

CHAPTER 7

Methods and Comparisons of Handling Missing Data in Student Data Analysis using Rough Set and Soft Set

7.1. Introduction

While conducting educational study, observations are repeatedly measure in the same subject and at each schedule visit. The educational trials may sometimes face the problems of missing observations and resultant information. The issue of missing data arises from bringing together of class test marks, percentage of the student, assignments, attendance, and grading of the students. It may further lead to serious problems like reducing the statistical power and the inability to rule out biases in the estimates. The problem is so serious that an immediate and complete solution can not be reached in statistical practice. The traditional statistical methods are marred by intentional faulty designs. Special attention should be directed to conduct the incomplete data details.

7.2. Background Literature

(Little and Rubin 2000)⁸¹ have classified missing data mechanisms into three different types based on the possible causes:

(i) missing completely at random (MCAR) if the messiness is not related to any observed and unobserved factors (such as domestic relocation, suffering an accident, or unrelated illness);
(ii) missing at random (MAR) if the missingness is conditional on observed factors and is independent of the unobserved data (such as lack of efficacy); and
(iii) missing not at random (MNAR) if the missingness depends on unobserved quantities as well as some observed factors. The MNAR missing mechanism is usually used to describe the students result. Researchers have pointed out that the MAR assumption may be more in practice than that of the MCAR (Collins, L.M., et.al 2001)¹⁵. In fact, by definition MCAR is only a special case of MAR. In other words, a MCAR missing mechanism is also a MAR one, but not every MAR is a MCAR. Actually, it is possible to formally test the MCAR assumption against its alternative hypothesis not MCAR (Little, R.J.A. 2000)⁸¹ and (Diggle, P.J., et.al 2002)²². However, it is not possible to test MAR or MNAR without using additional (external) information. MNAR is particularly useful in assessing the sensitivity of the results that are not MAR and it is highly recommended to be incorporated into the analysis.

In this literature a number of approaches have been applied in the statistical analysis of educational data with missing values. These approached should be applied basing on appropriate methods along with the data missing mechanism. The applications of different statistical approaches are valid only under selected situations (missing mechanism) in terms of specified missing rates to put it the other way. None of the methods is unique to be used for all situations. In the analysis of educational study it is not easy to test the missing mechanism as there is no clear rule to specify as to how much available missing data is in access to the required study. If the study is directed in its pragmatic and a more explanatory dimensions a particular method of handling missing data is adopted. Too much of missing data⁷ can stand as a problem for the study.

(Sprint and Dupin-Sprint 1993)¹⁴³ have prescribed an amount of missing data which can be tolerable excepted and it would not lead to an effect in the opposite direction from the performance of the

“*worst case*” analysis of missing data it can be determined whether the right level of missing data has been reached.

Despite these difficulties, several researchers have considered and constructed simulation studies for the proof of strong consistency of imputation methods to check the efficiency of the imputation methods. For example, (Myers 2000)⁹¹ compared the results of two imputation methods (that is, the complete case method and the multiple imputation method) based on simulated data sets with a dropout rate ranging from 20% to 60%, and they concluded that MI method provided results that are more closely mimicked the complete data set.

(Hening and Koonce 2014)⁴⁹ investigated five imputation methods (i.e., mean substitution, median substitution, zero value, hot-deck, and MI) and a first-year-student retention data with more than 20% missing values is used. The results shown that multiple and hot-deck imputations perform poorly in an accuracy comparison test, but they can slightly increase the predication accuracy rate compared with other methods.

(Ali, et al. 2011)⁵ performed a survival analysis in which missing data were simulated under MCAR and MAR to compare four imputation methods complete case analysis (CCA), means substitution (MS), and multiple imputation (MI) with the inclusion of the outcome (MI- and MI+). The simulation results suggested that in general MI+ is likely to be the best method. (Patrician) pointed out that MI is the best approach and should be considered to handle missing data compared with CCA and MS by an empirical investigation of AIDS care longitudinal data outcomes.

Recently, (Nakai, et al. 2014)⁹⁴ have shown that MI is the most effective imputation method in longitudinal data setting under MCAR via a simulation study. This indeed provides useful information about the performance of imputation methods under MCAR, but it is limited and restricted to clinical situations where MAR is more plausible. For example, (Lavori, et al. 1995)⁷⁷ have pointed out that the MCAR assumption is often not plausible in most clinical trial settings. The purpose of this paper is through a simulation approach to analytically evaluate the performance of four imputation methods for different missing mechanisms (MCAR and MAR) with various missing rates.

