

Association for Information Systems

## AIS Electronic Library (AISeL)

---

ICEB 2004 Proceedings

International Conference on Electronic Business  
(ICEB)

---

Winter 12-5-2004

### Privacy-Preserving Data Mining In Electronic Surveys

Justin Zhan

Stan Matwin

Follow this and additional works at: <https://aisel.aisnet.org/iceb2004>

---

This material is brought to you by the International Conference on Electronic Business (ICEB) at AIS Electronic Library (AISeL). It has been accepted for inclusion in ICEB 2004 Proceedings by an authorized administrator of AIS Electronic Library (AISeL). For more information, please contact [elibrary@aisnet.org](mailto:elibrary@aisnet.org).

# Privacy-Preserving Data Mining In Electronic Surveys

Justin Zhan, Stan Matwin

School of Information Technology & Engineering, University of Ottawa, Canada  
{zhizhan, stan}@site.uottawa.ca

## ABSTRACT

Electronic surveys are an important resource in data mining. However, how to protect respondents' data privacy during the survey is a challenge to the security and privacy community. In this paper, we develop a scheme to solve the problem of privacy-preserving data mining in electronic surveys. We propose a randomized response technique to collect the data from the respondents. We then demonstrate how to perform data mining computations on randomized data. Specifically, we apply our scheme to build a Naive Bayesian classifier from randomized data. Our experimental results indicate that accuracy of classification in our scheme, when private data is protected by randomization, is close to the accuracy of a classifier build from the same data with the total disclosure of private information. Finally, we develop a measure to quantify privacy achieved by our proposed scheme.

**Keywords:** Privacy, Data mining, Randomization

## 1. INTRODUCTION

Data mining has emerged as a means for identifying patterns and trends from large amounts of data. To conduct data mining computations, we need to collect data first. However, because of privacy concerns, people might decide to selectively divulge information, or give false information, or simply refuse to disclose any information at all. There is research evidence [4] that providing privacy protection measures is a key to the success of data collection.

There are many ways to collect data. For instance, data may be collected using transaction records. This can often be done without people's knowledge, and individuals have no control over what information can be collected. The evolving legal developments will hopefully soon preclude this questionable practice. Another way to collect data is to solicit respondents' responses via surveys, for example, respondents might be asked to rate certain products, or they might be asked whether they have a certain medical condition, etc. The collected data is entered into a database. Although answering survey questions gives respondents control over whether they want to disclose their information or not, privacy concerns might hinder the respondents from telling the truth or responding at all (we will refer to this problem as *respondent privacy in electronic surveys*). How can we improve the chance to collect more truthful data that are useful for data mining while preserving respondents' privacy? How can respondents contribute their personal information without compromising their privacy?

We propose to use *Randomized Response* techniques [9] to solve the problem of respondent privacy in electronic surveys. The basic idea of randomized response is to scramble the data in such a way that data collector cannot tell with probabilities better than a pre-defined threshold whether the data from a respondent contain truthful information about the sensitive, private

information. Although information from each individual respondent is scrambled, if the number of respondents is significantly large, the aggregate information of these respondents can be estimated with reasonable accuracy. Such property is useful for naive Bayesian classification since it is based on aggregate values of a data set, rather than individual data items.

The contributions of this paper are as follows: (1) We have modified naive Bayesian classification algorithm [7] to make it work with data disguised by randomized response techniques, and implemented the modified algorithm. (2) We then conducted a series of experiments to measure accuracy of our modified naive Bayesian algorithm on randomized data. Our results show that if we choose appropriate randomization parameters, the accuracy we have achieved is very close to the accuracy achieved by standard, unmodified naive Bayesian classifier on undisguised data. (3) We have developed a method to measure privacy achieved by our proposed approach.

There are some related works [1, 2, 3, 5, 6, 8, 10] in privacy-preserving data mining. We do not discuss them in details because of space.

