

# A BIG DATA APPROACH FOR DIMENSIONALITY REDUCTION TO DOLE OUT SCALABLE DEBASEMENT

<sup>1</sup>S.Neelakandan, <sup>2</sup>R.Annamalai, S.Velmurugan<sup>3</sup>

<sup>1&2</sup>Assistant Professor, Department of IT, Jeppiaar Institute of Technology,

<sup>3</sup>Assistant Professor, Department of CSE, VelTechMultitech Dr.Rangarajan  
Dr.Sakunthala Engineering college.

---

## ABSTRACT

In our day-day life, the tenders and geographies of information technology have been advanced in Big data. As big data and machine learning gross hold, storage technology frequently becomes automated, virtualized and more complex. Reducing Dimensionality big data wrenches attention as of late as an operative policy to detach the center evidence which is smaller to store and speedier to grip. With the exponential power enhancement of huge data, big data has put a wonderful weight on current groundwork. Reducing Dimensionality in big data fascinates a great deal of kindness as an efficient approach to abstract the core data that is minor to store and sooner to process. To tackle the basic complications closely related to distributed dimensionality reduction of big data, dimensionality reduction algorithm and construction of distributed computing platform. A chunk tensor method changes the unstructured, semi-structured and structured data as a unified model in which all appearances of the assorted data are appropriately agreed along the tensor orders. To reduce the dimensionality of the unified model is Lanczos-based High Order Singular Value Decomposition procedure is used. The result after analyzing the procedure are represented as storage scheme, convergence property, and calculation cost. The execution of dimensionality reduction task incurs the Transparent computing pattern to pattern a distributed computing stage as well as utilizes the linear predictive model to partition the data blocks. Final results establish that the planned approach is proficient.

*Index Terms - Big Data, Tensor, Dimensionality Reduction, Transparent Computing.*

---

## 1. INTRODUCTION

Big data consist of data sets with sizes exceeding the capability of regularly used software tool to accurate, manage, and process data within a appropriate time. Big data "size" is a continuously moving target, which has stretched from a limited dozen terabytes to many petabytes of data in 2012. Big data is an all-encompassing term for any collection of data sets so large and challenging that it becomes hard to exercise using traditional data processing applications. The amount of generated and deposited data. The mass of the data governs the value and potential insight- and whether it can really be well-thought-out big data or not. The type and nature of the data. This helps people who evaluate it to efficiently use the resulting perception. In this context, the faster at which the huge data is produced and analyzed come across the demands and dares that lie in the path of growth and development. The Inconsistency of the data set can hinder processes to knob and manage it. The quality of caught data can differ greatly, disturbing the accurate analysis. The issues include analysis, seizure, period, search, distribution, storage, handover, visualization, and privacy violations.

### 1.1 What Comes Under Big Data?

Big data involves the data produced by different devices and applications. Given below are some of the fields that come under the umbrella of Big Data.

#### 1.1.1 Black Box Data:

It is a component of helicopter, airplanes, and jets, etc. It captures voices of the flight crew, recordings of microphones and earphones, and the performance information of the aircraft.

**1.1.2 Social Media Data:**

Social media such as Facebook and Twitter hold information and the views posted by millions of people across the globe.

**1.1.3 Stock Exchange Data:**

The stock exchange data holds information about the buy and sell decisions made on a share of different companies made by the customers.

**1.1.4 Power Grid Data:**

The power grid data holds information consumed by a particular node with respect to a base station.

**1.1.5 Transport Data:**

Transport data includes model, capacity, distance and availability of a vehicle. Search Engine Data: Search engines retrieve lots of data from different databases.

Thus Big Data includes huge volume, high velocity, and an extensible variety of data. The data in it will be of three types.

- Structured data: Relational data.
- Semi-Structured data: XML data.
- Non-structured data: Word, PDF, Text, and Media Logs.

