



Society for Science and Education
United Kingdom

ISSN: 2055 - 1266
Volume 1. Issue 3

JOURNAL OF BIOMEDICAL ENGINEERING AND MEDICAL IMAGING



TABLE OF CONTENTS

EDITORIAL ADVISORY BOARD	I
DISCLAIMER	II
An Efficient Clustering based Segmentation Algorithm for Computer Tomography Image Segmentation V.V.Gomathi S.Karthikeyan	1
Multi-dimensional Multi-granularities Data Mining for Discovering Innovative Healthcare Services Johannes K. Chiang Chia-Chi Chu	12
Smart Obstacle Detector for Blind Person Daniyal Rajput Faheem Ahmed Habib Ahmed Engr Zakir Ahmed Shaikh Aamir Shamshad	31

EDITORIAL ADVISORY BOARD

Professor Kenji Suzuki

Department of Radiology, University of Chicago
United States

Professor Habib Zaidi

Dept. of Radiology, Div. of Nuclear Medicine, Geneva University Hospital,
Geneva, Swaziland

Professor Tzung-Pe

National University of Kaohsiung,, Taiwan
China

Professor Nicoladie Tam

Dept. of Biological Sciences, University of North Texas, Denton, Texas, United
States

Professor David J Yang

The University of Texas MD Anderson Cancer Center, Houston
United States

Professor Ge Wang

Biomedical Imaging Center, Rensselaer Polytechnic Institute. Troy, New York
United States

Dr Hafiz M. R. Khan

Department of Biostatistics, Florida International University
United States

Dr Saad Zakko

Director of Nuclear Medicine Dubai Hospital
UAE

Dr Abdul Basit

Malaysia School of Information Technology, Monash University
Malaysia

DISCLAIMER

All the contributions are published in good faith and intentions to promote and encourage research activities around the globe. The contributions are property of their respective authors/owners and the journal is not responsible for any content that hurts someone's views or feelings etc.

An Efficient Clustering based Segmentation Algorithm for Computer Tomography Image Segmentation

V.V.Gomathi¹, S.Karthikeyan²

¹Research and Development Centre, Bharathiar University, Coimbatore, India

²Department of Information Technology, College of Applied Sciences, Sohar, Oman

[1vv.gomathi@gmail.com](mailto:vv.gomathi@gmail.com), [2skaarthi@gmail.com](mailto:skaarthi@gmail.com)

ABSTRACT

Colossal amount of research has been done in creating many different approaches and algorithms for medical image segmentation, but it is still complicated to evaluate all the images. However the problem remains challenging, with no general and unique solution in computer-aided diagnosis. This paper provides medical image segmentation based on Clustering for computer tomography images. In this paper, we consider a mean shift segmentation and medoid shift segmentation method. We validate the mean shift and medoid shift medical image segmentation approach with the parameters in terms of sensitivity, specificity and accuracy. The Real time dataset is used to evaluate the performance of the proposed method. The experimental result shows that the medoid shift segmentation method gives more accurate and robust segmentation results than mean shift segmentation method.

Keywords: Clustering, Computer tomography, Segmentation, Mean shift segmentation, Medoid shift segmentation

1 INTRODUCTION

Image segmentation is the process of partitioning a digital image into segments. Segmentation refers to simplifying and/or change the representation of an image into more meaningful and easier to analyze [1]. One fundamental problem in medical image analysis is image segmentation, which identifies the boundaries of objects such as organs or abnormal regions in images. Medical image segmentation is becoming an increasingly important image processing step for a number of clinical applications.

With increasing use of imaging modalities for diagnosis, it has become almost compulsory to use computers to assist medical experts in clinical diagnosis. Various medical imaging techniques such as computed tomography (CT), magnetic resonance imaging (MRI), Ultrasound(US), Positron Emission Tomography (PET), etc provide different perspectives on the

DOI: 10.14738/jbemi.13.267

Publication Date: 30th June 2014

URL: <http://dx.doi.org/10.14738/jbemi.13.267>

human body. Computer tomography is very important imaging modalities to provide radiotherapy for tumor patient. Manual segmentation is time consuming task and be prone to errors, especially due to fatigue. Manual segmentation also gives inter and intra expert variability results. In this scenario reliable algorithms are essential for the delineation of anatomical structures and other regions of interest (ROI) to assist and automate the radiological tasks. Techniques for performing segmentations vary widely depending on the specific application, imaging modality, and other factors. There is no universal algorithm for segmentation of every medical image. Each imaging system has its own specific limitations [2].

Selection of a suitable approach to a segmentation problem can consequently be a complicated problem. In this paper the comparison has been made between Mean shift and Medoid shift segmentation algorithm for real time computer tomography images and proposed the best one. Mean and Medoid shift algorithm is a clustering based algorithm. Clustering based segmentation algorithm is more suitable for computer tomography image segmentation than other intensity based, region based and edge based segmentation algorithm. Mean shift and Medoid shift algorithm is a non-parametric and powerful clustering method that does not require a-priori Knowledge of the number of clusters. The Experimental result proves that the Medoid shift algorithm is more suitable and also gives robust results than mean shift algorithms.

The rest of the paper is organized as follows: Section 2 discusses Materials and Methods. Section 3 presents Computational results and discussion. Section 4 concludes the paper.

2 MATERIALS AND METHODS

2.1 Methods

2.1.1 Mean shift Segmentation Method

Mean shift is a popular mode seeking algorithm that is used in a large number of computer vision applications. The most important application is using Mean Shift for clustering. It is a non-parametric, unsupervised method that does not require a-priori Knowledge of the number of clusters, nor does it place any Limitations on the shape of the clusters [3].

Mean shift was first proposed by Fukunaga and Hostetler, later modified by Cheng for the purpose of image analysis and more recently extended by Comaniciu, Meer and Ramesh to low-level vision problems, including, segmentation, adaptive smoothing and tracking[4]. It represents a general non-parametric mode finding and clustering procedure. It considers feature space as an empirical probability density function. If the input is a set of points then Mean shift considers them as sampled from the underlying probability density function. If clusters are present in the feature space, then they correspond to the mode of the probability density function. For each data point, Mean shift associates it with the nearby peak of the data

set's probability density function. For each data point, mean shift defines a window around it and computes the mean of the data point. Then it shifts the center of the window to the mean and repeats the algorithm till it converges[5].

Generally the Mean shift method is computed as follows:

$$m(x) = \frac{\left[\sum_{i=1}^n x_i g\left(\frac{\|x - x_i\|^2}{h}\right) \right]}{\left[\sum_{i=1}^n g\left(\frac{\|x - x_i\|^2}{h}\right) \right]} - x \quad (1)$$

Here x_i is the initial estimate of this iterative method, $g\left(\frac{\|x - x_i\|}{h}\right)$ can be considered as the kernel function which determines the weight of nearby points for re-estimation of the mean. The center is altered to the new one unless the mode is found. This algorithm places $x \leftarrow m(x)$ and repeating is occurred until $m(x)$ is converged to x .

Benefits of Mean shift segmentation Algorithm

- 1) It is Suitable for real data analysis
- 2) This method does not assume any prior shape and number of clusters unlike k-means clustering method
- 3) It can handle arbitrary feature space
- 4) It is suitable for choosing Single parameter

Issues in Mean shift Segmentation Algorithm

- 1) Mean shift might not work well in higher dimensions. In higher dimensions, the number of local maxima is pretty high and it might converge to local optima soon.
- 2) The Window size is not trivial
- 3) Inappropriate window size can cause modes to be merged, or generate additional "shallow" modes->use adaptive window size
- 4) Does not scale well with dimension of feature space [6].
- 5) Initializing mean shift from every Data point is computationally expensive because each mean shift method iteration requires numerous nearest neighbor searches [3].
- 6) The choice of kernel bandwidth depends on the data and application [7].

Algorithm for Mean Shift Segmentation Method in Computer tomography Images

The Proposed Mean shift Algorithm is as follows:

Step I: Consider the Single DICOM image or slices of DICOM images

Step II: Apply the ECFT (Enhanced Curvelet Filter Technique) algorithm to get a noiseless image

Step III: Initialize the window radius

Step IV: Generate the window with the initialized radius value

Step V: Compute the convergence point

Step VI: Perform Superimposition of window over each pixel in the input image

Step VII: Check the convergence value and if it is not converged then it will move to the next pixel

Step VIII: If the value is converged means make it as the cluster value and move to the next data point

Step IX: Repeat the step VI, VII, VIII until obtain the convergence of all the pixels in the input image

Step X: Obtain the Clustered image after the convergence of all the pixels

2.1.2 Medoid shift Segmentation Method

Medoid shift algorithm is also a nonparametric clustering approach. It is a mode seeking method that computes shifts towards areas of greater data density using local weighted medoids. The use of medoids to discover structure in data is natural since, locally, the medoid can be considered a good representative of its neighborhood. Unlike means, medoids do not need an explicit feature space and require only a valid distance measure [8]. Medoid shift algorithm is proposed for seeking modes, based on manually selected bandwidth.

The medoid shift algorithms also automatically calculate the number of clusters during execution like mean shift. However, the medoid shift algorithm has major benefit over the mean shift algorithm. The medoid shift algorithm can operate directly on a distance matrix, irrespective of the original space in which the samples are distributed. This property allows medoid shift to find modes even when only a distance measure between samples is defined. In this sense, the relationship between the medoid shift algorithm and mean shift algorithm is similar. Medoid shift more related with the median operation.

Generally the Medoid shift method is computed as follows:

A weighted data point is calculated for every datapoint until the mode is obtained.

$$y_{K+1} = \arg \min_{y \in \{x\}} \sum_i \|x_i - y\|^2 G \left[\left\| \frac{x_i - y_k}{h} \right\|^2 \right] \quad (2)$$

Here y_{K+1} represents the weighted medoid that is considered to be best sample data point that minimizes the function.

Benefits of Medoid shift segmentation Algorithm

- 1) Automatically determines the number of clusters and does not need initialization [8]
- 2) Previous computations can be reused when new samples are added or old samples are deleted (good for incremental Clustering applications)
- 3) It can work in domains where only distances between samples are defined.
- 4) No need for heuristic terminating conditions.

Algorithm for Medoid shift Segmentation Method in Computer tomography Images:

Step I: Consider the Single DICOM image or slices of DICOM images

Step II: Apply the ECFT (Enhanced Curvelet Filtering Technique) algorithm to get a noiseless image

Step III: Find the Histogram of the input image

Step IV: Initialize the control parameter

Step V: Find the gray level cluster values based on an initialized control parameter

Step VI: Find no of pixel values present between each range of all gray level cluster values.

Step VII: The total number of pixels between each range are divided by two and Subtracted from the histogram value.

Step VIII: The minimum value from the subtracted data is found and the respective Intensity pixel value is calculated.

Step IX: This pixel value is substituted for the values between the respective ranges

Step X: Obtain the Clustered Image.

3 COMPUTATIONAL RESULTS AND DISCUSSION

Different type of Tumor patient dataset was collected by a SIEMENS SOMATOM EMOTION SPIRAL CT scanner located at Multi Speciality Hospital, Coimbatore. Experimentation was carried out on 100 numbers of different tumor patients contains 100 to 1000 slices of Computer Tomography images using Mean shift and Medoid shift Segmentation algorithms. The image format is DICOM (Digital Imaging Communications in Medicine). The algorithm has been implemented in Matlab environment. Manual Segmentation has been done by the radiation oncologist for comparing the performance of our algorithm. Experimental results of the images are illustrated here. In this paper we have demonstrated two patient DICOM images.

