

# International Journal of Applied Sciences and Smart Technologies

Volume 01, Issue 01, June 2019

**Indian Traffic Signboard Recognition and Driver Alert System Using Machine Learning**

Shubham Yadav, Anuj Patwa, Saiprasad Rane, Chhaya Narvekar

**Spur Gears Transmission Analysis on Countinous Passive Motion Machine Design for  
Shoulder Joint**

Felix Krisna Aji Nugraha, Antonius Hendro Noviyanto

**Influence of Annealing on the Electrical Properties of  $Ba_{0,5}Sr_{0,5}TiO_3$**

Dwi Nugraheni Rositawati

**Implementation of k-Medoids Clustering Algorithm to Cluster Crime Patterns in Yogyakarta**

Eduardus Hardika Sandy Atmaja

**Development Study of Deep Learning Facial Age Estimation**

Puspaningtyas S. Adi

**The Improvement of Watershed Algorithm Accuracy for Image Segmentation Handwritten  
Numbered Musical Notation**

Kartono Pinaryanto

**Factors Influencing the Difficulty Level of the Subject: Machine Learning  
Technique Approaches**

Hari Suparwito

ISSN 2655-8564

## **CONTENTS**

<b>CONTENTS</b>	i
<b>EDITORIAL BOARD</b>	ii
<b>PREFACE</b>	iii
<b>Indian Traffic Signboard Recognition and Driver Alert System Using Machine Learning</b> <i>Shubham Yadav*</i> , <i>Anuj Patwa</i> , <i>Saiprasad Rane</i> , <i>Chhaya Narvekar</i>	1–10
<b>Spur Gears Transmission Analysis on Countinous Passive Motion Machine Design for Shoulder Joint</b> <i>Felix Krisna Aji Nugraha<sup>1,*</sup></i> , <i>Antonius Hendro Noviyanto<sup>2</sup></i>	11–22
<b>Influences of Annealing on the Electrical Properties of Ba<sub>0,5</sub>Sr<sub>0,5</sub>TiO<sub>3</sub></b> <i>Dwi Nugraheni Rositawati</i>	23–32
<b>Implementation of k-Medoids Clustering Algorithm to Cluster Crime Patterns in Yogyakarta</b> <i>Eduardus Hardika Sandy Atmaja</i>	33–44
<b>Development Study of Deep Learning Facial Age Estimation</b> <i>Puspaningtyas Sanjoyo Adi</i>	45–50
<b>The Improvement of Watershed Algorithm Accuracy for Image Segmentation Handwritten Numbered Musical Notation</b> <i>Kartono Pinaryanto</i>	51–64
<b>Factors Influencing the Difficulty Level of the Subject: Machine Learning Technique Approaches</b> <i>Hari Suparwito</i>	65–82
<b>AUTHOR GUIDELINES</b>	83

**EDITORIAL BOARD**

**Editor in Chief**

Dr. I Made Wicaksana Ekaputra (*Sanata Dharma University, Yogyakarta, Indonesia*)

*Email: made@usd.ac.id*

**Associate Editor**

Dr. Pham Nhu Viet Ha (*Vietnam Atomic Energy Institute, Hanoi, Vietnam*)

Dr. Hendra Gunawan Harno (*Gyeongsang National University, Jinju, The Republic of Korea*)

Dr. Iswanjono (*Sanata Dharma University, Yogyakarta, Indonesia*)

Dr. Mukesh Jewariya (*National Physical Laboratory, New Delhi, India*)

Dr. Mongkolserj Lin (*Institute of Technology of Cambodia, Phnom Penh, Cambodia*)

Dr. Yohanes Baptista Lukiyanto (*Sanata Dharma University, Yogyakarta, Indonesia*)

Dr. Apichate Maneewong (*Thailand Institute of Nuclear Technology, Bangkok, Thailand*)

Dr. Sudi Mungkasi (*Sanata Dharma University, Yogyakarta, Indonesia*)

Dr. Pranowo (*Universitas Atma Jaya Yogyakarta, Yogyakarta, Indonesia*)

Dr. Mahardhika Pratama (*Nanyang Technological University, Singapore*)

Dr. Augustinus Bayu Primawan (*Sanata Dharma University, Yogyakarta, Indonesia*)

Prof. Dr. Leo Hari Wiryanto (*Bandung Institute of Technology, Bandung, Indonesia*)

**Editorial Proofreader**

Ir. Ignatius Aris Dwiatmoko, M.Sc. (*Sanata Dharma University, Yogyakarta, Indonesia*)

P. H. Prima Rosa, S.Si., M.Sc. (*Sanata Dharma University, Yogyakarta, Indonesia*)

**Editorial Assistant**

Eduardus Hardika Sandy Atmaja, M.Cs. (*Sanata Dharma University, Yogyakarta, Indonesia*)

Vittalis Ayu, M.Cs. (*Sanata Dharma University, Yogyakarta, Indonesia*)

**Administration**

Catharina Maria Sri Wijayanti, S.Pd. (*Sanata Dharma University, Yogyakarta, Indonesia*)

**Contact us**

International Journal of Applied Sciences and Smart Technologies

Faculty of Science and Technology

Sanata Dharma University

Kampus III Paingan, Maguwoharjo, Depok, Sleman

Yogyakarta, 55282

Phone : +62 274883037 ext. 523110, 52320

Fax : +62 272886529

Email : editorial.ijasst@usd.ac.id

Website : <http://e-journal.usd.ac.id/index.php/IJASST>

**IJASST** is an open-access peer-reviewed journal that mediates the dissemination of research and studies conducted by academicians, researchers, and practitioners in science, engineering, and technology.

## **PREFACE**

We are honored to announce the first issue of the journal *International Journal of Applied Sciences and Smart Technologies* (IJASST), which is managed and published by the Faculty of Science and Technology, Sanata Dharma University. IJASST is an open-access peer-reviewed journal that mediates the dissemination of research and studies conducted by academicians, researchers, and practitioners in science, engineering, and technology. Its scope also includes basic sciences which relate to technology, such as applied mathematics, physics, and chemistry. IJASST accepts submissions from all over the world, which of course, including Indonesia.

As a result of recent trends in research and development in the field of science and technology, it has become important to provide the most recent information and technology, not only on leading-edge research in specialist areas, but also on research and development to solve cross-cutting issues. IJASST will be published twice in a year (June and December). Submitted papers will be reviewed fairly and promptly using the open journal system (OJS) of IJASST. After the review process, accepted papers of the journal will be publicly available for free at the website of IJASST.

We are sure that IJASST will provide the opportunity to gain and present authentic as well as insightful scientific and technological information on the latest advances in science and technology. For future issues, we are looking forward to your submissions to IJASST.

Dr. I Made Wicaksana Ekaputra  
Editor in Chief  
IJASST

# **Indian Traffic Signboard Recognition and Driver Alert System Using Machine Learning**

Shubham Yadav\*, Anuj Patwa, Saiprasad Rane, Chhaya Narvekar

*Xavier Institute of Engineering, Mahim Causeway, Mahim (West), Mumbai,  
Maharashtra 400016, India*

*\*Corresponding Author: shubham.yadav.5497@gmail.com*

(Received 20-04-2019; Revised 14-05-2019; Accepted 14-05-2019)

## **Abstract**

Sign board recognition and driver alert system which has a number of important application areas that include advance driver assistance systems, road surveying and autonomous vehicles. This system uses image processing technique to isolate relevant data which is captured from the real time streaming video. The proposed method is broadly divided in five part data collection, data processing, data classification, training and testing. System uses variety of image processing techniques to enhance the image quality and to remove non-informational pixel, and detecting edges. Feature extractor are used to find the features of image. Machine learning algorithm Support Vector Machine (SVM) is used to classify the images based on their features. If features of sign that are captured from the video matches with the trained traffic signs then it will generate the voice signal to alert the driver. In India there are different traffic sign board and they are classified into three categories: Regulatory sign, Cautionary sign, informational sign. These Indian signs have four different shapes and eight different colors. The proposed system is trained for ten different types of sign. In each category more than a thousand sample images are used to train the network.

**Keywords:** image processing, signboard detection, svm algorithm, raspberry pi

## 1 Introduction

Recognition of signboard correctly at the right time and at the right place is very important for drivers to insure themselves and their passengers' safe journey. However, sometimes, due to the change of weather conditions or viewing angles, signs are difficult to be seen until it is too late. Now a days increases in computing power have brought computer vision to applications. On the other hand, the increase in traffic accidents accompanying the increasing amount of traffic has become a serious problem for society. The road accidents is particularly high under special road conditions, such as at the entrance to a one-way street, sharp curves, and intersections. One possible countermeasure is to install "STOP", "NO LEFT TURN" and other signs in order to notify the driver of the road conditions and other traffic information. Anyhow, there remains the possibility that the driver who is depending on his/her state of mind, fail to notice the sign while driving, a serious accident is possible if the driver fails to notice a sign such as "DO NOT ENTER", "STOP" etc [1].

It is possible that accidents can be prevented by utilizing an automatic sign board recognition system to provide traffic information to the driver, including information about the road in front of the vehicle. Signs also have distinct shapes like circles, triangles, rectangles and octagons. These systems assist drivers to drive safely. While driving the vehicle the driver gets the alert message like go slow, ahead is speed breaker. There are many detection techniques developed in recent days for traffic light and sign board detection. A system which involves detection process of traffic sign and sending the alert message does not exist. So keeping attention towards different traffic signs are difficult task for every drivers. So we proposed a system that can be used to detect traffic sign board. Traffic signs detection is an important part of driver assistant systems. These can be designed in different colors or shapes, in high contrast background. So in order to capture these images, traffic signs are oriented upright and facing camera. Hence there will be geometric and rotational distortions. In these cases accuracy is a key consideration. Any miss classified or undetected sign and lights will

produce adverse impacts on system. The basic idea of proposed system is to provide alertness to the driver about the presence of traffic sign at a particular distance apart. This system will be able to detect, recognize and infer the road traffic signs would be a prodigious help to the driver. The objective of an automatic road signs recognition system is to detect and classify one or more road signs from within live color images captured by a camera. The color of a traffic sign is easily distinguishable from the colors of the environment.

In this we provide alertness to the driver about the presence of signboard at a particular distance apart. The system provides the driver with real time information from road signs, which consist the most important and challenging tasks. Next generate an voice warning to the driver in advance of any danger. This warning then allows the driver to take appropriate corrective decisions in order to mitigate or completely avoid the event. First, it is necessary to select the hardware equipment to solve this problem. The second stage is based on color processing or object detection method based on rapid color changes. Image processing technology is mostly used for the identification of the signboards. The alertness to the driver is given as audio output.