Big data is challenging to work with most relational database management systems and desktop statistics and visualization packages, demanding instead "massively parallel software operating on tens, hundreds, or even thousands of servers. Information volume is ranging from Terabyte level or Petabyte level in different fields including biomedicine, Internet of Things (IoTs), online networking, and so on. Enormous information can bolster better comprehension of the certifiable and in addition, give great administrations to clients through proficient examination approaches. Be that as it may, the massively developing information are excessive and complex that they are exceeding the capacity of current figuring base and calculations to continue and process. Henceforth, new methodologies ought to be further investigated to address these difficulties brought about by enormous information. As the overabundance of dimensionality can achieve irregularity, excess, and perplexity in huge information applications, to remove the vital variables that emphatically catch the recognize qualities of heterogeneous information turns into an essential prerequisite. Dimensionality lessening is a procedure to diminish the quantity of variables in substantial scale information; it can maintain a strategic distance from the impact of condemnation of dimensionality. Unlike the previous decades, the investigation of dimensionality reduction has been accounted for writing. Essential Component Analysis (PCA) is a procedure that communicates the information so as to highlight their similitude and contrasts. PCA does the operations of changing the information to another coordinate framework in which the best change comes in main organize, the second in the following direction, etcetera. The technique to recognize the measurements along which information focuses show the most variety is Singular Value Decomposition (SVD). SVD locates the best estimation of the first information utilizing fewer measurements. Free Component Analysis (CA) is an efficient technique that looks only for straight projection, not for orthogonal to one another that are as about measurably free as could reasonably be expected. Element examination is a measurable technique to depict variability among watched, connected variables as far as a conceivably lower number of in secret variables called elements. Albeit numerous studies have been performed on dimensionality decrease, the said strategies endure from a few constraints. They mostly focus on similar information. These techniques are costly when they are connected to large scale information, therefore versatility issues may emerge if the decrease assignments are booked to heterogeneous registering gadgets. All the examinations neglect to perceive the significance of an efficient exploration for dimensionality decrease. The primary reason for this paper is to advance an all-encompassing way to deal with the location the three essential issues firmly identified with disseminated dimensionality diminishment of enormous information, that is huge information combination, dimensionality diminishment calculation, and development of a dispersed processing stage. We propose a lump tensor technique to intertwine the heterogeneous information from various sources as a brought together tensor model. The first information is spoken to as sub-piece Tensors in nearby customers, the removed center piece Tensors are transferred to servers for the mix. Ideas and operations of the piece tensor strategy are given and in addition, the combination procedure is represented in the refined element. We show a Lanczos-based High Order Singular Esteem Decomposition (L-HOSVD) calculation that can rapidly acquire the center Tensor and the truncated orthogonal bases from the bound together tensor model. The tensor developing is changed to symmetric frames and the Lancroz technique is utilized to figure the particular vectors and solitary standards. By analysis, the L-HOSVD calculation is effective and focused. We utilize the Transparent Computing worldview to build a conveyed registering stage with heterogeneous gadgets to perform the dimensionality decrease errands of enormous information. A straight prescient model is utilized to assess the measure of information squares that will be circulated to autonomic devices. Big data

is similar to small data but bigger in size, we use big data to solve new problems or old problems in a better way, it generates value from the storage and processing of large volumes of digital information that cannot be analyzed with usual computing techniques. In this proposed paper three types of important, smaller data are present here, Structured, Semi-Structured, NonStructured. Using big data fusion, we club the above data using chunk Tensor to be stored in HDFS, The clubbed data is processed by distributed computing platform with core data and task for mapping the data which the user requires. The data required by the user is reduced from the pre-processed and mapped data by dimensionally reduction algorithm. Dimensionality reduction is a process to decrease the number of variables in large-scale data; it can avoid the effect of the curse of dimensionality. It is also used to reduce the no of unplanned variables under consideration.

