^{abc} is shown in Fig. 3.

Deriving Pattern Support Consider two itemsets $I \subset J$, if the support of the lattice nodes of \mathcal{X}_I^J is available, one is able to derive the support of the pattern p of the form, $p = I(\overline{J \setminus I})$, according to the inclusion-exclusion principle:

$$T_{\mathcal{D}}(I(\overline{J \setminus I})) = \sum_{X \in \mathcal{X}_I^J} (-1)^{|X \setminus I|} T_{\mathcal{D}}(X)$$

Example 3: As illustrated in Fig. 3, given all the support of \mathcal{X}_c^{abc} w.r.t. $\mathcal{D}_s(12, 8)$, the support of pattern $p = \overline{abc}$ can be derived as 1.

Estimating Itemset Support Since the support of any pattern is non-negative, according to the inclusion-exclusion principle, if the support of the itemsets $\mathcal{X}_I^J \setminus \{J\}$ is available, one is able to bound the support of J as follows:

$$\begin{cases} T_{\mathcal{D}}(J) \leq \sum_{I \subset X \subset J} (-1)^{|J \setminus X|+1} T_{\mathcal{D}}(X) & |J \setminus I| \text{ odd} \\ T_{\mathcal{D}}(J) \geq \sum_{I \subset X \subset J} (-1)^{|J \setminus X|+1} T_{\mathcal{D}}(X) & |J \setminus I| \text{ even} \end{cases}$$

Example 4: In Fig. 3, given the support of c , ac and bc w.r.t. $\mathcal{D}_s(12, 8)$, one is able to establish the lower/upper bound for $T_{\mathcal{D}_s(12, 8)}(abc)$ as [2,5].

When the bound is tight, i.e., the lower bound equals to the upper bound, one can exactly determine the actual support. This technique is used in [16] to mine non-derivable frequent itemsets. In our context, an adversary can leverage this technique to exploit the privacy breaches existing in mining output.

B. Intra-Window Breaches

In a stream mining system without output privacy protection, the released frequent itemsets over one specific window may contain the intra-window breaches, which can be exploited by an adversary through the technique of deriving pattern support, as shown in Example 3.

Formally, if J is a frequent itemset, according to the Apriori rule, all $X \subseteq J$ are frequent, which are supposed to be reported with their support, so the information is available to compute the support of pattern $p = I(\overline{J \setminus I})$ for all $I \subset J$. This also implies that the number of breaches to be checked is potentially exponential in terms of the number of items.

Even if the information of J is unavailable, i.e., \mathcal{X}_I^J is incomplete to infer $p = I(\overline{J \setminus I})$, one could possibly apply the technique of estimating itemset support first to complete some missing ‘‘mosaics’’, then derive vulnerable pattern. In fact, the itemsets under estimation themselves could be vulnerable.

C. Inter-Window Breaches

In stream mining, the output of the previous window can be leveraged to infer the vulnerable patterns within the current one, and vice versa, even though no vulnerable patterns can be inferred from the output of each window per se.

Example 5: Consider the two windows $\mathcal{D}_s(11, 8)$ and $\mathcal{D}_s(12, 8)$ shown in Fig. 3. Assuming $C = 4$, and $K = 1$, then in $\mathcal{D}_s(11, 8)$, no \mathcal{P}_{hv} exists. In $\mathcal{D}_s(12, 8)$, the itemset abc is inaccessible (shown as a dashed box). From the available

information of $\mathcal{D}_s(12, 8)$, the best guess about abc is $[2, 5]$, as discussed in Example 4. Clearly, this bound is not tight enough to estimate that the pattern \overline{abc} is \mathcal{P}_{hv} . Thus both windows are currently immune to intra-window inference attack.

However, if one is able to derive that the support of abc decreases by 1 between $\mathcal{D}_s(11, 8)$ and $\mathcal{D}_s(12, 8)$, then based on the information released in $\mathcal{D}_s(11, 8)$, which is $T_{\mathcal{D}_s(11,8)}(abc) = 4$, the exact value of abc in $\mathcal{D}_s(12, 8)$ can be inferred, and the $\mathcal{P}_{hv} \overline{abc}$ is uncovered.

The main idea of inter-window inference is: (i) Estimating the transition of the support of certain itemsets from the previous window to the current one, using the technique of estimating itemset support; (ii) Uncovering vulnerable patterns, using the technique of deriving pattern support. Due to the space limit, a detailed discussion is referred to our technical report [17].

V. OUTPUT PRIVACY PROTECTION

A. Design Consideration of Solutions

Alternative to the reactive detecting-then-removing strategy, we propose to use a proactive approach to deal with these two types of attacks in a uniform way. Our approach is motivated by two key observations: (i) In many mining applications, users do not expect the exact support of frequent itemsets. Rather they care more about their semantic relationships in terms of the support, e.g., the ranking or the ratio of their support values. Thus it is tolerable to trade some precision of the support of frequent itemsets for the output privacy guarantee, provided that such desired output utility is maintained; (ii) Both intra- and inter-window inferences are based on the inclusion-exclusion principle, which involves multiple frequent itemsets. If we introduce some trivial uncertainty into each frequent itemset, the resulting inferred pattern can have considerable uncertainty, due to the accumulative property of uncertainty.

Based on these two observations, we propose **Butterfly**, a light-weighted output privacy preservation scheme based on random perturbation.

B. Mining Output Perturbation

Data perturbation refers to the process of modifying confidential data while preserving its utility for intended applications [9]. This is arguably the most important technique used so far for protecting original input data.

In our scheme however, we employ perturbation to inject uncertainty into the mining output. The perturbation over output pattern is significantly different from that over input data. In input perturbation, the data utility is defined on the overall statistical characteristics of the dataset. The distorted data is fed as input into the following mining algorithms. There is usually no utility constraints for individual data value. While in output perturbation, the perturbed results are directly presented to the end-user, and the data utility is defined over each individual value.

Specifically, there are two types of utility constraints for the perturbed results: (i) Each reported value should have enough accuracy, i.e., the perturbed value should not deviate from the actual value too far; (ii) The semantic relationships among the

results should be preserved to the maximum extent, e.g., the order or the ratio of the support values of frequent itemsets. There is non-trivial tradeoff among these utility metrics. To our best knowledge, no previous work has considered such multiple tradeoff in mining output perturbation.