For simplicity and also without loss of generality, a monotone pattern of missing data (meaning that once a patient has a missing response at an assessment visit, his or her data will be missing for all subsequent visits) is assumed. Under such assumptions, this paper primarily concentrates on the following four imputation methods:

- (i) Complete delete case (CDC)
- (ii) Last observation carried forward (LOCF)
- (iii) Rough Set Techniques (RST) and
- (iv) Soft Set.

7.3. Approaches to Handling Missing Data

There are so many techniques in handling missing data discussed in the literature. Especially, many methods have been proposed and developed to handle missing data in longitudinal clinical trials. However, there are few methods that are actually used in real trials with missing data. The purpose of this paper is to study four most frequently used methods for dealing with missing data and they will be described as follows.

To compare the performance of these methods, RST and Soft Sets are used as evaluation criteria.

Student No.	Class Test 1	Class Test 2	Assignment 1	Assignment 2	Lab test	Result
S01	First	Second	?	First	First	First
S02	Second	First	Second	First	Third	Second
S03	Third	First	Third	Third	Third	Third
S04	First	?	First	First	Second	First
S05	Second	Second	Second	Second	Second	Second
S06	Second	Second	?	Third	?	Second
S07	Third	?	Third	Second	Third	Third
S08	Second	Second	Second	First	First	Second
S09	Second	Third	Second	Second	Second	Second
S10	First	Second	First	?	First	First
S11	Second	?	Second	Third	Third	Second
S12	Second	Second	Third	?	Second	Second

Table-7.1: Student information system with missing values

7.3.1 Rough Set Analysis

Rough set theory introduced by (Pawlak in 1982)¹⁰³ is a mathematical tool to deal with vagueness and uncertainty of information. It has been proved to be very effective in many practical applications. The rough set uses the basic relation operations known as Equivalence relation, which is reflexive, symmetric and transitive. Using this relation the data table is classified according to the attributes and this classification is analysed using core and reduct of the relation.

- $S = (R, X, Y)$ is independent if all are indispensable in $x \in X$.
- The set of attributes is called a reduct of X , if $S' = (R, Z, Y)$ is independent and $POS_Z(Y) = POS_X(Y)$
- The set of all the condition attributes indispensable in S is denoted by $CORE(X)$.

$$CORE(X) = \cap RED(X)$$

where $RED(X)$ is the set of all reducts of X

However, in rough set theory, the deterministic mechanism for the description of error is very simple (Grzymala-Busse, J.W et.al 2001)⁴². Therefore, the rules generated by rough sets are often unstable and have low classification accuracy. Rough set methods can greatly accelerate the network training time and improve its prediction accuracy. In (Grzymala-Busse, J.W., 2003)⁴⁴ Rough set method was also applied to generating rules from trained neural networks. In these hybrid systems, rough sets were used only as a tool to speed up or simplify the process of mining knowledge from the databases and impute the missing values in the database.

For example, in (Grzymala-Busse, J. W., 2004)⁴⁵, a rule set, a part of knowledge, is first generated from a database by rough sets. In the prediction phase, a new object is first predicted by the rule set, if it does not match any of the rules, the model can get high classification accuracy. In this paper, from a new perspective we develop a hybrid system of rough sets and neural networks to mine classification rules from large databases. Compared with previous research works our study has the following contributions.

An information system is a 10-tuple $S = \{U, A, V, F \dots\}$ where U is a finite set of objects, called the universe, A is a finite set of attributes, $V = \bigcup_{a \in A} V_a$ is a domain of attribute a and $f : U \times A \rightarrow V$ is called an information function such that $f(x, a) \in V_a$ for $a \in A, x \in U$.

In the classification problems, an information system is also seen as a decision table assuming that $A = C \cup D$ and $C \cap D = \emptyset$ where C is a set of condition attributes and D is a set of decision attributes.

Let $S = \{U, A, V, F \dots\}$ be an information system, every $P \subseteq A$ generates a indiscernibility relation $IND(P)$ on U , which is defined as follows:

$$IND(P) = \{(x, y) \in U \times U : f(x, a) = f(y, a) \forall a \in P\}$$

Attribute rule generation (feature selection) is a process of finding an optimal subset of all attributes according to some criterion so that the attribute subset is good enough to represent the classification relation of data. A good choice of attribute subset provided to a classifier can increase its accuracy, save the computational time, and simplify its results.