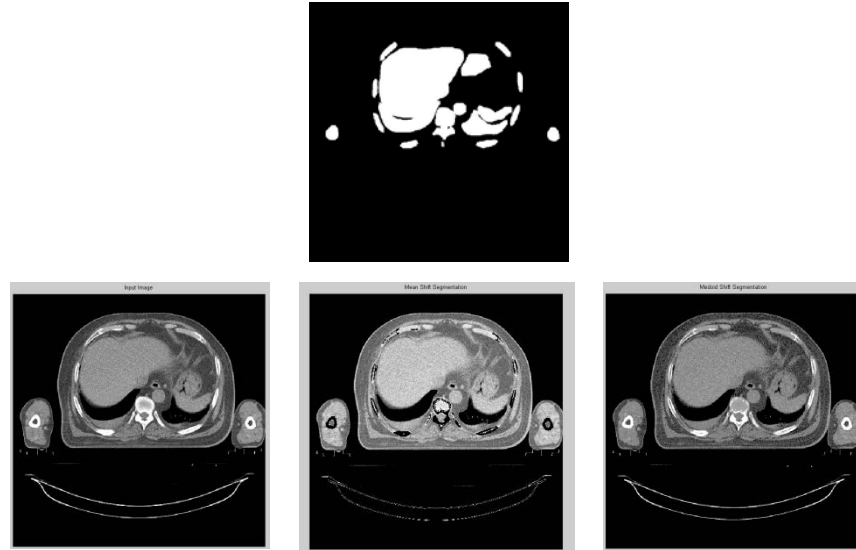


Figure 1: (a) Manual Segmentation done by the radiation oncologist (top row); (b) Input Image (bottom row), (c) Mean shift Segmentation, (d) Medoid shift Segmentation

3.1 Performance Analysis of Mean and Medoid shift segmentation method

Choosing an appropriate segmentation evaluation measure is a complicated task. A variety of performance measures to evaluate the medical image segmentation methods are available in current scenario. Generally sensitivity, specificity and accuracy are used to evaluate the segmentation methods. They are defined as

$$\text{Sensitivity} = \frac{TP}{TP + TN} \quad (3)$$

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (4)$$

$$\text{Accuracy} = \frac{TN + TP}{TP + TN + FP + FN} \quad (5)$$

where TP (True Positive) is the number of pixels of the foreground that are correctly classified, TN (True Negative) is the number of pixels of the background that are correctly classified, FP (False Positive) is the number of pixels of the background that are classified as foreground and FN (False Negative) is the number of pixels of the foreground that are Classified as background. Accuracy refers to the degree to which the segmentation results agree with the true segmentation i.e. Correct segmented pixels in the object. Fragments indicate the number of connected components in the required region to identify as organ. In this paper we also

consider the number of fragments parameter. If more number of fragments exists in the image, the segmentation task is also complicated.

Table 1: Performance Analysis of Mean shift and Medoid shift segmentation Algorithm

Methodology	Sensitivity	Specificity	Accuracy	Number of Fragments
Mean shift Segmentation	98.12	97.96	96.03	11400
Medoid shift Segmentation	99.71	97.23	98.02	8979

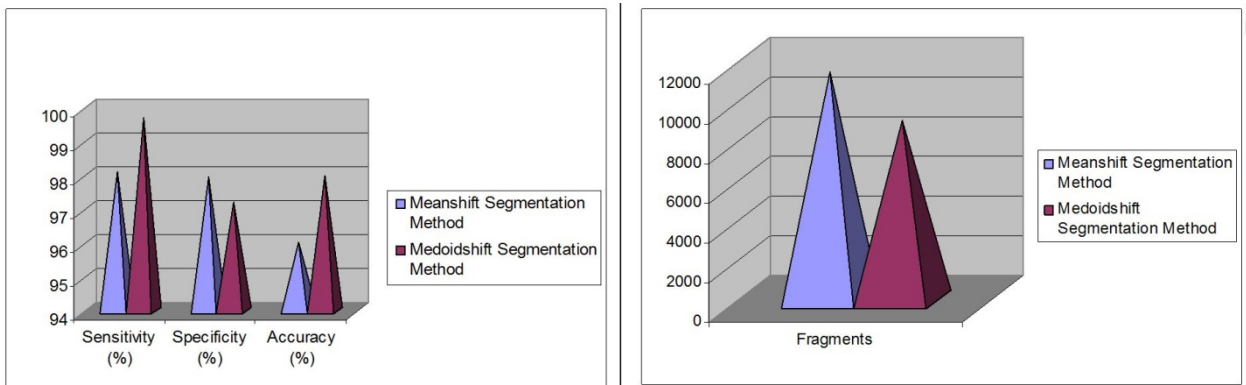


Figure 2: (a), (b) Graphical representation indicates the Performance analysis of Mean shift and Medoid shift Segmentation Algorithms in terms of Sensitivity, specificity, Accuracy, number of fragments.

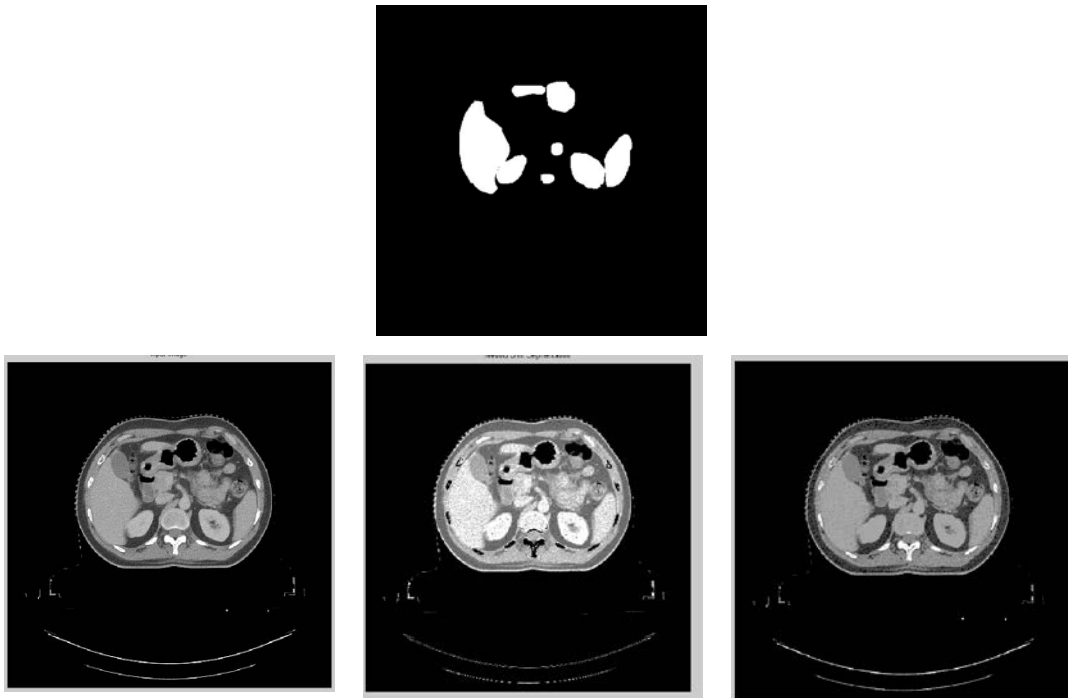


Figure 3: Patient Id: 77 Slice no: 101 (a) Manual Segmentation done by the radiation oncologist (top row); (b) Input Image (bottom row), (c) Mean shift Segmentation, (d) Medoid shift Segmentation

Table 2: Performance Analysis of Mean shift and Medoid shift segmentation Algorithm

Methodology	Sensitivity	Specificity	Accuracy	Number of Fragments
Mean shift Segmentation	97.06	96.75	95.98	5655
Medoid shift Segmentation	98.91	96.58	96.64	5398

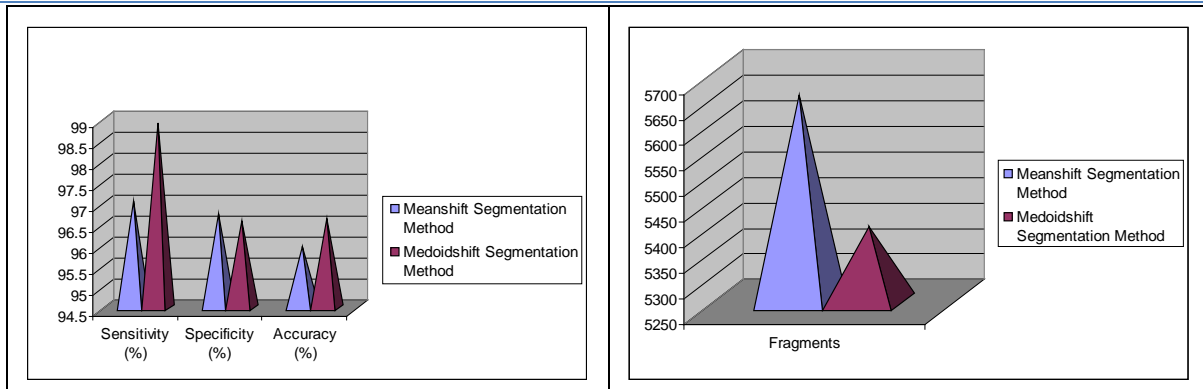


Figure 4: (a), (b) Graphical representation indicates the Performance analysis of Mean shift and Medoid shift Segmentation Algorithms in terms of Sensitivity, specificity, Accuracy, number of fragments.

3.2 Discussion

In this research work we have tried to implement a very good segmentation algorithm to segment the organs in Computer tomography images for effective Classification. Evaluating the results produced by segmentation algorithms is a challenging task. The segmentation is evaluated by assessing its consistency with the manual segmentation and their amounts of fragmentation. The value of Number of fragments indicates the number of connected components in the required region to identify as organ. Ultimately the value should be low for best segmentation method, since more connected components shall lead to inaccurate organ recognition and classification.

In this research, Table 1, Table 2 depicts the sensitivity, specificity, accuracy and numbers of Fragments are considered to compare the performance Mean shift and medoid shift Segmentation Method. Figure 1. (a). represent the Manual Segmentation done by the radiation oncologist. This slice consists of liver, heart, left lung, right lung, spleen, aorta, spinal cord and bones. Figure 1. (c), (d) shows the mean shift and medoid shift segmentation output. Figure 3. (a). also represent the Manual Segmentation done by the radiation oncologist. This slice consists of liver, spleen, stomach, aorta, spinal cord, left kidney and right kidney. Figure 3. (c), (d) shows the mean shift and medoid shift segmentation output. Medoid shift Segmentation methods produce higher accuracy and lower the number of fragments. It is clearly shows that segmentation results of Medoid shift Segmentation method are most promising than Mean shift segmentation method. But still overfragments exist in medoid shift segmentation method. In future we aimed to rectify this overfragmentation problem. However the medoid shift segmentation result is good for efficient organ classification.

4 CONCLUSION

In this paper we have evaluated the performance of mean shift and medoid shift segmentation algorithm which is suitable for real time computer tomography images. Medoid shift segmentation algorithm outperforms the mean shift segmentation in terms of sensitivity, specificity, accuracy and number of fragments. The Medoid shift segmentation algorithm

provides good quantitative metric values. An experimental results show that the medoid shift algorithm is more robust, effective, suitable and outperforms the mean shift segmentation algorithm. This segmentation results is very useful for effective organ classification.

REFERENCES

- [1]. Linda G. Shapiro and George C. Stockman., *Computer Vision*. 2001, New Jersey: Prentice-Hall. p.279-325.
- [2]. Neeraj Sharma and Lalit M. Aggarwal, *Automated medical image segmentation technique*. Journal of Medical Physics. 2010. 35(1):p.3–14.
- [3]. Lior Shapira, Shai Avidan , Ariel Shamir, *Mode-Detection via Median-Shift*, Computer Vision. IEEE 12th International conference on, 2009. p.1909-1916.
- [4]. Konstantinos G. Derpanis, *Mean Shift Clustering*, 2005.
- [5]. Miaoqing Huang, Liang Men, Chenggang Lai, *Accelerating Mean Shift Segmentation Algorithm on Hybrid CPU/GPU Platforms*. Ed: Xuan Shi, Volodymyr Kindratenko, Chaowei Yang. Modern Accelerator Technologies for Geographic Information Science, 2013. p.157-166.
- [6]. Dorin Comaniciu, Peter Meer, *Mean Shift: A Robust Approach toward Feature Space Analysis*. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 2002. 24(5): p.603-619.
- [7]. Tobias Weyand and Bastian Leibe, *Discovering Details and Scene Structure with Hierarchical Iconoid Shift*. Computer vision (ICCV), IEEE International conference on, 2013. p. 3479-3486.
- [8]. Yaser Ajmal Sheikh, Erum Arif Khan, Takeo Kanade, *Mode-seeking by Medoidshifts*. Computer Vision (ICCV), IEEE International conference on, 2007. P.1-8.
- [9]. Kaufman.L and P. J. Rousseeuw. Clustering by means of medoids. *Statistical Data Analysis Based on the L1 Norm*. Y.Dodge, Ed., Northi Holand /Elsevier.1987. p.405-416.
- [10]. MacQueen.J. *Some methods for classification and analysis of multivariate observations*. Mathematical Statistics and Probability, Proceedings of 5th Berkeley Symposium on, 1967, 1(1): p.281-297.
- [11]. Mohammad Talebi, Ahamd Ayatollahi, Ali Kermani, *Medical ultrasound image segmentation using genetic active contour*. Journal of Biomedical Science and Engineering, 2011. 4: p.105-109.
- [12]. Yizong Cheng, *Mean Shift, Mode Seeking, and Clustering*. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 1995.17(8): p.790-799.
- [13]. Arnaldo Mayer and Hayit Greenspan, *An Adaptive Mean-Shift Framework for MRI Brain Segmentation*. Medical Imaging IEEE Transactions on, 2009. 28(8): p.1238-1250.
- [14]. Jinghua Lu, Jie Chen, Juan Zhang, Lihui Zou, *Medical Image Segmentation Using Mean Shift Algorithm and General Edge Detection*. 18th IFAC World Congress Milano (Italy), 2011. p.9656-9661.

