

Enhanced Hybrid Privacy Preserving Data Mining Technique

Kundeti Naga Prasanthi¹, M V P Chandra Sekhara Rao², Ch Sudha Sree³, P Seshu Babu⁴

¹prasanthi.kundeti@gmail.com,

Dept. of CSE, Lakireddy Balireddy College of Engineering, Mylavaram,
Affiliated to JNTU Kakinada.

²manukondach@gmail.com

²Dept. of CSE, RVR & JC college of Engineering, Guntur.

³Dept. of CSE, RVR & JC college of Engineering, Guntur.

⁴Dept. of Mathematics and Statistics, Kakaraparti Bhavanarayana College, Vijayawada.

Abstract

Now a days, large volumes of data is accumulating in every field due to increase in capacity of storage devices. These large volumes of data can be applied with data mining for finding useful patterns which can be used for business growth, improving services, improving health conditions etc. Data from different sources can be combined before applying data mining. The data thus gathered can be misused for identity theft, fake credit/debit card transactions, etc. To overcome this, data mining techniques which provide privacy are required. There are several privacy preserving data mining techniques available in literature like randomization, perturbation, anonymization etc. This paper proposes an Enhanced Hybrid Privacy Preserving Data Mining (EHPPDM) technique. The proposed technique provides more privacy of data than existing techniques while providing better classification accuracy. The experimental results show that classification accuracies have increased using EHPPDM technique.

Keywords:

privacy, privacy preserving data mining, k-anonymization, perturbation, l-diversity.

1. Introduction

Modern machine learning models are applied on large volumes of data accumulated over the past few years. Different data analysis models are built using this humongous data. The data used for training or building models may contain personal data. Data owners may not want to share their personal data. This paper deals with providing privacy for the personal data as well as performing data analysis without revealing personal data of users.

Privacy definition is given by different persons in different manner. Westin(1968) gave privacy definition as “ the assertion of individuals, groups or institutions to specify when, how and to what extent their information can be shared to others”. Bertin⁰ et al.(2008) define privacy as “ the security of data about an individual contained in an electronic repository from unauthorized disclosure”.

Privacy threats can be categorized into three types, namely (a) Membership Disclosure, (b) Attribute Disclosure and (c) Identity Disclosure.

Membership Disclosure: In this type of attack, an attacker is able to check whether an individual's data is present in a data set or not and can infer some meta-information about an individual.

Attribute Disclosure: In this type of attack, some sensitive information about an individual can be inferred by the attacker by linking data entries with some data from other sources.

Identity Disclosure: A specific data entry in a data set can be directly related to a particular person revealing his identity. An attacker can identify all the sensitive data about an individual. This type of attack is illicit and may have legal consequences.

Privacy preservation methods protect the data from data leakage by altering the original data and protect owner's exposure. There are various privacy preservation techniques specified in literature. They are randomization, perturbation, suppression, generalization etc. Data utility is defined as the quantity of important data preserved after altering the data. Various data utility metrics are available in literature. Some of them are discernability metric, KL-divergence, entropy based information loss etc.

The data is present in tabular form for processing. Each and every row represents an entity in the real world. The attributes of the data table can be categorized into four types. They are (i) Identifier Attributes (Ids), (ii) Quasi-identifier Attributes (QIDs), (iii) Sensitive Attributes (SAs), (iv) Non-sensitive Attributes (NSAs). The attributes which are used to identify an individual from given data are called identifier attributes. For ex: SSN, Aadhar id etc. Generally these kind of attributes are removed from data before

sharing the data for data analysis as they reveal individuals' identity. Sensitive attributes contain sensitive information about individuals like type of disease, salary etc. Generally individuals don't want to share sensitive data about them. But removing sensitive data and using remaining data for data analysis may not yield good results. So, the sensitive data needs to be maintained but the identity of the person needs to be hidden. Quasi identifiers are the attributes which can be used by attacker to disclose identity of individual when combined with some background knowledge.