## 2. RELATED WORK

A Comparative Review,” refers to a variety of nonlinear reduction in dimensionality techniques have been proposed, Dimensionality reduction is the conversion of high-dimensional data into a meaningful depiction of reduced dimensionality is Reduction in Dimensionality . The reduced presentation should take a dimensionality that match up to the intrinsic dimensionality of the huge data. The minimum number of parameters needed to account for the perceived properties of the data many of which depend on the estimation of local properties of the data is the intrinsic dimensionality. The paper gives a review and systematic comparison of these methods. The efficiency of the methods are investigated on artificial and natural errands. A huge amount of nonlinear practices are perfectly able to discover this embedding, whereas linear practices fail to do so. The results of the experiments depicts that nonlinear techniques execute on carefully chosen artificial tasks, but cannot beat the traditional PCA on real-world tasks. The paper clarifies the results by detecting weaknesses of current non-linear techniques and propose a method to improve performance of non-linear dimensionality reduction practices . From the results acquired, we may conclude that nonlinear practices for dimensionality reduction are, despite the fact their large variance, not yet proficient of outperforming traditional PCA. In the forthcoming, we predict the improvement of new nonlinear practices for dimensionality reduction that do not rely on local properties of the data manifold. We also expect a shift in focus towards the increase of nonlocal practices for dimensionality reduction with objective functions that can be optimized well, such as (Kernel) PCA and autoencoders [1]. Eluding the bane of dimensionality in lively stochastic games Discrete-time stochastic games with a finite number of states have been broadly applied to study the strategic communication among onward-looking players in dynamic surroundings. The range of applications of finite-state stochastic games is limited by their high computational cost. The first source of the burden is the large size of the state space. Indeed, there can be a swear word of dimensionality in that the number of states upturns exponentially in the number of state variables, the dimension of the state vector. These games agonize from a “curse of dimensionality” when the price of computing players' anticipations over all likely future states increases exponentially in the number of state variables. Overall, the belief is that the benefits of continuous time are often substantial, and open the method to study more difficult and realistic stochastic games than presently feasible. We discover the alternative of continuous-time stochastic games with a limited number of states and claim that continuous time may have substantial benefits. In particular, under commonly used laws of motion, continuous time sidesteps the curse of dimensionality in calculating expectations, thereby speeding up the calculations by orders of magnitude in games with additional few state variables. This much reduced computational burden greatly spread out the range and richness of applications of stochastic games [2]. Principal component analysis (PCA) is a procedure that uses an orthogonal transformation to transform a set of remarks of possibly associated variables into a set of values of linearly uncorrelated variables so-called principal components. The number of components is less than or equal to the number of actual variables. This conversion is defined in a way that the first principal component takes the largest conceivable variance (that is, stands for as much of the variability in the data), and each consecutive component, in turn, has the utmost variance possible under the condition that it is orthogonal to the prior components. The causing vectors are an uncorrelated orthogonal base set. The principal components are orthogonal as they are the eigenvectors of the covariance matrix, which is symmetric. PCA is delicate to the comparative scaling of the actual variables[3]. “Singular Value Decomposition: Application to Exploration of Experimental Data Singular value decomposition (SVD) is a system normally used in the examination of spectroscopic data that acts as a noise filter and decreases the dimensionality of following least-squares fits. To begin the applicability of SVD to crystallographic data, we used SVD to calculate difference Fourier maps simulating those to be achieved in a time-resolved crystallographic reading of photoactive yellow protein. The atomic arrangement of one dark state and three intermediates were used in different kinetic mechanisms to produce time-dependent difference maps at particular time points. Random noise of changing stages in the variance structure factor amplitudes, dissimilar extents of reaction instigation, and dissimilar numbers of time points were all engaged to generate a range of realistic