C. Basic Butterfly Approach

On releasing the mining output of a stream window, one perturbs the support of each frequent itemset X , $T(X)$ ¹ by adding a random variable r_X drawn from a discrete uniform distribution over integers within an interval $[l_X, u_X]$. The sanitized support $T'(X) = T(X) + r_X$ is therefore a random variable, which can be specified by its bias $\beta(X)$ and variance $\sigma^2(X)$. Intuitively, the bias indicates the difference between the expected value $E[T'(X)]$ and $T(X)$, and the variance represents the average deviation of $T'(X)$ from $E[T'(X)]$.

While this operation is simple, the setting of $\beta(X)$ and $\sigma^2(X)$ is non-trivial, in order to achieve both sufficient privacy protection and utility guarantee, which is the focus of our following discussion. At this moment, just note that compared with $T(X)$, r_X is non-significant, i.e., $|r_X| \ll T(X)$.

Given the basic characteristics of the perturbation, we further define the metrics to measure the precision for outputted frequent itemsets, and the privacy guarantee for vulnerable patterns.

1) *Precision Measure*: Under the perturbation, the precision loss of each frequent itemset X can be measured by the *mean square error* (mse) of the perturbed support $T'(X)$: $\text{mse}(X) = E[(T'(X) - T(X))^2] = \sigma^2(X) + \beta^2(X)$.

Intuitively, $\text{mse}(X)$ indicates the average deviation of the perturbed support $T'(X)$ w.r.t. the actual value $T(X)$. A smaller mse implies higher precision of the frequent itemset. Also it is clear that the precision loss should depend on the actual support. A mse of 5 for frequent itemset I with $T(I) = 100$ may indicate sufficient accuracy, while the same mse for itemset J with $T(J) = 5$ may render the output of little value. Therefore, we have the following precision metric:

Definition 3 (Precision Degradation): For a frequent itemset X , its precision degradation $\text{pred}(X)$ is defined as the relative mean squared error of $T'(X)$:

$$\text{pred}(X) = \frac{\sigma^2(X) + \beta^2(X)}{T^2(X)}$$

2) *Privacy Measure*: Distorting the original support of frequent itemsets is only a part of the story, it is necessary to ensure that the distortion could not be filtered out. Therefore one needs to consider the power of the adversary in estimating the support of vulnerable patterns through the protection.

Without loss of generality, suppose that the adversary desires to estimate the support of pattern p of the form $I(\overline{J} \setminus \overline{I})$, and has full access to the sanitized support $T'(X)$ for all $X \in \mathcal{X}_I^J$. The privacy protection should be measured by the error of the adversary's estimation of the support of p , denoted as $T''(p)$. We will discuss this estimation from an adversary's perspective. Along the discussion, we will show

¹In the presentation below, when the context is clear, we omit the referred database \mathcal{D} in the notations.

how various prior knowledge the adversary possesses may impact the estimation.

From the adversary’s view, of each $X \in \mathcal{X}_f^J$, its actual support $T(X) = T'(X) - r_X$, is a variable with a discrete uniform distribution over the interval $[T'(X) - u_X, T'(X) - l_X]$ (since $|r_X| \ll T(X)$, this is a bounded distribution over positive integers), and has variance of $\sigma^2(X)$.

Prior Knowledge 1: The support values of different frequent itemsets are related by a set of inequations, derived from the inclusion-exclusion principle.

For the given frequent itemset-interval pairs, the adversary may attempt to apply these inequalities to tighten the intervals, therefore obtaining better estimation of the support. A key question she needs to answer is: for the modified interval(s), does there exist a database \mathcal{D} that satisfies the constraints of these itemset-interval pairs? It is called the itemset frequency satisfiability problem (FREQSAT), which however as shown in [18] is equivalent to the probabilistic satisfiability problem (pSAT), i.e., NP-Complete. This indicates that the inequality relationships among itemsets can hardly be leveraged to improve the estimation of a single itemset, and one can approximately consider frequent itemsets as independent in measuring the adversary’s power.

The actual support of p , $T(p)$ from the adversary’s view, is also a random variable. The mse of the estimation $T''(p)$ is defined as $mse(p) = E[(T''(p) - T(p))^2]$ w.r.t. $T(p)$. One has the following lemma (due to the space limit, the proofs of all the lemmas are referred to our technical report [17]):

Lemma 1: Given the distribution $f(x)$ of a random variable x , the mean square error of an estimation e of x , $mse(e)$ reaches its minimum value $Var[x]$, when $e = E[x]$.

In our scenario, $mse(p)$ is minimized when $T''(p) = E[T(p)]$, this is the best guess the adversary can achieve (note that the optimality is defined in terms of average estimation error, not the semantics, e.g., $E[T(p)]$ could possibly be negative). In this worst case (best case for the adversary), the $mse(p)$ equals to the variance of $T(p)$, which corresponds to the lower bound of the estimation error.

Considering the analysis regarding Prior Knowledge 1, and the fact that $T(p) = \sum_{X \in \mathcal{X}_f^J} (-1)^{|X \setminus I|} T(X)$ (a linear combination), the variance of $T(p)$ can be approximated by the sum of the variance of all involved $T(X)$, hence $mse(p) = \sum_{X \in \mathcal{X}_f^J} \sigma^2(X)$. Intuitively, a larger $mse(p)$ indicates a more significant error in estimating $T(p)$ by the adversary, and better privacy protection.

It is also noted that the privacy guarantee should depend on the actual value $T(p)$: if $T(p)$ is close to 0, trivial variance makes it hard for the adversary to estimate if $T(p)$ is zero or not, i.e., if such vulnerable pattern exists. Such “zero-indistinguishability” decreases as $T(p)$ grows. Therefore, we define the privacy metric for a vulnerable pattern p as the ratio of the variance of $T(p)$ from adversary’s view and the square of the actual support $T(p)$.

Definition 4 (Privacy Guarantee): For a vulnerable pattern p , its privacy guarantee $\text{prig}(p)$ is defined as its *relative estimation error*:

$$\text{prig}(p) = \frac{\sum_{X \in \mathcal{X}_f^J} \sigma^2(X)}{T^2(p)}$$

Prior Knowledge 2: The sanitized support of the same frequent itemsets may be published in consecutive windows.