These quasi-identifiers need to be modified to prevent identity disclosure by attackers. Non-sensitive attributes do not disclose any information about individuals. So, they are retained while sharing the data for data analysis purpose.

So, to provide privacy of data while sharing data for data analysis, several privacy preservation methods are proposed like randomization, perturbation etc.[1][2]. The data transformations are applied to provide privacy of data. But applying these data transformations may lead to inaccurate data mining results and also reducing utility of data. To balance both privacy preservation and accurate data mining results, Privacy Preserving Data Mining(PPDM) techniques are proposed. PPDM techniques ensure that the data is useful for data mining while preserving privacy of data and also utility of data is high. Utility of data can be defined as minimizing the divergence of data what the analysts see to the actual data. Several metrics are proposed to evaluate the privacy level and data utility of different PPDM techniques[3][4][5].

2. Review of PPDM Techniques

Data in a database can be anonymized by applying various privacy preserving techniques. Some of them are Generalization, Suppression, Anonymization and Perturbation.

- **Generalization:** In this method a data value is replaced with a more generalized one. For numerical attributes, a particular data value may be replaced with a range of values as a generalized one. For categorical attributes generalization is performed using a hierarchy. For example, engineer and lawyer are some of the data values for Occupation which can be replaced with a more generalized value of 'professional'.

- **Suppression:** This method prevents information disclosure by eliminating some attribute values. Generally replacing the original data value with("*").
- **Anonymization [5]:** In this, sensitive attributes and quasi identifiers are placed in two different Tables so that linking QIDs to sensitive attributes become very difficult.
- **Perturbation:** In this, original data values are replaced with synthetic values with the same statistical information.

Samarati and Sweeney [6], [7] proposed the most popular privacy model namely k-anonymization. According to [8] k-anonymity for a table is defined as follows [8]:

"Let $T(A_1, \dots, A_n)$ be a table.

Let QI be the set of quasi-identifiers corresponding to table T .

T fulfils k-anonymity property with respect to QI if and only if each sequence of values in $T[QI]$ appears at least with k occurrences in $T[QI]$ ".

Generalization and suppression techniques are applied on Quasi Identifiers(QIDs) as part of k-anonymization. All the QIDs in a group of size 'k' will have same values. This phenomenon ensures that the confidential data about individual users is not revealed when data is shared for analysis purpose. So, K-anonymized data provides privacy of data. An attacker can still infer sensitive information about individuals using K-anonymized table and some background knowledge, if the value of sensitive attribute is same for all individuals in a given k-group. For ex. Consider k-anonymized table shown below in table 1.

Table 1: 3-anonymized table

QI: Age	QI: city	Sensitive attribute: disease
20-30	mumbai	Flu
20-30	mumbai	Flu
20-30	mumbai	Flu
30-40	Delhi	Cancer
30-40	Delhi	Cancer
30-40	Delhi	Cancer

While k-anonymity is a promising approach to take for group based anonymization given its simplicity and wide array of algorithms that perform it, it is however susceptible to many attacks. When background knowledge is available to an attacker, such attacks become even more effective. Such attacks include:

- **Homogeneity Attack:** This attack leverages the case where all the values for a sensitive value within a set

Of k records are identical. In such cases, even though the data has been k -anonymized, the sensitive value for the set of k records may be exactly predicted.

-
- **Background Knowledge Attack:** This attack leverages an association between one or more quasi-identifier attributes with the sensitive attribute to reduce the set of possible values for the sensitive attribute. For example, Machanavajjhala, Kifer, Gehrke, and Venkatasubramanian (2007) showed that knowing that heart attacks occur at a reduced rate in Japanese patients could be used to narrow the range of values for a sensitive attribute of a patient's disease.