## 2. BUILDING NAÏVE BAYESIAN CLASSIFIERS USING MULTI-VARIANT RANDOMIZED RESPONSE TECHNIQUES

Randomized Response techniques were first introduced by Warner [9] to solve the following survey problem: to estimate the percentage of respondents in a population that has attribute A, queries are sent to a group of respondents. Since the attribute A is related to some confidential aspects of human life, respondents may decide not to reply at all or to reply with incorrect answers. For the purpose of this discussion, we will distinguish two types of questions in a survey: questions about the respondent's private information, and questions about the respondent's personal information. Both kinds of information refer to the attributes of the

respondent. The private information is an attribute the respondent would rather not disclose, including its probability distribution (e.g., whether the respondent has a certain medical condition; or whether she takes a given medication); personal information is also an attribute of the respondent, but unlike private information the respondents do not normally mind that data collector knows the probability distribution of personal information (e.g., what is the probability that the color of the respondent's hair being black, or what is the probability that she lives near a lake). We also assume that private and personal information are unrelated - e.g., taking a medication is unrelated to one's hair color. To enhance the level of cooperation, instead of asking each respondent whether she has the attribute A, data collector asks each respondent two unrelated questions. One of them asks private information, i.e., the one that data collector is interested in. The other refers to personal information. The answers to the two questions are unrelated to each other [9]. For example, the survey questions can be designed as follows:

1. Do you have the *private* attribute A?
2. Do you have the *personal* attribute Y?

In practice, the first question could be "Are you taking medicine A?", and the second question could be "Do you live near a lake?". Respondents answer one of these two questions. They use a randomization device to decide which question to answer, without letting data collector know which question is answered. Each randomization device tells the respondent which question she is to answer: the probability of choosing

The first question is  $q$ , and the probability of choosing the second question is  $1 - q$ . Although data collector learns a response (i.e., "yes" or "no"), he does not know which question was answered by the respondents. It is important to engineer the interaction between data collector and respondent in such a way that the respondent will trust the system, i.e., the respondent will clearly understand that data collector has no way of knowing which of the two questions is answered. Thus the respondent feels that her privacy is preserved. We further comment on this in Sec.4. Note that data collector only knows the probability distribution of the respondent's attribute Y. This is consistent with the interpretation of a personal attribute - data collector could know the distribution of the values (e.g., hair colors) of the personal attribute in the general population, without knowing the value of that attribute for a specific respondent.

The randomized response technique discussed above considers only one attribute. However, data sets usually consist of multiple attributes; finding the relationship among these attributes is one of the major goals for data mining. Therefore, we need techniques that can handle multiple attributes while supporting various data mining computations. In this paper, we provide multi-variant

randomized response technique (MRR) to address the problems of respondent privacy in electronic surveys.

## 2.1 Notations

In this work, we assume data are binary, but the techniques can be extended to categorical data. Suppose there are  $N$  *private* attributes ( $A_1, A_2, \dots, A_N$ ) in a data set A. We construct  $N$  *personal* attributes ( $Y_1, Y_2, \dots, Y_N$ ). We want one *private* attribute (question) to pair with one *personal* attribute (question), therefore we make the number of attributes of Y and the number of attributes of A be equal. Let A and Y represent any logical expression based on those attributes  $A_i (i \in [1, N])$  and  $Y_i (i \in [1, N])$ . For example, A can be  $(A_1 = 0) \wedge (A_2 = 1)$  and Y can be  $(Y_1 = 0) \wedge (Y_2 = 1)$ . Let  $P(Y)$  be the proportion of the records in the personal data that satisfy  $Y = \text{true}$ . Let  $P^*(A)$  be the proportion of the records in the whole randomized data set that satisfies  $A = \text{true}$ . Let  $P(A)$  be the proportion of the records in the whole non-randomized data set that satisfy  $A = \text{true}$  (the potential non-randomized data set which in reality does not exist).  $P^*(A)$  can be observed from the randomized data, but  $P(A)$ , the actual proportion that we are interested in, cannot be observed from the randomized data because the non-randomized data set is not available to anybody; we have to estimate  $P(A)$ . The goal of MRR is to find a way to estimate  $P(A)$  from  $P^*(A)$ .

