

# ENHANCING THE EFFECTIVENESS OF PRIVACY PRESERVING DATA MINING (PPDM) BY USING CORRELATION BASED TRANSFORMATION STRATEGY

Aryan Grover

## ABSTRACT

*Preservation of security is a huge part of information mining. The primary goal of PPDM is to cover up or give security to certain touchy data with the goal that they can be shielded from unapproved gatherings or unauthorised. Though security is accomplished by concealing the touchy or private information, it will influence the information mining calculations in information extraction, so a compelling strategy or methodology is required to give security to the information and at the same time ensuring the nature of information mining calculations. Rather than expelling or scrambling touchy or private information, we utilize information change methodologies that keep the factual, semantic and heuristic nature of information while securing the delicate or private information. In this paper, we considered the specialized possibility of acknowledging Privacy-Preserving Data Mining. In the proposed work, Correlation Based Transformation Strategy for Privacy Preserving Data Mining is utilized for ordinal information. We apply the technique on a few datasets to be specific soybean, Breast Cancer, Nursery dataset and Car dataset. We classify the final products applying the proposed system on both the first and the changed dataset and watch connection distinction, Information Entropy and Classification Accuracy with various AI calculations and Clustering Quality. As an improvement, the proposed work can be stretched out by utilization of vector stamping methods where these strategies help in expanding the productivity by maintaining a strategic distance from unapproved access to the data.*

## 1. INTRODUCTION

Data Mining is extensively used in varied areas like financial data analysis, retail industry, biological data analysis and many more. However, it has got its downsides. One of the key issues raised by data mining technology is not a business or technology one, but a social one. It is the privacy of an individual or a company. Data Mining makes it achievable to evaluate everyday business transactions and gather a considerable quantity of information about individuals buying habits and preferences. Many companies are making fortune aggregating petite pieces of information about people and putting scrap together to build a digital profile. Most of the times the information collected will be used to sell stuff, which is useful. However, the information extracted can be used for privacy violating purposes. Agencies, hospitals and other organizations often need to publish micro data for research and other purposes. However, the information extracted can be used for privacy violating purposes. As explained in 1 micro data is usually stored in the form of table where each row represents an individual.

Here the table has three types of attributes:

1. Identity attribute (To uniquely identify an individual like name),
2. Quasi identifier (which includes demographic attributes),
3. Sensitive attributes (which include confidential information like diseases).

Quasi identifiers attributes may be merged with other public databases to uniquely identify the individual and their sensitive data (Linking attack). Thus privacy is becoming a critical issue which led to a new research field called Privacy Preserving Data Mining (PPDM). PPDM comes into picture in the situations like the one described above. PPDM helps to perform data mining efficiently while preserving the private data or information about an individual or a company. Instead of hiding or encrypting, PPDM transforms the sensitive data to some other form while preserving the usefulness of the data. Many strategies have been proposed for PPDM, one of such is Correlation Based Transformation Strategy (CBTS) which is used on numerical data. The datasets likewise contain ordinal and ostensible information; the need is to change over the ordinal and ostensible information to numerical information by saving the information utility, with the goal that the calculation can be applied proficiently. In this paper, we propose a CBTS which can be applied to ordinal qualities. We depict a system to change over both ordinal and ostensible information to numerical information on which the CBTS can be applied. We measure the Information Entropy estimations of both Original information and Transformed information and the outcomes are practically identical and furthermore we measure Cluster Misclassification Error and demonstrate the mistake is less in our methodology. The paper is sorted out as follows: Section 2 depicts Related Work. Segment 3 clarifies Problem Definition. Design is introduced in Section 4. The result is explained in Section 5. We close this paper with future work in Section 6.

## 2. RELATED WORK

In [3], the creators have utilized a system called altered information transitive strategy in which the delicate numerical information thing is to be secured by changing the first information thing. There is an examination between the changed information transitive method and the perturbative concealing procedures, for example, added substance clamour, adjusting and miniaturized scale conglomeration and exhibitions are investigated and results are drawn by closing with the acceptable outcomes utilizing the transitive systems.