experimental conditions. Our results show that SVD supports for an unbiased differentiation between signals; a small subset of singular values and vectors denotes the signal well, decreasing the random noise in the data. Due to this, phase statistics of the difference structure aspects can be achieved. After recognizing and fitting a kinetic mechanism, the time-independent model of the intermediates could be improved. This shows that SVD will be a powerful tool in the reading of investigational time-resolved crystallographic data [4]. Factor analysis is a statistical method used to refer to variability among observed, correlated variables in terms of a potentially lesser amount of undetected variables termed as factors. Factor analysis is used as a data reduction or structure detection technique. Factor analysis is a useful tool for examining variable relationships for complex concepts like socioeconomic status, dietary patterns, or psychological scales [5]. Transparent Computing will be offered to solve this problem partially. Accordingly, we suggest and develop a pilot system that operates in a network environment and works at the assembler instruction level. This system lets users claim heterogeneous uses and applications on them from centered simple servers, similar to select different TV channels in day-to-day life. We also present some primitive actual and experimental outcomes to show that it is a viable and proficient solution for upcoming computing infrastructure. It is the characteristic of pervasive computing, the possible future state in which we will be bordered by computers universally that respond to our wants without our conscious use. TC is the implementation of pervasive or ubiquitous computing. The implementation of computer instruction and data is temporarily and spatially parted from their storage [6]. Impacts, Visions and Challenges, The rapid development of computer network technologies and social information has got many new chances and dares in information security. With improved information and service sharing relished by numerous people, how to fortify the statistics security has become an increasingly critical issue. In this paper, we put forward a new network security mechanism based on a novel computing hypothesis, i.e., transparent computing, which is based on the protracted von Neumann architecture. This paradigm separates the program storage and execution, which is applied in the network environment. It is recognized by a new generation server and client BIOS, namely EFI BIOS, and coordinated with the Meta OS management platform and associated switching and input/output devices of transparent computing. By the interface between hardware and software, it conducts in effect control, monitoring and management of data and information in a block-streaming mode for the operating system and the application programs above it. In Parallel, it accepts a security protection method to prevent and remove prevalent malicious software like worm and Trojan horses. Several examples are described in detail to demonstrate the attractive and hopeful security countryside and advantages [7].

### 3. PROBLEM IDENTIFICATION

Volume of Data is increasing towards Tera byte level or Peta byte level in various fields including biomedicine, Internet of Things (IoTs), social media, etc. Big data supports a better understanding of the real world as well as provide high-quality services to users through efficient analysis approaches. Although, the tremendously growing data are so large and complex that they are beyond the ability of current computing infrastructure and analysis algorithms to manage and process. Therefore, new approaches has to be explored to address the issues in big data. There are much analysis such as Principal Component Analysis PCA, Singular Value Decomposition & Independent Component Study helps to solve issues. These studies concentrate on similar data. The strategies are costly when they are connected to substantial scale information, and adaptability issues may emerge if the decrease assignments are planned to heterogeneous processing gadgets. Every one of the examinations neglects to perceive the significance of a deliberate exploration for dimensionality reduction.

#### *DIS ADVANTAGES*

Scalability issues like Computation are expensive to process on different ways of representing and storing the same data.

### 4. IMPLEMENTATION

This approach of addressing the three problems closely associated to distributed dimensionality reduction of big data such as big data fusion, dimensionality reduction algorithm, and construction of a distributed computing platform. This work proposes a method to fuse the heterogeneous data from multiple sources as a unified tensor model called chunk tensor. The original data are given as sub-chunk tensors in the local clients, core chunk tensors are extracted and uploaded to servers for integration. To reduce the dimensionality of the unified model Lanczos-

