



### AN EMPIRICAL APPROACH OF THE ACCUSATION OF UNREAL AND INCOMPLETE DATA USING LEARNING TECHNIQUES

S. Kanchana\*<sup>1</sup> and Dr. Antony Selvadoss Thanamani<sup>2</sup>

<sup>1</sup>Research Scholar, Research Department of Computer Science, NGM College, Pollachi 642001, Bharathiyar University, Coimbatore, India.

<sup>2</sup>Professor and Head, Research Department of Computer Science, NGM College, Pollachi 642001, Bharathiyar University, Coimbatore, India.

Article Received on 12/05/2016

Article Revised on 03/06/2016

Article Accepted on 24/06/2016

#### \*Corresponding Author

**S. Kanchana**

Research Scholar,  
Research Department of  
Computer Science, NGM  
College, Pollachi 642001,  
Bharathiyar University,  
Coimbatore, India.

#### ABSTRACT

Unreal and Incomplete data is a problem that focuses most important issue faced by researchers and practitioners who use industrial and research databases is incompleteness of data, usually in terms of missing or erroneous values. More or less of the data analysis algorithms can operate with incomplete data, a big share of the work require complete data. Therefore, variety of machine learning (ML) techniques are developed to reprocess the incomplete information. This

report centres on different imputation techniques and also proposes a supervised and unsupervised machine learning techniques Naïve Bayesian imputation method in MI model. The analysis is carried out employing a comprehensive range of databases, for which missing values were presented randomly. The goal of this report is to offer general guidelines for selection of suitable data imputation algorithms based on characteristics of the data.

**KEYWORDS:** Bayesian classifier, MI model, ML techniques, Supervised ML, Unsupervised ML.

#### I. INTRODUCTION

Missing data imputation is a realistic and challenging issue confronted by machine learning and information mining. Most of the real world datasets are characterized by an unavoidable problem of incompleteness, in terms of dropping values. Missing values may generate bias

and affect the caliber of the supervised learning procedure. Missing value imputation is an efficient means to detect or estimate the missing values based on other data in the data sets. Data mining consists of the various technical approaches, including machine learning, statistic and database system. The principal destination of the data mining process is to learn knowledge from large databases and transform into a human understandable format. This report concentrates on several algorithms such as missing data mechanisms, multiple imputation techniques and supervised machine learning algorithm. Experimental results are separately imputed in each real dataset and found out for accuracy.

The mechanism causing the missing information can determine the operation of both imputation and complete data methods. There are three different ways to categorize missing data as fixed in.<sup>[1]</sup> Missing Completely At Random (MCAR) lead to any particular data-item being missing are independent both of observable variables and of unobservable parameters. Missing At Random (MAR) is the alternative, suggesting that what caused the data to be missing does not depend upon the missing data itself. Not Missing At Random (NMAR) is data that is missing for a specific reason.

In the rest of this paper gives the background work or the literature review in section II, machine learning technique concepts in Section III, Section IV introduces new methods based on Naïve Bayesian Classifier to estimate and replace missing data. Experimental analyses of NBI model in Section V and the Conclusions are discussed in Section VI.

## II. LITERATURE REVIEW

Little and Rubin<sup>[1]</sup> summarize the mechanism of imputation method. Also introduces mean imputation<sup>[2]</sup> method to find out missing values. The drawbacks of mean imputation are sample size is overestimated, variance is underestimated, correlation is negatively biased. For median and standard deviation also replacing all missing records with a single value will deflate the variance and artificially inflate the significance of any statistical tests based on it. Different types of machine learning techniques are supervised and unsupervised machine learning techniques summarized in.<sup>[3]</sup> Classification of multiple imputation and experimental analysis are described in.<sup>[4]</sup> Min Pan et al.<sup>[5]</sup> summarize the new concept of machine learning techniques like NBI also analysis the experimental results which impute missing values. Comparisons of different unsupervised machine learning technique are referred from survey paper.<sup>[6]</sup> To overcome the unsupervised problem Peng Liu, Lei Lei et al.<sup>[7]</sup> applied the supervised machine learning techniques called Naïve Bayesian Classifier.

### III.MACHINE LEARNING APPROACH

In the data mining context, machine learning technique is generally classified as supervised and unsupervised learning technique both belong to machine learning technique.<sup>[8]</sup> Supervised classification focus on the prediction based on known properties and the classification of unsupervised focus on commonly used classification algorithm known as Naïve Bayesian imputation techniques.