In [4] authors proposed a new approach which involves in preserving sensitive information using fuzzy logic. Clustering is done, in which the original dataset i.e. numerical data is transformed into fuzzy data and then noise is added to the numeric data using an S shape fuzzy membership function. The Clusters which are generated using the fuzzified data is similar to the original cluster and privacy is also achieved. In [5] proposed a system which makes use of a perturbative system where encryption technique is applied to sensitive data items. The information has to be changed to a considerable extent before it is made available to the public for safe guarding the confidentiality of

the sensitive information. The proposed data transformation technique protects categorical sensitive data which is modified using advanced data transformation technique including cryptography technique which prevents sensitive items from public disclosure. This system gives greater results while preventing sensitive data from unauthorized disclosure and should not affect the importance of the original objective of data mining. In 6 and 7, the authors have proposed distortion based techniques to meet the privacy requirements. In the former randomized distortion technique is applied only on confidential categorical attribute. In latter probabilistic distortion method is used on original data before using frequent item set mining on the data. In 8 and 9, the authors have used correlation based techniques to achieve privacy in huge datasets. In paper 10 authors proposed a work which concentrates on finding an efficient solution for the classification problem over encrypted data in cloud. This work protects the privacy of sensitive data of users query and data access patterns. A k-NN classifier is developed firstly on a real world dataset for different parameters and the efficiency is resolved.

Authors of 11 proposed a new patient centric clinical decision support system, which is of a great help for a clinician complementary in diagnosing the risk of patient's disease without compromising its privacy. This method portrays correlation by spatial proximity. It involves the following methodologies which can handle categorical and numerical variables. Authors of 12 proposed various methods and possible risks by the method of Random Projection. It defines a number of reconstruction techniques over the data.

In paper 13 authors concentrate on decision tree learning, without accompanying loss of accuracy. This method strives at preserving the privacy of data which are partially lost. This deals with the production of a set of unreal datasets which can be obtained as a result of conversion of original dataset. Such that, redesigning of original samples without the entire group of unreal datasets becomes impossible. From these datasets the decision tree is built precisely. And also this method is congruent with that of the other approaches which preserve the privacy of sensitive data and thereby ensuring higher protection of data. In paper 14 the authors have proposed a method which provides an excellent spatial transformation method to protect the privacy concerns in cloud computing and this method also provides considerably good results with respect to the communication cost. In 15 presents an erratic system based chaotic signal generator. Due to the characteristics of chaotic signal, estimators find it hard to estimate original data since they work on noise Probability Distribution Function (PDF). The issue of maintaining data privacy while publishing is resolved. Data Perturbation level depends on trust on which the data is to be generated. Due to the different levels of trusts or same levels of trust of same data, a problem on security of data arises and may cause estimation of accurate data copies by Linear Least Squares Error (LLSE), which is an advance computational algorithm.

### 3. PROBLEM DEFINITION

Given large structured data constituting of sensitive information of ordinal nature, the objective is to preserve privacy by transforming the ordinal data into an equivalent numeric representation while retaining the original statistical nature with minimal entropy.

#### 4. ARCHITECTURE

Given a huge data containing ordinal sensitive information, our solution first converts the ordinal and nominal data to numerical data and transforms the resultant numeric data in such a way that it retains the correlation structure among the data values preserving its usefulness and maintaining the level of privacy. The conversion of ordinal data is done by taking input for each data value from the concerned user and the conversion of nominal data is done by assigning random numbers to each nominal data value. The numeric attributes are retained. We consider a dataset containing mixture of ordinal, nominal and numerical data attributes, in which many attributes are private and sensitive. The dataset is subjected to clustering method like Simple K Means to group the similar rows and classification algorithm like J48. The objective of this paper is to convert and transform the ordinal sensitive data such that the correctly classified instances and the decision trees of original data and transformed data are comparable. For the given dataset with numerical sensitive information, authors in paper [16] proposed CBTS for numerical data. Given a dataset comprising sensitive and private data, CBTS produces an outcome comprising of the subset of vectors correlated to sensitive data and produces equivalent components as substitutes. CBTS uses Pearson's correlation coefficient.

