



# COMPARISON ON TO HANDLING THE MISSING VALUES BY USING PRE- PROCESSING IN EMBEDDING TO APPLYING PCA TECHNIQUES

S. Muruganandam\* & S. Subbaiah\*\*

\* Research Scholar, PG and Research Department of Computer Science and Applications,  
Vivekanandha College of Arts and Sciences for Women (Autonomous),  
Tiruchengode, Tamilnadu.

\*\* Assistant Professor, Sri Krishna Arts and Science College, Coimbatore, Tamilnadu

## Abstract:

*As of now, information mining is one of the regions of incredible intrigue since it permits to find covered up and frequently fascinating examples with regards to enormous volumes of information. This present reality enormous datasets are gotten from numerous sources and contain information that will, in general, be deficient, boisterous and conflicting. In this unique circumstance, it is critical to get ready crude information to meet the prerequisites of information mining calculations. This is the job of information pre-handling stage, in which information cleaning, change, and coordination, or information dimensionality decrease are performed. Hence in this process has detected the missing values in data cleaning by using the data pre-processing techniques. It finds and detects and removes the tuples without noise and replaces to finds the value which helps of normalization which helps of data transformation. Data reduction used to perform dimensionality reduction in the pre-processing methods. It aggregates to the selective data sets through selective processing and reduces the size with supports of numerosity respectively. It overcomes the existing drawbacks which support embedding in data mining techniques to provide best solution to dimensionality reduction.*

## Keywords:

***Pre-processing, DR – Dimensionality Reduction, Normalization, Data Cleaning, Data Aggregation, Data Transformation.***

## 1. INTRODUCTION:

Information mining is the procedure of the extraction of helpful examples and models from a gigantic dataset. These models and examples have a successful job in a basic leadership task. Information mining essentially relies upon the nature of the information. The crude information normally used to identify the helpless to missing qualities, loud information, deficient information, conflicting information, and anomaly information. In this way, it is significant for this information to be handled before being mined. Pre-handling information is a fundamental advance to improve information productivity. Information pre-handling is perhaps the most datum mining step which manages information planning and change of the dataset and looks for simultaneously to make learning revelation increasingly

effective. Pre-preparing incorporates a few methods like cleaning, reconciliation, change, and decrease. It demonstrates a point by point portrayal of information pre-handling procedures which are utilized for information mining.

Information Pre-managing is required and it is a basic stage in Bioinformatics and Web Usage Mining. Data Cleaning and User Identification are the systems in Data Pre-planning. The inspiration for driving data cleaning is to discard insignificant things. This stream research is suffering from data pre-taking care of techniques that join data cleaning, data coordination, data change, and data decline. Different frameworks are obliged data cleaning anyway there are a couple of issues in data gathering and an exact estimation of customer unmistakable verification. Subsequently, Data Pre-getting ready is the key activity to complete Bioinformatics Mining structures and expect a basic occupation in choosing the idea of models. In Data Pre-setting up, the gathering of data differentiates in the kind of data available just as the data source site, the data source size and the way in which it is being completed. The Data Pre-treatment of Bioinformatics Mining is typically staggering.

## 2. RELATED WORKS:

Pilsung Kang, et.al execute Semi-managed learning (SSL) has been proposed to improve regular regulated taking in techniques via preparing from both unlabeled and marked information. To quantify naming vulnerability, the mark conveyance of the unlabeled information is evaluated with two probabilistic nearby recreations (PLR) models. At long last, the normal edge-based example determination (EMPS) is utilized to lessen preparing multifaceted nature [13].

Mathew Ngwae Maingi et.al proposed an alternate information pre-preparing method in information mining and laying out the significant reasons for learning disclosure [14].

Dongil Kim, et.al proposed a technique to VM dataset development strategy by distinguishing and expelling clamors. These commotions named flaw recognition and arrangement (FDC) clamors and metrology commotions [15].

Laila Benhlima et.al proposed an alternate strategy to build the exactness in huge information examination and maintain a strategic distance from repetitive information. These strategies need to deal with the missing information in the social database to give an ideal arrangement [6].

Wen-Yang Lin et.al decides another security model Closed MS ( $k, \theta^*$ ) - jumping and another anonymization strategy, Closed-MS segment, to process SRS information with missing qualities. This strategy averts the assailants of the learning security approach [7].