#### UNSUPERVISED MACHINE LEARNING TECHNIQUES

*Mean Imputation* is the process of replacing the missing data from the available data where the instance with missing attribute belongs.

*Median Imputation* is calculated by grouping up of data and finding average for the data. Median can be calculated by finding difference between upper and lower class boundaries of median class.

*Standard Deviation* measures the spread of the data about the mean value. It is useful in comparing sets of data which may have the same mean but a different range. Estimate standard deviation based on sample and entire population data.

#### SUPERVISED MACHINE LEARNING TECHNIQUES

Another way of learning technique is classified as supervised learning that focus on the prediction based on known properties. Naïve Bayes technique<sup>[9]</sup> is one of the most useful machine learning techniques based on computing probabilities. It analyzes relationship between each independent variable and the dependent variable to derive a conditional probability for each relationship. A prediction is made by combining the effects of the independent variables on the dependent variable which is the outcome that is predicted. It requires only one pass through the training set to generate a classification model, which makes it very efficient. The Naïve Bayesian generates data model which consists of set of conditional probabilities, and works only with discrete data.

### IV.EVALUATION OF MULTIPLE IMPUTATION METHOD

Multiple imputations for each missing values generated a set of possible values, each missing value is used to fill the data set, resulting in a number of representative sets of complete data set for statistical methods and statistical analysis. The main application of multiple imputation<sup>[10]</sup> process produces more intermediate interpolation values, can use the variation

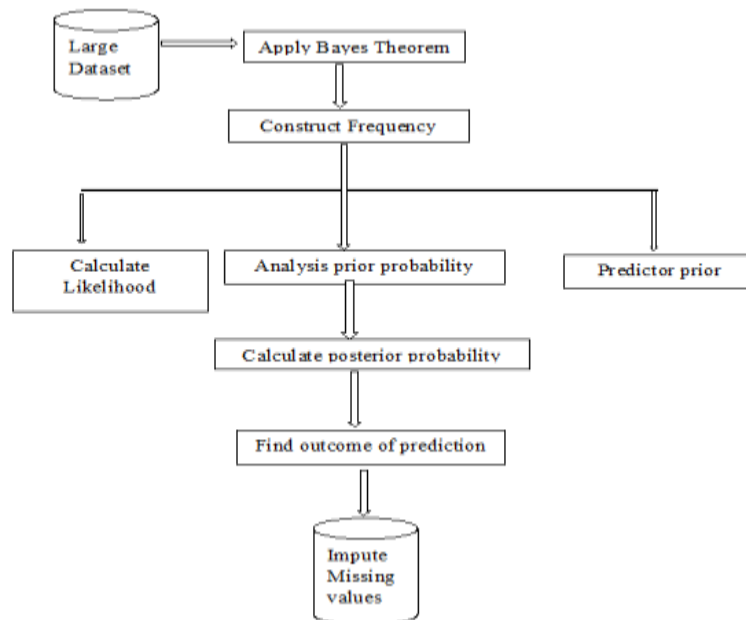
between the values interpolated reflects the uncertainty that no answer, including the case of no answer to the reasons given sampling variability and non- response of the reasons for the variability caused by uncertainty. Multiple imputation simulate the distribution that well preserve the relationship between variables. It can give a lot of information for uncertainty of measuring results of a single interpolation is relatively simple.

### **Naïve Bayesian Classifier(NBC)**

Naïve Bayesian Classifier is one of the most useful machine learning techniques based on computing probabilities.<sup>[11]</sup> It uses probability to represent each class and tends to find the most possible class for each sample. It analyzes relationship between each independent variable and the dependent variable to derive a conditional probability for each relationship. A prediction is made by combining the effects of the independent variables on the dependent variable which is the outcome that is predicted. NBC is a popular classifier, not only for its good performance, simple form and high calculation speed, but also for its insensitivity to missing data. It can build models on dataset with any amount of missing data. Naïve Bayesian Classifier generates full use of all the data in the present dataset. This paper focus a new method based on Naïve Bayesian classifier to handle missing data called Naïve Bayesian Imputation (NBI).