- [15]. Comaniciu. D and P. Meer, *Mean shift analysis and applications*. Computer Vision, International Conference on, 1999. pp. 1197–1203.

- [16]. Comaniciu. D, V. Ramesh, and P. Meer, *Real-time tracking of non-rigid objects using mean shift*. Computer Vision and Pattern Recognition, IEEE Conference on, 2000. 2: p.142–149.

- [17]. Beleznai.C, B. Frühstück, H. Bischof, *Human Tracking by Fast Mean Shift Mode Seeking*, Journal of Multimedia, 2006.1(1): p.1-8.

Multi-dimensional Multi-granularities Data Mining for Discovering Innovative Healthcare Services

Johannes K. Chiang¹, Chia-Chi Chu²

*Department of Management Information Systems, Cloud Computing and Operation
Innovation Center National Chengchi University, Taipei, Taiwan*

¹jkchiang@nccu.edu.tw, ²102356020@nccu.edu.tw

ABSTRACT

Data Mining is getting increasingly important for discovering association patterns for health service innovation and Customer Relationship Management (CRM) etc. Yet, there are deficits of existing data mining techniques. Since most of them perform a plain mining based on predefined schemata through the data warehouse as a whole, a re-scan must be done whenever new attributes are added. Secondly, an association rule may be true on a certain granularity but fail on a smaller one and vice versa. Last but not least, they are usually designed to find either frequent or infrequent rules.

After a survey of a category of significant health services, we propose a data mining algorithm along with a forest data structure to solve aforementioned weaknesses at the same time. At first, we construct a forest structure of concept taxonomies that can be used for representing the knowledge space. On top of it, the data mining is developed as a compound process to find the large-itemsets, to generate, to update and to output association rules that can represent services portfolio. After a set of benchmarks derived to measure the performance of data mining algorithms, we present the performance with respect to efficiency, scalability, information loss, etc. The results show that the proposed approach is better than existing methods with regard to the level of efficiency and effectiveness.

Keywords: Multidimensional Data Mining, Healthcare Services, Customer relationship Management (CRM), Association Pattern, Granular Computing.

1 INTRODUCTION

In the era of information economy, markets offer more variances of services and customers become demanding on more intensive information and better quality of services. While the term of Service Innovation becomes a focus in the scientific and business communities, data mining turns out to be increasingly important for knowledge discovery of innovative services.

DOI: 10.14738/jbemi.13.243

Publication Date: 30th June 2014

URL: <http://dx.doi.org/10.14738/jbemi.13.243>

As a whole, the conventional process of mass-marketing is being replaced by the customer-oriented view. As the second reason for seeking new way of services, healthcare institutions in many countries are facing a tail-off of healthcare assurance payments. Healthcare institutions need thus to target patients with new portfolios of service variances.

Under this condition, hospitals like to provide various new services such as prevention methods with education on patient with changing habits to prevent chronic illness and disease, treatment and physical check-up periodically to assist patients to improve their health quality. Moreover, hospitals like to improve their performance and to offer better quality of services. New tools and approaches such as CRM via data mining are needed to address this change.

Using association rules, we figure out simple yet useful insights on services [5, 13, 17]. Significant examples are finding new therapies and drugs for cancer cure as well as new portfolios of rationale services. For instance, “52% of the patients those take therapy X also take treatment Y”. With such association rules, we can reduce the costs of the therapy X, and raise the service level of the treatment Y to make more benefits.

However, most conventional data mining approaches only perform a plane scan over the databank based on a predefined schema for searching. Questions often arise such as: Should there be any other influencing factor like W for treatment Y taken into account? Since most association rules apply in a context of certain breadth, the knowledge usually exists in multidimensional insides [5]. In the in the meantime, adding attributes to the databank is meant to change the schema and lead to a full re-scan that consumes extra time.

The second problem of the conventional mining approaches lies in the assumption that the rules derived should be effective throughout a database as a whole. Nevertheless, this obviously is not true for real-life cases [5]. Different association rules can be found in different segments of the database. If the mining tool deals only with the database as a whole, meaningful rules that are partially true may be ignored.

The goal of this research is to invent an approach with novel data structure and efficient multi-dimensional data mining algorithm for association patterns in various granularities. The crucial issue here lies on a more efficient and accurate multidimensional mining approach to explore association patterns on different granularities. Last but not least, the data mining approach has to be very flexible and robust.

2 BASELINE OF THE RESEARCH

2.1 Data Mining for healthcare services

Data mining technology can contribute to hospitals with more understanding of patients' illness status and to improve quality of service (QoS). Hospitals use databases of patient's records, physical check-up, pathology etc. to analyze patients' status with aids of data mining and knowledge management. Based on the findings of above activities, hospitals can then

select different type of patient categories for different prevention, treatment services. Regarding data mining technology, they are now exploring five constructs for better service such as patient segmentations with respect to different type of service, different insurance reimbursement for varies type of patient, chronic illness, self-pay treatment and physical check-up services. Significant service categories can be summarized as follow:

- **Patient segmentations on different type of service:** By analyzing different types of patient illness, hospitals can provide various services for patients with their customization of treatment service, education, and wellness maintenance. Hospital notify patient to return back to the hospital for planning the best services for patient treatment.
- **Different insurance reimbursement for varies type of patients:** Hospital will analyze patients' insurance types of reimbursement, and also applies data mining to provide appropriate service to earn the maximum reimbursement. Furthermore, hospital will classify the contributions of different patient types to provide the best services to attract the higher level of patients to generate more revenue.
- **Chronic illness:** Hospital will analyze the patient's check-up results to define and predict different chronic illness types as well as different services for patients. Furthermore, hospital will notify patient back to hospital for routine check-up and treatments.
- **Self- pay services:** Hospital is capable of mining the patients' needs for self-pay services such as tumor/ cancer MRI check-up, skin disease for skin beauty treatments, hypertension for brain stoke check-up, cardiac disease for VCT cardiac service etc.
- **Physical check-up patient services:** Hospital applies data mining to retrieve patient illness status to notify patient for physical check-up.

2.2 Finding Association Rules

We are used to storing data in the transaction database containing simple items identified by the Transaction IDs (TID) as in Table 1. Let $I = \{i_1, i_2, \dots, i_n\}$ be the set of all n different items in D , each transaction in D is a subset of I . An itemset is defined as a subset of I . [4, 13, 17].

Table 1: An Example of Transaction Database

T_ID	Transaction content
001	Diagnosis-2, Therapy1.
002	Check-up-N, Therapy1.
003	MRI-Check, Diagnosis-3, Treatment-3

Table 2: An Example of Multidimensional Transaction Database

T_ID	Date	POS_No	Occupation	Sex	Age	Transaction content
001	05/03/01	003	Student	F	23	Diagnosis-2,Therapy1.
002	05/03/01	003	Student	M	14	Check-up-1,Diagnosis-2
003	05/03/01	003	Manager	M	47	MRI-Check,Diagnosis-3,Treatment-3

Rather than in an uni-dimensional transaction database, services and related information on customers are usually gathered in a relational database or data warehouse. Apart from keeping track of the item fed, a relational database may record other attributes associated with the transactions, and another table to record profile of patients, viz. a fact table. After joining several relational tables, a big data table can be obtained to store not only the items saved in the transaction [13, 17], but also 5W1H information corresponding to the transactions as Zachman Framework intends [4]. Table 2 illustrates an example of multidimensional transaction database MD, assuming each attribute is a dimension.

There are two important factors for association rules, viz. support, and confidence [13, 17]. Support means how often the rule applies, i.e. repeatability; Confidence means how often the rule is true, i.e. reliability [4]. Suppose we have a database MD as in Table 2, the support of an itemset X is the fraction of transactions containing X in MD. The confidence of $A \rightarrow B$ is the fraction of transactions containing A and B, and simultaneously also in transactions containing A. The formulas for support and confidence are as follows:

$$\text{Support}(X) = \frac{|\text{Transactions in D containing X}|}{|\text{Transactions in D}|} \quad (1)$$

$$\text{Confidence}(A \div B) = \frac{|\text{Transactions in D containing both A and B}|}{|\text{Transactions in D containing A}|} \quad (2)$$

Given a set of transaction MD and a threshold σ as minimum support, X is a large itemset in MD if the support of X in MD exceeds σ [13, 17]. The task for discovering association rules is to generate all association rules that own support and confidence greater than the user-specified minimum support (called minSup) and minimum confidence (called minConf) respectively [4, 12, 13, 17].

We are more likely to find association rules with high support and confidence, viz. frequent rules. Recently, the importance of vital few association rules is perceived, viz. infrequent rules [4].

2.3 Multidimensional Data Mining

Finding association rules involving various attributes efficiently is an important subject for data mining. Association Rule Clustering System (ARCS) was proposed in [], where association rule clustering is proposed for a 2-dimensional space. The restriction of ARCS is that it generates one rule in once of clustering. Hence, it takes massive redundant scans to find all rules.

The method proposed in [16] mines all large itemsets at first and then use a directed graph to assign attributes according the user given priorities of each attribute. Since the method is meant to discover the large itemsets over a database as the whole, it may loss some rules that hold only in specific segments of the database. Different priorities of the condition attributes will induce different rules so that user may need to try with all possible priorities to discover all possible rules.

2.4 Apriori Algorithm

2.4.1 Apriori Algorithm

The Apriori algorithm is a level-wise iterative search algorithm for mining frequent itemsets w.r.t association rules [1, 3, 5, 7, 13, 14, 17]. The key drawback of the Apriori algorithm is that it requires k passes of database scans when the cardinality of the longest frequent itemsets is k . In addition, the algorithm is computation intensive in generating the candidate itemsets and computing the support values, especially for applications with very low support threshold and/or vast amount of items. In this algorithm, if the number of first itemsets element is k , the database will be scanned k times at least. So, it is not efficient enough. The key point for improving the algorithm is to reduce the number of itemsets.

2.4.2 AprioriTID Algorithm [9]

The AprioriTID is a variant of the aforementioned Apriori algorithm which reduces the time needed for the frequency counting procedure by replacing every transaction in the database by the set of candidate sets that occur in that transaction [9]. This is done by iterating each candidate sets repeatedly.

While the AprioriTID algorithm is much faster in later iterations, it is much slower than original Apriori in early iterations. This is mainly due to the additional overhead that is created when the adapted transaction database C_k does not fit into main memory and has to be written into disk [4]. If a transaction does not contain any candidate k -sets, then C_k will not have an entry for the transaction. Hence, the number of entries in C_k may be smaller than the number of transactions in the database, especially at later iterations of the algorithm. Other drawbacks of AprioriTID are that the database modified by Apriori-Gen can be much larger than the initial database and only faster in the later stages of the scans.

2.5 Concept Description and Knowledge Taxonomy

The issues of data structures and concept description models for data mining when comparing works dealing with algorithms are less discussed till. The concept description task is problematic, since the term “concept description” is used in quite different ways in related discussions. In this situation, researchers argue for a de facto standard definition for the concept description [8, 18]. At this beginning stage, it is easier to deal with common criterion on higher abstraction level for the concept description, such as comprehension [8] and compatibility [4].

Researchers view concept description as a form of data generalization and define the concept description as a task that generates descriptions for the characterization and comparison of the data [8]. Similar concept appears in the development of ontology for Semantic Web/GRID. Semantic Web can be described as an extension of the existing Web where information is considered with priori well-defined meaning, enabling computer and people to work in cooperation centric to Internet [11]. The objective of such techniques is to enhance ill-structured content so that it can be interpreted universally by machines or humans.