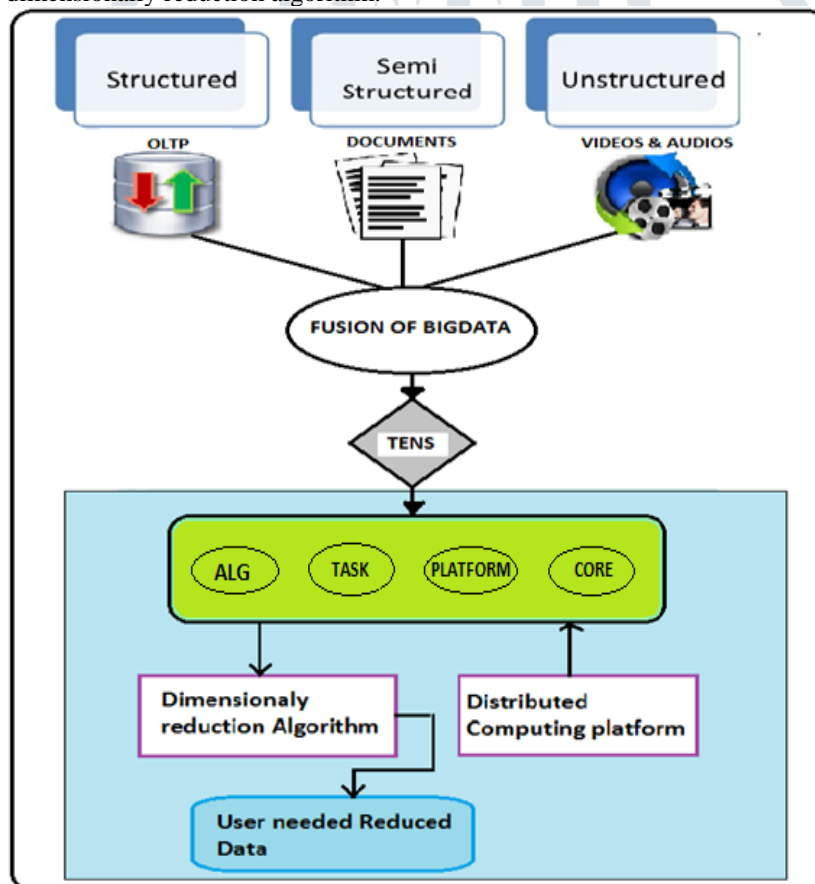
based High Order Singular Value Decomposition algorithm is proposed. Theoretical studies of the algorithm are provided in terms of storage scheme, convergence property, and computation cost.

#### ADVANTAGES

High scalability issues like Overhead computations are minimized; Computation cost is low and Easy to process on heterogeneous data.

#### 4.1 SYSTEM ARCHITECTURE

A chunk tensor approach is used to integrate the unstructured, semi-structured and structured data as a unified model in which features of the different data are appropriately arranged along the tensor orders. we club the above data to be stored in HDFS, The clubbed data is processed by distributed computing platform with core data and task for mapping the data which the user requires. The data required by the user is reduced from the pre-processed and mapped data by dimensionally reduction algorithm.



#### 4.2 MODULES

- 1) Data Fusion
- 2) Dimensional Reduction

### 3) Data Distribution

#### 4.2.1 DATA FUSION

Big data when used correctly brings benefits to the business, science, and humanity. The characteristics of big data like volume, velocity, variety, variation and veracity render the existing practices of data analysis ineffective. Big data analysis needs a combination of methods for data mining and machine learning. The k-means method is used in both the fields. This project depicts an appropriate algorithm based on k-means. It is an efficient way for big data analysis and is fast, scalable with high accuracy. It beats the drawback of k-means of an uncertain number of recurrences by setting the number of iterations, without dropping the precision.

#### 4.2.2 DIMENSIONALITY REDUCTION

Data columns with minute changes in the data contain little information. Thus each data column with variance lower than a given threshold is removed. A word of restraint: variance is range based; therefore normalization is mandatory before using this technique. Decision Tree Ensembles also known as random forests, are beneficial for feature choice in addition to being effective classifiers. A method for dimensionality reduction is to generate a large and carefully fabricated set of trees against a target feature and then make use of every attribute's usage statistics to discover the informative subset of features. In detail, we can produce a large set (2000) of narrow trees (2 levels), with each tree being trained on a lesser fraction (3) of the total number of attributes. If an attribute is chosen as best split, it is expected to be an informative feature to retain. A score computed on the attribute usage statistics in the random forest says relative to the added attributes which are the most extrapolative attributes.