Since our protection is based on independent random perturbation, if the same support value is repeatedly perturbed and published in multiple windows, the adversary can potentially improve the estimation by averaging the observed outcomes, according to the law of large numbers. To block this type of attack, once the support of a frequent itemset is perturbed, one keeps publishing this sanitized value if the actual support remains the same in consecutive windows.

Prior Knowledge 3: The adversary may have access to other forms of prior knowledge, e.g., the published statistics of the dataset, the support of the top- k frequent itemset or the ones near the threshold C , etc.

All these various forms of prior knowledge can be captured by the notion of *knowledge point*: a knowledge point is a specific frequent itemset X , for which the adversary has an average error less than $\sigma^2(X)$ in estimating $T(X)$. Our definition of privacy guarantee (prig) can readily incorporate this notion, by simply replacing the corresponding variance $\sigma^2(X)$ with the smaller estimation error.

3) *Effectiveness:* In summary, the effectiveness of our privacy protection method is evaluated in terms of its resilience against both intra- and inter-window inferences over stream mining output. We note three key implications.

First, the uncertainty of involved frequent itemsets are accumulated in the inferred vulnerable patterns. Moreover, more complicated inferences (i.e., harder to be detected) face higher uncertainty.

Second, the actual support of a vulnerable pattern is usually small (only a unique or less than K records match vulnerable patterns), hence adding trivial uncertainty can make it hard to tell the existence of such pattern in the dataset.

Third, the inter-window inferences follow a two-staged strategy, i.e., first deducing the transition between contingent windows, then inferring the vulnerable patterns. The uncertainty associated with both stages may result in estimation of even lower quality.

D. Tradeoff between Precision and Privacy

In our Butterfly approach, the tradeoff between privacy protection and output utility can be flexibly adjusted by the variance and bias setting of each frequent itemset. Specifically, the variance controls the overall balance between privacy and utility, while the bias gives a finer control over the balance between precision and other utility metrics, as we will show in the next section. Here we focus on the setting of variance. Intuitively, smaller variance leads to higher precision of the output, however also decreases the uncertainty of the inferred vulnerable patterns, therefore less privacy guarantee.

To ease the discussion, we assume that all the frequent itemsets are associated with the same variance σ^2 and bias β . In the next section when semantic constraints are taken into consideration, we will remove this simplified treatment, and develop more sophisticated setting scheme.

Let C denote the minimum support for frequent itemsets. From the definition of precision measure, it can be derived that for each frequent itemset X , its precision degradation $\text{pred}(X) \leq (\sigma^2 + \beta^2)/C^2$, because $T(X) \geq C$. Let $P_1(C) =$

$(\sigma^2 + \beta^2)/C^2$, i.e., an upper bound of the precision loss of the frequent itemsets. Meanwhile for a vulnerable pattern $p = I(\overline{J} \setminus I)$, it can be proved that its privacy guarantee $\text{prig}(p) \geq (\sum_{X \in \mathcal{X}_I^J} \sigma^2)/K^2 \geq (2\sigma^2)/K^2$, because $T(p) \leq K$ and the inference attack involves at least two frequent itemsets. Let $P_2(C, K) = (2\sigma^2)/K^2$, i.e., a lower bound of the privacy guarantee of the inferred vulnerable patterns.

P_1 and P_2 provide a convenient representation to control the tradeoff between precision and privacy protection. Specifically, setting an upper bound ϵ over P_1 guarantees sufficient accuracy of the reported frequent itemsets; While setting a lower bound δ over P_2 provides enough privacy protection for the inferred vulnerable patterns. Therefore one can specify the requirement of the output precision and privacy guarantee as a pair of parameters (ϵ, δ) , where $\epsilon, \delta > 0$. That is the setting of β and σ^2 should satisfy $P_1(C) \leq \epsilon$ and $P_2(C, K) \geq \delta$, as

$$\sigma^2 + \beta^2 \leq \epsilon C^2 \quad (1)$$

$$\sigma^2 \geq \delta K^2/2 \quad (2)$$

To make these two inequations compatible, it should be satisfied that $\epsilon/\delta \geq K^2/(2C^2)$. The term ϵ/δ is called *precision-privacy ratio (ppr)*, which tends to indicate the precision loss for user-specified privacy requirement. When the precision is a major concern, one can set *ppr* as its minimum value $K^2/(2C^2)$ for given K and C , resulting in the highest precision for the given privacy requirement. The minimum *ppr* also implies that $\beta = 0$ and the two parameters ϵ and δ are coupled. We refer to the perturbation scheme with the minimum *ppr* as our basic Butterfly approach. In the following section, we will relax this condition, and take into consideration other utility metrics in addition to precision.

VI. OPTIMIZED OUTPUT UTILITY

The basic Butterfly approach treats all the frequent itemsets uniformly (the same bias $\beta = 0$) without taking consideration of their semantic relationships. Though easy to implement and effective against inferences, this simple scheme may easily violate these semantic constraints directly related to the specific applications of the mining output. In this section, we refine the basic scheme by taking semantic constraints into our map, and develop constraint-aware Butterfly approach. Given the precision and privacy requirement (ϵ, δ) , our optimized version preserves the semantics to the maximum extent.

We specifically consider two types of utility-related semantic constraints, *absolute order* and *relative frequency*. By absolute order, we refer to the ranking of frequent itemsets according to their support. In many applications, users pay considerable attention to the absolute order, e.g., querying the top-ten popular purchase patterns. By relative frequency, we mean the ratio of the support of two frequent itemsets. In certain applications, users care much about the relative frequency, e.g., computing the confidence in mining association rules.

To help model the order and ratio of the support of itemsets, we first introduce the concept of *frequency equivalence class*:

Definition 5 (Frequency Equivalence Class): A frequency equivalence class (FEC) is a set of frequent itemsets, with the same support value. Given a FEC fec , we define its support $T(fec)$ as that of any of its members.

A set of frequent itemsets can be partitioned into a set of disjoint and strictly ordered FECs, based on their support. We say that two FECs, fec_i and fec_j , follow a partial order of $fec_i < fec_j$ if $T(fec_i) < T(fec_j)$. Without loss of generality, we assume that the given set of FECs are sorted according to their support, i.e., $T(fec_i) < T(fec_j)$ for $i < j$.