Li Tang et.al depicted a strategy to deal with the missing qualities utilizing Dynamic Bayesian Network (DBN). This method keeping up the connection between the properties of information. Bolster vector relapse (SVR) used to anticipate the missing qualities. And moreover, the approval completed by Symmetric Means Absolute Percentage Error (SMAPE) [8].

Sohail Sarwa et.al proposed a half and half methodology of Conditional Random Fields (CRF) and the Hidden Markov Model (HMM) is formulated for information purifying. A little volume of information is utilized in this methodology. The huge set utilizing a Pakistan Telecommunication Company (PTCL) for approval and train the information [3]. Weilu Chen

et.al depicts to plan a shortcoming location channel that is fulfilled by  $H \infty$  execution. Some adequate conditions are inferred regarding certain network imbalances and the unequivocal articulation of the required channel parameters [4].

Xianqiang Yang et.al decides to recognize the straight parameter shifting (LPV) framework within the sight of missing information. Autoregressive exogenous (ARX) with an obscure booking variable is a non-direct state-space model. The proposed calculations are at long last determined in the desire amplification calculation to appraise the obscure parameters of the LPV ARX model and booking variable [5].

Muammar Albayrak et.al clarify a grouping and most extreme probability estimation (MLE) based way to deal with the missing information issue is proposed. Three fundamental advances are pursued here called information decrease, grouping, and information fulfillment. The outcomes contrasted and the first informational collection with respect to a root mean square blunder (RMSE) [9].

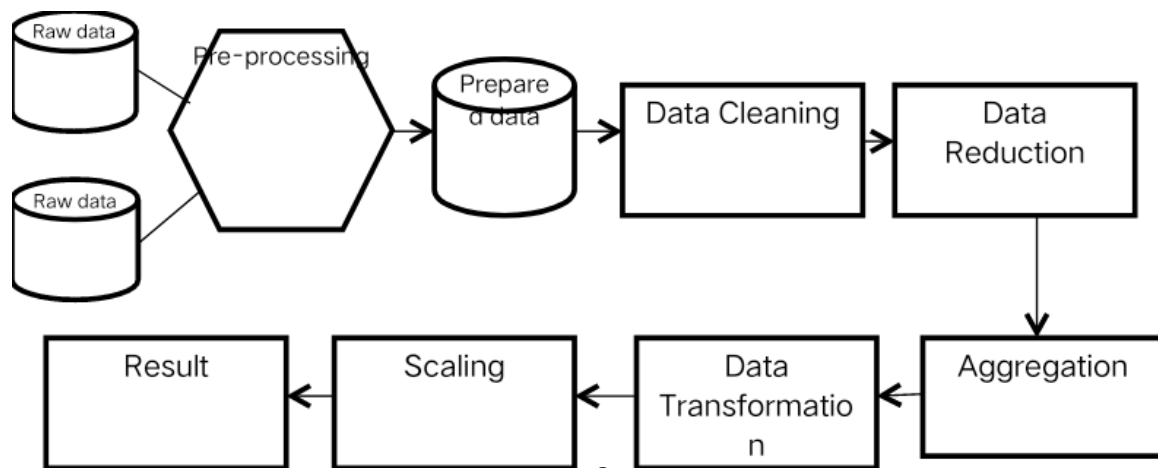
Jianke Chen, Pinghua Chen et.al proposed a technique to understand the missing information at high dimensional like over 3-dimensional by applying the tensor disintegration [10]. Sara Johansson Fernstad, Robert C. Glen et.al proposed visual examination to help the examination of missing information. It likewise underscores the absence of perception support in the territory and through this urges representation researchers to examine and address this profoundly pertinent issue [11].

Alina Lazar et.al proposed to research the missing qualities in straight out time arrangement groupings on normal information investigation errands. To conquer missing information issues is more diligently in the ostensible space versus the paired area. Furthermore, counterfeit bunches presented by the arrangement of driving missing qualities can be settled by tuning the missing worth substitution cost parameter [12].

Kim et.al going to propose a system to decrease the dimensionality incomplete least squares and least outright shrinkage and choice administrator relapse are used as expectation models [1].