**Bayes' Theorem** Bayes' Theorem is a statement from probability theory that allows for the computation of definite conditional probabilities. Conditional probabilities reflect the influence of one event on the probability of another event. **LogitBoost** The LogitBoost is an authoritative classifier. It can be successfully complemented with other existing powerful algorithms. **Random Forest** Random forests are very good in that it is an ensemble learning method used for taxonomy and regression. It makes use of many models for better performance that just uses a single tree model. **Decision Trees** Decision Trees are powerful and widespread tools for ordering and estimation. Decision Trees represent rules; *this* can be understood by humans and applied in knowledge system like Database.

#### 4.2.3 DATA DISTRIBUTION

The ability to control distributed computing and parallel processing practices dramatically changed the landscape and dramatically reduce latency. There are cases like High-Frequency Trading where low latency can only be accomplished by physically localizing servers in a single position.

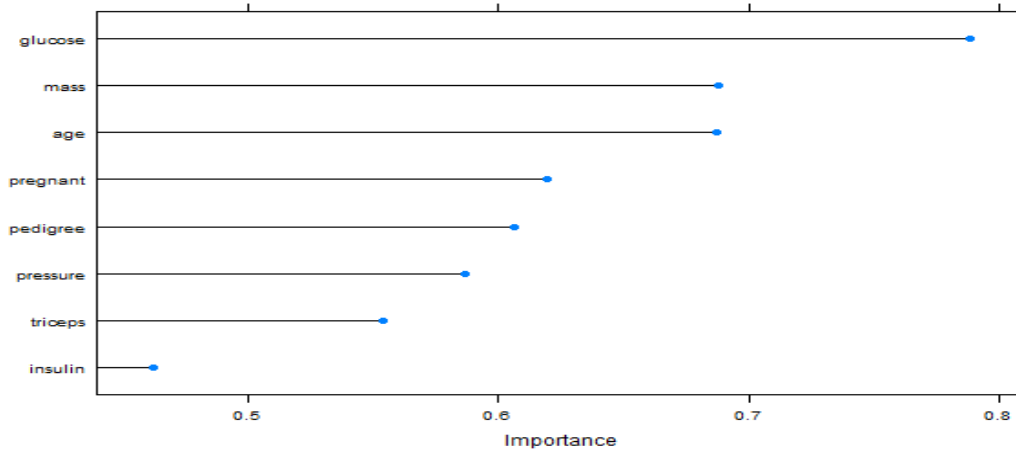
## 5. CONCLUSION

Theoretical examination and test results exhibit that the proposed far-reaching system is a strong way for dimensionality diminishment of huge data; it serves as a reason for considerable scale different data taking care of. Theoretical analysis and experimental test show that the proposed approach is a systematic approach for dimensionality reduction of big data; it serves as a basis for large-scale different data processing.