In general, rough set theory provides useful techniques to reduce irrelevant and redundant attributes from a large database with a lot of attributes (Schafer, J.L. 1997 and 2000)^{136 and 137}. However, it is not so satisfactory for the reduction of noisy attributes because the classification region defined by rough set theory is relatively simple and rough set based attribute rule generation criteria lack effective validation method.

There are different approaches to generate rules, direct and indirect methods. Direct methods generate rule from training data like sequential covering algorithms. Indirect methods build the classification model from which rule are extracted, e.g. decision tree, Neural Network, Genetic Algorithms etc.

From the Table-1 the missing attribute classes can be generated using rough set concept as follows.

R1: Class Test-1 = $\{\{S01, S04, S10\}, \{S02, S05, S06, S08, S09, S11, S12\}, \{S03, S07\}\}$

R2: Class Test-2 = $\{\{S02, S03\}, \{S01, S05, S06, S08, S10, S12\}, \{S09\}, \{S04, S07, S11\}\}$

R3: Assignment-1 = $\{\{S04, S10\}, \{S02, S05, S08, S09, S11\}, \{S03, S07, S12\}, \{S01, S06\}\}$

R4: Assignment-2 = $\{\{S01, S02, S04, S08\}, \{S05, S07, S09\}, \{S03, S06, S11\}, \{S10, S12\}\}$

R5: Lab Test = $\{\{S01, S08, S10\}, \{S4, S05, S09, S12\}, \{S02, S03, S07, S11\}, \{S6\}\}$

In the above classification the classes are First, Second, Third and the bold classes are the missing data, which will be filled in the common attribute rough set technique.

The bold classes are replacement of approximated information with the missing data in the Table-7.2. After replacing proper common values to the missing data, we get the different classes from the table-7.1 as follows.

Student No.	Class Test 1	Class Test 2	Assignment 1	Assignment 2	Lab test	Result
S01	First	Second	Second	First	First	First
S02	Second	First	Second	First	Third	Second
S03	Third	First	Third	Third	Third	Third
S04	First	Second	First	First	Second	First
S05	Second	Second	Second	Second	Second	Second
S06	Second	Second	Second	Third	Second	Second
S07	Third	Second	Third	Second	Third	Third

S08	Second	Second	Second	First	First	Second
S09	Second	Third	Second	Second	Second	Second
S10	First	Second	First	First	First	First
S11	Second	Second	Second	Third	Third	Second
S12	Second	Second	Third	First	Second	Second

Table-7.2: Student information system with RST analysis

7.3.2 Soft Set Analysis

The models such as theory of probability, Interval mathematics, fuzzy set theory, Intuitionistic fuzzy set theory are inadequate to handle some uncertainty problems. (D.A. Moldtsov 1999)⁹⁰ and (Majji et.al 2003)¹⁰⁰, noticed that, the problem might be due to inadequacy of parameterization tools in those models. So, Molodtsov initiated the concept of Soft Set theory which is a fusion of the notions of topology and set theory as a new mathematical tool to deal with uncertainties and free from some previous difficulties. Soft Sets can be called as (Binary, Basic, and Elementary) neighborhood systems. Soft set gives a general mathematical tool to deal with uncertain, fuzzy, not clearly defined (vague) objects.

Let U be an initial universal set and let E be a set of parameters. A pair (F, E) is called a *soft set* (over U) iff F is a mapping of E into the set U . The pair (U, E) is often regarded as a soft universe. Members of the universe and the parameter set are generally denoted by x and e respectively. Let A be the subset of E . A soft set over the soft universe (U, E) is denoted by (F, A) , where $F: A \rightarrow P(U)$. In other words, the soft set is a parameterized family of subsets of the set U . Every set $F(e)$, $e \in E$, from this family may be considered as the set of α -approximate elements of the soft set. The Sets of $F(e)$ may be arbitrary. Some of them may be empty, some may have non-empty intersection.

Analysis:

Minimum threshold ≥ 1 for missing data.