## 2 Review of Literature

There are many algorithms and methodologies have been proposed for road traffic sign detection [2-6]. Reza Azad proposed the system with Iranian Traffic signs with detection and recognition and the letters are segmented with SVM classifier. Another method has also been proposed by Gauri Tagunde based on color and shape Features by Detection and Recognition approaches have been proposed to deal with sign board detection and recognition. Most of these systems typically involve two tasks finding the locations and sizes of sign board in natural scene images (sign board detection) and recognizing the detected signs board to interpret its meaning (sign board recognition). Being designed with regular shapes and conspicuous colours, sign board attract human driver attention so as to be easily captured by human drivers. Mohammad Amen proposes the system with YCbCr colour space and shape based filtering the detected traffic signs are tracked and recognized using interest point descriptors. The algorithm is robust and can detect signs even when the traffic sign board is rotated. The traffic sign

template database can be updated easily. The method is aimed at achieving high accuracy in recognizing traffic signs at real-time, with a low computational cost. Reduced computational complexity of the algorithm enables the implementation of the proposed method in embedded systems for driver assistance. In the case of traffic sign detection majority of system make use of colour information as a method for segmenting images. The performance of colour based road sign detection is often reduced in scenes with strong illumination, poor lightning or adverse weather conditions. The vast majority of the existing systems consist of hand label real images which are repetitive time consuming and error prone process. Information about traffic symbols, such as shape and colour, can be used to place traffic symbols into specific groups; however there are several factors that can hinder effective detection and recognition of traffic signs. These factors include variations in illumination occlusion of signs, motion blur, and weather –worn deterioration of signs. Road scene is also generally much cluttered and contains many strong geometric shapes that could easily be misclassified as road signs. Accuracy is a key consideration because even one misclassified or detected sign could have an adverse impact on the driver.

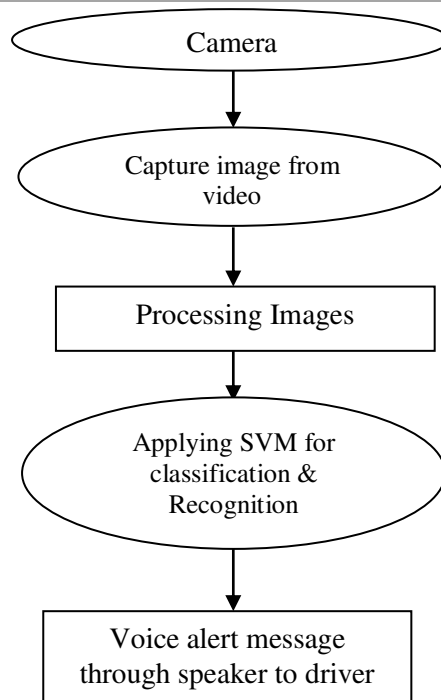
### 3 Design

Many times we see that many road accidents take place. This can be due to driver's ignorance of traffic sign board and road signs. As the road traffic is increasing day by day there is a necessity of following the traffic rules with proper discipline. Traffic signboard detection is an important part of driver assistant systems. The basic idea of proposed system is to provide real time voice signal to the driver about the presence of traffic sign board at a particular distance apart. The project is divided in to two part:

1. Training
2. Implementation

The system provides the driver with real time information from road sign board, which consist the most important and challenging tasks. It generates an voice signal to the driver in advance of any danger. This warning allows the driver to take some appropriate actions in order to avoid the accident.





**Figure 1.** Flow diagram of the Proposed System

The alertness to the driver is given as a voice signal through speaker as an output. There are two methods to classify the images in machine learning, convolution neural network (CNN) and Support Vector Machine (SVM). The proposed system uses support vector machine (SVM) for classification.

## **4 Working of the System**

Support Vector Machine is a supervised machine learning algorithm which is also known as the linear classifier mostly used for the classification purpose. The main advantage of SVM algorithm is its strong ability to classify any data. When the dataset has a clear classification boundary in that situation SVM is the best option than other available methods.

SVM is considered as one of the best classifiers and it is simple to use and understand than other classifiers. The proposed system uses 90% of sample data for training and 10% for testing. The working of the system is broadly divided into three phases:

1. Color Segmentation
2. Shape Classification
3. Recognition

## ***Phase 1: Color Segmentation***

In this phase candidate blob are extracted from the input image. Color segmentation phase is important phase because color every traffic sign are such that they appears different from the surrounding environment. HSI color space of image processing techniques used for segmenting the color. This is basically detection where region of interest is identified by using image processing techniques. Using the image processing technique system creates contours on each video frame and finds ellipses and circle among those contours. Detection strategy includes, increasing the contrast of video frame, removing unnecessary colors like green with HSV color range, using Laplacian of Gaussian to display the boarder of the candidate blob, making contour by binarization and detecting the ellips like and circle like contours.

## ***Phase 2: Shape Classification***

The candidate blob that are extracted from the video frame of the segmentation phase are now need to classify. The classification of these candidate are based on the shape. For classification of the candidate blob based on the shape linear SVM is used. There are two major task involved in shape classification.

### ***1. Shape Feature Extraction :***

First step in shape classification is to make feature vectors for the input to the linear SVM. Many methods have been proposed for obtaining the feature vectors (see [7, 8]). In this work, we have used distance to border vector (DtB) [8]. DtB stands for the distance of the blob from the external edge of the blob to its bounding box.

### ***2. Training and Testing Using Linear SVM :***

Once the feature Vector for the ROI is created then the classification is initiated. For Classification of the shape eight linear SVM is used. SVM is machine learning algorithm which can classify the data in different group. It is based on concept of decision plane where the training data is mapped to higher dimensional space and separated by plane defining two or more classes of data. The extensive introduction can befound in [9, 10].

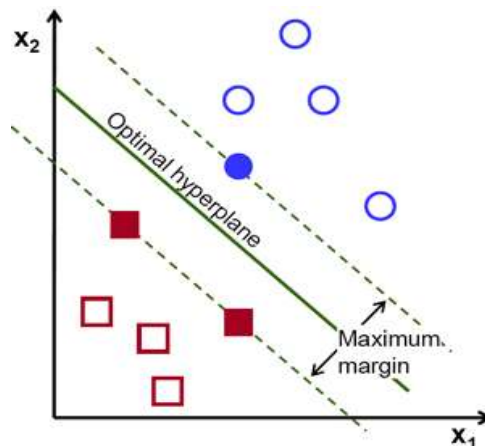


Figure 2. Shows possible hyper planes [11]

The proposed system is trained for ten traffic signs and the image of signs we can see in fig.3. In this 90% of the sample data is used to trained the system and 10% of the sample data is used for testing the system.



Figure 3. Trained traffic signs

### Phase 3: Recognition

Once the shape classification process is done then the next step is sending the blob to pattern recognition stage. To perform the recognition of pattern the radial basis function (RBF) is used. In this phase non linear SVM is used to recognized. In this classified blob is first converted into gray scale image then applying feature extractor to extract the features of the blob. Non linear SVM is used to recognition purpose in which extracted features are compared with all blob that are having the same shape and color. If the blob features matches with the trained signs features then system generates the

alert message calling the label of that class. The alert message is given in the form voice signal through speaker.

## 5 Working Result of Proposed System

The proposed system recognized almost all the traffic sign correctly when the traffic sign are stable while accuracy of the system is decreases while in motion. The environment and light also have the adverse effect on the system. sometimes the images which is captured from the real time streaming video have high contrast or low contrast in such cases system were not able to detect the traffic sign. So performance of the propose system in different environment not well but if the system have sample images with that environment then it works well. So we can say that accuracy of the system depend on number of sample images for a particular sign in that environment. Due to text to speech converter API some time the voice signal are delayed by few second.

## 6 Conclusions

The performance of the proposed system is quite good when system move slowly keeping signboard stationary but performance of the system while moved fast are not as per expectation. Environment and light also affect the system performance. Sometime due to text to speech converter API alert signal was getting delay. According to statistical report 3 death happens every 10 minutes due to road accident in India. On successful implementation of this project we expect to drastic reduction in road accident

## References

- [1] G. Revathi and G. Balakrishnan, “Indian sign board recognition using image processing techniques,” *International Journal of Advanced Research in Biology Engineering Science and Technology*, **2** (15), 326–330, 2016.
- [2] R. Azad, B. Azad, and I. T. Kazerooni, “Optimized method for Iranian road signs detection and recognition system,” *International Journal of Research in Computer Science*, **4** (1), 19–26, 2014.
- [3] K. M. Sumi and K. M. N. Arun, “Detection and recognition of road signs,” *International Journal of Computer Applications*, **160** (3), 1–5, 2017.

- [4] A. P. T. Agnes, C. A. Aiswarya, A. Augustine, A. S. Kumar, and N. Aswathy, “Real time traffic light and sign board detection,” *International Journal of Engineering Research and General Science*, **5** (3), 50–57, 2017.
- [5] W. Zhang, “Shift-invariant pattern recognition neural network and its optical architecture,” *Proceedings of Annual Conference of the Japan Society of Applied Physics*, p. 734, 1988.
- [6] <https://medium.com/@RaghavPrabhu/understanding-of-convolutional-neural-network-cnn-deep-learning-99760835f148> (Accessed on 14-05-2019).
- [7] P. G. Jiménez, H. G. Moreno, P. Siegmann, S. L. Arroyo, and S. M. Bascón, “Traffic sign shape classification based on support vector machines and the FFT of the signature of blobs,” *Proceedings of the 2007 IEEE Intelligent Vehicles Symposium*, 375–380, Istanbul, 13-15 June 2007.
- [8] S. L. Arroyo, P. G. Jiménez, R. M. Bascón, F. L. Ferreras, and S. M. Bascón, “Traffic Sign Shape Classification Evaluation I: SVM using Distance to Borders,” *Proceedings of IEEE Intelligent Vehicles Symposium*, 557–562, Las Vegas, June 2005.
- [9] S. Abe, *Support Vector Machines for Pattern Classification*, Springer, London, 2005.
- [10] C. C. Chang and C. J. Lin, *LIBSVM: A Library for Support Vector Machines*, 2001, <https://www.csie.ntu.edu.tw/~cjlin/papers/libsvm.pdf> (Accessed on 14-05-2019).
- [11] <https://towardsdatascience.com/support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47> (Accessed on 14-05-2019).

This page intentionally left blank

# Spur Gears Transmission Analysis on Countinous Passive Motion Machine Design for Shoulder Joint

Felix Krisna Aji Nugraha<sup>1,\*</sup>, Antonius Hendro Noviyanto<sup>2</sup>

<sup>1</sup>*Department of Mechatronic Product Design,  
Politeknik Mekatronika Sanata Dharma, Yogyakarta, Indonesia*

<sup>2</sup>*Department of Electromedical Technology,  
Politeknik Mekataronika Sanata Dharma, Yogyakarta, Indonesia*

*\*Corresponding Author: felix@pmsd.ac.id*

(Received 14-05-2019; Revised 15-05-2019; Accepted 15-05-2019)

## Abstract

An analysis of gear transmission on a continuous passive machine (CPM) from the 3-dimensional design has been carried out using SolidWorks software. Analysis of the strength of the gear structure is affected by the weight of the patient's arm. Analysis of gear transmission that is affected by the load of the passive arm uses static simulation, by entering the patient's arm load. The facilities used are static simulation with the condition of fixed geometry in the parts related of the shaft, the effect of gravity of  $10 \text{ m/s}^2$ , making mesh, and running simulation. The maximum stress that occurs in gear<sub>3</sub>  $z = 100$  is  $4.5524e + 006 \text{ N/m}^2$ , the maximum stress on gear<sub>2</sub>  $z = 80$  is  $4.81729e + 006 \text{ N/m}^2$ , the maximum stress on gear<sub>100</sub>  $z = 20$  is  $9.08982e + 006 \text{ N/m}^2$ .

**Keywords:** mesh, static, simulation, continuous passive machine, spur gear.

## **1 Introduction**

Continuous passive machine (cpm) is a therapeutic tool used to train patients in joint movements after joint surgery [1]. The CPM machine is designed to move flexion and horizontal adduction using spur gear transmission. The gears are used to reduce the speed from the actuator and increase torque, so that it can passively move the patient's shoulder to do therapy. With the development of computer aided design (CAD) technology is very helpful in designing a product or machine. The process of designing in manufacturing industries used a lot of time. An engineer who has experience in using CAD can use various tools/facilities in CAD software in various applications in mechanical engineering, so that the time spent designing can be done shorter, productivity and quality can be produced better. One CAD software for design and analyzing a 3-dimensional design is Solidworks.