In practical applications, ontology provides a vocabulary for specific domains and defines the meaning of the terms and relationships between them. In this article, ontology refers to the shared understanding (comprehension) of domains of interests which is often conceived as a set of concepts, relations, axioms etc. Hence, the term “Taxonomy” is hereby similar to “Ontology” and both terms can be used to denote the classification or categorization of concepts that describe entities and relations among them. This article applies the term Taxonomy rather than Ontology because the former is more flexible and even can cover the case with no semantic meaning.

3 METHODOLOGY

3.1 Representation schema and data structure

For the sakes of comprehension and compatibility, we use the forest structure consisting of Concept-Taxonomies to represent the overall searching space, i.e. the set of all the propositions of the concepts. On top of this structure, the sets of association patterns can be formed by selecting concepts from individual taxonomies. The notions can be clarified with examples as follows:

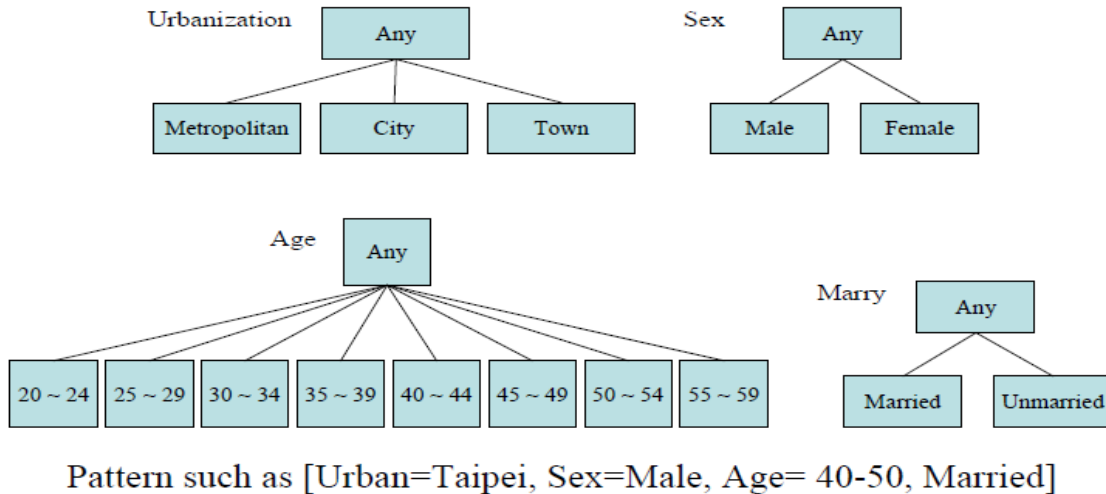


Figure 1: An Example for Forest of Concept Taxonomies

3.1.1 Taxonomy

A category consists of domain concepts in a latticed hierarchical structure, while each member per se can be in turn taxonomy. An Example (see Figure 1) for customer's characteristics can be [Age, Sex, Marry, Urbanization], while for instance the taxonomy of Sex can [Male, Female] and Marry can [Married, Unmarried] so on.

3.1.2 Forest of concept taxonomies:

A hyper-graph for representing the universe of discourse or the closed-world of interests is built with taxonomies under consideration. An example of forest of taxonomies with respect to the location and Sex of customers is shown in Figure 2 below:

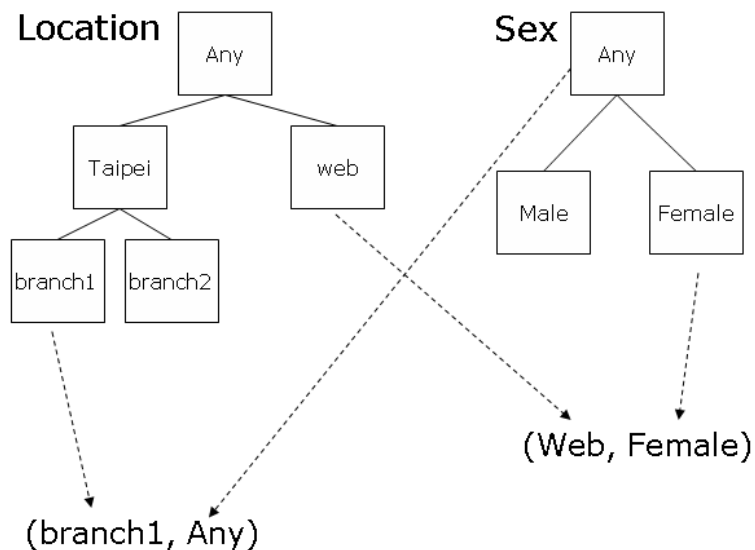


Figure 2: Examples of Forest Concept Taxonomies

3.1.3 Association Rule:

An association rule typically refers to a portfolio's pattern which consists of elements taken from various concept taxonomies such as [(Location=branch1), (Sex=female)]. It owns support and confidence greater than the user-specified minSup and minConf respectively [4].

3.1.4 Element patterns and generalized patterns:

An element pattern is composed of dimension atoms. On the other hand, if at least one of them is a dimension compound which combine several dimension atoms, we call this pattern a generalized pattern. For example, <web, Female> is an element pattern, <branch1, Any> is a generalized pattern, and both them are multi-dimension patterns. We use to denote the i -th element pattern, and use to denote the j -th generalized pattern.

By the proposed multidimensional data mining of association rules, the notion of relation will be implemented by the belonging relationship between elementary patterns and generalized patterns rather than the semantics [4]. Other notations to be used in the following text are shown in Table 3 below:

Table 3: Concepts and Notations

No tation	Meaning
CT	Concept Taxonomy
E_i	The i -th element segment
$T[E_i]$	an element segment over E_i in MD
G_j	The j -th generalized pattern
$T[G_j]$	The j -th combined segment over G
RE_i	Rules w.r.t the i -th element segment
RG_j	Rules w.r.t the j -th generalized pattern
(G_j, r)	association rules over G_j w.r.t to match ratio r

3.2 The Multidimensional Multi-granularity data mining algorithm

- 1) Input:
- 2) Multidimensional Transaction Database **MD**
- 3) Concept taxonomies for each dimension: $CT_x (X= 1-n)$
- 4) User given threshold: *minsup*, *minconf*, *match ratio m*
- 5) Procedure:
- 6) Phase0:
- 7) to generate all E_i and G_j by $CT_x (x = 1 to n)$;
- 8) build the pattern table;
- 9) Phase1:
- 10) For all $E_i \subset G$
- 11) to discover all association rules r in $T[E_i]$ as R_{E_i} ;
- 12) Phase2:
- 13) for all E_i
- 14) for all G_j that $E_i \subset G_j$
- 15) to update R_{G_j} using R_{E_i} ;
- 16) Phase3:
- 17) for all G_j
- 18) For all r (which satisfy m) in R_{G_j}
- 19) output (G_j, r) ;
- 20) Output:
- 21) all multidimensional association rules(p, r)

Figure 3: Outline of the proposed algorithm.

Outline of the proposed algorithm is shown in Figure 3. The input of the mining process involves 5 entities, namely (1) a multidimensional transaction database MD which is optional when a default MD is assigned, (2) a set of concept taxonomies for each dimension (CTs), (3) a minimal support, viz. minSup, (4) a minimal confidence, viz. minConf, and (5) a match ratio m for the relaxed match. The output of the algorithm encompasses all multi-dimensional associations with respect to the fully-relaxed match within the MD. The last three settings can help with finding frequent or infrequent rules.

The most significant feature of the algorithm is its capability to discover both frequent and infrequent associations rules R_{E_i} (based on different levels of granularities) in the element segment $T[E_i]$ for each element pattern E_i . After it, R_{E_i} is used to update R_{G_j} , i.e. the set of association patterns for every generalized pattern G_j which includes E_i . The heuristic regarding each element pattern is to find the large-itemsets per se and acknowledge its super generalized patterns with the result. The task of each generalized pattern is to decide which rules hold within it, according to the acknowledgements from the element patterns. The mining procedure needs only to work on each element segment to determine which rules hold in the compound segments. Thus, it is not necessary to scan all of the potential segments for finding the rules.

3.3 Pattern Generation and the Pattern Table

Being a pre-processing mechanism, the algorithm generates at first all elementary and generalized patterns with the given forest, where a pattern table for recording the belonging

relationship between the elementary and generalized patterns is built. Given a set of concept taxonomies, a multi-dimensional pattern can be generated by choosing a node from each of the taxonomy. The compound of different choices represents all the multidimensional patterns.

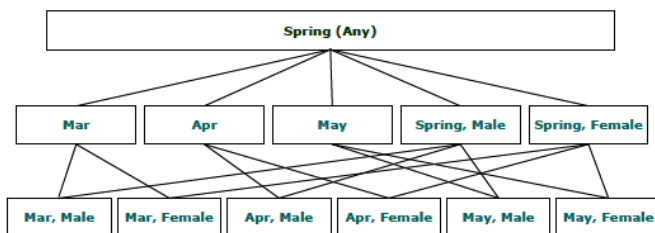


Figure 4: Belonging relationships between patterns

	(Mar)	(Apr)	(May)	(Spring, Male)	(Spring, Female)	(Spring)
(Mar, Male)	1	0	0	1	0	1
(Mar, Female)	1	0	0	0	1	1
(Apr, Male)	0	1	0	1	0	1
(Apr, Female)	0	1	0	0	1	1
(May, Male)	0	0	1	1	0	1
(May, Female)	0	0	1	0	1	1

Figure 5: The pattern table (for the relations in Figure 4) shows an example of the belonging relationship between 12 patterns in a lattice structure.

The relationships are recorded in the form of bit map as shown in Figure 5 which includes element patterns and generalized patterns. In the table, a “1” indicates that the element pattern belongs to the corresponding generalized pattern and “0” indicates the case vice versa.

3.4 Update process

- 1) **for** all R_{E_i}
- 2) **for** all $G_j \supset E_i$
- 3) **if** (R_{G_j} never be updated)
- 4) $R_{G_j} = R_{E_i}$;
- 5) **else**
- 6) $R_{G_j} = R_{G_j} \cap R_{E_i}$;

Figure 6: The “Update” algorithm for the full match

In order to be more optimization algorithm, we proposed full match and relaxed match method for update process. After all patterns and the pattern table have been generated, the procedure reads the transactions of each element segment and then discovers all the association rules. The output of this phase is all R_{E_i} for each element pattern E_i that will be fed as the input to the next phase for updating each R_{G_j} using R_{E_i} . For a full match illustrated in Figure 6, the update is done by intersection of the set R_{G_j} and the set R_{E_i} , where E_i belongs to G_j , let $R_{G_j} = R_{E_i}$ if R_{G_j} is updated for the first time. After all the intersections, the association pattern r left in R_{G_j} holds in all element segments covered by $T[G_j]$.

```

1) for all  $R_{E_i}$ 
2)   for all  $G_j \supset E_i$ 
3)     for all  $r$  in  $R_{E_i}$ 
4)       if ( $r \notin R_{G_j}$ )
5)         add  $r$  to  $R_{G_j}$ ;
6)          $R_{G_j}.r.count = 1$ ;
7)       else
8)          $R_{G_j}.r.count++$ ;

```

Figure 7: The “Update” procedure for the relaxed match

For the relaxed match as shown in Fig. 7, a counter for each rule in R_{G_j} is set. While using R_{E_i} for updating R_{G_j} , the counters of both R_{G_j} and R_{E_i} are incremented by one and the rules, those appear in R_{E_i} but not in R_{G_j} , will be added to R_{G_j} while setting the counter to one. After all the update process, the association rule r in R_{G_j} whose counts exceed $m|T[G_j]|$ holds in at least $m * 100\%$ of the element segments $T[E_i]$ that are covered by $T[G_j]$, and thus (G_j, r) is a multidimensional association rule for the relaxed match in MD.

Full match can ensure that all association rule be found in various granularities. But, it may be too restrictive to ignore some rules. On the other hand, relaxed match can solve “restrictive” problem and hold more association rules which may be our interesting rules. User can adjust the m ratio which ranges between 0 and 1.

For example, suppose we have a generalized segment <Spring> which covers three element segment <March>, <April>, and <May>. Finding patterns of each element segment <March>{A},{B},{C} ∙ <April>{B},{C} and <May>{B},{E}. As we above-mentioned algorithm that update each R_{G_j} using the R_{E_i} come from previous phase. For the full match case, we just can hold rule B in <Spring> generalized segment R_{G_j} because only rule B exists every element segment R_{E_i} . For the relaxed match case, we suppose $m = 0.6$ (result of count numbers should greater than 1.5 times) and count numbers of all rules in each element segment R_{E_i} : {A=1} ∙ {B=3} ∙ {C=2} ∙ {E=1}. Hence, we hold rule B and C in <Spring> generalized segment R_{G_j} .