An attacker who has access to this 3-anonymous table can use background knowledge from other data sources and identify that all patients in Mumbai have disease 'Flu'. So, sensitive information about an individual residing Mumbai is revealed. To overcome this security breach l -diversity principle is applied on sensitive attribute. [9] defines l -diversity as being:

"Let a q^* -block be a set of tuples such that its non-sensitive values generalize to q^* . A q^* -block is l -diverse if it contains l 'well represented' values for the sensitive attribute S . A table is l -diverse, if every q^* -block in it is l -diverse."

Li et al [10] define l -diversity as being:

The l -diversity Principle – "An equivalence class is said to have l -diversity if there are at least l "well-represented" values for the sensitive attribute. A table is said to have l -diversity if every equivalence class of the table has l -diversity".

Machanavajjhala et. al.[11] define "well-represented" in three possible ways:

1. **Distinct l -diversity** – The simplest definition ensures that at least l distinct values for the sensitive field in each equivalence class exist.
2. **Entropy l -diversity** – The most complex definition defines *Entropy* of an equivalent class E to be the negative of summation of s across the domain of the sensitive attribute of $p(E,s) \log(p(E,s))$ where $p(E,s)$ is the fraction of records in E that have the sensitive value s . A table has entropy l -diversity when for every equivalent class E , $Entropy(E) \geq \log(l)$.
3. **Recursive (c - l -diversity)** – A comprehensive definition that ensures the most common value does not appear too often while less common values are ensured to not appear too infrequently.

Aggarwal and Yu (2008) note that when there is more than one sensitive field the l -diversity problem becomes more difficult due to added dimensions.

3. Methodology

Kundeti N et.al[12] proposed a hybrid privacy preserving data mining (HPPDM) technique which provides more privacy and lesser attacks. The technique can be extended with more privacy by applying l -diversity principle. l -diversity provides more privacy against different background attacks.

Algorithm (Enhanced Hybrid Privacy Preserving Data Mining(EHPPDM) Technique)

Input:- Adult Dataset D

Output:- Privacy enabled Adult Data set D'

Step1: Categorize attributes of Adult Data set into Identifiers, Quasi Identifiers, Sensitive and Non-Sensitive Attributes.

Step2: Consider the Quasi Identifiers and create value generalization hierarchies for quasi identifiers.

Step3: For numerical quasi identifiers apply geometric perturbation technique to obtain perturbed numerical quasi identifier.

Step4: For categorical quasi identifiers create generalization hierarchies and choose different levels in generalization hierarchy based on k -value chosen for anonymization.

Step5: For sensitive attributes apply l -diversity based on number of different values for class present.

Step 6: Obtain the privacy preserved Adult data set D'.

4. Implementation

Enhanced Hybrid Privacy Preserving Data Mining(EHPPDM) technique is implemented using R language. ARX anonymization tool is used for performing K -Anonymization.

UCI machine learning repository's Adult Dataset is used for evaluating EHPPDM technique. The dataset consists of 15 attributes including the Class attribute. The attributes are age(numerical), work-class(categorical), fnlwgt(numerical), education(categorical), education-num(numerical), marital-status(categorical), occupation(categorical), relationship(categorical), race(categorical), sex(categorical), capital-gain(numerical), capital-loss(numerical), hours-per-week(numerical), native-country(categorical) and class variable. These attributes can be divided into quasi-identifiers, sensitive

attributes and Insensitive attributes. The quasi identifiers in this data set are age, work class, education and native-country. Class attribute is sensitive attribute. Remaining attributes are classified as Insensitive attributes.

Among the quasi identifiers, age is the numerical attribute. Geometric data perturbation technique[13] is applied on numerical quasi identifier i.e. age. Value generalization hierarchies are created for categorical quasi identifiers. K-anonymization algorithm is applied to these categorical quasi identifiers. For different values of K, different anonymization levels are obtained, which provide privacy at different levels. The k-values considered are 50,100, 150,200, 250, 300, 350, 400, 450, 500. After anonymization, the anonymized data sets are applied with classification algorithms like naive bayes, J48 and decision tree. The accuracies of classification are noted down.