In complying with the constraints of order or ratio, the equivalence of itemsets in a FEC should be maximally preserved in the perturbed output. Therefore in the optimized Butterfly schemes, the perturbation is performed over each of the FECs, instead of each specific itemset.

Clearly, this revision does not affect the privacy guarantee, considering the fact that the inference of a vulnerable pattern involves at least two frequent itemsets with different support, i.e., at least two FECs.

A. Order Preservation

When the absolute order of itemset frequency is an important concern, the perturbation over all FECs can not be uniform, since that would easily render the inversion of the orders of two FECs, especially when their support values are close. To model the probability of inversion, we first introduce the concept of *uncertainty region* of a FEC.

Definition 6 (Uncertainty Region): The uncertainty region of a FEC fec_i is defined as the set of values that its perturbed support $T'(fec_i)$ can take: $\{x | Pr[T'(fec_i) = x] > 0\}$.

For instance, when adding a random variable drawn from the integers from the interval $[l, u]$ to fec_i , the uncertainty region of fec_i is all the integers in the interval $[T(fec_i) + l, T(fec_i) + u]$.

1) *Minimizing Inversion Probability:* Below we formally define the problem of preserving absolute order. To simplify the notations, we use the following short version: $t_i = T(fec_i)$, $t'_i = T'(fec_i)$, and let β_i denote the bias setting for fec_i .

Without loss of generality, consider two FECs fec_i, fec_j with $t_i < t_j$. The order of fec_i and fec_j can be possibly inverted if their uncertainty regions overlap, that is $Pr[t'_i \geq t'_j] > 0$, and larger $Pr[t'_i \geq t'_j]$ indicates higher probability that this inversion occurs.

One can minimize this inversion probability $Pr[t'_i \geq t'_j]$ by adjusting β_i and β_j . However, this adjustment is bounded by the requirement of precision and privacy specified in Inequations 1 and 2 as introduced in Section 5.3. We therefore define the concept of *maximum adjustable bias*:

Definition 7 (Maximum Adjustable Bias): Given a FEC fec_i , its bias is allowed to be adjusted within the interval of $[-\beta_i^m, \beta_i^m]$, β_i^m is called the maximum adjustable bias, given ϵ and δ , which is defined as $\beta_i^m = \lfloor \sqrt{\epsilon t_i^2 - \delta K^2/2} \rfloor$, derived from Inequations 1 and 2.

The problem of preserving absolute order can therefore be formalized as: Given a set of FECs $\mathcal{FEC} = \{fec_1, \dots, fec_n\}$ one finds the optimal bias setting for each FEC fec_i within its maximum adjustable bias $[-\beta_i^m, \beta_i^m]$, to minimize the sum of pairwise inversion probability: $\min \sum_{i < j} Pr[t'_i \geq t'_j]$.

We now show how to compute $Pr[t'_i \geq t'_j]$ in our setting. For a discrete uniform distribution over the interval $[l, u]$, $\alpha = u - l$ is called the length of the region. The variance of this distribution is $\sigma^2 = ((\alpha + 1)^2 - 1)/12$. According to

Inequalities 2 in Section V, one has $\alpha = \lceil \sqrt{1 + 6\delta K^2} \rceil - 1$. Let d_{ij} be the distance between the estimators $e_i = t_i + \beta_i$ and $e_j = t_j + \beta_j$ of fec_i and fec_j : $d_{ij} = e_j - e_i$. The optimization problem above can be simplified as: $\sum_{i < j} (\alpha + 1 - d_{ij})^2$ for $d_{ij} \geq 0$ (the details are referred to [17]).

Note that the discussion so far has not considered the characteristics of each FEC, such as the number of its members. The inversion of two FECs containing five frequent itemsets each, is much more serious than that of two FECs with only one member respectively. Quantitatively, let s_i be the number of members in fec_i , the inversion of fec_i and fec_j means the ordering of $s_i + s_j$ itemsets is disturbed.

Therefore, our aim is to solve the following weighted optimization problem:

$$\begin{aligned} \min \quad & \sum_{i < j} (s_i + s_j)(\alpha + 1 - d_{ij})^2 \\ \text{s.t.} \quad & d_{ij} = \begin{cases} \alpha + 1 & e_j - e_i \geq \alpha + 1 \\ e_j - e_i & e_j - e_i < \alpha + 1 \end{cases} \\ & \forall i < j, e_i \leq e_j \quad \forall i, e_i \in \mathbb{Z}^+, |e_i - t_i| \leq \beta_i^m \end{aligned}$$

This is a quadratic integer programming (QIP) problem, with piecewise cost function. In general, quadratic programming is NP-hard, even without integer constraints [19]. Instead of applying off-the-shelf quadratic optimization tools, we are more interested in on-line algorithms that can flexibly trade accuracy for efficiency. Following we present such a solution based on dynamic programming.

2) *A Near Optimal Solution*: Although this optimization problem in the general setting is NP-hard, by relaxing the constraint that $\forall i < j, e_i \leq e_j$ to $e_i < e_j$, one can construct the following optimal substructure property, leading to an efficient dynamic programming solution.

Lemma 2: Assume that the bias setting of the last α FECs $\{fec_{n-\alpha+1}:fec_n\}^2$ are fixed as $\{\beta_{n-\alpha+1}^*:\beta_n^*\}$ respectively, and $\{\beta_1^+:\beta_{n-\alpha}^+\}$ are optimal w.r.t. $\{fec_1:fec_n\}$, then for given $\{\beta_{n-\alpha}^+,\beta_{n-\alpha+1}^*:\beta_{n-1}^*\}$, $\{\beta_1^+:\beta_{n-\alpha-1}^+\}$ must be optimal w.r.t. $\{fec_1:fec_{n-1}\}$.