3.5 The Output Function

For a full match, the algorithm outputs all the (G_j, r) pairs for every r left in each R_{G_j} . For a relaxed match, it outputs all the (G_j, r) pair for every r in each R_{G_j} where the count exceed $|mT[G_j]|$. By means of this approach, loss of finding the rules that only hold in some segments can be prevented. And, pickup of multidimensional association rules that do not hold over all the range of the domain can also be avoided. For example, the full match can guarantee that the corresponding rules, those hold only in two months of spring but fail in the rest one, will never be counted as an association rules with respect to whole spring.

3.6 The Breakthroughs for Incremental Data Mining

A breakthrough hereby is that the incremental data mining can be realized with the proposed approach. By keeping out the rules deduced in each element segment, we only need to search the new data. That is, using the proposed approach, we can produce the new association rules by combining the rules discovered from the new data with existing rules to reduce redundant scan on the old data. The following section will present our experimentation results.

3.7 Design of metrics for measuring data mining

In order to assure the performance, we need to design metrics for measuring the mining performance, at least to measure whether it is better than the prior algorithms. By cascade evaluating the results of a hypothetical measurement, we can evaluate the consequence from any sequence of measurements to determine the optimal next measure. For this reason, a one-step look-ahead strategy based on Shannon's Entropy Function is adopted and the capacity of ICT systems can be described in the following form [4, 15]:

$$C = B * [\log_2 (1 + S/N)] \quad (3)$$

where B is the bandwidth, (S/N) is Signal-to-Noise(S/N) ratio.

Drawing on this equation, the function for the performance of data mining can be formulated as follow:

$C = |D| [\log_2 (1 + \text{information lost ratio})]$, where |D| is the number of transactions in whole transaction database [4].

While WSE_i denotes each element segment in the measure, the WSE_i of an element segment $T[E_i]$ can be generated by a uniform distribution between 0 and SM. Suppose there are N element segments, the number of transactions in the element segment $T[E_i]$ is:

$$|D_{E_i}| = \frac{|D|}{\sum_{a=0}^n WSE_a} WSE_i \quad (4)$$

Thereafter, the definitions of information loss are:

$$\text{discrete ratio} = \frac{|\{r \mid r \text{ holds in } T[G_j] \text{ \<Gj,r> doesn't hold in MD } \}|}{|\{r \mid r \text{ holds in } T[G_j] \}|} \quad (5)$$

Definition 1: discrete ratio is the ratio of the number of rules pruned by the improved algorithm to the number of rules discovered by prior mining approaches.

$$\text{lost ratio} = \frac{|\{ \langle G_j, r \rangle \mid \langle G_j, r \rangle \text{ holds in MD } \text{ \&Gj,r> doesn't hold in } T[G_j] \}|}{|\{ \langle G_j, r \rangle \mid \langle G_j, r \rangle \text{ holds in MD } \}|} \quad (6)$$

Definition 2: lost ratio is the ratio of the number of rules discovered by the improved algorithm but lost in the previous mining approaches to the number of rules discovered by the improved algorithm.

4 EXPERIMENT AND EVALUATION

4.1 Experiment scenario on a case of hospital

A scenario for a medical center and related data were established to evaluate the performance of the proposed approach. The center contains various departments and a website for e-services. The testbench is implemented with Java on a PC Server with an AMD processor and the data mining software is implemented with Java.

Data from different departments of the medical center and the website are gathered for the experiment (ref. Figure 1). There are various attributes in the database of patients' records that may influence the healthcare behaviors. We take five of them, *viz.* (Address, Sex, Occupation, Age, and Marriage) as the dimensions for the test. Adding with the therapy/service catalog and the cost records, there are 7 dimensions, *i.e.* 7 concept-taxonomies for each dimension.

4.2 Experiment Data

The medial center provided basic patterns resulted from their mining tool and ca. 50K basic data. We then generated with Apriori-Generator three types of synthetic data sets respectively, as shown in table 4. There are 110 multidimensional patterns with respect to these taxonomies, where 40 of them are element patterns and the other 70 of them are generalized patterns. The proposed mining tool should find all large itemsets for the 70 generalized patterns.

Table 4: Three Types of Experimental Data Set

Type 1	To generate a single set of maximal potentially large itemsets and then generate transactions for each element pattern E_i following apriori-gen.[3]
Type2	Diagnosis-2, Therapy1. Beside a set of common maximal potential large itemsets, to generate maximal potentially large itemsets for each element pattern E_i . and then generate transactions for each element pattern E_i and the common maximal potentially large itemsets respectively following the apriori-gen[3]
Type3	generating a set of maximal potentially large item-sets for each element pattern E_i , and then generating transactions for each element pattern E_i from its own maximal potentially large itemsets following the apriori-gen.[3]

The first task for the evaluation is to determine the size of the transactions, where the size is picked from a Poisson distribution with the mean value μ equal to the average transaction size $|T|$. As the second step, each transaction is assigned a series of potentially large-itemsets. If the large-itemset on hand does not fit in the transaction, the itemset is put in the transaction randomly in half of the cases, and the itemset is fed to the next transaction of the rest. The number of maximal potentially large itemsets is set to the maximal size of potential large itemsets $|L|$. A maximal potentially large itemsets is generated by picking the size of the itemset from a Poisson distribution with mean μ equal to its average size $|I|$.

4.3 The Results of Experiment

At first, the 74 generalized patterns are successfully found. The key feature of the algorithm as illustrated in Figure 9 is that it is linear (and hence highly scalable) to the number of records and that it is flexible in terms of reading various data types. The test result w.r.t scalability in Figure 9 illustrates that the algorithm takes execution time linear to the number of transactions of all three data types. The experiment results of both the test (see Figure 8 and 9) illustrates that the new algorithm is superior to conventional methods in several areas:

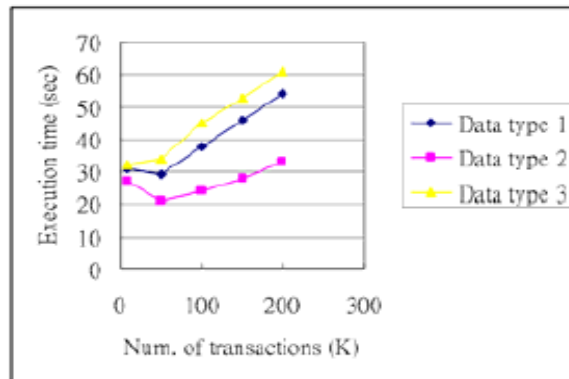


Figure 8: Scalability test w.r.t. the no. of transactions

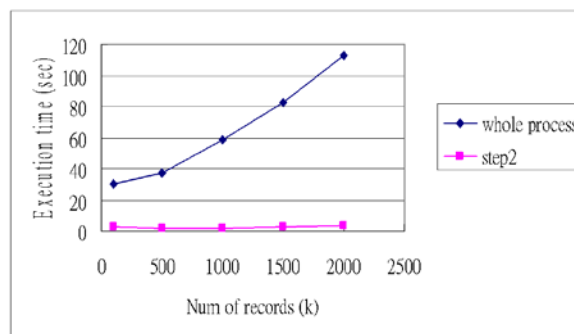


Figure 9: Scalability experiment w.r.t. the no. of records

Execution time with regards to number of transactions is linear for the data types tested for the whole process. This means that the time and space cost of executing our algorithm do not increase exponentially as compared to conventional methods.

Phase 2 (the update phase) of our algorithm is an important space and time saver as illustrated by the Figure 8; execution time is also linear and time taken to read up to 2000k records took less than 5 seconds. This means that data patterns from new data can be quickly extracted and used to update the existing pattern table for immediate use.

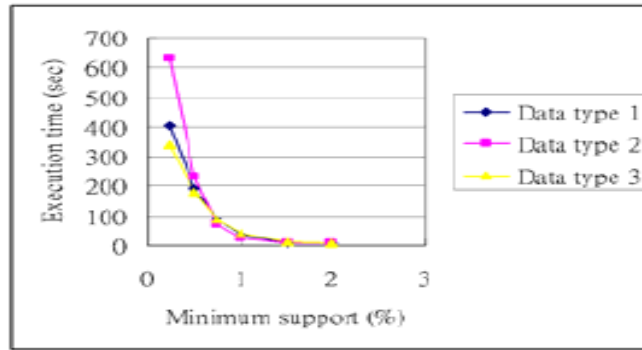


Figure 10: Efficiency in Relation to Minimum Support

In general, an increase of element patterns with result in an increase in execution time; the key to scalability is having the execution time increasing in a linear manner with an increase in element patterns. In Figure 11, all three data types experienced an increase of execution time with an increase of element pattern in a linear fashion, thus making our algorithm efficient.

Most importantly, an increase in element patterns leads to a less than proportion increase in execution time, making out the algorithm highly scalable. Reading off Figure 10, a 4 time increase of 30 element patterns from 10 to 40 will result in:

- 75 times increase in execution time for data type 1 from 20 seconds to 35 seconds.
- 1.67 times increase in execution time for data type 2 from 15 seconds to 25 seconds.
- 2.05 times increase in execution time for data type 3 from approximately 22 seconds to 45 seconds.

The impact of minSup on the algorithm can be categorized in terms of efficiency, discrete ratio and lost ratio. All of such algorithms are sensitive to the minimum support; the smaller the minimum support, the longer the execution time. However, we have shown that the real execution time of the step 2 (the update) in the proposed algorithm is relatively much shorter than the whole process (see Figure 8).

The test results proved that an increase in minSup will lead to greater returns of investment in terms of time efficiency; this is in line with one of the core objectives of building an efficient algorithm. Our algorithm is more efficient than conventional methods in terms of execution time over data. For instance in Figure 11, a 10 time increase (from 0.1 to 1) in minSup leads to a more than proportionate decrease in execution time across all data types:

- Execution time for data type 1 decreased by approximately 10 times, from approximately 400 seconds to approximately 40 seconds in terms of execution time.
- Execution time for data type 2 decreased by more than 30 times, from more than 600 seconds to approximately 20 seconds in terms of execution time.

- Execution time for data type 3 decreased by more than 11 times, from approximately 350 seconds to approximately 30 seconds in terms of execution time.

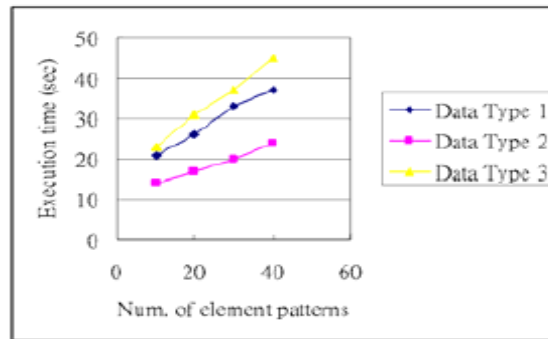


Figure 11: Efficiency in Relation to Minimum Support

The discrete ratio is the ratio of the number of rules pruned by the proposed algorithm to the number of rules discovered by prior mining approaches. Figure 12 illustrates the ratio of rules pruned by the proposed algorithm against minSup. In general, all three data types (except for data type 1) exhibited an increase of ratio with an increase of minSup from approximately 0.2% to 2%.

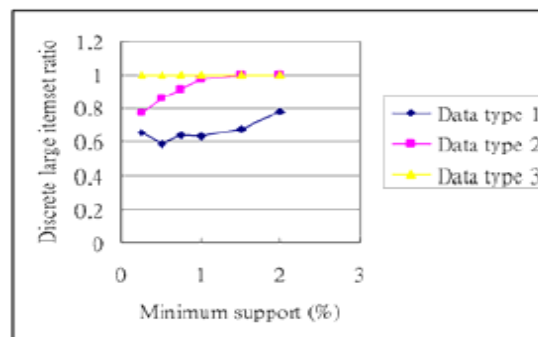


Figure 12: Effects of MinSup on discrete large itemsets ratio

The test results point the fact that the proposed algorithm can effectively decrease unwanted generalized patterns in which elemental data patterns is not true. This greatly helps users to focus on data patterns that are useful for their organizations while uncovering niche data patterns. For instance with a higher setting value, only <Female, Age 30-50, buy SK-II > will be found instead of <Age 30-50, buy SK-II>.