The purpose of this study was to analyze the strength of material from the gear transmission to the patient's arm load. The strength of material of the gears are analyzed by the stress and strain on the gears. All analyzes of this study use SolidWorks software. Our main references are [2-5].

## **2 Research Methodology**

To complete the analysis in this study, the steps taken before doing the simulation are as follows.

### ***2.1. Research methods***

The research methods is done by designing 3 Dimensional CPM machine, from the results of the design will be analyzed the transmission of straight gears. By using software, the stress and strain experienced by each gear that is exposed to the specified load will be carried out. The research method is carried out by the following steps:

- a. Collect the geometry of the CPM machine that will be designed.
- b. Design 3 Dimensional CPM machine with spur gear transmission using Solidworks software.



- c. Analyzing the strength of material of the spur gear transmission of the CPM machine designed with the influence of the load determined using Solidworks software.
- d. Analyzing the strength of material of the spur gear of the CPM machine using Solidworks software.

## **2.2. Design Methods of 3-Dimensional CPM Machine**

In the CPM machine that will be designed there are 2 gearboxes that are used to reduce speed and increase the torque of a DC motor actuator. Two gearboxes are used for flexion movement and horizontal adduction has the same transmission pair. The method for carrying out the analysis of CPM gears transmission using the SolidWorks software is as follows:

- a. Design spur gear
  - 1. Determine the patient's arm load
  - 2. Determine the torque that is required for the movement of the CPM machine
  - 3. Determine the gear ratio used in design of the CPM machine
  - 4. Determine the level of spur gear transmission that used at CPM machine
  - 5. Determine dimensions, modules, number of teeth of spur gears based on the ratio of each gear transmission level
  - 6. Determine the material of spur gear in the design of CPM machine

- b. Design of spur gear transmission.

It is assumed that the patient's arm weight is 5 kg and the patient's arm is 80 cm long, so that the torque produced by the arm is 20 Nm. Shown in the equation below.

$$T = \frac{1}{2} \text{ arm length} \times \text{ arm weight} \times \text{ gravity}$$

$$T = 40 \text{ cm} \times 5 \text{ kg} \times 10 \text{ m/s}^2$$

$$T = 20 \text{ Nm}$$

The Dc motor actuator used has a torque characteristic of 1.5 Nm and a rotational speed of 70 rpm, so the gear transmission ratio is obtained by the equation

$$\tau_{nc} = \tau_d \times i \times \mu$$

Symbol and description:

$\tau_{nc}$  : torque needed

$\tau_d$  : DC motor torque

$i$  : ratio

$\mu$  : efficiency 75%

so obtained :

$$\tau_{nc} = \tau_d \times i \times \mu$$

$$i = \frac{\tau_{nc}}{\tau_d \times \mu}$$

$$i = \frac{20}{1,5 \times 0,75}$$

$$i = 17,75$$

From the results of the calculation the ratio is obtained at 17.75. The ratio value is increased for the safety factor to 20. The total ratio can be divided into 2 levels of gear transmission which are equal to 4 and 5. The gears are directly connected to the actuator gear<sub>1</sub>, then transmitted to gear<sub>2</sub>. Gear<sub>2</sub> and gear<sub>3</sub> are on the same axis. Next gear<sub>3</sub> is transmitted to gear<sub>4</sub>. All gears are determined using the same module, which is equal to 1mm. Gear<sub>1</sub> is determined to have 20 teeth. So that in the first transmission the number of teeth obtained in gear<sub>2</sub> is equal to:

$$i_1 = \frac{z_2}{z_1}$$

$$z_2 = i_1 \times z_1$$

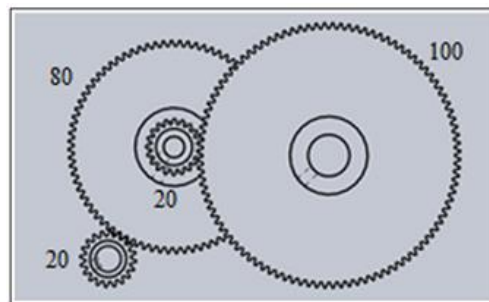
$$z_2 = 4 \times 20$$

$$z_2 = 80$$

For second level spur gear transmissions obtained:

$$i_2 = \frac{z_2}{z_1}$$
$$z_2 = i_2 \times z_1$$
$$z_2 = 5 \times 20$$
$$z_2 = 100$$

From these calculations of the level of the spur gears transmissions ore obtained as shown in figure 1.



**Figure 1.** Spur gears transmission in CPM machine

c. Simulation Method of structural strength in Solidworks is as follows

1. Use Solidworks Simulation facility.
2. Select the static test form.
3. Insert the material used for spur gear.
4. Determine the fixed of part design.
5. Determine the gravity on the spur gear.
6. Determine the part of spur gear that affected by patients arm weight/load and enter the value of the.
7. Create a mesh on the spur gears.
8. Run the simulation.

### **3 Results and Discussions**

In this section will explain the results and discussion of the simulation of the spur gears that used in CPM machine. The analysis that will be carried out in this study includes:

1. Analysis of spur gear 1 ( $z = 20$ )
2. Analysis of spur gear 2 ( $z = 80$ )
3. Analysis of spur gear 3 ( $z = 100$ )

### **3.1. Results**

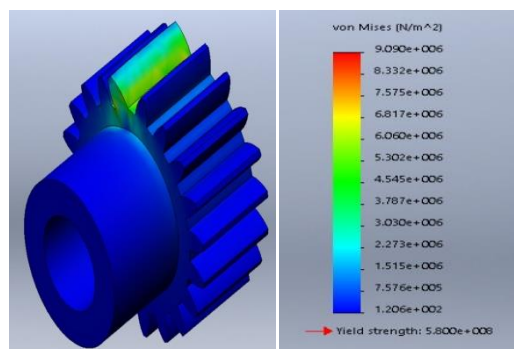
From the design of CPM machine is shown as shown in the figure 2.



**Figure 2.** Design of CPM machine

#### **1. Simulation analysis of load on gear1 $z = 20$**

In the analysis of the results of the design of gear1 with  $z = 20$  are shown as shown in Figure 3, Figure 4, and Figure 5.



**Figure 3.** Stress on gear1  $z = 20$  due to the patient's arm load

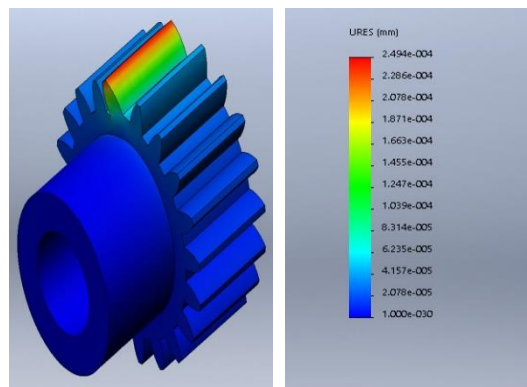


Figure 4. Displacement on gear1  $z = 20$  due to the patient's arm load

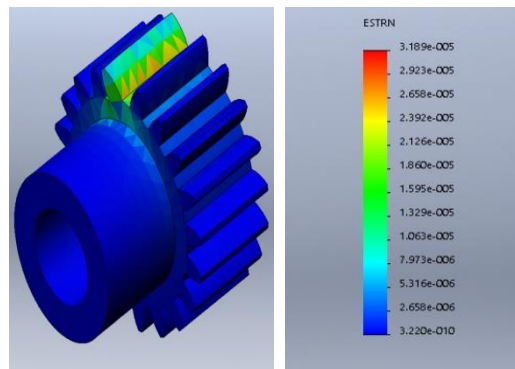


Figure 5. Strain on gear1  $z = 20$  due to the patient's arm load

Description spur gear1  $z = 20$  :

Material	: 1.0503 (C45)
Mass	: 0.0277809 kg
Volume	: $3.56165e - 006 m^3$
Density	: $7200 kg/m^3$
Weight	: 0.277809 N
Resultant Forces	: 51.7429 N
Total Nodes	: 22126
Total elements	: 13140

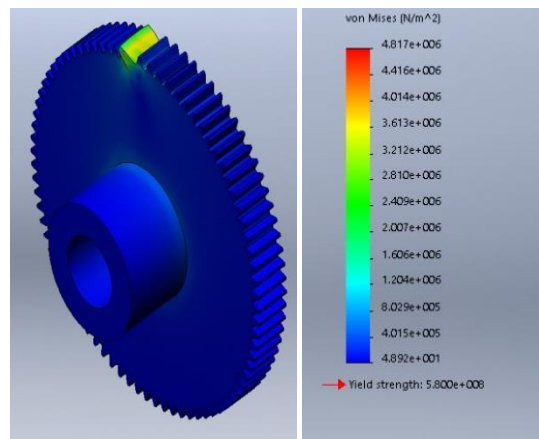
**Table 1.** Table of results of analysis due to of patient arm load on spur gear1  $z = 20$

	Type	Min	Max
Stress	VON: von Mises Stress	120.602 $N/m^2$ Node: 4339	9.08982e + 006 $N/m^2$ Node: 5019
Displacement	URES: Resultant Displacement	0 m Node: 838	0.000249418 mm Node: 227
Strain	ESTRN: Equivalent Strain	3,2199e – 010 Element: 3550	3.18916e – 005 Element: 6033

In Table 1. It shows the results of the analysis of the structure due to the patient's arm load on the gear1  $z = 20$ , the displacement/deflection on gear1 is 0.000249418 mm. And the strain that happened on gear1 is 3.18916e – 005. The gear1 is still relatively safe because the maximum Yield Strength is 9.08982e + 006  $N/m^2$ , and far below the allowable Yield Strength which is equal to 5.800e + 008  $N/m^2$ .

**2. Simulation analysis of load on gear2  $z = 80$**

In the analysis of the results of the design of gear2 with  $z = 80$  are shown in Figure 6, Figure 7, and Figure 8.