Based on this optimal sub-structure, we propose a dynamic programming solution, which adds FECs sequentially according to their order. Let $C_{n-1}(\beta_{n-\alpha}:\beta_{n-1})$ represent the achievable minimum cost by adjusting $\{fec_1:fec_{n-\alpha-1}\}$ with the last α FECs fixed as $\{\beta_{n-\alpha}:\beta_{n-1}\}$, and let c_{ij} denote $(s_i + s_j)(\alpha + 1 - d_{ij})^2$. When adding fec_n , the minimum cost $C_n(\beta_{n-\alpha+1}:\beta_n)$ is computed using the rule:

$$C_n(\beta_{n-\alpha+1}:\beta_n) = \min_{\beta_{n-\alpha}} C_{n-1}(\beta_{n-\alpha}:\beta_{n-1}) + \sum_{i=n-\alpha}^{n-1} c_{in}$$

The optimal setting is the one with the global minimum value among all the combinations of $\{\beta_{n-\alpha+1}:\beta_n\}$.

Computation Complexity Let β_{max} be the maximum value of maximum adjustable bias of all FECs: $\beta_{max} = \max_i \beta_i^m$, the time and space complexity of this scheme are both bounded by $O(\beta_{max}^\alpha n)$, i.e., $O(n)$ in terms of total number of FECs. The empirical time/space complexity is usually much lower than this bound, given the facts that under the constraint $\forall i <$

$j, e_i < e_j$, a number of combinations are invalid, and β_{max} is an over-estimation of the average maximum adjustable bias.

To even lower the complexity, on adding a FEC, one can approximately consider its intersection with its previous γ FECs, instead of α ones ($\gamma < \alpha$). This approximation is close to the exact solution when the distribution of FECs is not extremely dense, which is usually the case, and verified by our experiments. The complexity is then bounded by $O(\beta_{max}^\gamma n)$. By adjusting γ , one can flexibly control the tradeoff between accuracy and efficiency.

Algorithm 1: Order Preserving Bias Setting

```

Input:  $\{t_i, \beta_i^m\}$  for each  $fec_i \in \mathcal{FEC}$ ,  $\alpha, \gamma$ .
Output:  $\beta_i$  for each  $fec_i \in \mathcal{FEC}$ .
begin
  /* initialization */;
  for  $\beta_1 = -\beta_1^m : \beta_1^m$  do
     $C_1(\beta_1) = 0$ ;
  for  $i = 2 : \gamma$  do
    for  $\beta_i = -\beta_i^m : \beta_i^m$  do
      /*  $e_i < e_j$  */;
      if  $\beta_i + t_i > \beta_{i-1} + t_{i-1}$  then
         $C_i(\beta_1 : \beta_i) = C_{i-1}(\beta_1 : \beta_{i-1}) + \sum_{j=1}^{i-1} c_{ji}$ ;
  /* dynamic programming */;
  for  $i = \gamma + 1 : n$  do
    for  $\beta_i = -\beta_i^m : \beta_i^m$  do
      if  $\beta_i + t_i > \beta_{i-1} + t_{i-1}$  then
         $C_i(\beta_{i-\gamma+1} : \beta_i) = \min_{\beta_{i-\gamma}} C_{i-1}(\beta_{i-\gamma} : \beta_{i-1}) + \sum_{j=i-\gamma}^{i-1} c_{ji}$ ;
  /* find the optimal setting */;
  find the minimum  $C_n(\beta_{n-\gamma+1} : \beta_n)$ ;
  backtrack and output  $\beta_i$  for each  $fec_i \in \mathcal{FEC}$ ;
end

```

The order-preserving bias setting scheme is shown in Algorithm 1: After initializing the cost function for the first γ FECs, one computes the cost function for each newly added FEC, via the dynamic programming scheme. The optimal configuration is the one with the global minimum value $C_n(\beta_{n-\gamma+1}:\beta_n)$.

B. Ratio Preservation

In a number of applications, it is desirable to maximally preserve the relative frequency (ratio) information during the output perturbation process. Similar to the order preservation optimization, one can achieve this by carefully adjust the bias setting of FECs.

1) *Maximizing $(k, 1/k)$ Probability of Ratio*: Without loss of generality, consider two FECs fec_i and fec_j . To preserve the relative frequency of fec_i and fec_j , one is interested in making the ratio of the perturbed support t'_i/t'_j appear in a tight neighborhood of the actual value t_i/t_j with high probability.

Definition 8 (($k, 1/k$) Probability): For the ratio $\frac{t'_i}{t'_j}$ of two random variables t'_i and t'_j , its $(k, 1/k)$ probability $Pr_{(k,1/k)} \left[\frac{t'_i}{t'_j} \right]$ is defined as $Pr \left[k \frac{t_i}{t_j} \leq \frac{t'_i}{t'_j} \leq \frac{1}{k} \frac{t_i}{t_j} \right]$, where $k \in (0, 1)$, is an indicator of the tightness of the approximation area, the larger k , the tighter the interval.

This $(k, 1/k)$ probability can quantitatively measure the impact of the perturbation over the ratio. Therefore the problem of ratio preservation is formalized as follows:

²Without ambiguity, following we use $\{x_i : x_j\}$ as a short version of $\{x_i, x_{i+1}, \dots, x_{j-1}, x_j\}$.

$$\begin{aligned} \max \quad & \sum_{i < j} Pr^{(k, 1/k)} \left[\frac{t'_i}{t'_j} \right] \\ \text{s.t.} \quad & \forall i, e_i \in \mathbb{Z}^+, |e_i - t_i| \leq \beta_{mi} \end{aligned}$$

In the case of discrete uniform distribution, the $(k, 1/k)$ probability of the ratio of two random variables is a non-linear piecewise function, resulting in a non-linear integer optimization problem, which in general sense, is **NP-hard**, even without integer constraints. Below we present an efficient approximate solution that can find the near-optimal configuration with complexity linear in terms of the number of FECs.

2) *A Near Optimal Solution:* In order to maximize the $(k, 1/k)$ probability, one can alternatively minimize the probability $Pr \left[\frac{t'_i}{t'_j} \geq \frac{1}{k} \frac{t_i}{t_j} \right] + Pr \left[\frac{t'_j}{t'_i} \geq \frac{1}{k} \frac{t_j}{t_i} \right]$. With Markov's Inequality, one knows that the probability $Pr \left[\frac{t'_i}{t'_j} \geq \frac{1}{k} \frac{t_i}{t_j} \right]$ is bounded as $Pr \left[\frac{t'_i}{t'_j} \geq \frac{1}{k} \frac{t_i}{t_j} \right] \leq k \frac{t_i}{t_j} E \left[\frac{t'_i}{t'_j} \right]$, which can be further simplified as (the details are referred to [17]):

$$\min \quad \frac{t_j e_i}{t_i e_j} + \frac{t_i e_j}{t_j e_i}$$

Bottom-up Bias Setting Assuming that e_i is fixed, by differentiating the expression w.r.t. e_j , and setting the derivative as zero, one gets the solution of e_j as $e_j/e_i = t_j/t_i$, i.e., $\beta_j/\beta_i = t_j/t_i$.