Case 1: The no. of missing data =1 Classify according to attributes,

Class Test1= $\{\{1,4,10\}, \{6,11,12\}, \{7\}\}$ Result= $\{\{1,4,10\}, \{6, 11,12\}, \{7\}\}$

Lab Test= $\{\{6\}, \{1,10\}, \{4,12\}, \{7,11\}\}$

Max possibility for S06, is in Result and Class Test 1= $\{6, 11,12\}$. So the missing value is **Second**.

Case 2: The no. of missing data =2 Classify according to attributes,

Class Test1= $\{\{1,4,10\}, \{6,11,12\}, \{7\}\}$ Result= $\{\{1,4,10\}, \{6, 11,12\}, \{7\}\}$

Assignment 1= $\{\{1,6\}, \{4,10\}, \{7,12\}, \{11\}\}$

Max possibility for S01, is in Result and Class Test 1= $\{1,4,10\}$. So the missing value is **First**

Max possibility for S06, is in Result and Class Test 1= $\{6,11,12\}$. So the missing value is **Second**

Assignment 2= $\{\{10,12\}, \{1,4\}, \{6,11\}, \{7\}\}$

Max possibility for S10, is in Result and Class Test 1= $\{1,4,10\}$. So the missing value is **First**

Max possibility for S12, is in Result and Class Test 1= $\{6,11,12\}$. So the missing value is **Second**

Case 3: The no. of missing data =3 Classify according to attributes,

Class Test1= $\{\{1,4,10\}, \{6,11,12\}, \{7\}\}$ Result= $\{\{1,4,10\}, \{6, 11,12\}, \{7\}\}$

Class Test 2= $\{\{4, 7,11\}, \{1,6,10,12\}\}$

Max possibility for S04, is in Result and Class Test 1= $\{1,4,10\}$. So the missing value is **First**.

Max possibility for S07, is in Result and Class Test 1= $\{7\}$. So the missing value is **Third**.

Max possibility for S11, is in Result and Class Test 1= $\{6,11,12\}$. So the missing value is **Second**.

Result: From the above analysis using Rough Set and Soft Set, it is found that there is no certainty to impute the correct missing value and we have use our own intuitionistic knowledge to select the value. So it is necessary to use some other techniques to replace the missing data with an appropriate value. Some of these techniques are discussed in the following sections.

Student No.	Class Test 1	Class Test 2	Assignment 1	Assignment 2	Lab test	Result
S01	First	Second	First	First	First	First
S02	Second	First	Second	First	Third	Second
S03	Third	First	Third	Third	Third	Third
S04	First	First	First	First	Second	First
S05	Second	Second	Second	Second	Second	Second
S06	Second	Second	Second	Third	Second	Second
S07	Third	Third	Third	Second	Third	Third
S08	Second	Second	Second	First	First	Second
S09	Second	Third	Second	Second	Second	Second
S10	First	Second	First	First	First	First
S11	Second	Second	Second	Third	Third	Second
S12	Second	Second	Third	Second	Second	Second

Table-7.3: Analysis using Soft Set

7.3.3. Complete Delete Case (CDC) Analysis

This method deletes all cases with missing data and then performs statistical analyses on the remaining complete data set (which has a smaller sample size). Since all cases containing missing data have been removed, there is no missing data problem to handle. Therefore, all statistical methods can be used to analyze the smaller data set. Obviously, one major advantage of this method is its ease of use. In fact, virtually all statistical programs incorporate this method as a default method because it accommodates any type of statistical analysis (Allison, P.D. 2001)⁶. The method may be preferred under the situation in which the sample size is large, the proportion of missing data is small, and the missing data mechanism is MCAR (Kim, J.O. and Curry, J. 1977)⁶⁸. For MCAR missing data, the method will yield unbiased parameter estimates and larger standard errors due to the smaller sample size. However, even when data are MCAR, loss of data will result in loss of precision (larger standard errors), particularly in multivariate data analyses.

In general, the major disadvantage of the method is that it could possibly lead to losing statistical power due to the reduction of the sample size (Allison, P.D. 2001)⁶.

Student No.	Class Test 1	Class Test 2	Assignment 1	Assignment 2	Lab test	Result
S02	Second	First	Second	First	Third	Second
S03	Third	First	Third	Third	Third	Third
S05	Second	Second	Second	Second	Second	Second
S08	Second	Second	Second	First	First	Second
S09	Second	Third	Second	Second	Second	Second

Table-7.4: Student information system with CDC analysis

7.3.4. Last Observation Carried Forward (LOCF)

The simplest imputation approach is the LOCF method that replaces every missing value with its corresponding last observed value. LOCF method is often used in longitudinal studies of continuous outcomes under MCAR. Conceptually, this method assumes that the outcome would not change after the last observed value. Therefore, there is no time effect since the last observed data. In fact, LOCF has been a popular method that is frequently used in handling missing data problems because it is easy to understand and can be implemented easily as well. Also, unlike the CDC method, the sample size does not change. For example, in an educational trial (see the data below), the bold words are newly imputed in missing attributes.