**Figure 6.** Stress on gear2  $z = 80$  due to the patients arm load

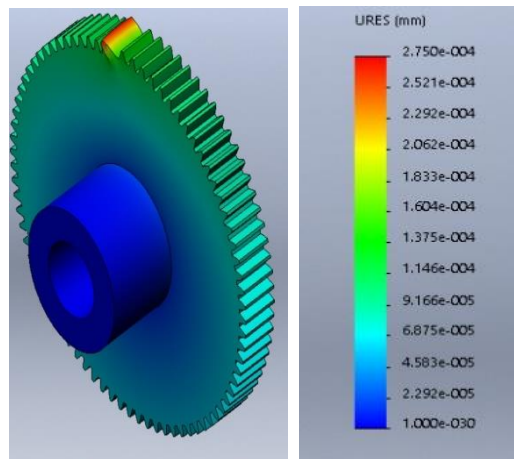


Figure 7. Displacement on gear<sub>2</sub> z = 80 due to the patient's arm load]

Description spur gear<sub>2</sub> z = 80:

Material	:	1.0503 (C45)
Mass	:	0.45213 kg
Volume	:	5.79653e – 005 m <sup>3</sup>
Density	:	7200 kg/m <sup>3</sup>
Weight	:	4.43087 N
Resultant Forces	:	49.4847 N
Total Nodes	:	20410
Total elements	:	12580

Table 2. Table of results of analysis due to of patient arm load on spur gear<sub>2</sub> z = 80

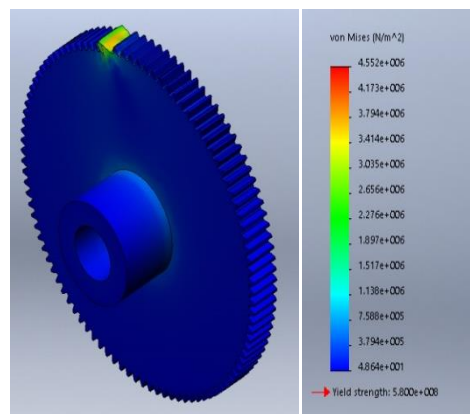
Type		Min	Max
Stress	VON: von Mises Stress	48.9205 N/m <sup>2</sup> Node: 15026	4.81729 + 006 N/m <sup>2</sup> Node: 19255
	URES: Resultant Displacement	0 m Node: 110	0.000274988 mm Node: 1763
Strain	ESTRN: Equivalent Strain	2.44294e – 010 Element: 5820	1.77017e – 005 Element: 1609

Table 2 shows the results of the analysis of the structure due to the patient's arm load on gear<sub>2</sub> z = 80, the displacement/deflection on gear<sub>2</sub> is 0.000274988 mm and the strain that happened on gear<sub>2</sub> is 1.77017e – 005. The gear<sub>2</sub> is still relatively safe because the

maximum Yield Strength is  $4.81729 + 006 N/m^2$  and far below the allowable Yield Strength which is equal to  $5.800e + 008 N/m^2$ .

### 3. Simulation analysis of load on gear<sub>3</sub> z = 100

In the analysis of the results of the design of gear<sub>3</sub> with z = 100 are shown in Figure 9, Figure 10, and Figure 11.



**Figure 9.** Stress on gear<sub>3</sub> z = 100 due to the patients arm load safe because the maximum Yield Strength is  $9.08982e + 006 N/m^2$ , and far below the allowable Yield Strength which is equal to  $5.800e + 008 N/m^2$ .

### 3.2. Discussions

Analysis of the strength of spur gear material due to patient's arm load, to analyze the gearbox transmission of the affected part of the transmission gear to transmit torque to the spur gear of each gears. This can be simulated because when the rotating gears of each gear will experience the same load and direction .

From the design and simulation the maximum stress on the gear<sub>1</sub> z = 20 is  $9.08982e + 006 N/m^2$ . On gear<sub>2</sub> z = 80 the maximum stress that occurs is equal to  $4.81729 + 006 N/m^2$ . The maximum stress that occurs in gear<sub>3</sub> z = 100 is equal to  $9.08982e + 006 N/m^2$ . The value of maximum stress indicates that teeth of the spur gear is affected by load of patients arm. The results of the simulations show that material 1.0503 (C45) still safe. This is indicated by the maximum stress occurs in each gear is still far below Yield Strength. Material yield strength is 1.0503 (C45) which is equal to  $5.8e + 008 N/m^2$ .



## **4 Conclusion**

For designing and simulating three dimensions, it is very helpful in evaluating the design of a device. Loading on the simetris section can be assumed on one part that is exposed to the load of the tool. The analysis in this study was carried out with static loading. The possibility of analysis through simulation can still be done using the dynamic loading method

## **References**

- [1] S. Miyaguchi, N. Matsunaga, K. Nojiri, and S. Kawaji, “Impedance control of cpm device with flex-/extension an pro-/supination of upper limbs,” *IEEE/ASME International Conference on Advanced Intelligent Mechatronics*, 2007.
- [2] A. C. Lad and A. S. Rao, “Design and Drawing Automation Using Solid Works Application Programming Interface,” *International Journal of Emerging Engineering Research and Technology*, **2** (7), 157–167, 2014.
- [3] M. H. H. Razali, M. A. H. A. Ssomad, S. M. Sapuan, M. N. A. Noordin, M. Hasbullah, and R. Syazili, “Simulation and analysis of innovative hand tool harvester,” *Scientific Research and Essays*, **7** (19), 1864–187, 2012.
- [4] SolidWorks, *Tutorial SolidWorks Simulation*, Dassault Systèmes SolidWorks Corporation, 2016.
- [5] SolidWorks, *An Introduction to Flow Analysis Applications with Solid Works Flow Simulation*, Student Guide, Dassault Systèmes SolidWorks Corporation, 2010.

This page intentionally left blank

## **Influences of Annealing on the Electrical Properties of $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$**

Dwi Nugraheni Rositawati

*Department of Physics Education, Faculty of Teacher Training and Education, Sanata Dharma University, Yogyakarta, Indonesia*

*\*Corresponding Author: wiwikfis@gmail.com*

(Received 05-05-2019; Revised 14-05-2019; Accepted 14-05-2019)

### **Abstract**

The research aims to investigate influences of annealing on the electrical properties of  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$ .  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  material which was annealed at  $900^\circ\text{C}$  for 1, 2 and 4 hours has better mechanical properties. It needs investigation for its electrical contribution, namely the correlation between grain and grain boundaries to values of resistance and capacitance. The changing of electrical properties was controlled by grain, grain boundary and the area between the sample and contact. The electrical properties of  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  were investigated by impedance spectroscopy in the room temperature. This method is able to separate the electrical and dielectric properties of the grain, grain boundary and the area between contact with the sample. ZsimpWin software was used to find out the equivalent electrical circuit, the resistance and capacitance value. It was observed that with the increase in annealing time the small grains resistance, the grain boundaries resistance, and the large grain capacitance value also increases. The resistance values of small grains and large grains were smaller than the grain boundaries resistance. The value of capacitance-resistance of the small grains and large grains were obtained values that tend to be smaller.

**Keywords:** Ba<sub>0,5</sub>Sr<sub>0,5</sub>TiO<sub>3</sub>, annealing, grain, grain boundary.

## **1 Introduction**

The rapid development and advancement of technology were influenced by the development of material as its basic material. The development of the material certainly is inseparable from the development of discoveries of properties superior of a material as its basic material [1]. Solid materials have been conveniently grouped into three basic categories: metals, ceramics, and polymers, a scheme based primarily on chemical makeup and atomic structure. Most materials fall into one distinct grouping or another. In addition, there are the composites that are engineered combinations of two or more different materials. A brief explanation of these material classifications and representative characteristics is offered next. Another category is advanced materials—those used in high-technology applications, such as semiconductors, biomaterials, smart materials, and nanoengineered materials [1].

Barium Titanate material (BaTiO<sub>3</sub>) was originally discovered in 1941. This material was ferroelectric [2]. The continues research were carried out in line with the discovery of interesting properties on Barium Titanate (BaTiO<sub>3</sub>) material, namely the discovery of various attractive properties including the material is very practical because of its very stable chemical and mechanical properties. It has ferroelectric properties [3]. The application of Barium Titanate material (BaTiO<sub>3</sub>) includes the fields of thermal, electricity, electromechanics, and electro-optics, namely as multilayer capacitors (MLCs), PTC thermistors, electro-optical equipment, dynamic random access memories (DRAM) and tunable capacitors for microwave technology [3-5]. Barium Strontium Titanate which has the chemical formula BaSrTiO<sub>3</sub> or better known as BST is one type of material in the ceramic group. BST is a ferroelectric material which belongs to the type of perovskite formed from Barium Titanate (BaTiO<sub>3</sub>) doped with Strontium (Sr). Addition of Strontium to Barium Titanate is able to change the nature of Barium Titanate because the nature of a material can be changed by heat treatment and by the addition of other substances [1].

This research was intended for the application of Ba<sub>0,5</sub>Sr<sub>0,5</sub>TiO<sub>3</sub> as a thermistor PTC. One of the characteristics of PTC is the resistance of material will rise significantly if

the temperature of material increased [6]. The influences temperature on the material change the size of the grain which will cause a shift in the Curie point - the transition point from ferroelectric to paraelectric in the material  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  and phase transition. It shows that changing in electrical properties and transport mechanisms at room temperature and low temperatures are controlled by grain and grain boundaries. The addition of Sr to  $\text{BaTiO}_3$  will reduce the Curie temperature to room temperature [5] therefore it is important to examine the electrical conductivity of  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  material at room temperature.

The stimulus of electrical properties such as electrical conductivity and dielectric constant is the electric field [1]. The method that was used in this study is Impedance Spectroscopy. Impedance spectroscopy is an analytical method that is popular in the research and development of material science. This method provides relatively simple electrical measurements and the results can be related to complex material variables: starting from mass transport, chemical reaction rate, corrosion, amorphous and polycrystalline dielectric behavior, microstructure and the influences of composition on the conductance of solids. The Impedance Spectroscopic Method could separate the electrical and dielectric properties of the grain, grain boundary and the area between contact with the sample. The measurement impedance parameters helps identify the physical process and determines the types of electrical parameters that represent the system [7]. It is important to have an equivalent model that can provide electrical properties. The electrical properties of the material are determined by a series combination between grain and grain boundaries, each of them is represented by a parallel RC element. It can be said that material electrical circuits are equivalent to a series of two parallel RC elements [8].  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  material which was annealed at  $900^\circ\text{C}$  for 1, 2 and 4 hours has better mechanical properties than  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  material which was sintered [9]. In order to obtain a complete understanding of  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  material properties, it is necessary to examine the electrical properties of  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$  which was annealed  $900^\circ\text{C}$  for 1, 2 and 4 hours. For this reason, it is necessary to examine the contribution of electricity, namely the correlation between grain and grain boundary to the value of resistance and capacitance. The research aims to investigate influences of annealing on the electrical properties of  $\text{Ba}_{0,5}\text{Sr}_{0,5}\text{TiO}_3$ .

## **2 Research Methodology**

The materials which were used in this study were  $\text{Ba}_{0.5}\text{Sr}_{0.5}\text{TiO}_3$  samples. It were obtained from the annealing temperature at  $900^\circ\text{C}$  for 1, 2 and 4 hours using Ney Vulcan furnaces 3-550. The sample  $\text{Ba}_{0.5}\text{Sr}_{0.5}\text{TiO}_3$  are shaped like a piece that have a diameter of 10 mm, 2 mm thick and have a mass of 0.5 gr. Samples of  $\text{Ba}_{0.5}\text{Sr}_{0.5}\text{TiO}_3$  which had been washed were given contact from fiber wire. Sample preparation is done by heating samples that have been given silver glue at temperature of  $120^\circ\text{C}$  for 1 hour using Memmert 1534 Furnace. The heating function is to quicken the glue drying process and to further glue the contact wire to the sample.