Following this solution is our bottom-up bias setting scheme: for each FEC fec_i , its bias β_i should be set in proportion to its support t_i . Note that the larger $t_i + \beta_i$ compared with α , the more accurate the approximation applied here [17], therefore β_i should be set as its maximum possible value. The whole scheme is shown in Algorithm 2.

Algorithm 2: Ratio Preserving Bias Setting

```

Input:  $\{t_i\}$  for each  $fec_i \in \mathcal{FEC}$ ,  $\epsilon, \delta, K$ .
Output:  $\beta_i$  for each  $fec_i \in \mathcal{FEC}$ .
begin
  /* setting of the minimum FEC */;
  set  $\beta_1 = \lfloor \sqrt{\epsilon t_1^2 - \delta K^2 / 2} \rfloor$ ;
  /* bottom-up setting */;
  for  $i = 2 : n$  do
    [ set  $\beta_i = \lfloor \beta_{i-1} \frac{t_i}{t_{i-1}} \rfloor$ ;
  end

```

Further, we have the following lemma to show that for each FEC fec_i , β_i falls within the interval $[-\beta_i^m, \beta_i^m]$.

Lemma 3: For two FECs fec_i and fec_j with $t_i < t_j$, if the setting of β_i is feasible for fec_i , namely within the interval $[-\beta_i^m, \beta_i^m]$, then the setting $\beta_j = \beta_i \frac{t_j}{t_i}$ is feasible for fec_j .

C. A Hybrid Scheme

While order preserving (OP) and ratio preserving (RP) bias settings achieve the maximum utility at their ends, in some applications where both semantic relationships are important, it is desired to balance the two factors in order to achieve the overall optimal quality.

We therefore develop a hybrid scheme that takes advantage of the two approaches, and can flexibly trade between order and ratio preservation quality.

Specifically, for each FEC fec , let β_{op} and β_{fp} denote its bias setting obtained by the order preserving and frequency

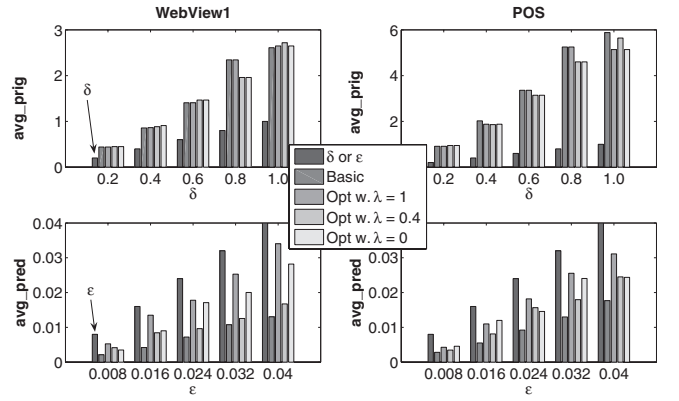


Fig. 4. Average privacy guarantee (avg_prig) and precision degradation (avg_pred).

preserving approaches respectively. We incorporate them using a linear combination: $\beta = \lambda \beta_{op} + (1 - \lambda) \beta_{fp}$.

The parameter $\lambda \in [0, 1]$, controls the balance between the two quality metrics. Intuitively, larger λ tends to assign more importance over absolute order, but less over relative frequency, and vice versa. Especially, the order-preserving and ratio-preserving schemes are the special cases when $\lambda = 1$ and 0 respectively.

VII. EXPERIMENTAL ANALYSIS

In this section, we investigate the effectiveness and efficiency of the proposed Butterfly approaches. Specifically, the experiments are designed to measure the following three properties: (i) privacy guarantee: the effectiveness against both intra-/inter-window inferences; (ii) result utility: the degradation of the output precision, the order and ratio preservation, and the tradeoff among these utility metrics; (iii) efficiency: The time taken to perform our approaches. We start by describing the datasets and the setup of the experiments.

A. Experimental Setting

We tested our approaches over two real-life datasets. The first one is BMS-WebView-1, which contains a few months of clickstream data from an e-commerce web site, and the second one is BMS-POS, which contains several years of point-of-sale from a large number of electronic retailers. Both datasets have been used in stream frequent pattern mining [20].

We built our Butterfly prototype on top of *Moment* [20], a stream frequent pattern mining framework, which finds closed frequent itemsets over a sliding window model. By default, the minimum support C and vulnerable support K are set as 25 and 5 respectively, and the window size is set as $2K$. Note that the setting here is designed to test the effectiveness of our approach with high ratio of vulnerable/minimum threshold (K/C). All the experiments are performed over a workstation with Intel Xeon 3.20GHz and 4GB main memory, running Red Hat Linux 9.0 operating system. The algorithm is implemented in C++ and compiled using g++ 3.4.

B. Experimental Results

To provide an in-depth understanding of our output privacy preservation schemes, we evaluated four different versions of our approach: the basic version, the optimized version with λ

= 0, 0.4 and 1 respectively, over both real datasets. Note that $\lambda = 0$ corresponds to the ratio preserving scheme, while $\lambda = 1$ corresponds to the order preserving one.

Privacy and Precision To evaluate the effectiveness of output privacy protection of our approach, one needs to find all potential privacy breaches in the mining output. This is done by running an analysis program over the results returned by the mining algorithm, and finding all possible vulnerable patterns that can be inferred through either intra-window or inter-window inferences.