### **Naïve Bayesian Classifier Model**

The Naïve Bayesian classifier is based on Bayes' theorem with independence assumptions between predictors. Bayes theorem<sup>[12]</sup> provides a way of calculating the posterior probability  $P(C/X)$  of class from  $P(C)$  is the prior probability of class,  $P(X)$  is the prior probability of predictor and  $P(X/C)$  is the likelihood which is the probability of predictor given class. Naïve Bayes classifier assumes that the effect of the value of a predictor (X) on a given class (C) is independent of the values of other predictors called conditional independence. Figure 1 shows the pictorial representation of proposed system.



**Figure1: Flowchart of the Proposed System.**

### 1) *Algorithm for posterior probability*

- Construct a frequency table for each attribute against the target.
  - Transform frequency table to likelihood tables
  - Finally use the Naïve Bayesian equation to calculate the posterior probability for each class. The class with the highest posterior probability is the outcome of prediction.
- 2) *Zero-Frequency Problem*: When an attribute value doesn't occur with every class value adds 1 to the count for every attribute value class combination.
- 3) *Numerical Predictors*: Numerical variables need to be transformed to their categorical counterparts before constructing their frequency tables.

## V. EXPERIMENTAL RESULTS

### Design

Experimental datasets were carried out from the Machine Learning Database UCI Repository. Table1. describes the dataset with electrical impedance measurements in samples of freshly excised tissue dataset contains number of instances and number of attributes about the datasets used in this paper. The main objective of the experiments conducted in this work is to analyze the algorithm. Datasets without missing values are taken and few values are removed from it randomly. The rates of the missing values removed are from 5% to 25%. In these experiments, missing values are artificially imputed in attribute.

**Table1: Datasets used for Analysis.**

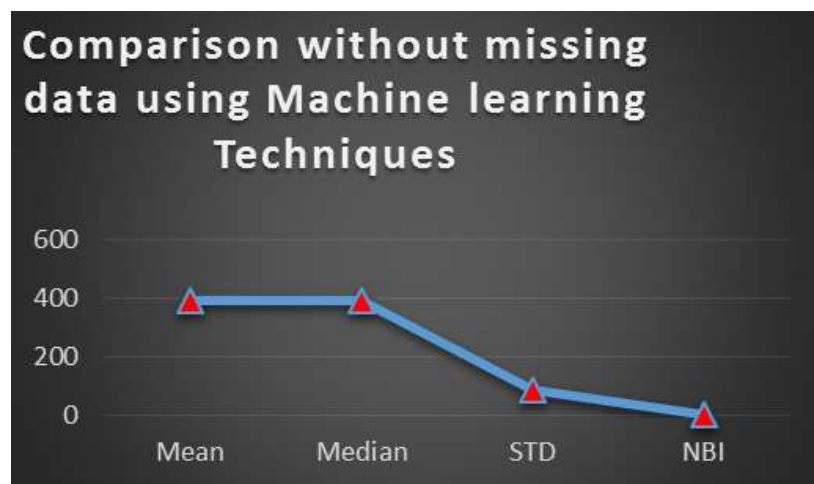
Datasets	Breast Tissue
Instances	106
Attributes	10 (9features + 1 classes)
Missing rates	5% to 25%
Unsupervised	Mean, Median, Standard Deviation
Supervised	Naïve Bayesian

**EXPERIMENTAL EVALUATION**

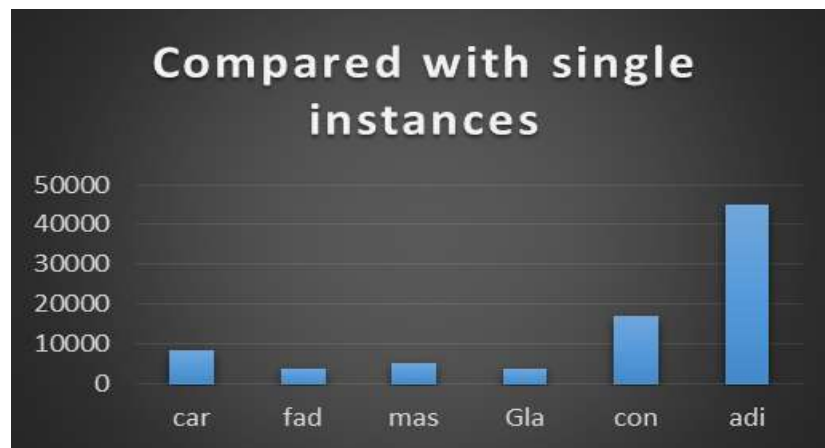
**Table 2:** describe the complete structure of all the attributes and classes without any missing values.