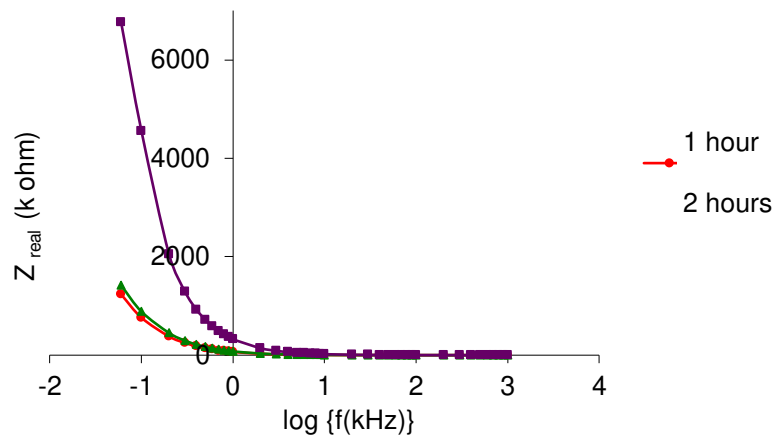
The prepared sample is ready to measured the value of the RLC with the RCL meter which has been calibrated. Since the measurements with RCL meter were very sensitive, it were done in a stable sample state. The measuring values are the impedance and phase angle of each sample at a frequency of 50 Hz-1 MHz. The measurement starts from a high frequency of 1 MHz to a low frequency of 50 Hz to maintain the stability of the reading of the impedance value.

Data of impedance and phase angles are used to determine the real part of impedance ( $Z_{\text{real}}$ ) and imaginary part of impedance ( $Z_{\text{imaginary}}$ ). The data are used as input data for processing data with ZsimpWin program to obtain the impedance spectrum in the Nyquist plot. The ZsimpWin software is used to obtain the equivalent electrical circuit according to the physical state of the sample. By providing input impedance data and selecting the desired electrical circuit, the software will automatically match the impedance curve. The equivalent circuit depends on the character of the measured sample. The value of each electrical component characterizes the electrical properties of the sample.

## **3 Results and Discussions**

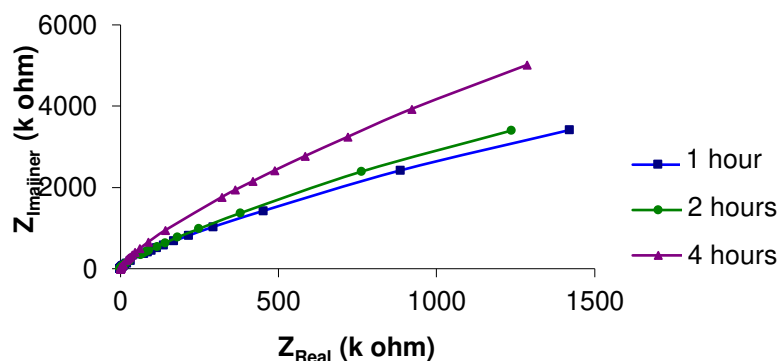
Variation of  $Z_{\text{real}}$  vs log frequency on samples annealed at  $900^\circ\text{C}$  for 1, 2 and 4 hours is presented Figure 1. The  $Z_{\text{real}}$  value on samples annealed at  $900^\circ\text{C}$  for 4 hours increase significantly compared to the samples annealed at  $900^\circ\text{C}$  for 1 hour and 2 hours. The

increasing in annealing time will increase the  $Z_{real}$  value. It shows decreasing in AC conductivity.



**Figure 1.** Graph of  $Z_{real}$  vs log frequency (annealed at 900°C for 1 hour, 2 hours and 4 hours)

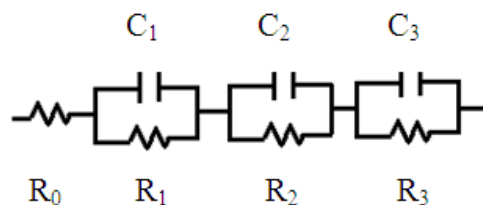
Figure 1 also shows that the  $Z_{real}$  value of the high-frequency log is much smaller than the low-frequency log. It indicates that  $Z_{real}$  at high frequencies, grains contribute to electrical conduction, while at low frequencies the role is at the grain boundary. The changing in impedance spectrum occurs more often at low frequencies, namely at the grain boundary because it is a less stable area.



**Figure 2.** Comparison of Nyquist plots (annealed at 900°C for 1 hour, 2 hours and 4 hours)

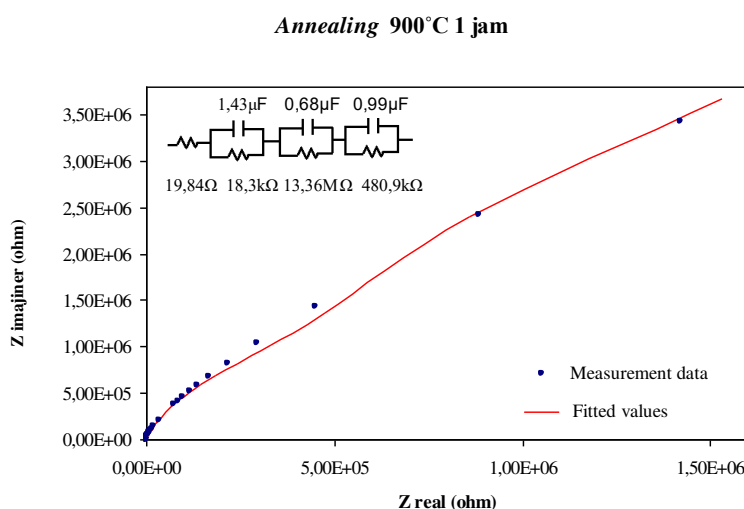
Figure 2 shows comparison of Nyquist plots (annealed at 900°C for 1 hour, 2 hours and 4 hours). It shows that the increasing time of annealing causes the Impedance Spectrum curve to become larger and higher. It indicates the changing in the resistance and capacitance values in the sample.

The value of electrical elements can be obtained by modeling the electrical circuit from the Nyquist curve with fittings using the ZsimpWin program. After fitting in several electrical circuits, the most suitable electrical circuit obtained which is the model circuit in the sample is as shown in Figure 3.



**Figure 3.** Equivalent electrical circuit;  $R_0$  = interface resistance;  $R_1$  &  $C_1$  = resistance & grain capacitance (small size);  $R_2$  &  $C_2$  = grain boundary resistance & capacitance;  $R_3$  &  $C_3$  = grain resistance & capacitance (large size)

Comparison of measurement data with the results of fittings on samples annealed at 900°C for 1 hour, 2 hours and 4 hours are shown in Figure 4, Figure 5 and Figure 6..



**Figure 4.** Comparison of measurement data & the results of fitting on samples annealed at 900°C for 1 hour



Annealing 900°C 2 jam

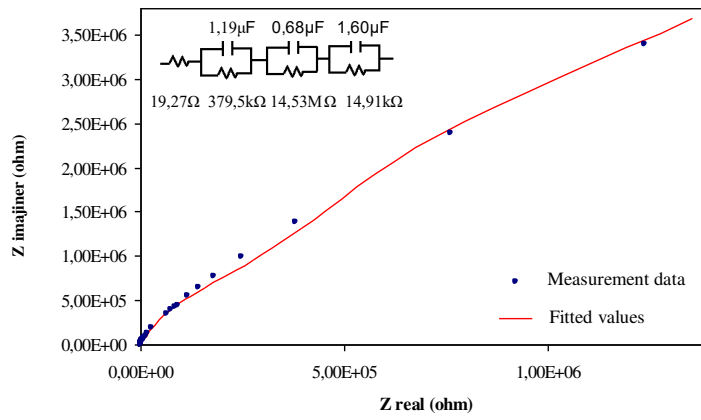


Figure 5. Comparison of measurement data & the results of fitting on samples annealed at 900°C for 2 hours

Annealing 900°C 4 jam

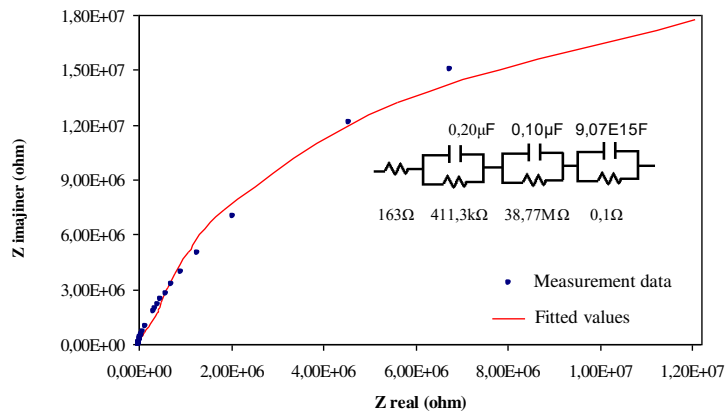


Figure 6. Comparison of measurement data & the results of fitting on samples annealed at 900°C for 4 hours

The resistance and capacitance values which were obtained from the results of fitting by ZsimpWin were summarized in Table 1.

**Table 1.** Value of resistance and capacitance

Circuit element	Annealed at 900°C		
	1 hour	2 hours	4 hours
R <sub>0</sub> (Ω)	19,84	19,27	163
C <sub>1</sub> (μF)	1,43	1,19	0,20
R <sub>1</sub> (kΩ)	18,3	379,5	411,3
C <sub>2</sub> (μF)	0,68	0,68	0,10
R <sub>2</sub> (MΩ)	13,36	14,53	38,77
C <sub>3</sub> (F)	0,99E – 6	1,60E – 6	9,07E15
R <sub>3</sub> (kΩ)	480,9	14,91	0,1E – 3

Table 1 shows that the resistance values of small grains and large grains were smaller than the grain boundaries resistance. Grains were more stable than grain boundaries so that the resistance value were smaller. It was observed that with the increase in annealing time the small grains resistance, the grain boundaries resistance, and the large grain capacitance value also increases. The annealing process can be used to reduce the thermal stress that is known from cracks which are found to be smaller according to the increase in annealing time. Annealing is also able to improve microstructure which is to produce smaller, homogeneous grains and not to find porosity [9]. The contact area between one item and another becomes more and more so that it produces resistance and the capacitance gets bigger. The value of capacitance-resistance of the small grains and large grains were obtained values that tend to be smaller, this is probably due to the distribution of grains that have not been homogeneous in the sample layer [9].

## 4 Conclusions

Influences of annealing on the electrical properties of Ba<sub>0,5</sub>Sr<sub>0,5</sub>TiO<sub>3</sub> can be written as follows:

- a. It was observed that with the increase in annealing time the small grains resistance, the grain boundaries resistance, and the large grain capacitance value also increases.

- b. The resistance values of small grains and large grains were smaller than the grain boundaries resistance.
- c. The value of capacitance-resistance of the small grains and large grains were obtained values that tend to be smaller.

## References

- [1] W. D. Callister, “Material science and engineering-third edition”, *John Willey and Sons*, New York, 1994.
- [2] W. Heywang and H. Thomann, “Tailoring of piezoelectric ceramics,” *Annual Review of Materials Science*, **14**, 27–47, 1984.
- [3] H. Lin and Wang, “Structure and dielectric properties of perovskite–barium titanate (BaTiO<sub>3</sub>),” San Jose State University, 2002.
- [4] H. Y. Tian, W. G. Luo, A. L. Ding, J. Choi, C. Lee, and K. S. No, “Influences of annealing temperature on the optical and structural properties of (Ba,Sr)TiO<sub>3</sub> thin films derived from sol-gel technique,” *Thin Solid Films*, 200–205, 2002.
- [5] T. Hungría, M. Algueró, A. B. Hungría, and A. Castro, “Dense, fine-grained Ba<sub>1-x</sub>Sr<sub>x</sub>TiO<sub>3</sub> ceramics prepared by the combination of mechanosynthesized nanopowders and spark plasma sintering,” *Chemistry of Materials*, 6205–6212, 2005.
- [6] W. Cao, H. H. Cudney, and R. Waser, “Smart materials and structures,” *Proceedings of the National Academy of Sciences of the United States of America*, 8330–8331, July 1999.
- [7] S. Sen, R. N. P. Choudhary, and P. Pramanik, “Impedance spectroscopy of Ba<sub>1-x</sub>Sr<sub>x</sub>Sn<sub>0.15</sub>Ti<sub>0.85</sub>O<sub>3</sub> ceramics,” *British Ceramic Transactions*, 250–256, 2004.
- [8] K. Prabakar, S. K. Narayandass, and D. Mangalaraj, “Impedance and electric modulus analysis of Cd<sub>0.6</sub>Zn<sub>0.4</sub>Te thin films,” *Crystal Research and Technology*, **37** (10), 1094–1103, 2002.
- [9] D. N. Rositawati, “Studi pengaruh annealing terhadap keramik barium strontium titanate,” *Prosiding Seminar Nasional Sains dan Pendidikan Sains VI*, 210–215, Juli 2011.