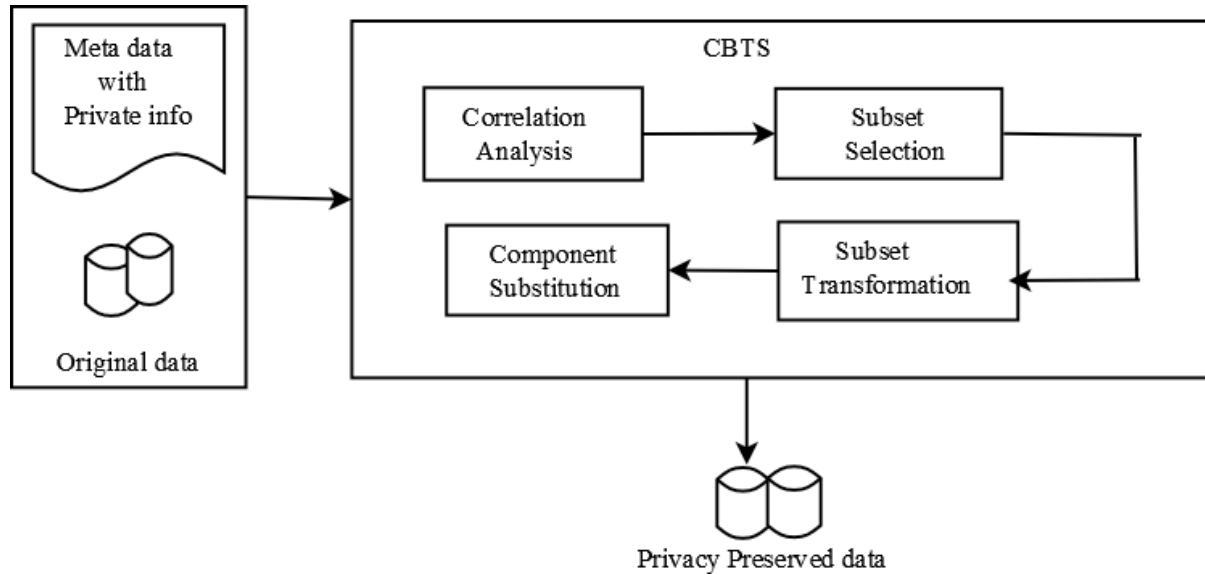
$$X^2 = \sum_{K=1}^n \frac{(O_k - E_k)^2}{E_k}$$

$O_k$  - Observed frequency.

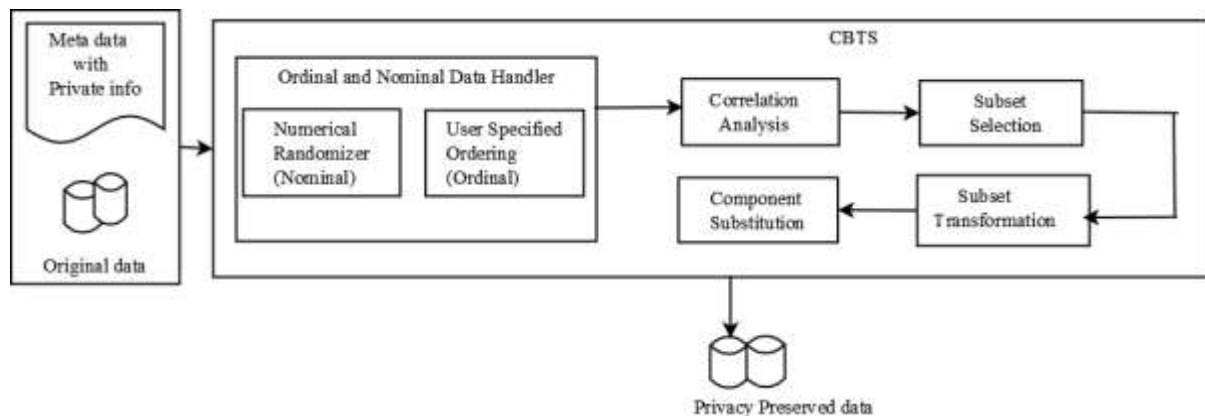
$E_k$  - Expected frequency.