Concretely, given a stream window, let \mathcal{P}_{hv} denote all the hard vulnerable patterns that are inferable from the mining output. After the perturbation, we evaluate the square relative deviation of the inferred value and the estimator for each $p \in \mathcal{P}_{hv}$ for 100 continuous windows. we use the following average privacy (avg_prig) metric to measure the effectiveness of the privacy preservation:

$$\text{avg_prig} = \sum_{p \in \mathcal{P}_{hv}} \frac{(T'(p) - E[T'(p)])^2}{T^2(p)|\mathcal{P}_{hv}|}$$

The decrease in precision of the output is measured by the average precision degradation (avg_pred) of all frequent itemsets \mathcal{I} :

$$\text{avg_pred} = \sum_{I \in \mathcal{I}} \frac{(T'(I) - T(I))^2}{T^2(I)|\mathcal{I}|}$$

In this set of experiments, we fix the precision-privacy ratio $\epsilon/\delta = 0.04$, and measure avg_prig and avg_pred for different setting of ϵ (δ).

Specifically, the two plots in the top tier of Fig. 4 show that as the value of δ increases, all four versions of the Butterfly approaches provide similar amount of average privacy protection for both datasets, all above the minimum privacy guarantee δ . The two plots in the lower tier show that as σ increases from 0 to 0.04, the output precision decreases, however all four versions of the Butterfly approaches have average precision degradation below the the system-supplied maximum threshold ϵ . Among them, the basic Butterfly offers the lowest precision loss, which can be explained by the fact that the basic version trades the minimum precision loss for the privacy guarantee, without considering semantic constraints.

Order and Ratio For given privacy and precision requirement (ϵ , δ), we measure the effectiveness of our approaches in preserving order and ratio of frequent itemsets.

The quality of order preservation is evaluated by the proportion of the order preserved pairs over all possible pairs, referred to as the rate of order preserved pairs (ropp):

$$\text{ropp} = \frac{\sum_{I, J \in \mathcal{I} \cap T(I) \leq T(J)} \mathbf{1}_{T'(I) \leq T'(J)}}{C_{|\mathcal{I}|}^2}$$

where $\mathbf{1}$ is the indicator function, returning 1 if the condition is met, and 0 otherwise.

Analogously, the quality of ratio preservation is evaluated by the fraction of the number of $(k, 1/k)$ probability preserved pair over the number of possible pairs, referred to as the rate of ratio preserved pairs (rpp) (k is set 0.95 in all the experiments):

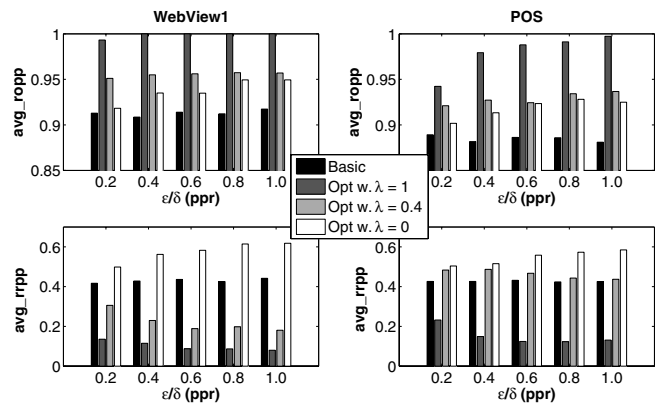


Fig. 5. Average order preservation (avg_ropp) and ratio preservation (avg_rrpp).

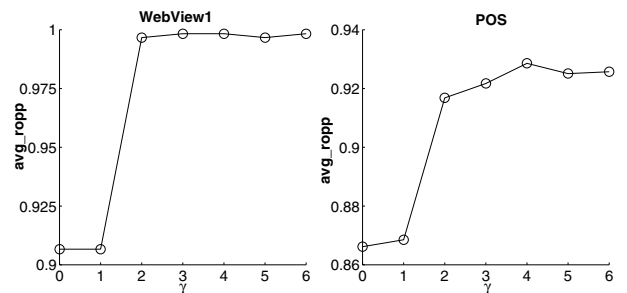


Fig. 6. Average rate of order-preserved pairs versus the setting of γ .

$$\text{rpp} = \frac{\sum_{I, J \in \mathcal{I} \cap T(I) \leq T(J)} \mathbf{1}_{k \frac{T'(I)}{T'(J)} \leq \frac{T(I)}{T(J)} \leq \frac{1}{k} \frac{T(I)}{T(J)}}}{C_{|\mathcal{I}|}^2}$$

In this set of experiments, we vary the precision-privacy ratio ϵ/δ for fixed $\delta = 0.4$, and measure the ropp and rpp for four versions of the Butterfly approaches (the parameter $\gamma = 2$ in all the experiments), as shown in Fig. 5.

As predicted by our theoretical analysis, the order-preserving ($\lambda = 1$) and ratio-preserving ($\lambda = 0$) bias settings are quite effective, both outperform all other approaches at their ends. The ropp and rpp increase as the ratio of ϵ/δ grows, due to the fact that larger ϵ/δ offers more adjustable bias therefore leading to better quality.

It is also noticed that order-preserving scheme has the worst performance in the term of avg_rrpp, even worse than the basic approach. This can be explained by that in order to distinguish overlapping FECs, the order-preserving scheme may significantly disturb the ratio of pairs of FECs. In all these cases, the hybrid scheme $\lambda = 0.4$ achieves the second best in terms of avg_rrpp and avg_ropp, and an overall best quality when order and ratio preservation are equally important.

Tuning of Parameters γ and λ Here we give a detailed discussion over the setting of the parameters γ and λ .

Specifically, γ controls the depth of the dynamic programming in the order-preserving bias setting. Intuitively, larger γ leads to better order preservation, but also higher time/space complexity. We desire to characterize the increase of order-preservation quality in term of γ to find the setting that balances the computation complexity and the quality.

For both datasets, we measure the ropp versus the setting

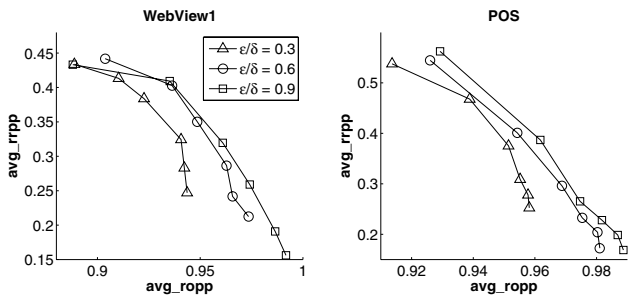


Fig. 7. Tradeoff between order-preservation and ratio-preservation.

of γ , with result shown in Fig. 6. It is noted that the quality of order-preservation increases sharply at certain points $\gamma = 2$ or 3, and the trend becomes much flatter after that. This is explained by the fact that in most real datasets, the distribution of FECs is not extremely dense, hence under proper setting of (ϵ, δ) , a FEC can intersect with only 2 or 3 neighboring FECs on average. Therefore, a setting of small γ is usually sufficient for most datasets.