Miriam Seoane Santos et.al proposed a system to survey an alternate methodology for the manufactured missing information age and their pragmatic subtleties. In this methodology containing subjective highlights is the most testing [2].

### 3. METHODOLOGY:



*Fig 3.1: Overall Framework*

Data course of action is a huge issue for both Bioinformatics and Web Usage Mining, as evident data will, by and large, be lacking, tumultuous, and clashing. Information planning incorporates information cleaning, information coordination, information change, and information decrease. Pre-handling improves the exhibition of bioinformatics and web mining information. Information cleaning schedules can be utilized to fill in missing qualities, smooth loud information, recognizing anomalies, and the right information irregularities. Information coordination joins information from numerous sources to shape an intelligent information store.

Metadata, relationship investigation, information struggle identification, and the goals of semantic heterogeneity contribute towards smooth information joining. Information change schedules affirm the information into proper structures for mining. Information decrease methods, for example, information block accumulation, measurement decrease, information pressure, various decreases, and discretization can be utilized to get a diminished portrayal of the information while limiting the loss of useful substance. Albeit a few strategies for information readiness are created, information arrangement remains a functioning and significant zone of research. To whole up, in the blink of an eye, the ideas of information mining and improving the exhibition of bioinformatics and web mining information by utilizing pre-preparing procedures are broke down and introduced.

Parameter	Data cleaning	Data Integration	Data Transformati on	Data reduction	Data Discretization
Input	Incomplete Noise  Inconsistent information	• Different wellspring of same information/i nformational index • Entity Identification	• Large records set  • High dimensional informational index	• Different sort and type of information  • Multiple angularities information	• Large set Continuous information and scope of the quality  • non-arranged

		issue •Redundancy information	• Large information volume	• wrong data(noise information)	information
<b>Output</b>	<ul style="list-style-type: none"> <li>• Quality information</li> <li>• Reliable information</li> <li>•Completed information</li> </ul>	<ul style="list-style-type: none"> <li>• Provide a Metadata</li> <li>• Data strife recognition</li> <li>• Quality information with consideration taken</li> </ul>	<ul style="list-style-type: none"> <li>• Normalized information</li> <li>• Summarized information</li> <li>• Compressed information</li> <li>• Correct wrong information</li> <li>• Feature development.</li> </ul>	<ul style="list-style-type: none"> <li>• Reduce information and evacuate undesirable things in the information al index</li> <li>•Minimizing the loss of data • Content. Lessen information volume</li> <li>• Replacing estimations of a nonstop characteristic by an interimnames</li> </ul>	<ul style="list-style-type: none"> <li>• Categorized information</li> <li>• Reduce information volume</li> <li>• Simplified informational collection.</li> </ul>
	<ul style="list-style-type: none"> <li>• Ignore the record</li> <li>• Fill in the missing woth</li> </ul>	<ul style="list-style-type: none"> <li>• Data coordination</li> </ul>	<ul style="list-style-type: none"> <li>• Smoothing</li> </ul>	<ul style="list-style-type: none"> <li>• Data solid shape total</li> <li>• Attribute subset</li> </ul>	<ul style="list-style-type: none"> <li>• Entropy-based Discretization</li> <li>• Binning</li> </ul>



<b>Technique</b>	physically • Use worldwide steady • Attribute Mean • Binning • Regression • bunching	• Schema coordination • Correlation investigation of absolute information utilizing X2	• Aggregation • Generalization • Normalization • Attribute development	determination • Dimensionality decrease • Sampling • Binning • Clustering	• Histogram investigation • X <sup>2</sup> - based Discretization • Concept Hierarchy Generation
<b>The complexity of pre-processing (Challenges)</b>	• Identifying anomalies. • exactness and nature of the information • filling on information	• semantic heterogeneity and structure of information • redundancies and irregularities • exactness, speed, and nature of the information	• exactness and comprehension of structure in high dimensional information • connections between information qualities	• The respectability of the first information • Complex information investigation • Mining on a lot of information • limit the disappointment content	• Replacing various estimations of a nonstop trait by a fixed number of interim marks • Generalized information important and simpler to translate • Classification exactness and quality
<b>Extra Memory</b>	Yes	No	Yes	Yes	Yes

**Table 3.1: Process of Pre-processing table**

For exact information mining PC based systems of information pre-handling offer arrangements that help the information under preparing to adjust typical structures which thusly significantly improve the exhibition of AI calculations. In this procedure, the exact assurance of exceptions, extraordinary qualities and topping off holes presents considerable difficulties. Different strategies have in this manner been created to recognize these strayed or conflicting qualities called anomalies.