Figure 13 illustrates the test result on lost ratio, i.e. the influence of minSup values on the lost rules by other mining tools in comparison to this approach. All three data types experienced an increase in lost ratio over an increase in minSup from 0.25% to 2%, with the greatest increase in data type 2, followed by data type 3 and finally data type 1.

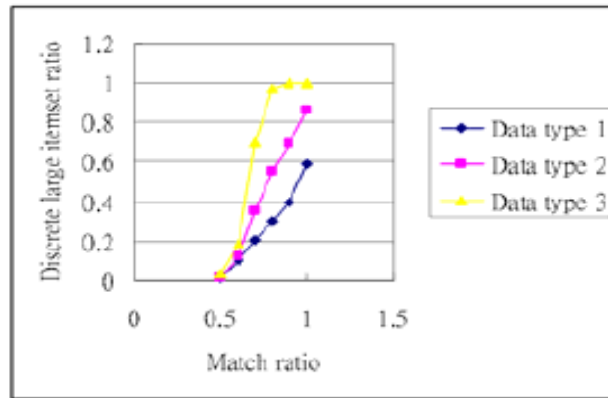


Figure 13: Effects of match ratio on discrete large itemsets ratio

The test results prove that the proposed algorithm will help users uncover useful data patterns which otherwise would be uncovered by traditional approaches. Thus, our objective of uncovering niche data patterns that would otherwise be left out is met and proved by this test result.

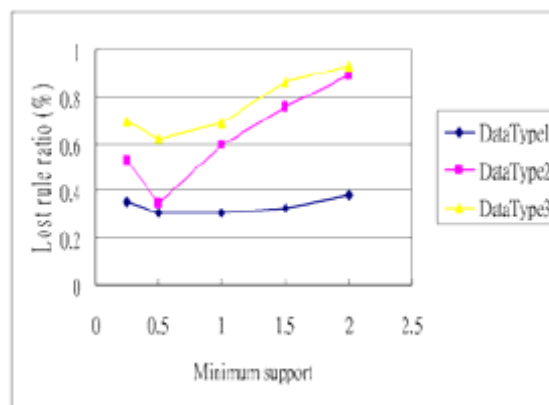


Figure 14: Effects of match ratio on lost itemsets.

Increasing the match ratio would decrease unwanted data patterns in general. Figure 13 illustrates the effect of match ratio (r) on discrete ratio. Similar to the above test results, an increase of m from 0.5 to 1 results in a more than proportional increase in discrete ratio across all three forms of data types. The significance of this test result is congruent with the test results above; the algorithm is efficient and scalable without losing flexibility and helps uncover niche data patterns.

5 SUMMARY

This paper presents at first the categories of innovative healthcare services as well as the way to find new service patterns. Then, we propose a data mining approach for managing such new healthcare services, including a novel data structure and an effective algorithm for multi-

dimensional mining association rules on various granularities. It is proved to be very useful for discovering new service patterns. The advantages of this approach over existing approaches include (1) more comprehensive and easy-to-use (2) more efficient with limited scans (3) more effective with finding rules hold in different granularity levels (4) capable of finding frequent patterns and infrequent patterns while users can choose the full match and the relaxed match (5) low information loss rate (6) capable of incremental mining of association rules to avoid unnecessary re-scan.

The design and evaluation of the multidimensional multi-granularity data mining approach were discussed in this paper. Since there is in our knowledge no metrics serving as the base for the measuring the data mining methods, we derive new metrics from Shannon's Entropy Function. The evaluation results prove the performances of the proposed approach, including efficiency, scalability and information loss rate, are better than existing approaches we know. The results show that we can use the proposed approach to find frequent and infrequent rules on different granularities by user-defined minSup value and match ratio.

Beyond the research so far, the effects of perceived issues and potential development of data mining without thresholds as well as concept description are worthy of further investigation.

REFERENCES

- [1]. R. Agrawal and J. C. Shafer (1996). "Parallel Mining of Association Rules," IEEE Transactions on Knowledge and Data Engineering, vol. 8, no. 6, pp. 962-969.
- [2]. R. Agrawal and R. Srikant (1994). "Fast Algorithms for Mining Association Rules in Large Databases," in Proceedings of the 20th International Conference on Very Large Data Bases.
- [3]. R. Agrawal, T. Imielinski and A. N. Swami (1993). "Mining Association Rules between Sets of Items in Large Databases," in Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data.
- [4]. J. K. Chiang (2007). "Developing an Approach for Multidimensional Data Mining on various Granularities ~ on Example of Financial Portfolio Discovery," in ISIS 2007 Proceedings of the 8th Symposium on Advanced Intelligent Systems, Sokcho City, Korea.
- [5]. J. K. Chiang and J. C. Wu (2005). "Mining Multi-Dimension Rules in Multiple Database Segmentation-on Examples of Cross Selling," in Proceedings of the 16th International Conference on Information Management, Taipei, Taiwan.
- [6]. T. M. Cover and J. A. Thomas (2006). *Elements of Information Theory*, 2nd ed., Wiley.
- [7]. R. Feldman and J. Sanger (2007). *The Text Mining Handbook – Advanced Approaches in Analyzing Unstructured Data*, Cambridge University Press.
- [8]. J. Han and M. Kamber (2006). *Data Mining - Concepts and Techniques*, 2nd ed., Morgan Kaufman.

- [9]. L. J. He, L. C. Chen and S. Y. Liu (2003) "*Improvement of AprioriTid Algorithm for Mining Association Rules,*" Journal of Yantai University(Natural Science and Engineering Edition), vol. 16, no. 4.
- [10]. B. Lent, A. Swami and J. Widom (1997). "*Clustering Association Rules,*" in Proceedings of the 13th International Conference on Data Engineering.
- [11]. M. Li and M. Baker (2005). *The GRID – Core Technologies*, Wiley.
- [12]. B. Liu, W. Hsu and Y. Ma (1999), "*Mining Association Rules with Multiple Minimum Supports,*" in Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.
- [13]. G. Shmueli, N. R. Patel and P. C. Bruce (2007)."*Association Rules,*" in *Data Mining for Business Intelligence, Concepts, Techniques, and Applications*, Wiley, pp. 203-215.
- [14]. R. Srikant and R. Agrawal (1995). "*Mining Generalized Association Rules,*" in Proceedings of the 21th International Conference on Very Large Data Bases, Zurich, Switzerland.
- [15]. W. Stallings (2004). "*Channel Capacity,*" in *Business Data Communications*, 6th ed., Prentice Hall, pp. 470-471.
- [16]. P. S. Tsai and C. M. Chen (2004). "*Mining interesting association rules from customer databases and transaction databases,*" *Information Systems*, vol. 29, no. 8, p. 685–696.
- [17]. C. Vercellis (2009). "*Association Rules,*" in *Business Intelligence, Data Mining and optimization for Decision Making*, Wiley, pp. 277-290.
- [18]. The CRISP-DM Consortium, CRISP-DM 1.0 (2000), www.crisp-dm.org.

The Author

Prof. Dr.-Ing. Johannes K. Chiang is now a faculty member of the Department of MIS and the Deputy Director of the Center for Cloud Computing and Operation Innovation at National Chengchi University Taipei. He received his academic degree of Doctor in Engineering Science (*Dr.-Ing., Summa Cum laude*) from the RWTH University of Aachen Germany. His current research interests include Cloud Computing, Semantic Web, Business Intelligence, Data Mining, e-Business and ebXML. He also serves as a consultant for several government agencies in Taiwan and as an active member of various international affiliations, such as IEEE, ACM, CSIM and ITMA etc. before 1995, he has been a research fellow at RWTH of Aachen and a Manager of EU/CEC ESPRIT Programmes.

Smart Obstacle Detector for Blind Person

Daniyal Rajput, Faheem Ahmed, Habib Ahmed, Engr Zakir Ahmed Shaikh, Aamir Shamshad
*Institute Of Biomedical Technology, Liaquat University Of Medical & Health Sciences Jamshoro,
and Mehran University of Engineering and Technology, Jamshoro, Pakistan;*
12bme11@gmail.com, faheemaffandi@gmail.com, h.habibahmedpirzada@yahoo.com,
urchoice_zakir@yahoo.com, aamir_shamshad01@yahoo.com

ABSTRACT

Smart obstacle detector helps blind people in moving and allowing them to perform their work easily and comfortably. In normal cane/stick, detection is done by the sensor. However, it is not much efficient because the blind person does not know what type of things or objects come in front of him, what is the size of that object and how far is he from the object? So it is difficult for blind person to move here and there. But SOD's output comes in two forms i.e. sound and vibration. In SOD, we detect the object by video processing method with the help of camera. For this we have used MATLAB software and then converted the video processing output into sound. Here we have used video processing for efficient and fast detection of objects, Stick measure distance between objects and SOD stick by Ultrasonic sensor. Moreover we have connected the vibrating motor with Ultrasonic sensor. When objects come in range of ultrasonic sensor then handle of the stick will vibrates. The vibration increases if the object comes toward the stick and vibration decreases when object goes far from the stick. We can also use solar chip for operating the stick which is definitely included in the future enhancement of our research based project.

KEYWORDS-: SOD (Smart Obstacle Detector), Video processing by MATLAB, Camera, Ultrasonic sensor, Arduino Board, vibration alarm.

1 INTRODUCTION

Goal of our research paper is to help the blind person. [1]More than 180 million people are visually disabled throughout the World. Of this group, 45 million people are completely blind. And of this group of 45 million blind people, 90% of them live in the developing world.[2] "In Pakistan, around two million people are blind, out of which 340,000 are in Sindh province, while cataract and diabetes are the major causes of blindness", said community

ophthalmologist and patron-in-Chief of the Disabled Welfare Association Karachi, Dr. M Shahnawaz Munami.



Figure 1: The condition of blind person, structural design, working principle of stick

When someone becomes blind in the developing world: Firstly; 90% of these individuals can no longer work. Secondly; Life expectancy drops down to 1/3 that of a matched peer, in age and health. Thirdly; 50% of the blind report a loss of social standing and decision-making authority. Furthermore; 80% of all women note a loss of authority within their families. (Report by Javitt, Int. Congr. Ophthalmol., 1983).

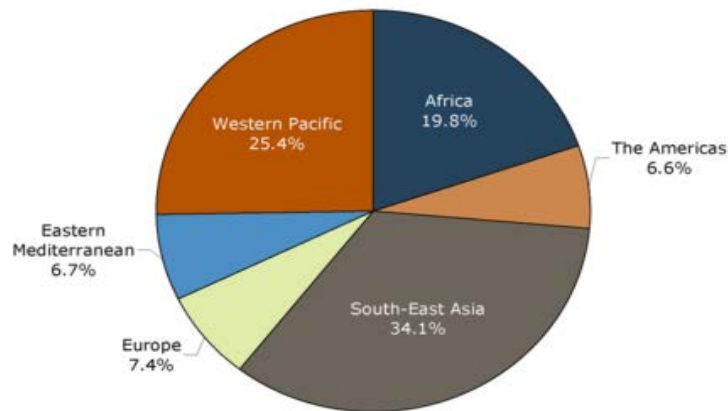


Figure 2 Geographical Distribution of Global Blindness

2 WORKING PRINCIPLE

Working principle of Smart Obstacle Detector for blind person consist of

- Video processing by Matlab Software
- Conversion of video processing into sound

- Ultrasonic sensor
- Ear phone
- Solar cell



Figure 3: Block Diagram of Smart Obstacle Detector for Blind Person

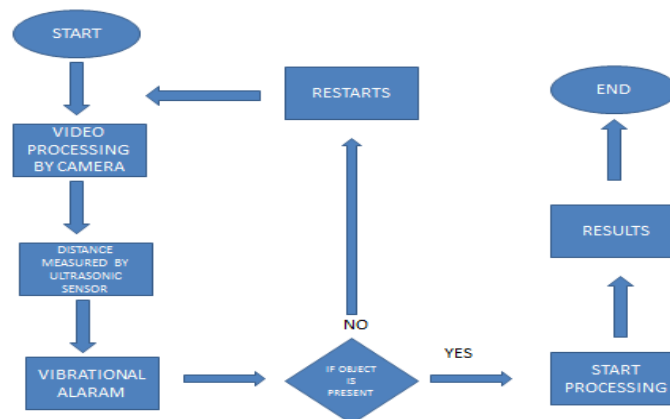


Figure 4: Flow Chart of Smart Obstacle Detector for Blind Person

2.1 Video processing by Matlab Software

2.1.1 Video Processing

Video processing systems [3] require a stream processing architecture, in which video frames from a continuous stream are processed one (or more) at a time. We use video processing for detection of objects because it gives fast and efficient response in which we use camera for detection purpose.