This page intentionally left blank

# **Implementation of k-Medoids Clustering Algorithm to Cluster Crime Patterns in Yogyakarta**

Eduardus Hardika Sandy Atmaja

*Department of Informatics, Faculty of Science and Technology,*

*Sanata Dharma University, Yogyakarta, Indonesia*

*Corresponding Author: edo@usd.ac.id*

(Received 02-05-2019; Revised 21-05-2019; Accepted 21-05-2019)

## **Abstract**

The increase in crime from day to day needs to be a concern for the police, as the party responsible for security in the community. Crime prevention effort must be done seriously with all knowledge that they have. To increase police performance of crime prevention effort, it is necessary to analyze crime data so that relevant information can be obtained. This study tried to analyze crime data to obtain relevant information using clustering in data mining. Clustering is a data mining method that can be used to extract valuable information by grouping data into groups that have similar characters. The data used in this study were crime patterns which were then grouped using K-medoids clustering algorithm. The obtained results in this study were three crime groups, namely high crime level with 4 members, medium crime level with 6 members and low crime level with 8 members. It is expected that this information can be used as material for consideration in crime prevention effort.

**Keywords:** clustering, k-medoids, criminality

## 1 Introduction

Crime is any act that is prohibited by public law to protect the public and given punishment by the state. These acts are punished because they are violating the social norms such as act that conflict with legal norms, social norms and religious norms that applied in the society [1]. The existence of punishment applied by law enforcement does not make the criminals undermine their intentions, and in fact criminal in Yogyakarta are increasing widespread.

The increase of criminal cases in the society can result in losses both materially and immaterially. For this reason, efforts are needed from law enforcement to reduce crime in the society. Such efforts can be done by finding relevant information related to crime. Such information can be obtained by processing and analyzing crime data owned by the police.

The crime data owned by the Yogyakarta Police is still stored in the manual form such as register books and excel. The data is only stored and is not used to produce any information. Where the data can be processed and analyzed to produce valuable information in efforts to prevent crime. Data mining is a proper technique to extract important information from a data set.

Crime data owned by the police can be processed using data mining to become crime patterns that represent relationship between crimes. The research was successfully done by Atmaja [2], the result was crime patterns presented in graph form. The weakness in that study is that there is no clear grouping on crime level from generated crime patterns. This study tried to refine previous research by grouping crime patterns into three categories, namely high crime level, medium crime level and low crime level.

Clustering is one of the data mining techniques that aims to group data based on information found in the data [3]. The grouping is based on the similarity between data so the data in the same cluster is homogeneous. Thus clustering is a very appropriate method for classifying crime patterns into high, medium and low crime level.

Researches on implementation of clustering method have been done, as done by Singh et. al. [4]. They tried to implement K-means clustering algorithm by using three different distance measurements namely Euclidean, Manhattan and Chebychev. The

result is that the implementation of K-means algorithm using Euclidean distance measurements can produce the best group from the other distance measurements. So it can be concluded that the best pair for K-means algorithm is the Euclidean distance measurement.

Research on the use of Euclidean distance in K-means algorithm has been successfully done by Atmaja [5]. The aim of his study was to cluster crime data into three categories, namely high, medium and low crime level. Although the objective of the research was achieved, K-means algorithm is classified as an ineffective algorithm because it involves too much noise and outliers caused by the average selection of clusters [6].

This study tried to improve previous study by replacing K-means algorithm with K-medoids algorithm. K-medoids algorithm is one of the clustering algorithms that are not influenced by outliers or other extreme variables [6]. K-medoids work by determining the center point of existing data without performing an average calculation as in K-means. The following is the K-medoids algorithm [6]:

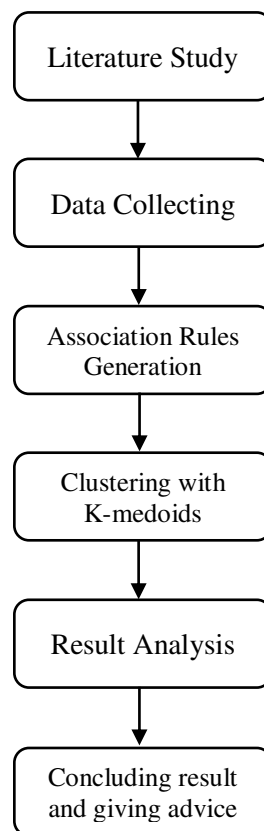
- (1) arbitrarily choose  $k$  objects in  $D$  as the initial representative objects or seeds;
- (2) **repeat**
- (3) assign each remaining object to the cluster with the nearest representative object;
- (4) randomly select a nonrepresentative object,  $o_{\text{random}}$ ;
- (5) compute the total cost,  $S$ , of swapping representative object,  $o_j$ , with  $o_{\text{random}}$ ;
- (6) **if**  $S < 0$  **then** swap  $o_j$  with  $o_{\text{random}}$  to form the new set of  $k$  representative objects;
- (7) **until** no change;

**Figure 1.** K-medoids algorithm

The result of this study is crime patterns that have been divided into three groups, namely high, medium and low crime level. It is expected that the police can use this information to improve crime prevention efforts in the society.

## 2 Research Methodology

Research methodology done by this research is activity steps to implement K-means algorithm to cluster crime patterns from Yogyakarta Police data which are presented in Figure 2.



**Figure 2.** Research Methodology

Figure 2 shows research methodology which began with literature study to study relevant theories related to solve problems. The next step was data collecting related to research, in this case the processed data was crime data from Yogyakarta Police. The crime data that has been collected then processed using association techniques in data mining to produce association rules that described crime patterns. Generated rules was used as input to K-medoids algorithm to produce crime patterns accompanied by grouping based on low, medium and high crime level. The next step was result analyzing that has been obtained to find out whether the objective achieved or not. Finally, the result analysis will draw conclusions from the research that has been done. Suggestions were also given to correct existing disadvantages to be applied in the future research.



### 3 Results and Discussions

#### 3.1 Crime Patterns

There are 18 samples of crime patterns as results of association technique processing accompanied by support and confidence. The data will be grouped using the K-medoids algorithm based on variable support and confidence. These data are presented in Table 1.

**Table 1.** Crime patterns

No.	Rules	Support	Confidence
1	IF Embezzlement THEN Theft	0.02	0.03
2	IF Theft THEN Embezzlement	0.02	0.29
3	IF Embezzlement THEN Deception	0.54	0.81
4	IF Deception THEN Embezzlement	0.54	0.82
5	IF Embezzlement THEN Document Forgery	0.02	0.03
...	...	...	...
18	IF Unpleasant Act THEN Defamation	0.02	0.38

#### 3.2 Determining Initial Medoids

In the first stage, three medoids were randomly selected from data sample in Table as shown in Table 2.

**Table 2.** Three initial medoids

	Medoid		
	C1	C2	C3
Support	0.54	0.08	0.03
Confidence	0.81	0.12	0.30

#### 3.3 Calculating Euclidean Distance Iteration 1

The next step is euclidean distance calculation from each data to the three selected medoids. Euclidean distance is calculated based on the following formula [6]:

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2}$$

Here,  $d(i, j)$  represents distance between data and medoid,  $x_{i1}$  denotes support value in each data,  $x_{j1}$  is medoid (c) for support,  $x_{i2}$  denotes confidence value in each data and  $x_{j2}$  is medoid (c) for confidence. Table 3 presents results from euclidean distance

calculation on each data along with medoid information which has the shortest distance to the data.

**Table 3.** Rules with euclidean distance

Rules	Support	Confidence	Distance to Medoid			Shortest Distance
			C1	C2	C3	
1	0.02	0.03	0.937	0.108	0.270	0.108
2	0.02	0.29	0.735	0.180	0.014	0.014
3	0.54	0.81	0.000	0.829	0.721	0.000
4	0.54	0.82	0.010	0.838	0.728	0.010
5	0.02	0.03	0.937	0.108	0.270	0.108
6	0.02	0.41	0.656	0.296	0.110	0.110
7	0.09	0.13	0.815	0.014	0.180	0.014
8	0.09	0.97	0.478	0.850	0.673	0.478
9	0.02	0.03	0.937	0.108	0.270	0.108
10	0.02	0.41	0.656	0.296	0.110	0.110
11	0.08	0.12	0.829	0.000	0.187	0.000
12	0.08	0.94	0.478	0.820	0.642	0.478
13	0.03	0.30	0.721	0.187	0.000	0.000
14	0.03	0.86	0.512	0.742	0.560	0.512
15	0.02	0.23	0.779	0.125	0.071	0.071
16	0.02	0.69	0.534	0.573	0.390	0.390
17	0.02	0.44	0.638	0.326	0.140	0.140
18	0.02	0.38	0.675	0.267	0.081	0.081

From Table 3, it can be seen that medoid C1 has 5 members rules {3,4,8,12,14}, medoid C2 has 5 members rules {1,5,7,9,11} and medoid C3 has 8 members rules {2,6,10,13,15,16,17,18}.

### 3.4 Calculating Total Cost Iteration 1

Calculating total cost is the final step from iteration 1, by summing the shortest distance from data in Table 3, so the total cost is 2.734.

### 3.5 Determining Random Medoids Iteration 2

The process continues to iteration 2 by selecting a new random medoid from the data to replace the medoid C3 temporarily. The selection of a new medoid should not be the

same as one of the medoids that has been selected. Table 4 shows three medoids for iteration 2.

**Table 4.** Three medoids iteration 2

	Medoid		
	C1	C2	C Random
Support	0.54	0.08	0.03
Confidence	0.81	0.12	0.86

### 3.6 Calculating Euclidean Distance Iteration 2

After a new medoid has been determined, the next step is to recalculate the euclidean distance for each data based on three medoids from Table 4. The results is shown in Table 5.

**Table 5.** Rules with euclidean distance iteration 2

Rules	Support	Confidence	Distance to Medoid			Shortest Distance
			C1	C2	C3	
1	0.02	0.03	0.937	0.108	0.830	0.108
2	0.02	0.29	0.735	0.180	0.570	0.180
3	0.54	0.81	0.000	0.829	0.512	0.000
4	0.54	0.82	0.010	0.838	0.512	0.010
5	0.02	0.03	0.937	0.108	0.830	0.108
6	0.02	0.41	0.656	0.296	0.450	0.296
7	0.09	0.13	0.815	0.014	0.732	0.014
8	0.09	0.97	0.478	0.850	0.125	0.125
9	0.02	0.03	0.937	0.108	0.830	0.108
10	0.02	0.41	0.656	0.296	0.450	0.296
11	0.08	0.12	0.829	0.000	0.742	0.000
12	0.08	0.94	0.478	0.820	0.094	0.094
13	0.03	0.30	0.721	0.187	0.560	0.187
14	0.03	0.86	0.512	0.742	0.000	0.000
15	0.02	0.23	0.779	0.125	0.630	0.125
16	0.02	0.69	0.534	0.573	0.170	0.170
17	0.02	0.44	0.638	0.326	0.420	0.326
18	0.02	0.38	0.675	0.267	0.480	0.267

From Table 5, it can be seen that medoid C1 has 2 members rules {3,4}, medoid C2 has 12 members rules {1,2,5,6,7,9,10,11,13,15,17,18} and medoid C3 has 4 members rules {8,12,14,16}.