The subsets generated are subjected to transformation strategies that tend to converge on the obtained similarity forming new components. Hence the components obtained are a mathematical representation of the sensitive data and used instead of sensitive data for data mining. Figure 1 gives the Architecture of CBTS for Numerical data. Existing transformation methods PCA, SVD and NNMF have been used prior in PPDM by [17–19] demonstrating the required property of convergence. The method was able to remove the highly correlated sensitive data and transform the non-correlated sensitive data. CBTS is applied to datasets which has numerical values, the information entropy values are compared for the original data and the transformed data and the results are obtained. Thorough experiment analysis proved the proposed dataset transformation method has low clustering misplacement error and minimal deviation in classifier accuracies. In this paper we are extending CBTS to support ordinal data. The proposed architecture is shown in Figure 2. Our method first converts both ordinal and nominal data to equivalent numerical data. The conversion step has two sub-steps. Initially, the dataset is parsed to extract the unique data values in each column which is given to next step. In the next step, based on the type of the data values of the column, conversion is done. When the column has ordinal data values, they are converted to numerical values based on the user provided ordering. In this work we have assumed all the nominal data to have some ordinal nature. Nominal data are substituted by unsupervised statistical methods. Correlation coefficient is calculated for the respective values against the data

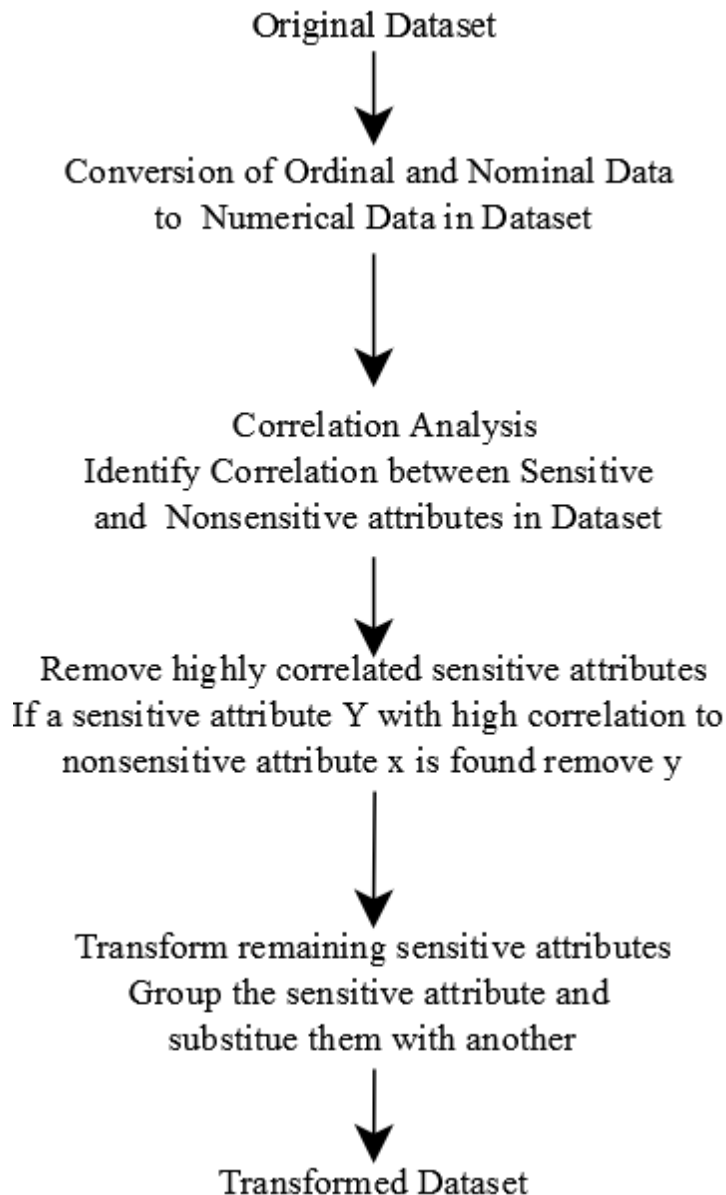
vectors. If there exists a strong correlation, then they are converted to random numbers. If the correlation is weak, then the conversion is done by substituting categories with closeranged numbers to avoid and minimize error bias. Chisquaredtest is done to determine the correlation between nominal data values.