If there are more attributes missing then this method might give a biased conclusion about the effect of the student group. In our example, the measurement of student 1,4,6,7,10,11 and 12 are missing randomly. After impute the missing data the following table 7.5 shows the bold letters.

Student No.	Class Test 1	Class Test 2	Assignment 1	Assignment 2	Lab test	Result
S01	First	Second	Second	First	First	First
S02	Second	First	Second	First	Third	Second
S03	Third	First	Third	Third	Third	Third
S04	First	First	First	First	Second	First
S05	Second	Second	Second	Second	Second	Second
S06	Second	Second	Second	Third	Third	Second
S07	Third	Third	Third	Second	Third	Third
S08	Second	Second	Second	First	First	Second
S09	Second	Third	Second	Second	Second	Second
S10	First	Second	First	First	First	First
S11	Second	Second	Second	Third	Third	Second
S12	Second	Second	Third	Third	Second	Second

Table-7.5: Student information system with LOCF analysis

Rigorously speaking, LOCF is not an analytic approach, but it is a method that is very easy to impute missing values. Analytic proofs (Molenberghs, G., 2004)⁸⁹ and studies in simulated data (Shao, J. and Zhong, B. 2003)¹³⁸ and (Mallinckrodt, C.H., 2001)⁸⁷ have been clearly shown that LOCF can bias results and lead to either overestimation or underestimation of the parameter estimates.

7.4. Missing Data Generation

After the original data sets were created, the measurements at different time points for different subjects were set to missing, according to the MCAR or MAR missing mechanism. However, the measurement at the first time point of each subject was assumed always observed. In the MCAR setting, missing data were generated randomly at visits 2 through 5 based on the missing probabilities listed in Table 7.1. Therefore, the missing probabilities do not depend on either observed or unobserved data. Furthermore, Little's MCAR test was performed to make sure the missing mechanism is indeed MCAR otherwise that data set was discarded and another data set was generated anew.

7.5. Simulation Results

The simulation results are summarized in Table-7.6. In this table, Original data set, RST analysis, SS analysis and LOCF analysis are shown.

In the following Comparison, the numbers 1, 2 and 3 are represented as First, Second and Third respectively.

Method	Student No.	Class Test 1	Class Test 2	Assignment 1	Assignment 2	Lab test	Result	Maching %
Original Data	S01	1	2	1	1	1	1	
	S04	1	1	1	1	2	1	
	S06	2	2	2	3	3	2	
	S07	3	2	3	2	3	3	
	S10	1	2	1	1	1	1	
	S11	2	3	2	3	3	2	
	S12	2	2	3	2	2	2	
Rough Set Analysis	S01	1	2	2	1	1	1	38%
	S04	1	2	1	1	2	1	
	S06	2	2	2	3	2	2	
	S07	3	2	3	2	3	3	
	S10	1	2	1	1	1	1	
	S11	2	2	2	3	3	2	
	S12	2	2	3	1	2	2	
Soft Set Analysis	S01	1	2	1	1	1	1	62%
	S04	1	1	1	1	2	1	
	S06	2	2	2	3	2	2	
	S07	3	3	3	2	3	3	
	S10	1	2	1	1	1	1	
	S11	2	2	2	3	3	2	
	S12	2	2	3	2	2	2	
LOCF Analysis	S01	1	2	2	1	1	1	50%
	S04	1	1	1	1	2	1	
	S06	2	2	2	3	3	2	
	S07	3	3	3	2	3	3	
	S10	1	2	1	1	1	1	
	S11	2	2	2	3	3	2	
	S12	2	2	3	3	2	2	
CDC Analysis	No Results came as it is removing all the tuples with missing values retaining the non-missing tuples.							0%