## 2.2 Multi-variant Randomized Response Scheme

In this scheme, all the attributes including the class label will be treated as a group. They either keep the same values or obtain the values from personal data. In other words, when sending the private data to data collector, respondents either tell their answers to the private questions or tell their answers to the personal questions. The probability for the first event is  $q$ , and the probability for the second event is  $1 - q$ . For example, assume a respondent's attribute values  $A_1$  and  $A_2$  are 11 for private data; and the respondent's attribute values  $Y_1$  and  $Y_2$  are 01. The respondent generates a random number between 0 and 1; if the number is less than  $q$  she sends 11 to data collector; if the number is bigger than  $q$ , she sends 01 to data collector. Since data collector only knows  $q$  which is the same for all respondents and does not know the random number generated by each respondent, he cannot know whether the respondent tells the values from private data or personal data. To simplify our presentation, we use  $P(A(11))$  to represent  $P(A_1 = 1 \wedge A_2 = 1)$ ,  $P(Y(11))$  to represent  $P(Y_1 = 1 \wedge Y_2 = 1)$  where " $\wedge$ " is the logical *and* operator. Because the contributions to  $P^*(A(11))$  partially come from  $P(A(11))$ , and partially

come from  $P(Y(11))$ , we can derive the following equation:

$$P^*(A(11)) = P(A(11)) \cdot \mathbf{q} + P(Y(11)) \cdot (1 - \mathbf{q})$$

Since  $P(Y(11))$  is known as  $Y$  is personal data,  $\mathbf{q}$  is determined before collecting the data, and  $P^*(A(11))$  can be directly computed on the disguised (randomized) data set. By solving the above equation, we can obtain  $P(A(11))$ , the information needed to build a naive Bayesian classifier. The general model is described in the following:

$$P^*(A) = P(A) \cdot \mathbf{q} + P(Y) \cdot (1 - \mathbf{q}) \quad (1)$$

### 2.3 Building Naïve Bayesian Classifiers

The naive Bayesian classifier is one of the most successful algorithms in many classification domains. Despite of its simplicity, it is shown to be competitive with other complex approaches, especially in text categorization and content based filtering. The naive Bayesian classifier applies to learning tasks where each instance  $x$  is described by a conjunction of attribute values and where the target function  $f(x)$  can take on any value from some finite set  $V$ . A set of training examples of the target function is provided, and a new instance is presented, described by the tuple of attribute values  $\langle a_1, a_2, \dots, a_n \rangle$ . The learner is asked to predict the target value for this new instance. Under a conditional independence assumption, i.e.,

$P(a_1, a_2, \dots, a_n | y_j) = \prod_{i=1}^n P(a_i | y_j)$ , a naive Bayesian classifier can be derived as follows:

$$\begin{aligned} v_{NB} &= \arg \max_{v_j \in V} P(v_j) \prod_{i=1}^n P(a_i | v_j) \\ &= \arg \max_{v_j \in V} P(v_j) \prod_{i=1}^n \frac{P(a_i \wedge v_j)}{P(v_j)} \end{aligned}$$

To build a NB classifier, we need to compute  $P(v_j)$  and  $P(a_j \wedge v_j)$ . To compute  $P(v_j)$ , we can use the general model (Eq.(1)) with  $A$  being ( $C = v_j$ ) and  $Y$  being ( $CY = v_j$ ) where  $C$  is the class label for the private data  $A$  and  $CY$  is the class label of personal data  $Y$ .  $P^*(A)$  can be computed directly from the (whole) randomized data set.  $P(Y)$  is known since it is personal and  $\mathbf{q}$  is known as well. By knowing  $\mathbf{q}$ , data collector, who conducts the training, only knows the probability of the training data being private, but does not exactly know if each value is private data or not. By solving the above equation, we can get  $P(A)$  which is  $P(C = v_j)$  in this case. Similarly, we can compute  $P(a_i \wedge v_j)$  using the general model (Eq.(1)) with  $A$  being ( $A_i = a_i \wedge C = v_j$ ) and  $Y$  being ( $Y_i = a_i \wedge CY = v_j$ ).