2.1.2 Procedure of Object Detection by Video Processing

Detection of moving objects and motion-based tracking are important components of many computer vision applications, including activity recognition, traffic monitoring, and automotive safety. The motion-based object tracking can be divided into two parts:

1. Detecting moving objects in each frame

2. Associating the detections corresponding to the same object over time

2.1.3 Create System Objects:

Create System objects used for reading the video frames, detecting foreground objects, and displaying results.

2.1.4 Initialize Tracks

The initialize Tracks function creates an array of tracks, where each track is a structure representing a moving object in the video. The purpose of the structure is to maintain the state of a tracked object. The state consists of information used for detection to track assignment, track termination, and display.

2.1.5 Detect Objects

The detect Objects function returns the centroids and the bounding boxes of the detected objects. It also returns the binary mask, which has the same size as the input frame. Pixels with a value of 1 correspond to the foreground, and pixels with a value of 0 correspond to the background. The function performs motion segmentation using the foreground detector. It then performs morphological operations on the resulting binary mask to remove noisy pixels and to fill the holes in the remaining blobs.

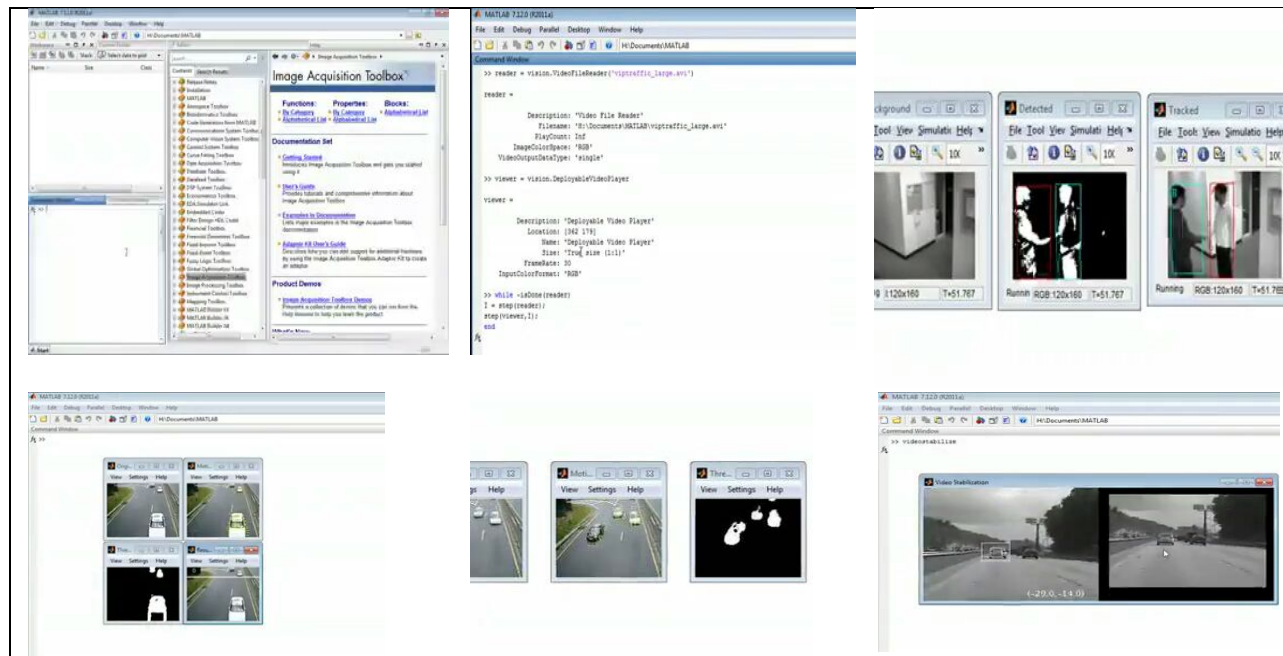


Figure 5: The video processing method

2.1.6 Predict New Locations of Existing Tracks

Use the Kalman filter to predict the centroid of each track in the current frame, and update its bounding box accordingly.

2.1.7 Assign Detections to Tracks

Assigning object detections in the current frame to existing tracks is done by minimizing cost. The cost is defined as the negative log-likelihood of a detection corresponding to a track.

2.1.8 Create New Tracks

Create new tracks from unassigned detections. Assume that any unassigned detection is a start of a new track. In practice, you can use other cues to eliminate noisy detections, such as size, location, or appearance.

2.1.9 Display Tracking Results

The display Tracking Results function draws a bounding box and label ID for each track on the video frame and the foreground mask. It then displays the frame and the mask in their respective video players.

2.2 Ultrasonic Sensor

2.2.1 Ultrasonic principle

In SOD the working principle of Ultrasonic sensor is similar to sonar which evaluate qualities of a target by interpreting the echoes from sound waves respectively. Ultrasonic sensors generate high frequency sound waves and evaluate the echo which is received back by the sensor. The time interval between the sent signal and received signal is determined to measure the distance from an object. [7] Ultrasonic sensors emit short, high-frequency sound pulses at regular intervals. These propagate in the air at the velocity of sound. If they strike an object, then they are reflected back as echo signals to the sensor, which itself computes the distance to the target based on the time-span between emitting the signal and receiving the echo. [8] Independent to target materials, surface and color. Detect small objects over long operating distance. Work very well by dust, dirt or high moisture environments. Resistant to external disturbances such as vibration, light, noises...[9]An ultrasonic sensor consists of a transmitter and receiver which are available as separate units or inserted together as single unit.

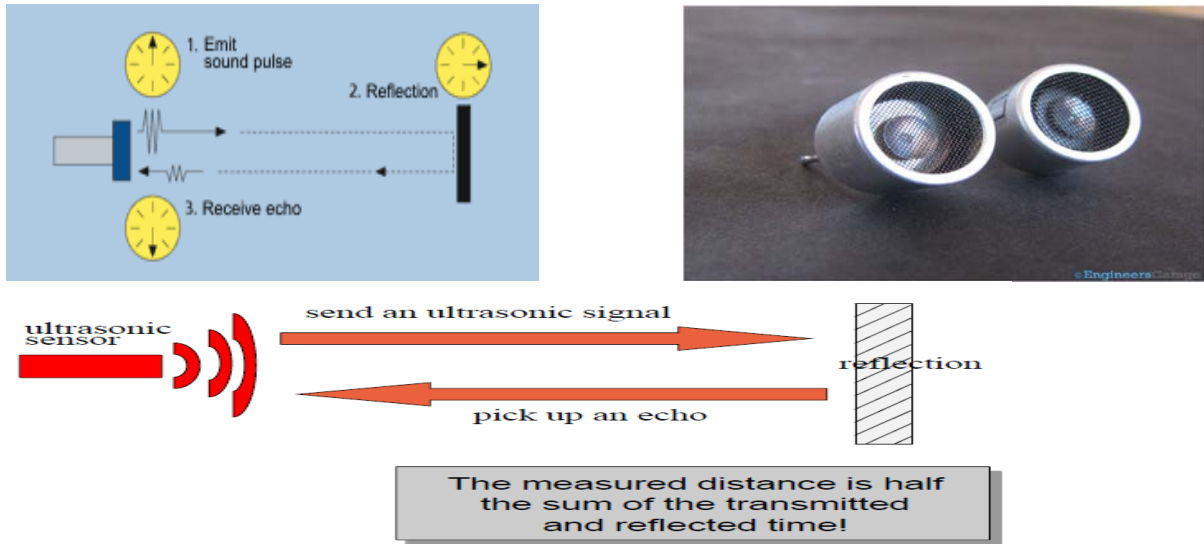


Figure 6: Working principle of ultrasonic sensor in which the sound signals are generated from the Ultrasonic Sensor which strike to the object and bounce back and receiver receives these signals.

The Timing diagram is shown below. You only need to supply a short 10uS pulse to the trigger input to start the ranging, and then the module will send out an 8 cycle burst of ultrasound at 40 kHz and raise its echo. The Echo is a distance object that is pulse width and the range in proportion .You can calculate the range through the time interval between sending trigger signal and receiving echo signal.

$$\text{Formula: } \mu\text{S} / 58 = \text{centimeters or } \mu\text{S} / 148 = \text{inch;}$$

$$\text{or: range} = \text{high level time} * \text{velocity (340M/S)} / 2;$$

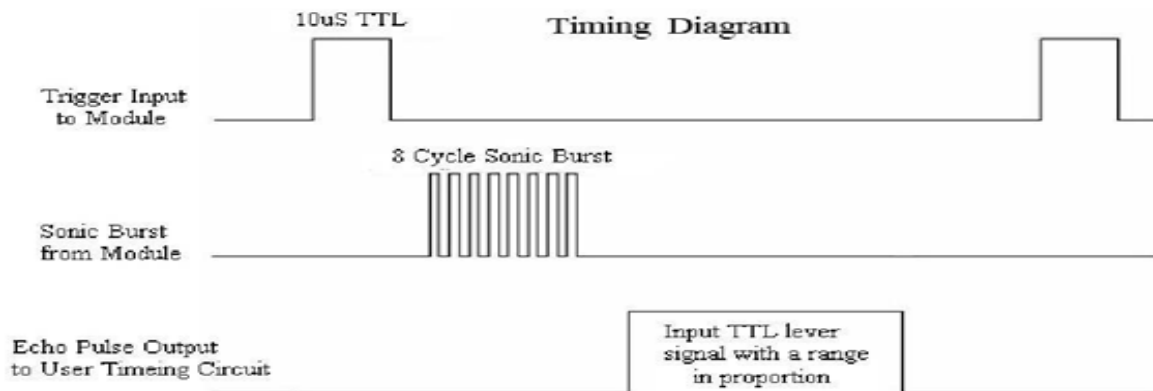


Figure 7 : The above figure shows the timing diagram of sound signals.

2.2.2 Arduino Board

We use Arduino board for operating Ultrasonic Sensor as it detects the presence of object in easy way and for this we do the programming in Arduino Board. It also reduces the size of Smart Obstacle Detector.

Schematic Diagram

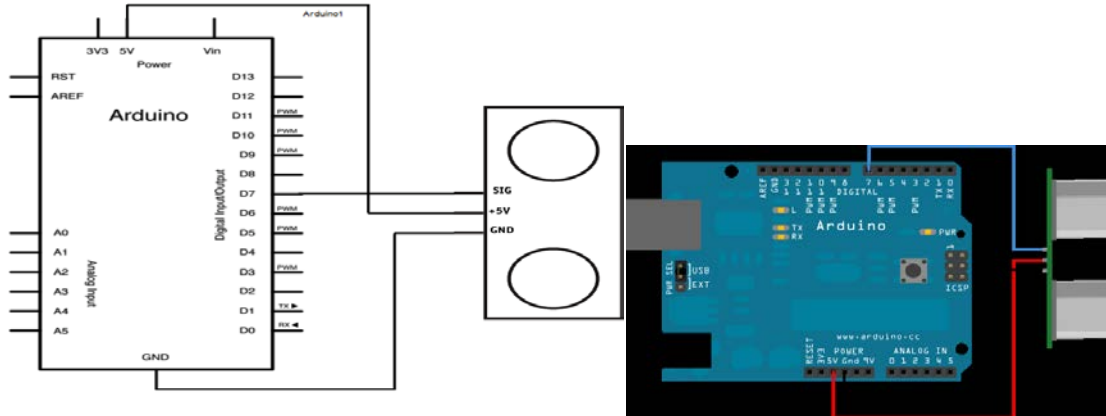


Figure 8: Schematic and layout diagram of Arduino connected with Ultrasonic sensor.

2.3 Vibration Alarm

Device currently uses Arduino Micro [29], Ping Ultrasonic sensor and vibrating motor. Micro-controller is used to measure the distance of the forth coming object by measuring time of flight of Ultrasonic waves emitted by Ping sensor. Depending on the distance, controller changes the vibration of the vibrating motors. Hence, Vibrating motor is attached with the sensor which vibrates faster if the distance between the obstacle and device decreases and vibration dampens as distance increases.