### ***3.7 Calculating Total Cost Iteration 2***

Calculating total cost is the final step from iteration 2, by summing the shortest distance from data in Table 5, so the total cost is 2.416. To determine the next iteration, total cost from iteration 2 is compared with iteration 1, which is  $2,416 > 2,734$ . Because the total cost of iteration 2 is not greater than iteration 1, the iteration is continued to iteration 3 and the medoid C Random replaces medoid C3.

### ***3.8 Determining Random Medoids Iteration 3***

The process continues to iteration 3 by selecting a new random medoid from the data to replace the medoid C3 temporarily (C Random from iteration 2). The selection of a new medoid should not be the same as one of the medoids that has been selected. Table 6 shows three medoids for iteration 3.

**Table 6.** Three medoid iteration 3

	<b>Medoid</b>		
	<b>C1</b>	<b>C2</b>	<b>C Random</b>
Support	0.54	0.08	0.02
Confidence	0.81	0.12	0.44

### ***3.9 Calculating Euclidean Distance Iteration 3***

After a new medoid has been determined, the next step is to recalculate the euclidean distance for each data based on three medoids from Table 6. The results is shown in Table 7.

**Table 7.** Rules with euclidean distance iteration 3

Rules	Support	Confidence	Distance to Medoid			Shortest Distance
			C1	C2	C3	
1	0.02	0.03	0.937	0.108	0.410	0.108
2	0.02	0.29	0.735	0.180	0.150	0.150
3	0.54	0.81	0.000	0.829	0.638	0.000
4	0.54	0.82	0.010	0.838	0.644	0.010
5	0.02	0.03	0.937	0.108	0.410	0.108
6	0.02	0.41	0.656	0.296	0.030	0.030
7	0.09	0.13	0.815	0.014	0.318	0.014
8	0.09	0.97	0.478	0.850	0.535	0.478
9	0.02	0.03	0.937	0.108	0.410	0.108
10	0.02	0.41	0.656	0.296	0.030	0.030
11	0.08	0.12	0.829	0.000	0.326	0.000
12	0.08	0.94	0.478	0.820	0.504	0.478
13	0.03	0.30	0.721	0.187	0.140	0.140
14	0.03	0.86	0.512	0.742	0.420	0.420
15	0.02	0.23	0.779	0.125	0.210	0.125
16	0.02	0.69	0.534	0.573	0.250	0.250
17	0.02	0.44	0.638	0.326	0.000	0.000
18	0.02	0.38	0.675	0.267	0.060	0.060

From Table 7, it can be seen that medoid C1 has 4 members rules {3,4,8,12}, medoid C2 has 6 members rules {1,5,7,9,11,15} and medoid C3 has 8 members rules {2,6,10,13,14,16}.

**3.10 Calculating Total Cost Iteration 3**

Calculating total cost is the final step from iteration 3, by summing the shortest distance from data in Table 7. So the total cost is 2.510. To determine the next iteration, total cost from iteration 3 is compared with iteration 2, which is  $2.510 > 2.416$ . Because the total cost of iteration 3 is greater than iteration 2, the iteration stops.

**3.11 Results**

Each medoid represents 1 group of crime level based on support and confidence. C1 represents high crime level, C2 represents medium crime level and C3 represents low crime level. The results of crime patterns grouping are shown in Tables 8, 9 and 10.

**Table 8.** High level crime patterns

<b>No.</b>	<b>Rules</b>	<b>Support</b>	<b>Confidence</b>
1	IF Embezzlement THEN Deception	0.54	0.81
2	IF Deception THEN Embezzlement	0.54	0.82
3	IF Fiduciary THEN Embezzlement	0.09	0.97
4	IF Information violation and electronic transaction THEN Deception	0.08	0.94

**Table 9.** Medium level crime patterns

<b>No.</b>	<b>Rules</b>	<b>Support</b>	<b>Confidence</b>
1	IF Embezzlement THEN Theft	0.02	0.03
2	IF Embezzlement THEN Document Forgery	0.02	0.03
3	IF Embezzlement THEN Fiduciary	0.09	0.13
4	IF Deception THEN Document Forgery	0.02	0.03
5	IF Deception THEN Information violation and electronic transaction	0.08	0.12
6	IF Persecution THEN Beating	0.02	0.23

**Table 10.** Low level crime patterns

<b>No.</b>	<b>Rules</b>	<b>Support</b>	<b>Confidence</b>
1	IF Theft THEN Embezzlement	0.02	0.29
2	IF Document Forgery THEN Embezzlement	0.02	0.41
3	IF Document Forgery THEN Deception	0.02	0.41
4	IF Persecution THEN Domestic Violence	0.03	0.3
5	IF Domestic Violence THEN Persecution	0.03	0.86
6	IF Beating THEN Persecution	0.02	0.69

Tables 8, 9 and 10, show that some crime patterns are classified as high and some others are classified as low. Information about high level crime can be used by the police to prevent potential crime in the society.

## 4 Conclusions

It can be concluded that K-medoids algorithm can be used to cluster crime patterns into three crime levels namely, 4 rules classified as high level crime, 6 rules classified

as medium level crime and 8 rules classified as low level crime. Suggestions that can be given based on the results of this study are:

- a) There is a need to compare some distance method for K-medoid algorithm. Thus, it can be known the most appropriate distance calculation method for K-medoid algorithm.
- b) There is a need to apply weighting mechanism for each variable because not all variables have the same interests and priorities.

## Acknowledgements

The author thanks Polisi Resor Kota Yogyakarta (POLRESTA Yogyakarta) who provided crime data in this research by hiding some sensitive variables regarding victims and criminals.

## References

- [1] J. E. Sahetapy and B. M. Reksodiputro, “Paradoks dalam Kriminologi,” Rajawali : Jakarta, 1982.
- [2] E. H. S. Atmaja, “Visualisasi aturan asosiasi berbasis graph untuk data tindak kejahatan,” *Media Teknika*, **12** (1), 46–57, 2017.
- [3] P. Tan, M. Steinbach, and V. Kumar, “Introduction to data mining,” *Addison-Wesley*, Boston, 2006.
- [4] A. Singh, A. Yadav, and A. Rana, “K-means with three different distance metrics,” *International Journal of Computer Applications*, **67** (10), 13–17, 2013.
- [5] E. H. S. Atmaja, “Pengelompokan tingkat kriminalitas di kota yogyakarta dengan menggunakan metode k-means clustering,” *Seminar Nasional Riset dan Teknologi Terapan*, Agustus 2018.
- [6] J. Han, “Data mining: concepts and techniques second edition,” *Morgan Kaufmann*, San Francisco, 2006.

This page intentionally left blank



# **Development Study of Deep Learning Facial Age Estimation**

Puspaningtyas Sanjoyo Adi

*Department of Informatics, Faculty of Science and Technology,*

*Sanata Dharma University, Yogyakarta, Indonesia*

*Corresponding Author: pusp@usd.ac.id*

(Received 28-05-2019; Revised 29-05-2019; Accepted 31-05-2019)

## **Abstract**

Human age estimation is one of the most challenging problem because it can be used in many applications relating to age such as age-specific movies, age-specific computer applications or website, etc. This paper will contribute to give brief information about development of age estimation researches using deep learning. We explore three recent journal papers that give significant contribution in age estimation using deep learning. From these papers, they selected classification methods and there is gradual improvement in result and also in selected loss function. The best result gives MAE (mean average error) 2.8 years and VGG-16 is the most selected CNN architecture.

**Keywords:** age estimation, facial analysis

## **1 Introduction**

Human age can be estimated by facial appearance. Our faces show a special pattern in every lifetime so that our faces will have a huge difference every lifetime such as in childhood and adulthood. For the same person, the photo taken at different years indicate the aging process on their faces. The longer the interval is, the more obvious

changes there will be. Facial age estimation has potential application such as age-specific movies, age-specific products vending machine like tobacco, alcohol, and other age-specific computer applications or websites.

Estimating age from images is one of the most challenging work in facial analysis. It is hard to accurately predict human age because human facial aging is a slow and complicated process effected by many factors. With rapid advances in computer vision and pattern recognition, this problem becomes an interesting topic.

A typical pipeline of the existing methods for age estimation usually consists of two modules: age image representation and age estimation techniques [1]. Recently, deep learning schemes, especially Convolutional Neural Networks (CNNs), have been successfully employed for many tasks related to facial analysis. This paper aims to provide a brief description about some papers that have done age estimation research using CNN or deep learning. We will limit discussion to only a few paper published in journals or conferences in the last 5 years and became an important milestone of age estimating work. This paper is organized as follows: in section 2, age estimation algorithm will be explained and in section 3, we will explain about CNN architecture.

## 2 Age Estimation Algorithm

There has been a significant volume of research done for age estimates. This paper will focus on some papers that contributed significant development. We will explain these researches together with the estimation algorithm used. For age estimation, there are three methods that have been worked on, namely, classification, regression and ranking. In classification method, human age is assumed to be classified according to age-groups. The weakness of classification method is the sharing of important information between adjacent age groups. This is addressed by regression methods which appear to perform better. A different approach to deal with this challenge is to adopt ranking methods.

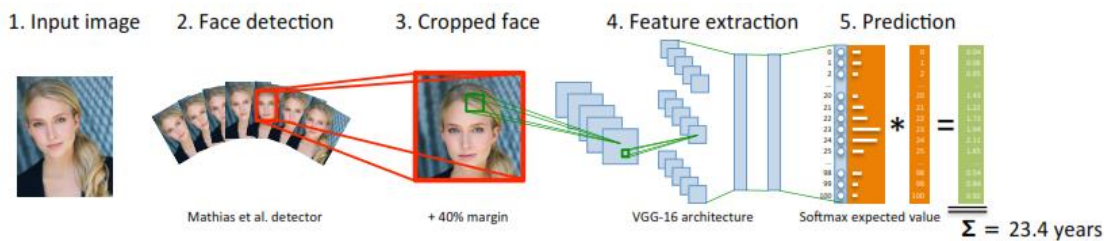


Figure 1. Pipeline of DEX method [2]

We choose Rothe’s work [2] as first paper examined and the winner of the LAP 2015 challenge [3] on apparent age estimation. Age estimation done by Rothe is a classification method. They use VGG-16 [4] as base CNN architecture called DEX (Deep Expectation). Fig. 1 shows pipeline of DEX method. System will get face image and then, it will be classified using CNN into 101 classes. These classes describe possible age groups from face image samples. They train CNN for classification and at test time, they compute expecting value over the softmax-normalized output probabilities of  $|Y|$  neurons.

$$E(O) = \sum_{i=1}^{|Y|} y_i o_i, \tag{1}$$

where  $O = \{1, 2, \dots, |Y|\}$  is the  $|Y|$ -dimensional output layer and  $O_i \in O$  is the softmax-normalized output probability of neuron  $i$ . Their research result a MAE (mean average error) 3.09 years with using IMDB-WIKI [2] as training dataset and FG-NET as testing dataset [5].