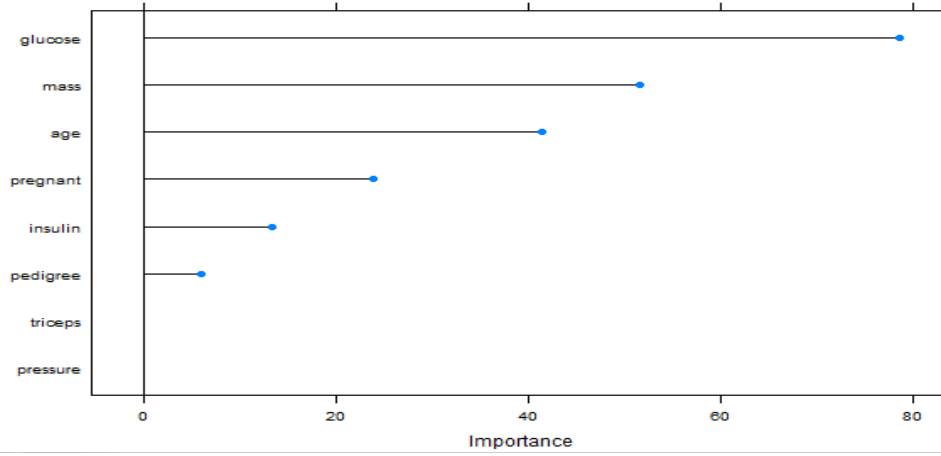
## 6. RESULT

Chunk tensor model is suggested to consolidate the heterogeneous data from different sources as a united tensor model. It goes for giving a widely inclusive approach to managing distributed dimensionality decreasing of immense data. Thoughts and operations of the knot tensor model are developed in this paper. Additionally, to gain the inside data which are little yet contain imperative information a Lanczos-based High Order Singular Value Decomposition (L-HOSVD) computation is proposed. At the present time, space (from hospital to residence and carry) and moment are no longer a hesitant stone for modern healthcare by means of more powerful investigation technologies. Nevertheless, with all advantages outlined, there is more that needs to be prepared especially concerning the service environment in which these systems will be offloaded. Chunk Tensor is mainly used to combine all kinds of data. Dimensionality reduction algorithm and distributed parallel computing are used to analyze the patient's data quickly, in addition to this, we added 5 algorithms to improve their performance and get the accurate result. From this, we can be able to know the medical history of the patient at a particular period.