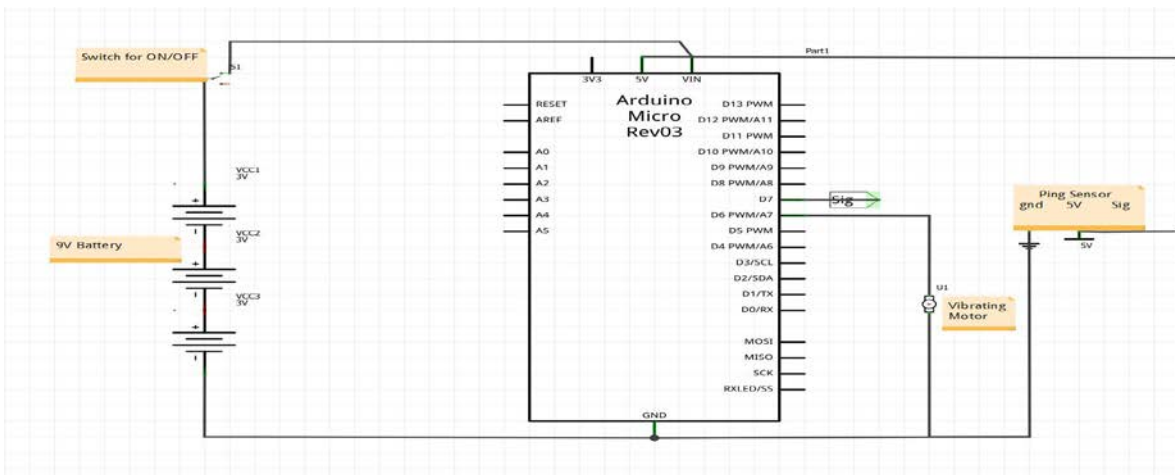


Figure 9: Schematic Diagram

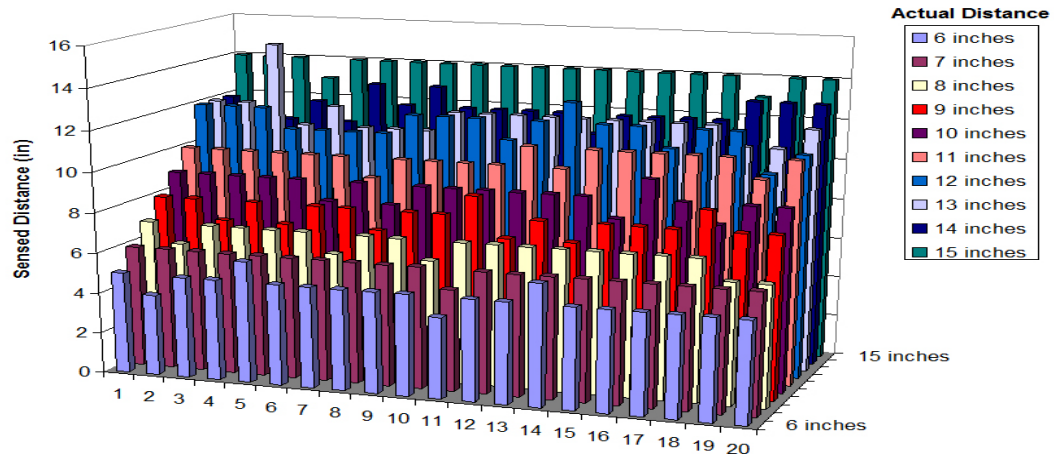


Figure 10: Sensed distance for 20 data points at a range of actual distances

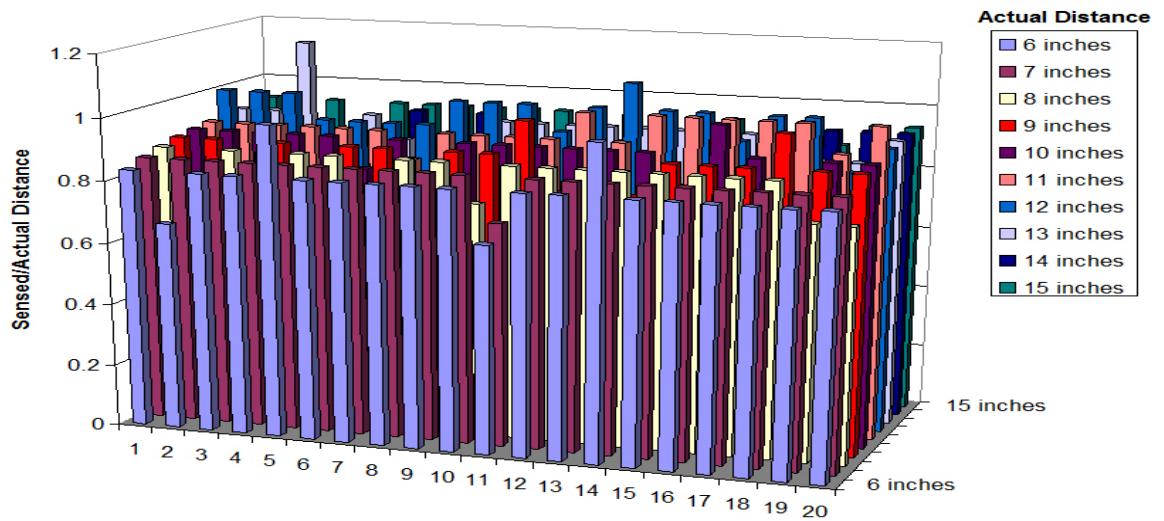


Figure 11: Sensed distance normalized by actual distance for 20 data points at a range of actual distances

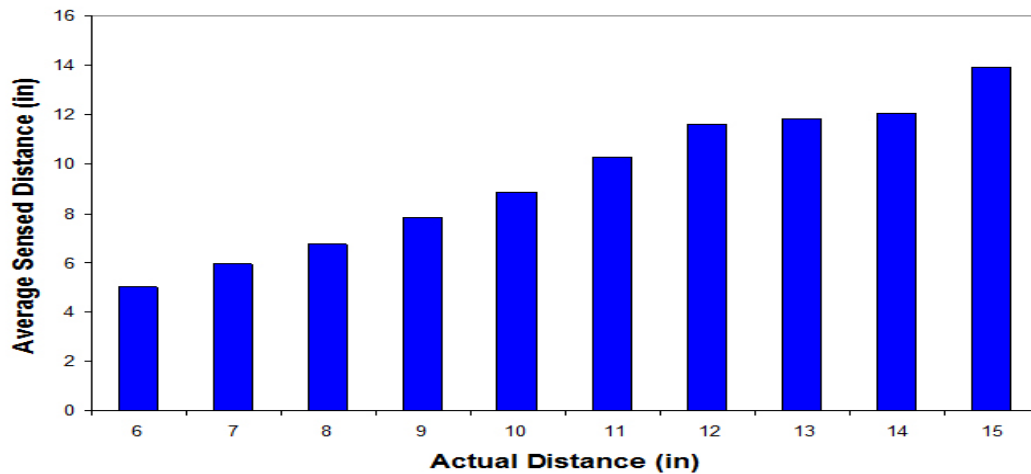


Figure 12: Average sensed distance of a large number of data points compared to actual distance

3 CONCLUSION

This paper has described the system and research mechanism which provides an immense help and support to blind persons. It is true that every organ of the body is very important and has its own and specific objectives. Similarly; eye is also a very important organ of the body. Unfortunately, blind persons' life is really colourless and is void of many happiest moments of life. The project will help the blind persons to detect the obstacles through the video processing mechanism by carrying just this small piece of stick (Smart obstacle detector). Ultimately; this research based project will result in serving the humanity which is indeed a greatest act.

REFERENCE

- [1]. <http://www.cureblindness.org/world-blindness/>
- [2]. <http://www.thenews.com.pk/Todays-News-4-173929-Two-million-blind-people-in-Pakistan>
- [3]. <http://www.mathworks.de/products/computer-vision/description4.html>
- [4]. P.Charbonnier, F.Diebolt, Y.Guillard, and F.Peyret. Road markings recognition using image processing. In IEEE Conf. on Intelligent Transportation System, pages 912-917, 1997.
- [5]. <http://www.mathworks.de/image-video-processing/video-processing.html>
- [6]. <http://www.mathworks.com/help/vision/examples/motion-based-multiple-object-tracking.html?prodcode=VP&language=en>
- [7]. Benjamin J.M., Ali N.A. laser cane for the blind in proceedings of San Diego Biomedical Symposium, volume 12, pages 53-57. 1973
- [8]. John Boren and Yoram Koren, The Guide Cane A computerized travel aid for the Active Guidance for the blind Pedestrian, proceedings of the IEEE international conference on Robotics and Automation, Albuquerque, 1997, page 1283-1288, April 21-27.
- [9]. C.Gearhart, A.Herold, B.Self, C.Birdsong, L.Slivovsky, Use of ultrasonic sensor in the development of an Electronic Travel Aid, Sensor Application Symposium, 2009, SAS 2009, IEEE, pp. 275-280, 17-19 Feb.
- [10]. R.Nagarjan, S.Yaqoob, and G.Sainarayanan, Role of Object Identification in Sonification System for Visually Impaired, In IEEE Tencon (IEEE Region 10 Conference On Convergent Technologies for the Asia Pacific), Bangalore, India, 2003, October 15-17.
- [11]. L.Kay, A Sonar Aid to Enhance Spatial Perception of the Blind: Engineering Design and Evaluation, Radio and Electronic Engineer, 1974, 44(11), pp 605-627.
- [12]. Z. W. Kim. Robust lane detection and tracking in challenging scenarios. *IEEE Trans. on Intelligent Transportation Systems*, 9(1):16–26, 2008.

- [13]. D. Schreiber, B. Alefs, and M. Clabian. Single camera lane detection and tracking. In *IEEE Conf. on Intelligent Transportation Systems*, pages 302–307, 2005.
- [14]. Felipe Jiménez, José Eugenio Naranjo, "Improving the obstacle detection and identification algorithms of a laserscanner-based collision avoidance system," *Transportation Research Part C*, Article in press.
- [15]. Kunsoo Huh, Jaehak Park, Junyeon Hwang, Daegun Hong, "A stereo vision-based obstacle detection system in vehicles," *Optics and Lasers in Engineering*, Vol. 46, pp. 168–178, 2008.
- [16]. Bruno Andò, Salvatore Graziani "Multisensor strategies to assist blind people: a clear-path indicator" *IEEE Trans. on Instrumentation and Measurement*, Vol. 58, No. 8, pp. 2488-2494, 2009.
- [17]. John F. Canny "A computational approach to edge detection" *IEEE Trans. on Pattern Analysis and Machine Intelligence*, pp. 679-698, 1986
- [18]. Richard Duda, Peter Hart "Use of the Hough transform to detect lines and curves in pictures" *Comm. ACM*, Vol. 15, pp. 11-15, 1972
- [19]. Zhengyin Zhou, Tianding Chen, Di Wu, Changhong Yu "Corridor navigation and obstacle distance estimation for monocular vision mobile robots" *JDCTA: Int. J. of Digital Content Technology and its Applications*, Vol. 5, No. 3, pp. 192- 202, 2011.
- [20]. Jeff Wilson, Bruce Walker, Jeffrey Lindsay, Craig Cambias, Frank Dellaert "SWAN: System for wearable audio navigation" In *Proc. 11th IEEE Int. Symp. on Wearable Computers*, pp. 91-98, 2007.
- [21]. <http://www.microsonic.de/en/Interesting-facts/Ultrasonic-principle.htm>
- [22]. Ultrasonic sensor distance measuring in almost any conditions-arteos GmbH/Germany-www.arteos.com
- [23]. http://www.seeedstudio.com/wiki/Ultra_Sonic_range_measurement_module.
- [24]. <http://arduino.cc>.
- [25]. Arduino Programming Notebook" by - by Brian W. Evans.
- [26]. Beginning Android ADK with Arduino by Mario Bohmer
- [27]. [http://www.robosoftsystems.co.in/wikidocs/index.php?title=Ultrasonic_Sensor_\(HC-SR04\)](http://www.robosoftsystems.co.in/wikidocs/index.php?title=Ultrasonic_Sensor_(HC-SR04))
- [28]. http://www.maxbotix.com/Ultrasonic_Sensors/People_Sensors.htm
- [29]. http://en.wikipedia.org/wiki/Proximity_sensor
- [30]. <http://www.dhavalmalaviya.com/my-research-art/eye-for-the-blind/>
- [31]. <https://wiki.engr.illinois.edu/display/ae498mpa/Ultrasonic>