### 2.4 Testing

Conducting the testing is straightforward when data are

not randomized, but it is a non-trivial task when the testing data set is randomized. When we choose a record from the testing data set, compute a predicted class label using the naive Bayesian classifier, and find out that the predicated label does not match the record's actual label, an we say this record fails the testing? If we knew whether the record represents the private or the personal data, and if we knew the true class for each data, we could easily answer this question. But how can we compute the accuracy score of a NB classifier when data are randomized? Our answer is to apply the multi-variant randomized response technique once again to compute the accuracy. Let us use an example to illustrate how to compute the accuracy. Assume the number of attributes is 2. To test a record ( $A_1 = 1, A_2 = 0$ ) denoted by  $A(10)$ , we feed  $A(10)$  and  $Y(10)$ , where  $Y = (Y_1 = 1, Y_2 = 0)$  to the NB classifier built in Sec.2.3. Let  $P^*(A(cc))$  be the proportion of correct predictions using the disguised (randomized) testing data set,  $P(Y(cc))$  be the proportion of correct predictions in the personal data, and let  $P(A(cc))$  be the proportion of correct predictions in the private data.  $P(A(cc))$  is what we want to estimate. Because  $P^*(A(cc))$  consist of contributions from  $P(A(cc))$  and  $P(Y(cc))$ , we have the following equation:

$$P^*(A(cc)) = P(A(cc)) \cdot \mathbf{q} + P(Y(cc)) \cdot (1 - \mathbf{q})$$

Where  $P^*(A(cc))$  can be obtained from disguised testing data set.  $\mathbf{q}$  is known and by knowing  $\mathbf{q}$ , data collector, who conducts the testing, only knows the probability of the testing data being private, but does not exactly know if each value is private data or not. How does data collector know  $P(Y(cc))$ ? One implementation is as follows: each respondent is given the same classifier by data collector. The classifier is constructed during the training (Sec.2.3). Each respondent applies this classifier on her personal data  $Y$  and communicates the number of correct predictions (0 or 1) to data collector, who then computes ( $Y(cc)$ ). Note that data collector does not know the values of the  $Y$  attributes, only the result of the classifier. Data collector can solve the above equation and get  $P(A(cc))$ , the accuracy score of testing.

## 3. EXPERIMENTAL RESULTS

To evaluate the effectiveness of our proposed scheme, We conducted experiments on two real life data sets *Adult* and *Breast Cancer* which were obtained from the UCI Machine Learning Repository.

### 3.1 Experimental Steps

We modified naive Bayesian classification algorithm to handle randomized data based on our proposed scheme. We applied our scheme to obtain a privacy-oriented classifier. We also ran naive Bayesian classification algorithm on original data set, and obtained a base classifier. We then applied the same testing data to both

classifiers. Our goal is to compare classification accuracy of these two classifiers. Obviously we want accuracy of privacy-oriented classifier to be close to accuracy of the base classifier. Our experiments consist of the following steps:

#### Step I: Preprocessing

Since we assume that data set contains only binary data, we first discretize original non-binary data to become binary. We split the value of each attribute from the median point of the range of the attribute. After preprocessing, we randomly divided data sets into a training data set D (80%) and a testing data set B (20%). Note that B will be used for comparing our results with benchmark results.

#### Step II: Benchmark

We use D and the original NB classification algorithm to build a classifier  $T_D$ ; we use data set B to test the classifier, and get an accuracy score. We call this score original accuracy (or benchmark score).

#### Step III: $q$ Selection

For  $q = 0, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$  and 1.0, we conduct the following 4 steps:

##### Sub-step 1: Randomization

We create a disguised data set G. For each record in training data set D, we generate a random number  $r$  from 0 to 1 using uniform distribution. If  $r < q$ , we copy the record of D to G without any change; if  $r \geq q$ , we randomly generate the values for a record of Y according to the pre-defined probability and copy the record values to G. In this paper, each record of Y is randomly generated such that each logical expression (Y) appears with the probability of 0.5. That is  $W_y = 0.5$  (ref.Sec.4). We perform this randomization step for all the records in the training data set D, then generate the new data set G.

##### Sub-step 2: Classifier Construction

We use data set G and our modified NB classification algorithm to build a naive Bayesian classifier  $T_G$ .

##### Sub-step 3: Testing

We use data set B to test and get an accuracy score S.

##### Sub-step 4: Repeating

We repeat steps 1-3 for 1000 times, and get  $S_1, S_2, \dots, S_{1000}$ . We then compute mean and variance of these 1000 accuracy scores.