**Figure 1.** Architecture of CBTS for numerical data.



**Figure 2.** Architecture of CBTS for ordinal and nominal data.



**Figure 3.** Transformation method.

## 5. RESULT

The datasets used in this paper are Soybean and BreastCancer. Both the datasets are taken from UCI MachineLearning Repository. Soybean dataset is a dataset with 307 instances and 35 attributes. Among 35 attributes some are ordinal and some are nominal. Breast Cancer is another dataset with ordinal, nominal and numerical attributes. There are 286 instances and 10 attributes in this dataset. This dataset contains two classes and among 286 instances, 201 belong to one class and the other 85 belong to another class. Data Entropy of unique information against bothered information utilizing CBTS for Ordinal information with change techniques is condensed in Table 1. We can gather from the table that deviation in Information Entropy is least utilizing the proposed CBTS strategy against utilizing change strategies alone. Table 2 gives the examination of classifier

correctness's for different AI calculations utilizing CBTS against unique information. It is plainly detectable from the outcomes the classifier execution is similar to the first information. Table 3 shows the Misclassification Error ME esteems with k-implies grouping.

**Table 1.** Comparison of information entropy

Types of Data	Original Entropy	Information Entropy( $I_E$ ) Using CBTS Method/Using existing Methods		
		PCA	SVD	NNMF
Soybean (683x36)	3.317	3.30/10.25	3.41/5.12	3.25/9.26
Car (1729x6)	2.31	2.28/5.16	2.37/2.58	2.22/9.14
Nursery Dataset (12960x7)	1.88	1.88/4.6	1.95/2.8	1.86/11.0
Breast Cancer (286x9)	3.02	3.7/6.39	3.5/3.8	3.39/8.04

## 6. CONCLUSION AND FUTURE WORK

CBTS achieves accountable privacy by applying correlationtransformation based methods. CBTS hasapplications over varied areas involving huge data.Joined with the CBTS we have introduced a method for change by changing over delicate ordinal and ostensible information to numerical information of a considered dataset all the while saving the protection and the information utility of the equivalent. The proposed work can be stretched out by utilization of vector checking systems where these methods help in expanding the productivity by maintaining a strategic distance from unapproved access to the data.

**Table 2.** Comparison of various machine learning algorithms using CBTS (ME)

Types of Data	Machine Learning Algorithms	Observed Classifier Accuracy (%)				
		Ordinal and Nominal Data	Numerical Data	Transformation using CBTS		
				PCA	SVD	NNMF
Soybean (683x36)	Decision Tree	97.0	96.3	97.6	97.0	97.0
	Multilayer Perceptron	99.8	93.3	94.8	95.0	95.0
	Naïve Bayes	93.7	82.1	82.5	81.8	81.8
Breast Cancer(286x9)	Decision Tree	81.4	81.4	81.4	81.4	81.4
	Multilayer Perceptron	84.6	84.6	84.2	84.6	84.6
	Naïve Bayes	73.4	73.4	73.4	73.4	73.4

Table 3. Cluster misclassification error (ME)

Types of Data	Clusters (k)	$M_k$ (with CBTS)			$M_k$ (without CBTS)		
		PCA	SVD	NNMF	PCA	SVD	NNMF
Soybean(683x36)	2	0.253	0.455	0.248	0.999	0.999	1.0
	3	1.22	0.88	0.74	1.09	2.6	0.9
Breast Cancer (286 x 9)	2	0.017	0.7	0.7	1.3	1.60	1.50
	3	0.7	0.74	0.74	0.5	1.91	1.54