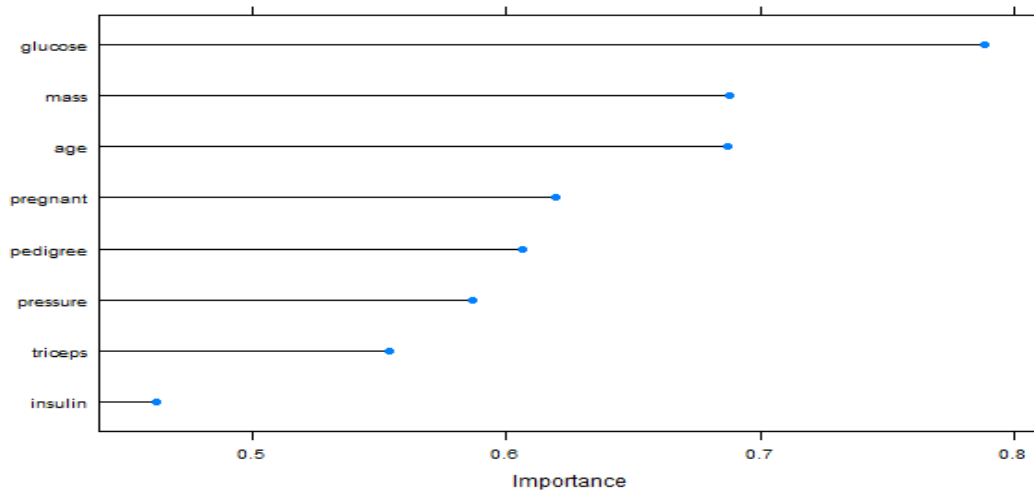
Variable importance plot



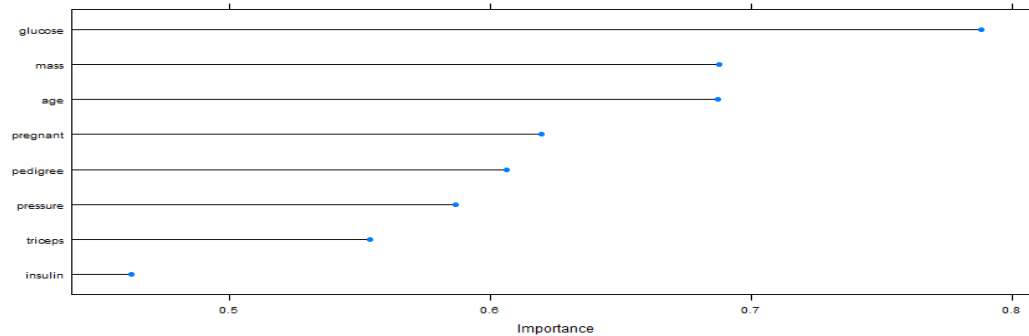
Variable importance plot



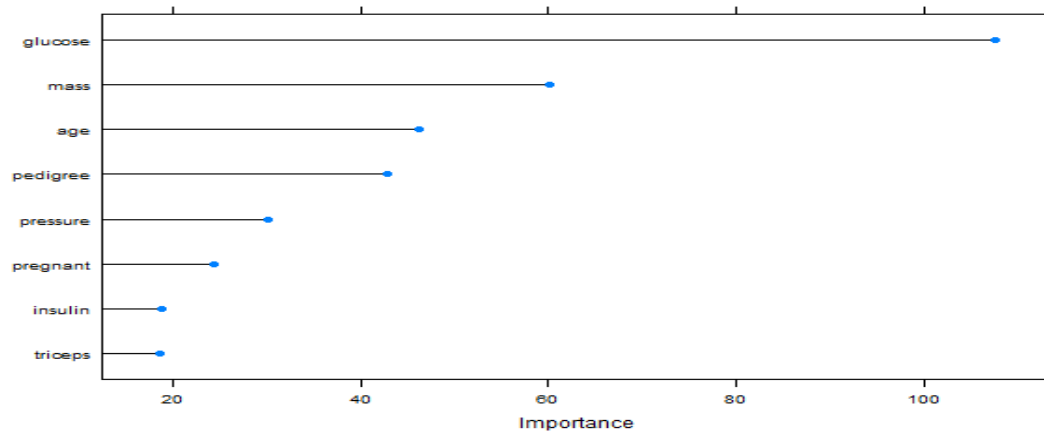
Variable importance plot



Variable importance plot



Variable importance plot



## 7. REFERENCES

- [1] L. J. van der Maaten, E. O. Postma, and H. J. van den Herik, "Dimensionality Reduction: A Comparative Review," *Journal of Machine Learning Research*, vol. 10 (2009):66–71.
- [2] U. Doraszelski and K. L. Judd, "Avoiding the Curse of Dimensionality in Dynamic Stochastic Games," *Quantitative Economics*, vol. 3(2012): 53–93.
- [3] H. Abdi and L. J. Williams, "Principal Component Analysis," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2(2010):433–459.
- [4] E. Henry, J. Hofrichter et al., "Singular Value Decomposition: Application to Analysis of Experimental Data," *Essential Numerical Computer Methods*, vol. 210 (2010):81–138.
- [5] E. E. Cureton and R. B. D'Agostino, *Factor Analysis: An Applied Approach*. Psychology Press, 2013.
- [6] Y. Zhang and Y. Zhou, "Transparent Computing: A New Paradigm for Pervasive Computing," in *Ubiquitous Intelligence and Computing*. Springer, (2006) :1–11.
- [7] Y. Zhang, L. T. Yang, Y. Zhou, and W. Kuang, "Information Security Underlying Transparent Computing: Impacts, Visions and Challenges," *Web Intelligence and Agent Systems*, vol. 8 (2010):203–217.
- [8] P. Comon and C. Jutten, *Handbook of Blind Source Separation: Independent Component Analysis and Applications*. Academic Press, 2010.
- [9] A. A. Abbasi and M. Younis, "A survey on clustering algorithms for wireless sensor networks," *Comput. Commun.*, vol. 30, nos. 14\_15, pp. 2826\_2841, Oct. 2007.
- [10] C. C. Aggarwal and C. Zhai, "A survey of text clustering algorithms," in *Mining Text Data*. New York, NY, USA: Springer-Verlag, 2012, pp. 77\_128.
- [11] R. Xu and D. Wunsch, "Survey of clustering algorithms," *IEEE Trans. Neural Netw.*, vol. 16, no. 3, pp. 645\_678, May 2005.
- [12] S. Guha, R. Rastogi, and K. Shim, "CURE: An efficient clustering algorithm for large databases," in *Proc. ACM SIGMOD Rec.*, Jun. 1998, vol. 27, no. 2, pp. 73\_84.
- [13] Bu, Y. and Howe, B. and Balazinska, M. and Ernst, M. HaLoop: Efficient iterative data processing on large clusters *Proceedings of the VLDB Endowment*, 3(1- 2), 285-296, 2010.