To enhance the privacy of data further, l-diversity is applied on sensitive attribute i.e. Class attribute. L-diversity is applied to reduce background attacks and linkage attacks. As l-diversity ensures that the class attribute value in a given anonymized group does not have single value, then the attacker can not identify an individual's sensitive attribute value. The anonymized and l-diversity applied dataset is obtained. Classification algorithms are applied on the anonymized data. Classification accuracies are noted down. Risk analysis for various types of attacks is given in following figures.

Fig.1 shows the classification accuracies for Adult data set when applied with k-anonymization. K-anonymization for different values of k is applied. Fig.2 shows the classification accuracies for Adult data set when l-diversity is applied to decrease background attacks. It is observed from results that classification accuracies have remained same and privacy is increased when l-diversity principle is applied.

Fig.3 shows the classification accuracies for Adult data set when applied with Hybrid Privacy Preserving Data Mining(HPPDM)[12] technique. It is observed that classification accuracies have increased when HPPDM technique is applied than k-anonymization.

Fig.4 shows the classification accuracies for Adult data set when Enhanced Hybrid Privacy Preserving Data Mining(EHPPDM) is applied.

Fig.5-8 show the risk analysis for Adult data set. Fig.5 shows the risk analysis against various types of attacks, after applying k-anonymization on Adult data set. Fig.6 shows the risk analysis against various types of attacks, after applying k-anonymization and l-diversity on Adult data set. Fig.7 shows the risk analysis against various types of attacks, after applying Hybrid Privacy Preserving Data Mining(HPPDM) technique on Adult data set. Fig.8 shows the risk analysis against various types of attacks,

after applying Enhanced Hybrid Privacy Preserving Data Mining(EHPPDM) technique on Adult data set. It is observed from results that risks have reduced to negligible levels when HPPDM and EHPPDM techniques are applied.

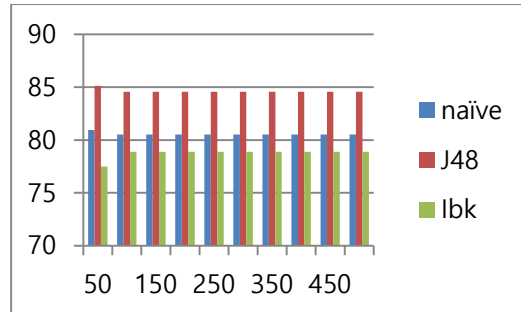


Fig.1 : Classification Accuracies for Adult K-anonymized Data for different k-values

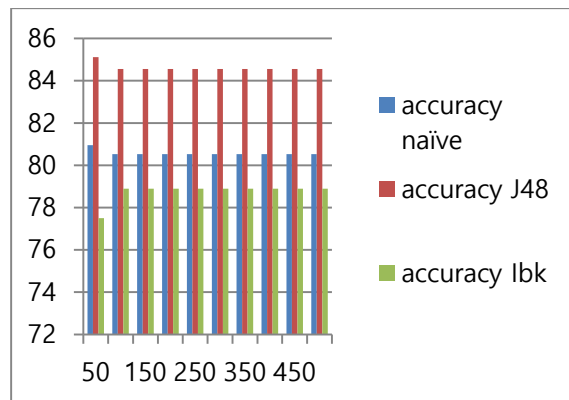


Fig.2: Classification accuracies for Adult K-anonymized and l-diversity(l-value=2) applied.