**Table 2: Original datasets without missing values.**

Class	I0	PA500	HFS	DA	Area	A/DA	Max II	DR	P
car	8279	4.62	3.87	3534	120186	673	1355	3213	10079.42
fad	3688	1.43	1.06	815.9	9152.7	150	345	717.7	4033.11
mas	5226	2.22	2	1319	19483	226	566	1145	5668.581
Gla	3813	1.87	1.53	645.6	6586.6	126	422	440	4184.044
con	16980	0.98	0.73	5152	74544	196	1021	5012	14909.9
adi	45145	1.62	2.96	8734	547574	1117	4281	7144	47052.58

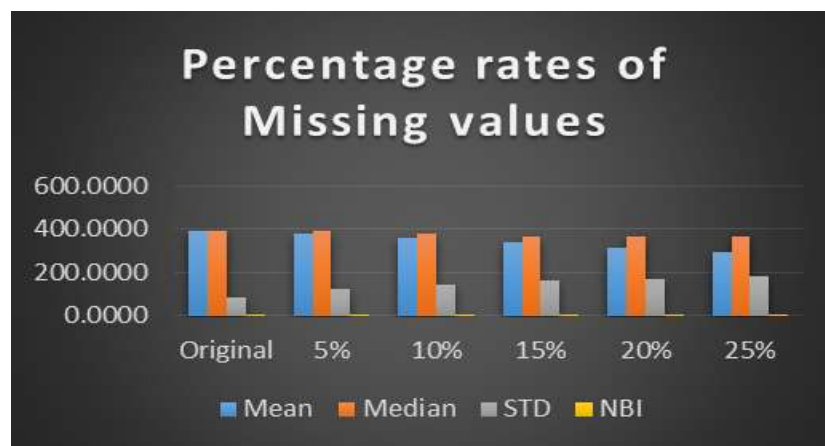
**Figure2: Original Datasets without missing values.**

The above Figure 2 represents the classification of all attribute of original dataset using both the machine learning techniques without missing values.



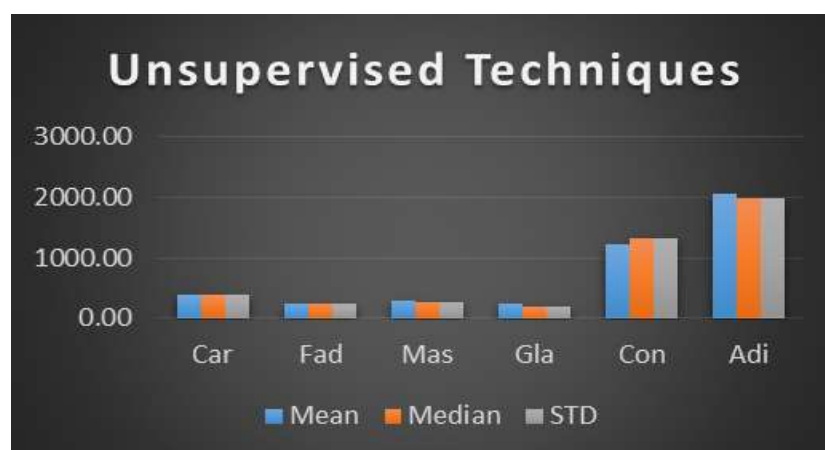
**Figure3: Single instance of original datasets.**

Figure 3: describes the single instance of Breast tissue dataset without missing values.



**Figure4: Missing value rates for experimental analysis.**

Figure 4 specifies the different percentage rates of missing values for experimental analysis.