S.NO	TECHNIQUES	MERITS	DEMERITS	PERFORMANCE
1	SVM – Support vector machine	• Regularization capabilities • Handles non-linear data efficiently	• Choosing an appropriate kernel function is difficult. • Required memory.	90.56%

2	Classification	<ul style="list-style-type: none"> <li>• Simplification</li> <li>• Briefness</li> <li>• Utility</li> <li>• Comparability</li> <li>• Attractive and effective</li> <li>• Scientific arrangement</li> </ul>	<ul style="list-style-type: none"> <li>• Based on subjective judgment may or may not be shared.</li> </ul>	86%
3	Clustering	<ul style="list-style-type: none"> <li>• Automatic recovery from failure</li> <li>• Recovery without user intervention</li> </ul>	<ul style="list-style-type: none"> <li>• Complexity and Inability to recover from Database corruption</li> </ul>	90%
4	Naive Bayes	<ul style="list-style-type: none"> <li>• Better performance compared to other models</li> <li>• Require small amount of training data to estimate the test data</li> </ul>	<ul style="list-style-type: none"> <li>• Assumption of independent predictors</li> <li>• Zero Frequency</li> </ul>	88.15%
5	Random Forest	<ul style="list-style-type: none"> <li>• Reduces Overfitting</li> <li>• Improves the accuracy</li> </ul>	<ul style="list-style-type: none"> <li>• Complexity</li> <li>• Longer training period</li> </ul>	86%
6	Decision tree	<ul style="list-style-type: none"> <li>• Path through possibilities</li> <li>• Desirable outcomes</li> </ul>	<ul style="list-style-type: none"> <li>• Loss of innovation</li> <li>• Complex process into discrete steps</li> </ul>	80%

7	ANN – Artificial Neural Networks	<ul style="list-style-type: none"> <li>• Storing information on the entire network</li> <li>• Ability to work with incomplete knowledge</li> <li>• Fault tolerance</li> </ul>	<ul style="list-style-type: none"> <li>• Hardware Dependence</li> <li>• Unexplained behaviour of the network</li> </ul>	93.17%
---	----------------------------------	---	---	--------

8	PCA – Principal Component Analysis	<ul style="list-style-type: none"> <li>Removes co-related features</li> <li>Improves algorithm performance</li> </ul>	<ul style="list-style-type: none"> <li>Independent variables become less interpretable.</li> </ul>	96.75%
9	T - PCA	<ul style="list-style-type: none"> <li>Efficient</li> <li>Easy to use</li> <li>Reliable</li> </ul>	<ul style="list-style-type: none"> <li>Required More memory</li> </ul>	98%
10	SICA - Subjectively Interesting Component Analysis	<ul style="list-style-type: none"> <li>Easy to understand</li> <li>Easy to test</li> <li>Easy to Upgrade</li> <li>Easy to identifying constraints</li> <li>Easy to diagnose the problems</li> </ul>	<ul style="list-style-type: none"> <li>Discrete – event model</li> </ul>	98.78%
11	HSIC - Hilbert-Schmidt independence criterion	<ul style="list-style-type: none"> <li>To identify the cross-covariance operator.</li> </ul>	<ul style="list-style-type: none"> <li>Memory Required</li> </ul>	98.75%
12	RKHS-reproducing Kernel Hilbert spaces	<ul style="list-style-type: none"> <li>Completeness</li> <li>Good performance.</li> </ul>	<ul style="list-style-type: none"> <li>Variable determined is difficult.</li> </ul>	98.9%
13	HSIC-NDR	<ul style="list-style-type: none"> <li>Network data representation</li> </ul>	<ul style="list-style-type: none"> <li>Required memory</li> <li>Non-delivery report</li> </ul>	99.01%