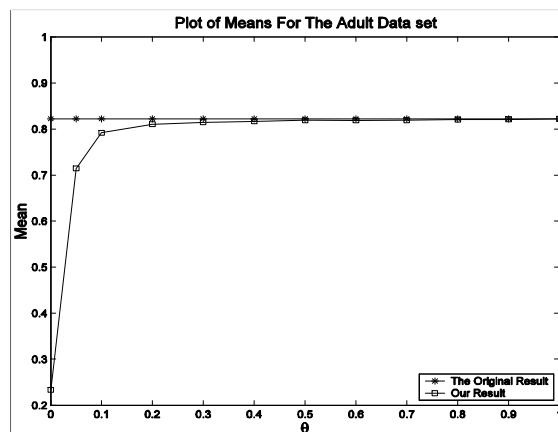


Figure 1

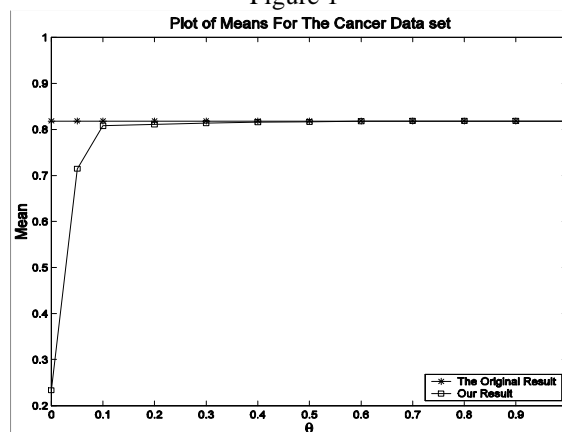


Figure 2

## 3.2 Accuracy Analysis

### 3.2.1 Analysis of Mean

Fig.1 and Fig. 2 shows mean values of accuracy scores for *Adult* and *Breast-Cancer* data sets respectively. We can see from figures that when  $q=1$ , the results are exactly the same as the results when standard, unmodified classification algorithm is applied. This is because when  $q=1$ , randomized data sets are all from private data D. When  $q$  approaches 1, contribution of private data is enhanced; with  $q$  deviating from 1, the contribution of private data is decreasing (when  $q=0$ , collected data set is all from personal data). Therefore, when  $q$  moves from 1 towards 0, the mean of accuracy has the trend of decreasing.

### 3.2.2 Analysis of Variance

Fig. 3 and Fig.4 shows variances of accuracy scores. When  $q$  moves from 1 towards 0, the degree of randomness in disguised data is increasing, variance of estimation used in our method becomes larger. Variance changes with different randomization levels  $q$ . When  $q$  is near 0, randomization level is much higher and private data is better disguised. We do not show variance when  $q=0$ . In this case, since collected data set is actually personal data and the probability

distribution for it is always the same for each iteration, variance is 0.