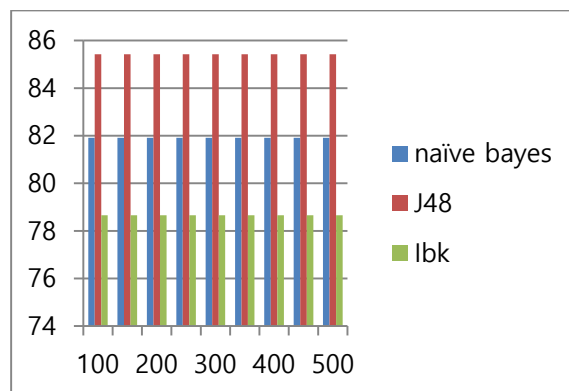


Fig.3: classification accuracies for adult after applying Hybrid Privacy Preserving Data Mining technique

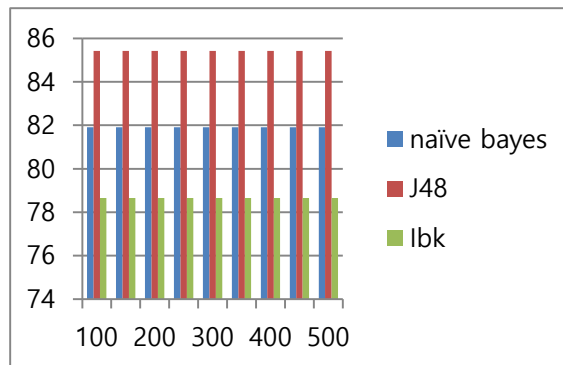


Fig.4 : Classification accuracies for Enhanced Hybrid Privacy Preserving Data Mining Technique.

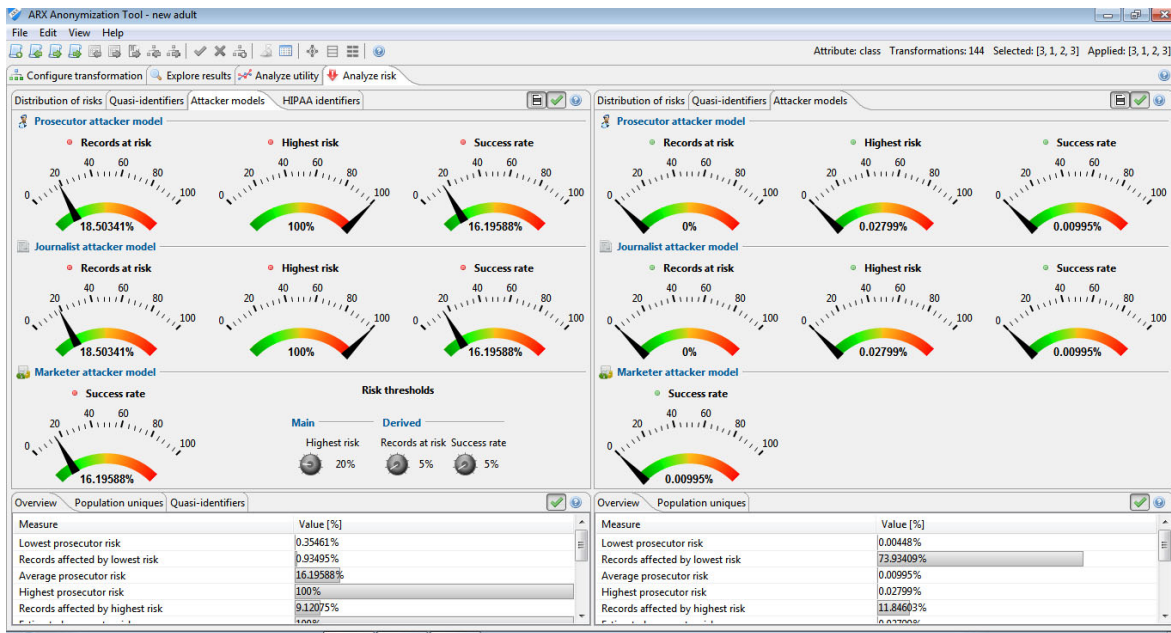


Fig.5 : Risk analysis for various types of attacks after applying k-anonymization(k-value=100)