Table-7.6: Comparison of different Methods

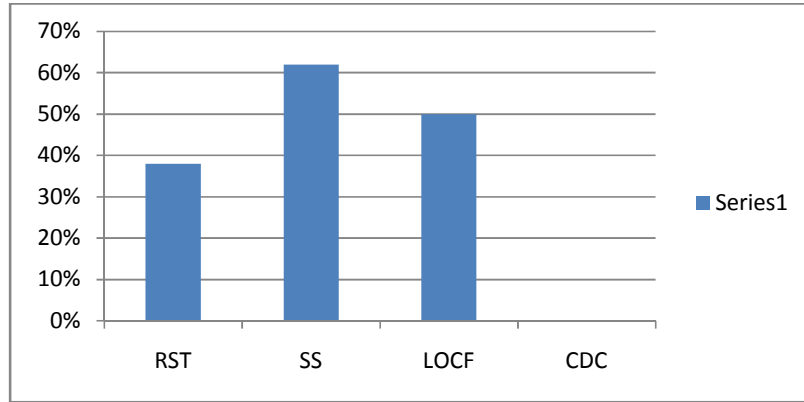


Fig-7.1. Comparison of different Methods

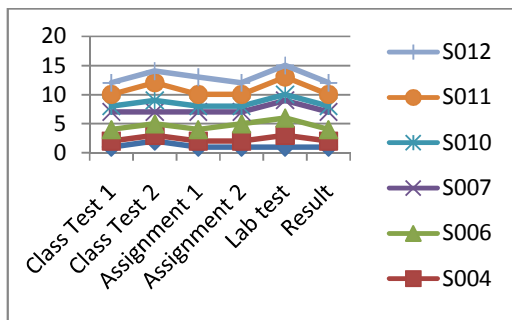


Fig-7.2. Original Data Set

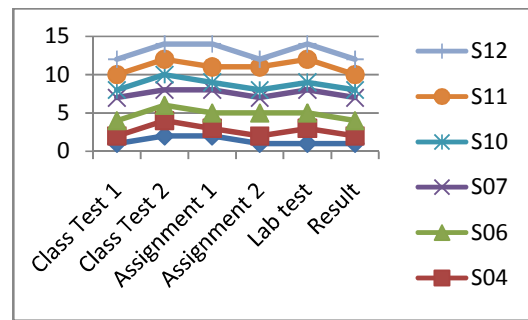


Fig-7.4. RST Analysis

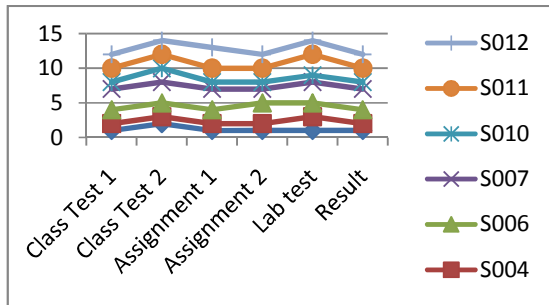


Fig-7.3. LOCF Analysis

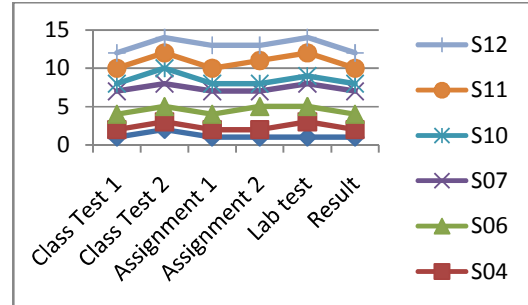


Fig-7.5. Soft Set Analysis

7.6. Discussion and Conclusions

The simulation results suggested that there is no one single method that is the best under all situations. In the above experiments the original student data set contains 12 rows and 7 attributes; in this the data set randomly missed 8 attributes. In CDC method the missing data are completely deleted and result data set will have 5 rows. While SS method was superior to LOCF and RST also the LOCF method is better than RST because in SS analysis the missing data rate matched 62.5% whereas in LOCF analysis the missing data matching is 50% and RST analysis matching is only 38%.

This paper discusses different methods to impute the missing values. Missing values are replaced by probability distributions over possible values for the missing feature, which allows the corresponding transaction to support all item sets that could possibly match the data. Transactions which do not

exactly match the candidate item set may also contribute a partial amount of support this behavior is beneficial for databases with many missing values or containing numeric data. Handling missing values using the most probable information for all the samples belonging to the same class gives better result as compare to other techniques. Missing values filled with better accuracy leads to better results, this phenomenon is also observed.