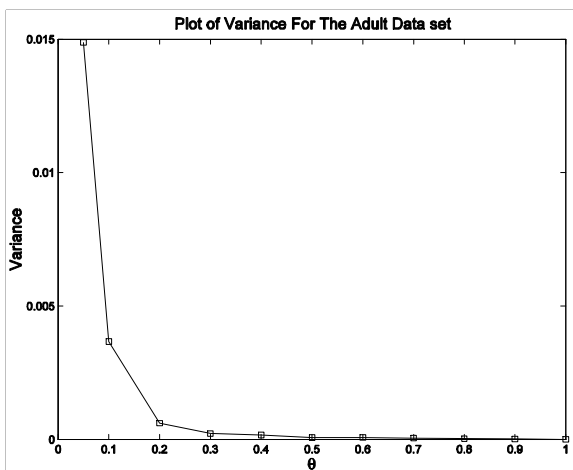


Figure 3

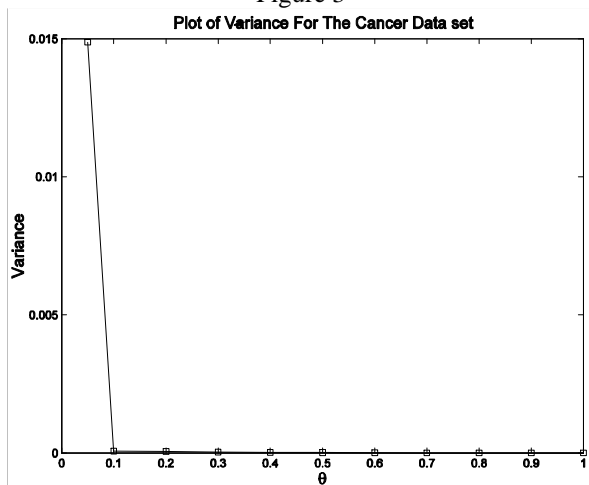


Figure 4

#### 4. MEASURING PRIVACY

Various privacy enhancing techniques [5, 8, 11, 14, 15, 16] have been developed to protect data privacy. To trust a privacy enhancing technology, we need to know how much privacy that a particular technique protects. A challenge faced by privacy-preserving data mining community is how to measure privacy. A general privacy measure which can quantify privacy for any privacy protection scheme is still an open question. In this section, we develop a privacy measure for our proposed multi-variant randomized response technique. Our measure contains two steps: First, we measure privacy for a single entry. Second, we select the minimal privacy value and treat it as the privacy level for the group. The reason why we choose the minimal value for the group is that, when the entries are randomized together, finding the original value for one entry will cause disclosing the original values for other entries in the group.

For a single entry, original value can be 1 or 0; randomized value can be 1 or 0 as well. Privacy comes from uncertainty of original value given a randomized

value. In other words, if original value is 1, given randomized value 1 or 0, privacy will be the probability that data collector guesses the original value being 0. There are four possible randomization results:

- Original value is 1, value after randomization is 1;
- Original value is 1 but value after randomization is 0;
- Original value is 0 but value after randomization is 1;
- Original value is 0, value after randomization is 0.

Let's use the following denotations:

- Let's  $O_m$  be the original value;
- Let's  $R_m$  be the value after randomization;
- Let's  $G_m$  be the guessed value.
- Let's  $W_a$  be the probability that a value is 1 in data set A, and the probability that a value is 0 in data set A is  $(1 - W_a)$ ;
- Let's  $W_y$  be the probability that a value is 1 in data set Y, and the probability that a value is 0 in data set Y is  $(1 - W_y)$ ;

Privacy for a single entry, denoted by  $PSE$ , can be derived as follows:

$$\begin{aligned}
 PSE &= \Pr(O_m = 1) \cdot \Pr(R_m = 1 | O_m = 1) \cdot \Pr(G_m = 0 | R_m = 1) \\
 &\quad + \\
 &\Pr(O_m = 1) \cdot \Pr(R_m = 0 | O_m = 1) \cdot \Pr(G_m = 0 | R_m = 0) \\
 &\quad + \\
 &\Pr(O_m = 0) \cdot \Pr(R_m = 1 | O_m = 0) \cdot \Pr(G_m = 1 | R_m = 1) \\
 &\quad + \\
 &\Pr(O_m = 0) \cdot \Pr(R_m = 0 | O_m = 0) \cdot \Pr(G_m = 1 | R_m = 0) \\
 &= Comp_1 + Comp_2 + Comp_3 + Comp_4
 \end{aligned}$$

The first component contains three parts:

1.  $\Pr(O_m = 1)$  is the probability that a value is 1 in private data set (A), which is  $W_a$ .
2.  $\Pr(R_m = 1 | O_m = 1)$  is the probability that a randomized value is 1 given the original value is 1. There are two possibilities: (1) a randomized value comes from data set A, and the probability for the case is  $q$ ; (2) a randomized value comes from data set Y, and the probability for this case is  $(1 - q) \cdot W_y$ .
3. Let  $(G_m | R_m)$  be the guessed value given a randomized value and  $(O_m | R_m)$  be the original value given the same randomized value. For the same  $R_m$  value, there are two possibilities:  $G_m = O_m$  or  $G_m \neq O_m$ . We should notice when  $(G_m | R_m) \neq (O_m | R_m)$ , there is zero contribution to  $PSE$ . For instance, in the first component, the original value is 1, the guessed value is 0 given the randomized value is 1. When  $(G_m | R_m) \neq (O_m | R_m)$ , then the guessed will be 1. That is the guessed value and the original value is the same (both of them are 1s) and it contributes zero to the  $PSE$ . Therefore, we only consider the case where