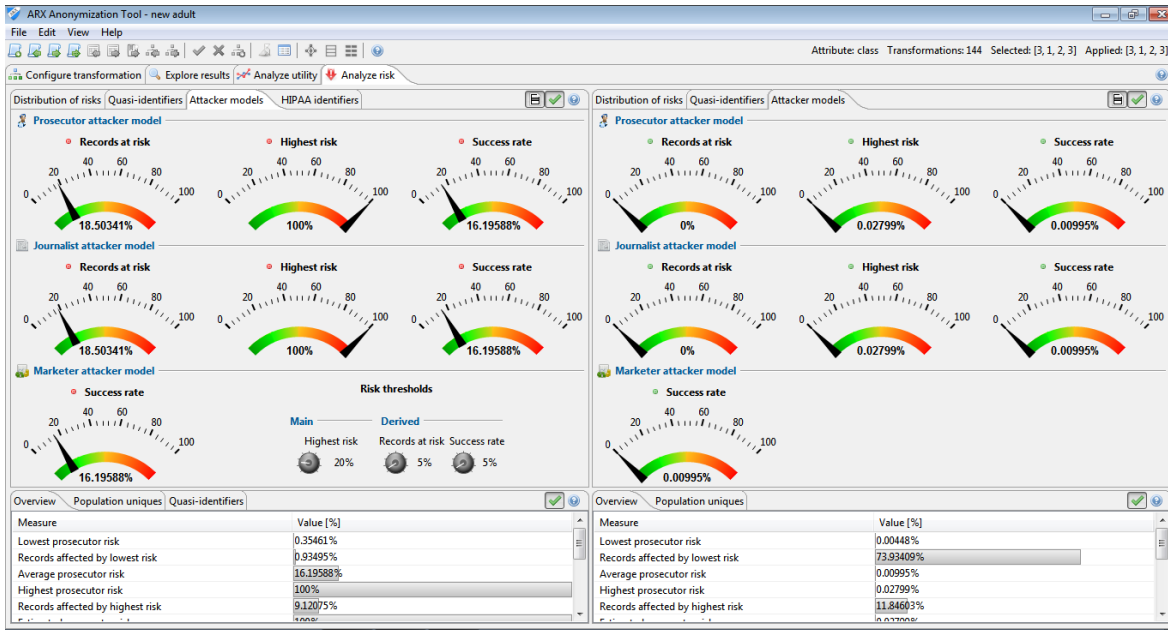


Fig.6: Risk analysis for various types of attacks after applying k-anonymization(k-value=100) and l-diversity(l-value=2).



Fig.7 : Risk analysis for various types of attacks after applying Hybrid Privacy Preserving Data Mining(HPPDM) technique for kvalue=100



Fig.8: Risk analysis for various types of attacks after applying Enhanced Hybrid Privacy Preserving Data Mining(EHPPDM) technique for kvalue=100, 1-2 diversity

5. Conclusion

This proposed Enhanced Hybrid Privacy Preserving Data Mining(EHPPDM) technique is applied on datasets from UCI machine learning repository. EHPPDM technique combines two privacy preservation techniques namely perturbation and k-anonymization. The numerical quasi identifiers are applied with geometric data perturbation and categorical quasi identifiers are applied with k-anonymization technique. To enhance privacy and reduce attacks l-diversity(lvalue=2) is applied to sensitive attribute. The experimental results show that classification accuracy has increased by applying EHPPDM technique. EHPPDM technique can be extended with t-closeness property in future works.