14	MLGODR - multiple locality-constrained graph optimization for dimensionality reduction	<ul style="list-style-type: none"> <li>Many fold validation function.</li> <li>Difficult to interpret</li> <li>Feature Extraction</li> <li>Pattern recognition</li> </ul>	<ul style="list-style-type: none"> <li>Linear methods dependent</li> </ul>	99.5%
15	MPCA - multilinear principal component analysis	<ul style="list-style-type: none"> <li>Analyzing tensor structured data</li> <li>Reduce the dimensions for both real data.</li> </ul>	<ul style="list-style-type: none"> <li>In-efficient unstable prediction</li> </ul>	95.90%
16	OMPCA - online multilinear principal component analysis	<ul style="list-style-type: none"> <li>Supervised feature selection</li> <li>Robust</li> </ul>	<ul style="list-style-type: none"> <li>Uncorrelated multi-linear feature</li> </ul>	98.1%
17	CART - Classification and Regression Tree	<ul style="list-style-type: none"> <li>It is not significantly impacted by outliers.</li> <li>Data belongs to particular types of distribution.</li> </ul>	<ul style="list-style-type: none"> <li>Take a large tree to get good lift</li> <li>Instability of model structure</li> </ul>	85%
18	Gaussian Naive Bayes	<ul style="list-style-type: none"> <li>Very easy for implementing</li> <li>Needs few amount of training data</li> </ul>	<ul style="list-style-type: none"> <li>Chance to loss of accuracy</li> <li>Cannot modify dependencies</li> </ul>	88.90%
19	C4.5	<ul style="list-style-type: none"> <li>Build models that can be easily interpreted</li> <li>Easy to implement</li> </ul>	<ul style="list-style-type: none"> <li>Data may be over fitted</li> <li>Only one attribute at a time tested</li> </ul>	90.70%

20	ELM - extreme learning machine	<ul style="list-style-type: none"> <li>• Faster to train</li> <li>• Fairly robust to mapping function</li> </ul>	<ul style="list-style-type: none"> <li>• Ignore deep learning Architectures</li> <li>• Cannot encode more than one layer of abstraction</li> </ul>	91%
21	Fuzzy	<ul style="list-style-type: none"> <li>• Based on linguistic model</li> <li>• High precision</li> <li>• Rapid operation</li> </ul>	<ul style="list-style-type: none"> <li>• Lower speed with longer run time.</li> </ul>	91.5%
22	K-means	<ul style="list-style-type: none"> <li>• Computationally faster</li> <li>• Produce tighter cluster</li> </ul>	<ul style="list-style-type: none"> <li>• Difficult to predict k values.</li> <li>• Global cluster didn't work well</li> </ul>	96.17%
23	T-SNE	<ul style="list-style-type: none"> <li>• Handles Non-linear data Efficiently</li> <li>• Preserves Local and Global Structure</li> </ul>	<ul style="list-style-type: none"> <li>• Computationally Complex</li> <li>• Pairwise conditional probabilities for each data point</li> </ul>	98.95%
24	K-means +Fuzzy	<ul style="list-style-type: none"> <li>• Easy to use</li> <li>• Easy to diagnose the problem</li> </ul>	<ul style="list-style-type: none"> <li>• Depends the numerical methods</li> </ul>	98.57%
25	PCA + T-SNE	<ul style="list-style-type: none"> <li>• Probability distribution</li> <li>• Strongly chosen data parameterization</li> <li>• To detect the false findings.</li> <li>• Approximately simple form of special clustering.</li> </ul>	<ul style="list-style-type: none"> <li>• Access depends on kull-back Leibler divergence</li> </ul>	99.95%

**Tab. Comparison Table**

In this collective raw data stored in the temporary database then it will move on Pre- processing techniques to transformation started to training data set respectively. First it removes the noise to make consistency data. Hence, the data has to compress and getting without loss of data while compressing techniques in the given problem. Then choose to apply the concept of data transformation in aggregation techniques. Similarly those techniques involved finding the relations between the two objects in the embedding methods.