$(G_m | R_m) = (O_m | R_m)$  for the third part. We can apply Bayes rule to tackle this part. After applying the Bayes rule, we obtain  $\Pr(R_m = 1 | O_m = 0) \cdot \Pr(O_m = 0) / \Pr(R_m = 1) \cdot \Pr(O_m = 0)$  is the probability that a value is 0 in private data set (A), which is  $1 - W_a$ .  $\Pr(R_m = 1 | O_m = 0)$  is the probability that a randomized value is 1 given the original value is 0. In this case, randomized value cannot come from private data since original value is 0 and randomized value is 1. The only possibility is that randomized value is from personal data set Y, and the probability is  $(1 - q) \cdot (1 - W_y)$ . As for  $\Pr(R_m = 1)$ , we can extend this term and details are shown as follow:

$$\begin{aligned} Comp_1 &= W_a \cdot [q + (1 - q) \cdot W_y] \cdot \frac{\Pr(R_m = 1 | O_m = 0) \cdot \Pr(O_m = 0)}{\Pr(R_m = 1)} \\ &= \frac{W_a \cdot [q + (1 - q) \cdot W_y] \cdot (1 - q) \cdot (1 - W_y) \cdot (1 - W_a)}{\Pr(R_m = 1 | O_m = 1) \cdot \Pr(O_m = 1) + \Pr(R_m = 1 | O_m = 0) \cdot \Pr(O_m = 0)} \\ &= \frac{W_a \cdot [q + (1 - q) \cdot W_y] \cdot (1 - q) \cdot (1 - W_y) \cdot (1 - W_a)}{[q + (1 - q) \cdot W_y] \cdot W_a + (1 - q) \cdot (1 - W_y) \cdot (1 - W_a)} \end{aligned}$$

Similarly, we can obtain other components, and we then get

$$\begin{aligned} PSE &= \frac{(1 - W_a) \cdot W_a \cdot (1 - q) \cdot [q + (1 - q) \cdot W_y]}{[q + (1 - q) \cdot W_y] \cdot W_a + (1 - q) \cdot W_y \cdot (1 - W_a)} \\ &\quad + \\ &\quad \frac{2W_a \cdot (1 - W_a) \cdot (1 - q) \cdot (1 - W_y) \cdot [q + (1 - q) \cdot (1 - W_y)]}{[q + (1 - q) \cdot (1 - W_y)] \cdot (1 - W_a) + (1 - q) \cdot (1 - W_y) \cdot W_a} \end{aligned} \quad \text{Eq. (2)}$$

We compute  $PSE$  for each single entry. We then select the smallest value  $PSE(\text{Min})$  as the privacy value for the group. We can see from Eq.(2) that,  $PSE$  is determined by three parameters: (1) control parameter  $q$ ; (2) private data (ref. to as data set A in Sec.2) distribution  $W_a$ ; (3) personal data (ref. to as data set Y in Sec.2) distribution  $W_y$ . What we can see from  $PSE$  equation is that, when  $W_a = 0.9, 0.8, 0.7, 0.6$  privacy is symmetric with respect to privacy when  $W_a = 0.1, 0.2, 0.3, 0.4$ . In other words, given a certain  $q$ ,  $PSE$  value will equal to the  $PSE$  value when control parameter is  $1 - q$ .

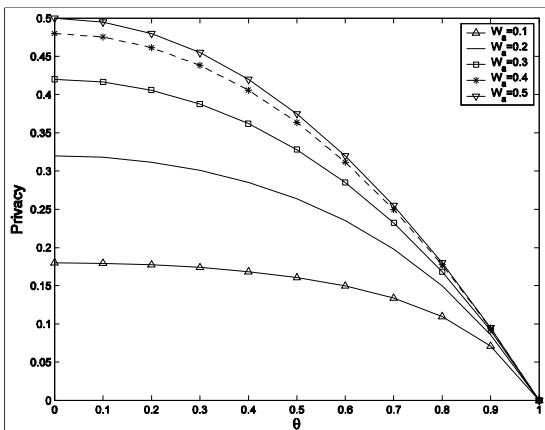


Figure 5

To get better sense of our proposed privacy measure, we conducted a set of experiments on private data sets with

various distributions and personal data set with  $W_y = 0.5$ . Specifically, we conduct experiments when  $W_a = 0.1, 0.2, 0.3, 0.4, \text{ and } 0.5$ . For each data distribution, we compute privacy value for the cases where  $q = 0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1$ . As we see from results in Fig. 5:

- When  $q = 1$ , private data is fully disclosed. Privacy value is 0;
- When  $q = 0$ , data collector gets no private data, and the data obtained are all personal data. In this case, privacy level of private data is the highest.
- When  $q$  is away from 1 and approaches 0, the elements of private data contribute less to the classification, and the probability of disclosing private data is decreasing. Therefore privacy level of private data increases.
- When private data (A) distribution approaches to uniform ( $W_a = 0.5$ ), privacy level is increasing. Since uniform distribution will make original data recoverability be the lowest.

Empirical results from Sec.3 and Sec.4 confirm that recoverability and privacy are complementary goals. Given  $W_a$  and  $W_y$ , the best privacy is achieved when control parameter  $q$  is 0; however, the accuracy will be the worst in this case. The best accuracy is attained when  $q = 1$  but privacy is the worst. Trade-offs are also applied when  $q$  has a value between 0 and 1. In practice, how to select  $q$  is dependent upon our primary goals. If we want the results to be very precise, we need choose the values near 1; in contrast, if privacy is the primary goal, we choose the values near 0.

## 6. CONCLUDING REMARKS

In this paper, we have presented a method to build naïve Bayesian classifiers using multi-variant randomized response technique. Experimental results show that when we select an appropriate randomization parameter  $q$ , we can get fairly accurate classifiers comparing to the classifiers built from undisguised data. A privacy measure was developed and privacy analysis was also conducted. Trade-offs between privacy and accuracy are discussed. The proposed multi-variant unrelated question model has a broader impact in the sense that it can be used not only for naïve Bayesian classification, but also can be utilized in many other privacy-preserving data mining computations, such as decision tree induction, Bayesian classification, probabilistic-based clustering. As future work, we will apply the proposed scheme to other data mining problems.

## REFERENCES

- [1] H. Kargupta, S. Datta, Q. Wang, and K. Sivakumar, "Random Data Perturbation Techniques and Privacy Preserving Data Mining", the IEEE International

- Conference on Data Mining, 2003, Florida, USA.
- [2] A. Evfimievski, J. E. Gehrke, and R. Srikant, "Limiting Privacy Breaches in Privacy Preserving Data Mining", Proceedings of the 22nd ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, 2003, San Diego, CA.
- [3] R. Agrawal and R. Srikant, "Privacy-Preserving Data Mining", Proceedings of The ACM SIGMOD Conference On Management of Data, 2000, Dallas, Texas, USA.
- [4] L. F. Cranor and J. Reagle and M. S. Ackerman, "Beyond concern: Understanding net users' attitudes about online privacy", AT&T Labs-Research, 1999, April, Available from <http://www.research.att.com/library/trs/TRs/99/99.4.3/report.htm>
- [5] W. Du and Z. Zhan, "Using Randomized Response Techniques For Privacy-Preserving Data Mining", Proceedings of The 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2003, Washington, DC, USA, August 24-27
- [6] Y. Lindell and B. Pinkas, "Privacy Preserving Data Mining", Advances in Cryptology - CRYPTO '00, 2000, 1880 of Lecture Notes in Computer Science. Springer-Verlag, 36-54
- [7] Langley, P., Iba, W., and Thompson, K., "An Analysis of Bayesian Classifiers, National Conference on Artificial Intelligence, 223-228, 1992, url = "citeseer.nj.nec.com/langley92analysis.html"
- [8] S. Rizvi and J.R. Haritsa, "Maintaining Data Privacy in Association Rule Mining", Proceedings of the 28th VLDB Conference, 2002, Hong Kong, China
- [9] S. Warner, "Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias", The American Statistical Association, Volume:60, number:309, pages:63-69, March,1965
- [10] J. Vaidya and C. Clifton, "Privacy-Preserving K-Means Clustering over Vertically Partitioned Data", Proceedings of The 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2003, Washington, DC, USA