The setting of λ balances the order and ratio-preservation. For each dataset, we evaluate ropp and rppp with different setting of λ (0.2, 0.4, 0.6, 0.8 and 1) and precision-privacy ratio ϵ/δ (0.3, 0.6 and 0.9), as shown in Fig. 7.

These plots give good estimation of the gain of order preservation, given the ratio preservation one is willing to sacrifice. A larger ϵ/δ gives more room for this adjustment. In most cases, the setting of $\lambda = 0.4$ offers a good balance between the two metrics. The trade-off plots could be made more accurate by choosing more settings of λ and ϵ/δ to explore more points in the space.

Efficiency In evaluating the efficiency, we measure the computation overhead of our Butterfly approaches over the original mining algorithm, given different settings of minimum support C . We divide the execution time into three parts: the mining algorithm (Mining alg), the optimization part seeking the optimal setting (Opt) and the basic perturbation part (Basic). The window size is set 5K for both datasets.

The result plotted in Fig. 8 shows clearly that the overhead of our approaches, especially the basic perturbation operation is almost unnoticeable, therefore, it can be readily implemented in current stream mining systems. While the current version of our methods are window-based, in the future work we aim at developing incremental version, and expect even lower overhead.

Note that in most cases, the running time of both mining algorithm and optimization part grow significantly as C decreases, however the growth of the overhead of Butterfly is much less evident than that of the mining algorithm. This is explained by that as the minimum support decreases, the number of frequent itemsets increases super-linearly, but the number of FECs has much lower growth rate, which dominates the performance of the Butterfly approaches.

VIII. CONCLUSIONS

In this work, we investigated the problem of protecting output privacy in data stream mining. We presented the inference and disclosure scenarios where the adversary performs attack

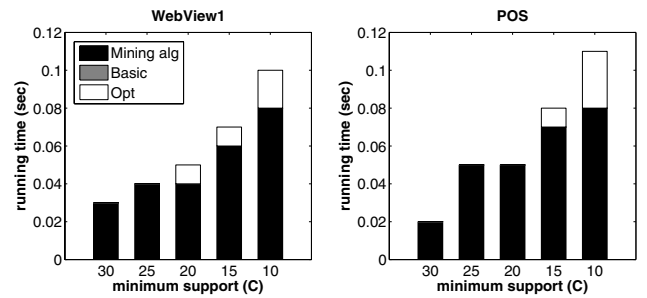


Fig. 8. Overhead of Butterfly approaches in mining systems.

over the mining output. Motivated by the basis of the attack model, we proposed Butterfly, a two-tier countermeasure: In the first tier, it counters the malicious inferences by amplifying the uncertainty of vulnerable patterns, at the cost of trivial decrease in the output precision; In the second tier, for given privacy and precision requirement, it maximally preserves the utility-related semantics of output, therefore achieving the optimal balance between privacy guarantee and output quality. The effectiveness and efficiency of our approaches are validated by extensive experiments over real-life datasets.

ACKNOWLEDGMENT

This research is partially sponsored by grants from NSF CyberTrust, an IBM SUR grant, and an IBM faculty award.

REFERENCES

- [1] R. Agrawal and R. Srikant, "Privacy preserving data mining," in *SIGMOD*, 2000.
- [2] K. Chen and L. Liu, "Privacy preserving data classification with rotation perturbation," in *ICDM*, 2005.
- [3] A. Evfimevski, R. Srikant, R. Agarwal, and J. Gehrke, "Privacy preserving mining of association rules," in *SIGKDD*, 2002.
- [4] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramanian, "l-diversity: Privacy beyond k-anonymity," in *ICDE*, 2006.
- [5] L. Sweeney, "k-anonymity: a model for protecting privacy," *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*, vol. 10, no. 5, 2002.
- [6] S. Bu, L. Lakshmanan, R. Ng, and G. Ramesh, "Preservation of pattern and input-output privacy," in *ICDE*, 2007.
- [7] F. Chin and G. Ozsoyoglu, "Statistical database design," *ACM Trans. on Database Systems*, vol. 6, no. 1, 1981.
- [8] M. Atzori, F. Bonchi, F. Giannotti, and D. Pedreschi, "Anonymity preserving pattern discovery," *VLDB Journal*, vol. 16, no. 4, 2006.
- [9] N. Adam and J. Wortman, "Security-control methods for statistical databases," *ACM Computing Surveys*, vol. 21, no. 4, 1989.
- [10] A. Shoshani, "Statistical databases: Characteristics, problems and some solutions," in *VLDB*, 2002.
- [11] D. Agrawal and C. Aggarwal, "On the design and quantification of privacy preserving data mining algorithms," in *PODS*, 2001.
- [12] Y. Lindell and B. Pinkas, "Privacy preserving data mining," in *CRYPTO*, 2000.
- [13] F. Li *et al.*, "Hiding in the crowd: Privacy preservation on evolving streams through correlation tracking," in *ICDE*, 2007.
- [14] M. Kantarcioglu, J. Jin, and C. Clifton, "When do data mining results violate privacy?" in *SIGKDD*, 2004.
- [15] K. Wang, B. Fang, and F. Yu, "Handicapping attacker's confidence: An alternative to k-anonymization," *Knowl. and Info. Syst.*, vol. 11, no. 3, 2007.
- [16] T. Calders and B. Goethals, "Mining all non-derivable frequent itemsets," in *PKDD*, 2002.
- [17] T. Wang and L. Liu, "Output privacy protection in stream mining," <http://www.cc.gatech.edu/~twang/>, 2007.
- [18] T. Calders, "Computational complexity of itemset frequency satisfiability," in *PODS*, 2004.
- [19] S. Vavasis, "Quadratic programming is in np," *Information Processing Letter*, vol. 36, no. 2, 1990.
- [20] Y. Chi, H. Wang, P. Yu, and R. Muntz, "Moment: Maintaining closed frequent itemsets over a stream sliding window," in *ICDM*, 2004.