It detects to determine the process of,

$$SD = \sigma = \sqrt{\sum ((x - \bar{x})^2 / n)} \dots\dots\dots(1)$$

$$\text{Mean} = \sum ((x_i - \bar{x}) / n) \dots\dots\dots(2)$$

Hence,

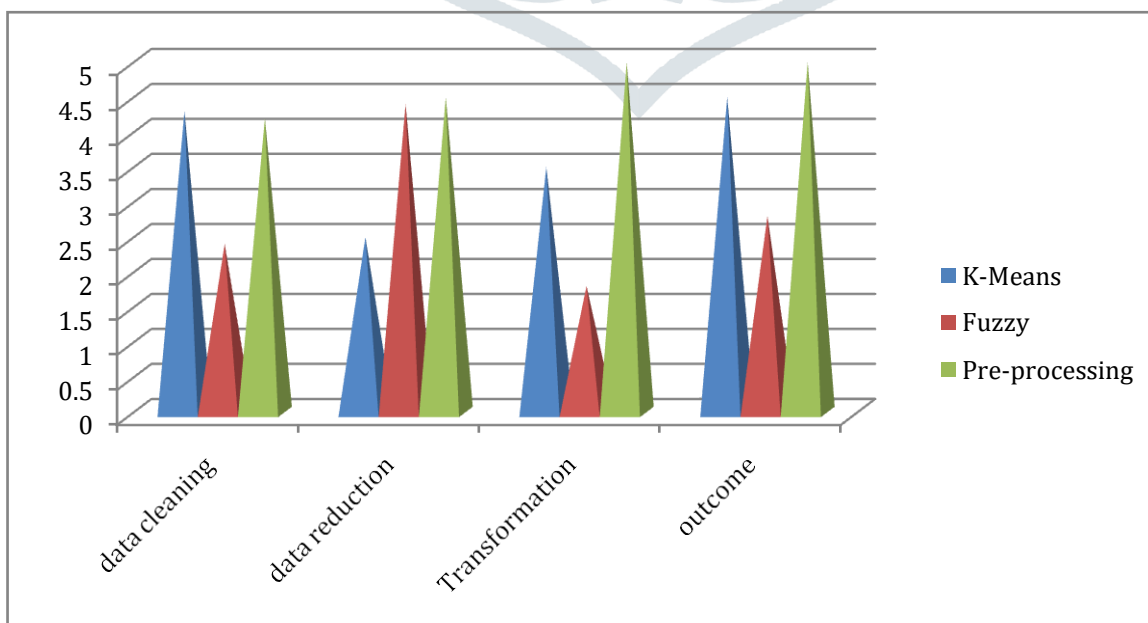
The mean absolute deviation has to produce the absolute likelihood in the given data points specifically.

$$\text{Average absolute deviation} = |x - \bar{x}| \dots\dots\dots(3)$$

It has been used to detect the process of ranges appears to access the pre-processing in the missing values to applying the data mining techniques. Scaling which helps to determine the mode of value has to be found in different variations in the pre-processing techniques. It will help to the different fields using to mess-up the data entered in the directory. So pre- processing used to find the missing values which help of standardization and normalization techniques to produce the outcome exactly. It is very simple and efficient to produce the outcome as best way without manpower.

#### 4. RESULT AND DISCUSSION:

Information cleaning and planning are essential to advance in the information mining process. We initially recognize various sorts of missing information and afterward talk about two ways to deal with arrangements with missing information in various situations. It is utilized to tend the issues of dealing with missing characteristics in datasets and methods in which missing characteristics can be taken care of. We initially talk about the various kinds of missing information and break down their effect on the dataset. We currently investigate the issue of missing qualities in dreary datasets. We recommend a basic pre-handling strategy which when utilized with different procedures helps in killing missing qualities and helps in keeping up the dataset dull.



**Fig 4.1: Accuracy Report**

Thus the above methods have been proved to produce the outcome with best for the given data in missing values handled by the pre-processing techniques. It conducts basic investigations to test the calculation and locate that taking the most continuous worth and supplanting it instead of missing qualities give better outcomes. Missing information some of the time additionally mask themselves as substantial information and are hard to recognize. It consequently, proposes a heuristic method to manage to handle a sensible and testing issue of cleaning conceal missing data. With the assistance of this methodology, we distinguish a suspicious example of information and after that build up a fair-minded example heuristic way to deal with find missing qualities.