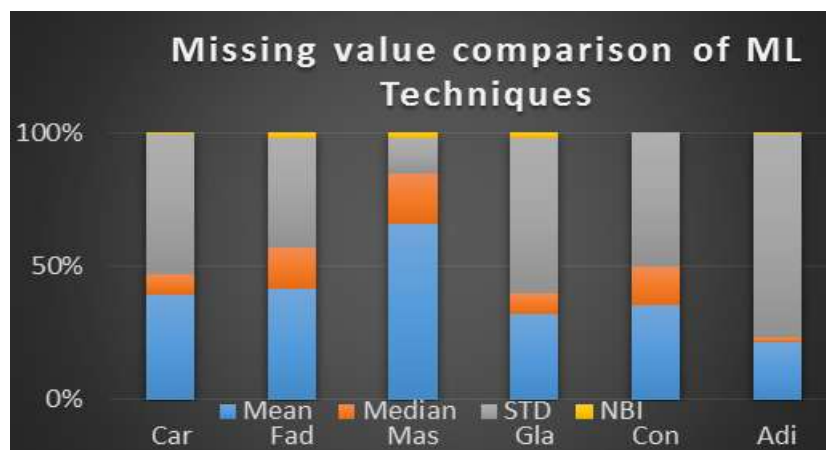
**Figure 5: Experimental results for Mean, Median and STD.**

Figure 5 & 6 represent the experimental results of both supervised and unsupervised machine learning techniques using missing value with the rate of 5%, 10%, 15%, 20% & 25% respectively.





**Figure 6: Experimental results for Supervised Techniques.**



**Figure 7: Comparative results using missing values for both ML Techniques.**

Figure 7 specifies the comparison of both ML techniques using missing value and Table 3 describes the percentage of missing value occur in the original dataset.

## VI. CONCLUSIONS

We did an observational evaluation of machine learning styles for alleging the non-existent data. Our investigation proves the complete view about the multiple imputation of missing values in large dataset. Single imputation technique generates bias result and affects the calibre of the execution. This paper focused multiple imputation using machine learning techniques of both supervised and unsupervised algorithms. The comparative study of mean, median, standard deviation in which standard deviation generates stable result in unsupervised algorithm. As well this report presents the observational resolution of standard deviation and Naïve Bayesian using less parameter for their analysis and the performance evaluation express among the other missing value imputation techniques the proposed method performs best. In time to come, it can be expanded to handle categorical attributes and it can be substituted by other supervised machine learning techniques.



**REFERENCES**

1. R. J. Little and D. B. Rubin. Statistical Analysis with missing Data, John Wiley and Sons, New York, 1997.
2. S. Kanchana, Dr. Antony Selvadoss Thanamani, "Classification of Efficient Imputation Method for Analyzing Missing values", International Journal of Computer Trends and Technology, 12(I): 2349-0829.
3. R.S. Somasundaram, R. Nedunchezian, "Evaluation on Three simple Imputation Methods for Enhancing Preprocessing of Data with Missing Values", International Journal of Computer Applications, May 2011; 21(10): 14-19.
4. S. Kanchana, Dr. Antony Selvadoss Thanamani, "Multiple Imputation Of Missing Data Using Efficient Machine Learning Approach", International Journal of Applied Engineering Research, 2015; 10(1): 0973-4562, 1473-1482.
5. Jeffrey C. Wayman, "Multiple Imputation for Missing Data: What is it and How Can I Use It?" Paper presented at the. Annual Meeting of the American Educational Research Association, Chicago, IL, 2003; 2-16.
6. Mrs.R. Malarvizhi, Dr. Antony Selvadoss Thanamani, "K-Nearest Neighbor in Missing Data Imputation", International Journal of Engineering Research and Development, November 2012; 5(1).
7. Alireza Farhangfar, Lukasz Kurgan and Witold Pedrycz, "Experimental Analysis of Methods for Imputation of Missing Values in Databases.
8. K. Lakshminarayan, S. A. Harp, and T. Samad, "Imputation of Missing Data in Industrial Databases", Applied Intelligence, 1999; 11: 259-275.
9. Peng Liu, Lei Lei, "Missing Data Treatment Methods and NBI Model", Sixth International Conference on Intelligent Systems Design and Applications, 0-7695-2528-8/06.
10. S.Hichao Zhang, Jilian Zhang, Xiaofeng Zhu, Yongsong Qin, Chengqi Zhang, "Missing Value Imputation Based on Data Clustering", Springer-Verlag Berlin, Heidelberg, 2008.
11. Blessie, C. E., Karthikeyan, E, Selvaraj. B. NAD – A Discretization approach for improving interdependency, Journal of Advanced Research in Computer Science, 2010; 9-17.
12. R. Kavitha Kumar and Dr. R. M. Chandrasekar, "Missing Data Imputation in Cardiac data set".
13. Ingunn Myrtveit, Erik Stensrud, "IEEE Transactions on Software Engineering", November 2001; 27(11).