References

- [1] C. C. Aggarwal and P. S. Yu, "A general survey of privacy-preserving data mining models and algorithms," in *Privacy-Preserving Data Mining*. New York, NY, USA: Springer, 2008, pp. 11-52.
- [2] C. C. Aggarwal, "Data Mining: The Textbook", New York, NY, USA: Springer, 2015.
- [3] E. Bertino, D. Lin, and W. Jiang, "A survey of quantification of privacy preserving data mining algorithms," in *Privacy-Preserving Data Mining*. New York, NY, USA: Springer, 2008, pp. 183-205.
- [4] E. Bertino and I. N. Fovino, "Information driven evaluation of data hiding algorithms," in *Proc. Int. Conf. Data Warehousing Knowl. Discovery*, 2005, pp. 418-427.
- [5] S. Fletcher and M. Z. Islam, "Measuring information quality for privacy preserving data mining," *Int. J. Comput. Theory Eng.*, vol. 7, no. 1, pp. 21-28, 2015.
- [6] Pierangela Samarati and Latanya Sweeney, "Protecting Privacy when Disclosing Information: k- anonymity and its Enforcement Through Generalization and Suppression", *Proc. of the IEEE Symposium on Research in Security and Privacy*, pp. 384-393, 1998.
- [7] Pierangela Samarati and Latanya Sweeney, "Generalizing data to provide anonymity when disclosing information", in *PODS*, vol. 98, p. 188, 1998.
- [8] Pierangela Samarati, "Protecting respondents identities in microdata release", *IEEE Transactions on Knowledge and Data Engineering*, Volume: 13, issue 6, Nov/Dec 2001.
- [9] Aggarwal, Charu C.; Yu, Philip S. (2008). "A General Survey of Privacy-Preserving Data Mining Models and Algorithms" (PDF). *Privacy-Preserving Data Mining – Models and Algorithms*. Springer. pp. 11–52. ISBN 978-0-387-70991-8
- [10] Li, Ninghui; Li, Tiancheng; Venkatasubramanian, S. (April 2007). t-Closeness: Privacy Beyond k-Anonymity and l-Diversity. *IEEE 23rd International Conference on Data Engineering, 2007. ICDE 2007*. pp.106-115. [CiteSeerX 10.1.1.158.6171](https://doi.org/10.1.1.158.6171). doi:10.1109/ICDE.2007.367856. ISBN 978-1-4244-0802-3. S2CID 2949246.

- [11] Machanavajjhala, Ashwin; Kifer, Daniel; Gehrke, Johannes; Venkatasubramanian, Muthuramakrishnan (March 2007). "L-diversity: Privacy Beyond K-anonymity". *ACM Transactions on Knowledge Discovery from Data*. **1** (1):3–es. doi:[10.1145/1217299.1217302](https://doi.org/10.1145/1217299.1217302). ISSN 1556-4681. S2CID 679934
- [12] Kundeti Naga Prasanthi, Chandra Sekhara Rao MVP, "A Novel Method of Privacy Preserving Classification Mining Balancing Utility and Accuracy", *Journal of Advanced Research in Dynamical & Control Systems*, Vol. 11, Issue-05, 2019
- [13] Keke Chen, Ling Liu, "Geometric data perturbation for privacy preserving outsourced data mining", *Knowledge Information and Systems*, 2010
- [14] Bertino E., D. Lin and W. Jiang, "A survey of quantification of privacy preserving data mining algorithms", *Privacy-Preserving Data Mining*, Springer, ISBN: 978-0-387-70991-8, pp:183-205, 2008.
- [15] Langheinrich M., "Privacy in Ubiquitous Computing", *Ubiquitous Computing Fundamentals*, CRC Press, ISBN:9781420093605, pp:95-159, 2009.
- [16] Prasanthi, K.N., Chandra Sekhara Rao, M.V.P., "A Comprehensive Assessment of Privacy Preserving Data Mining Techniques", *Lecture Notes in Networks and Systems*, 351, pp. 833–842, 2022.



Kundeti Naga Prasanthi received her B.Tech. degree from Acharya Nagarjuna University in 2000 and M.Tech. degree from JawaharLal Nehru technological University in 2003. She received doctoral degree from Acharya Nagarjuna University in 2020. Her research interests include data mining, privacy preserving data mining, cybersecurity, Image processing.