## 5. CONCLUSION AND FUTURE ENHANCEMENT:

An immense measure of information produced in the medicinal services area which is deficient with regards to potential data to settle on choice and examination. Even though rich data is historical in healthcare, researchers finding difficulties while data collected through primary and secondary data sources. Real-world healthcare data contains inconsistency, noisy data which leads to wrong decisions and treatment. Healthcare data enforced to pre-process. After information accumulation, different information pre-handling techniques can be connected to clean the information. For handling missing value, imputation method is best for the healthcare data, since every attribute plays an important role in decision making. Integrating the medical data needs a very strong knowledge-based system and transformation

requires mapping of data present in each format. Finally, data reduction is necessary to cut the cost of data management and predict the accurate values from large scale medical data.

For, further feature research planned to use for the best embedding techniques applying to get best efficiency and accuracy of the pre-processing handling the missing values. It will help to improve the performance without lack of data entered in the given field such as college management etc.

## 6. REFERENCE:

- 1) Variable Selection Under Missing Values and Unlabeled Data in Semiconductor Processes, Kyung-Jun Kim, Kyu-Jin Kim, Chi-Hyuck Jun, Il-Gyo Chong, Geun-young Song, IEEE Transactions on Semiconductor Manufacturing, 2019.
- 2) Generating Synthetic Missing Data: A Review by Missing Mechanism, Miriam Seoane Santos, Ricardo Cardoso Pereira, Adriana Fonseca Costa, Jastin Pompeu Soares, João Santos, Pedro Henriques Abreu, IEEE Access, 2019.
- 3) Machine learning-based intelligent framework for data pre-processing, Sohail Sarwar,Z. Ul Qayyum,A. Kaleem, International Arab Journal of Information Technology, 2019.
- 4) Protocol-Based Fault Detection for Discrete Delayed Systems With Missing Measurements: The Uncertain Missing Probability Case, Weilu Chen, Jun Hu, Xiaoyang Yu, Dongyan Chen, IEEE Access, 2018,
- 5) LPV model identification with an unknown scheduling variable in the presence of missing observations – a robust global approach, Xianqiang Yang | Xin Liu | Boxuan Han, IET Control Theory

&Applications, 2018,

- 6) A Study of Handling Missing Data Methods for Big Data, Imane Ezzine, Laila Benhlila, IEEE 5th International Congress on Information Science and Technology, 2018.
- 7) Privacy-Preserving SRS Data Anonymization by Incorporating Missing Values, Wen- Yang Lin, Kuang-Yung Hsu, Zih-Xun Shen, Conference on Technologies and Applications of Artificial Intelligence (TAAI), 2018.
- 8) Imputation of missing value using dynamic Bayesian network for multivariate time series data, Li Tang, Chengke Wu, H.T. Tsui, Shigang Liu, International Conference on Data and Software Engineering, 2017.
- 9) A missing data imputation approach using clustering and maximum likelihood estimation, Muammar Albayrak, Kemal Turhan, Burçin Kurt, Medical Technologies National Congress, 2017.
- 10) A method based on tensor decomposition for missing multi-dimensional data completion, Jianke Chen, Pinghua Chen, IEEE 2nd International Conference on Big Data Analysis, 2017.
- 11) Visual analysis of missing data — To see what isn't there, Sara Johansson Fernstad, Robert C. Glen, IEEE Conference on Visual Analytics Science and Technology (VAST), 2017.
- 12) Data quality challenges with missing values and mixed types in joint sequence analysis, Alina Lazar, Ling Jin, C. Anna Spurlock, Kesheng Wu, Alex Sim, IEEE International Conference on Big Data, 2017.
- 13) Semi-supervised support vector regression based on self-training with label uncertainty: An application to virtual metrology in semiconductor manufacturing, Pilsung Kang, Dongil Kim, Sungzoon Cho, Expert Systems with Applications, 2016.
- 14) Survey on Data Preprocessing Concept Applicable in Data Mining, Mathew Ngwae Maingi, Neurocomputing 239, 2015.
- 15) Improvement of virtual metrology performance by removing metrology noises in a training dataset, Dongil Kim, Pilsung Kang, Seung-Kyung Lee, Seok-ho Kang, Seungyong Doh, Sungzoon Cho, Pattern Analysis and Applications, 2015.