The same research was also conducted by Antipov [6]. They also use VGG-16 as base CNN architecture. They did the research with 3 kind age encoding, Fig. 2,: (1) pure year classification, called 0/1 Classification Age Encoding (0/1 CAE), (2) pure regression, called Real Value Age Encoding (RVAE), (3) soft classification, called Label Distribution Age Encoding (LDAE). Each encoding has its loss function but LDAE gives the best result.

$$L_{LDAE} = -\frac{1}{N} \sum_{k=1}^N \sum_{i=1}^{100} (t_i^k \log p_i^k + (1 - t_i^k) \log(1 - p_i^k)) \tag{2}$$

where  $N$  denotes number of images in batch,  $k$  denotes number of age class,  $t$  denotes targets,  $p$  denotes to prediction. Loss function refer to Gaussian distribution.

Lost function become differentiator between Rothe [2] and Antipov [6], but they are still in classification method. Antipov research result MAE 2.84 years using FG-NET as testing dataset.

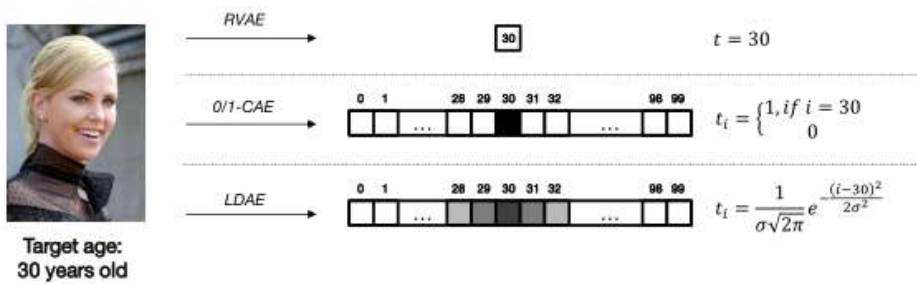


Figure 2. Example of encoding [6].  $t$  denotes encoding result and  $\sigma$  is a hyper parameter of LDAE.

Two papers before showed evolution of classification methods especially in loss function development. Hu et al [7] made improvement with adding age difference estimator. This component is built to overcome limitation of ground-truth age label dataset. Making ground-truth age label dataset is a costly effort so Hu et al propose a new dataset consisted of pair face images of same person with different taken time. This new dataset is used to make age difference loss function.

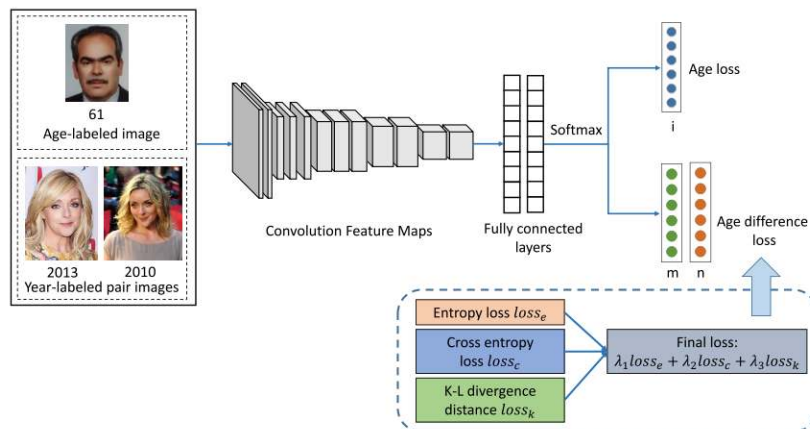


Figure 3. An overview of proposed method by Hu [7]

Fig. 3 shows an overview of proposed method by Hu. The system will give 2 outputs: age estimation and age difference. After training phase, CNN architecture will have values that will be tested with testing dataset. Initial probability distribution of age classes is set to Gaussian distribution. The age difference information with three kinds of loss functions, i.e. entropy loss, cross entropy loss and Kullback-Leibler (K-L) divergence distance. These loss functions can not only force the probability distribution of age classes to have one single peak value but also make the probability distribution locate within the correct range. This research result MAE 2.8 years using FG-NET as testing dataset.

### 3 Discussions

From three recently age estimation researches [2], [6], [7], we know that CNN architecture give good result in age estimation. Estimation method using classification approach gives good result. The challenge in CNN architecture is to find the best loss function which Gaussian distribution is the most choice used by researchers.

Based on three papers above, CNN architecture giving best prediction is VGG-16. This architecture is basically designed for face recognition but based on these papers, this architecture give good result for age estimation.

The other challenge is to find effective and efficient training dataset. From 3 papers, 2 papers [2], [7] contribute new dataset that is used by another similar paper for its training. IMDB-WIKI dataset [2] is not only used by [2] but it is also used by [6] for training phase. The big challenge is to make age label automatically. Further implication is hard to label age for its facial image.

### 4 Conclusions

In this paper, we have reviewed a few milestone papers in age estimation using deep learning. All papers result significant improvement in age estimation. Significant development component for these problem solving are loss function and dataset. From the recent researches, the task still open for improvement especially using deep learning.

In the future, this problem still gives challenging because aging process is a complex process influenced by many internal and external factor such as gene, environment, etc.

## **Acknowledgements**

Authors wishing to acknowledge assistance or encouragement from colleagues, special work by technical staff or financial support from organizations should do so in an unnumbered. Acknowledgments section immediately following the last numbered section of the paper.

## **References**

- [1] Y. Fu, G. Guo, and T. S. Huang, “Age synthesis and estimation via faces: a survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **32** (11), 1955–1976, 2010.
- [2] R. Rothe, R. Timofte, and L. V. Gool, “Deep expectation of real and apparent age from a single image without facial landmarks,” *International Journal of Computer Vision*, **126** (2-4), 144–157, 2018.
- [3] S. Escalera, J. González, X. Baró, P. Pardo, J. Fabian, M. Oliu, H. J. Escalante, I. Huerta, and I. Guyon, “ChaLearn looking at people 2015: apparent age and cultural event recognition datasets and results,” *Proceedings of the IEEE International Conference on Computer Vision*, 243–251, 2015.
- [4] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv*, 1–10, 2014.
- [5] G. Panis, A. Lanitis, N. Tsapatsoulis, and T. F. Cootes, “Overview of research on facial ageing using the FG-NET ageing database,” *IET Biometrics*, **5** (2), 37–46, 2015.
- [6] G. Antipov, M. Baccouche, S. A. Berrani, and J. L. Dugelay, “Effective training of convolutional neural networks for face-based gender and age prediction,” *Pattern Recognition*, **72**, 15–26, 2017.
- [7] Z. Hu, Y. Wen, J. Wang, M. Wang, R. Hong, and S. Yan, “Facial age estimation with age difference,” *IEEE Transactions on Image Process*, **26** (7), 3087–3097, 2017.

# **The Improvement of Watershed Algorithm Accuracy for Image Segmentation Handwritten Numbered Musical Notation**

Kartono Pinaryanto

*Department of Informatics, Faculty of Science and Technology,*

*Sanata Dharma University, Yogyakarta, Indonesia*

*Corresponding Author: kartono@usd.ac.id*

(Received 13-05-2019; Revised 21-05-2019; Accepted 21-05-2019)

## **Abstract**

In the implementation of image processing to translate the image of the numbered musical notation into a numerical character requires some initial process that must be passed like image segmentation process. The advantage of successful segmentation process is that it can reduce the failure rate in the object recognition process. Segmentation process determines the success of object recognition process, it takes segmentation algorithm that can perform accurate object separation. The combination segmentation process developed in this research used projection profile algorithm, watershed and non-object filtering. Profile projection algorithm is used to crop the image of the musical horizontally and vertically. The watershed algorithm is used to segment the numerical object of numerical notation generated from the projection profile process. Non-object filtering is a continuation of the watershed algorithm that includes the non-object reduction process and the process of combining objects so that the original object segment will be generated. Based on the results of the research, the accuracy of the segment on watershed segmentation is 99.74% higher than watershed segmentation without combination of 94.82%.

**Keywords:** watershed algorithm, profile projection algorithm, non-object filtering, image number notation musical

## **1 Introduction**

Indonesian regional songs are a valuable cultural heritage of Indonesia and have important characteristics in each region. Each regional song provides valuable advice or knowledge to the younger generation. Regional songs in addition to consisting of sounds also have a tone or music in the form of score numbers notation [1] so that we as young people are expected to be able to understand and preserve regional songs. One way to understand folk songs is to learn folk songs, especially on scores of number notations by performing image processing.

In the application of image processing to translate the image number notation into numerical characters requires several initial processes (pre-processing) that must be passed like the image segmentation process. The process of image segmentation is a process to separate one object from another object. One advantage of the success of the segmentation process is that it can reduce the level of failure in the object recognition process. In research [2] a method of combining intensity filtering (high pass filtering and low pass filtering) was developed as image pre-processing for noise reduction and watershed transformation to produce better quality segmentation. The result of combining this method is able to reduce excessive over segmentation. Based on research [3] river segmentation is an important process in the river tracking system using unmanned aerial vehicles. This process greatly determines the success of the river tracking system, this is because the output of image segmentation is an input that determines the outcome of the next process.

Based on the research [4] profile projection segmentation algorithm on the Javanese literary text document image Hamong Tani is a relatively good algorithm for document image segmentation with an accuracy of 84.255% and standard deviation of 14.093%. Based on the research [5] it was tested on Batak characters and the results of segmentation yielded a percentage value of truth segmentation ranging from 65% to 87.68% with a confidence level of 95%. Based on the research [6] which examined the image segmentation of braille documents using the projection profile method. The



results of the study were able to distinguish front side dots (Recto) and back side dots (Verso). Based on research [7] which examines the text image segmentation on soil maps using radon transformation based projection profile can be used to cluster text blocks. Based on research [8] which examines intelligent systems with the introduction of numerical music notation (NMRPIS) automatically using Optical Music Recognition (OMR). Based on the results of the experiment the level of introduction of NMRPIS reached 99% and the performance rate was 99.06%.

The watershed segmentation algorithm is one of the segmentation algorithms based on topographic forms [9]. Based on the research [10] the combination of medical image segmentation algorithms, namely reconstruct gradient, float-point active-image, watershed algorithm and the Grab Cut algorithm that functions to cut the image smoothly. In quality the combination of algorithms produces a good image, so that in the process of analysis and analysis will give better results. Based on research [11] watershed segmentation algorithms on Javanese literary text document images produce many objects, so the results of watershed segmentation are not good with an accuracy of 57.433% which can increase failure in the object recognition process. The results of watershed segmentation contain many objects and not objects so it is necessary to filter out objects and non-objects by throwing away not objects so that the expected level of accuracy is more accurate. Given the importance of the results of the segmentation process as an initial process of object recognition, a segmentation algorithm is needed that can accurately separate objects.

To correct this problem, it is necessary to create a score segmentation algorithm, number notation on a regional song which is a combination of a watershed algorithm, profile projection and filtering not an object. This research is expected to help attract the interest of the younger generation in preserving Indonesian regional culture and applying combinations can help improve the performance of the proposed algorithm in terms of the accuracy of segmentation results.

## 2 Research Methodology

The research method discusses the design of research methods, implementations and datasets. In the design of the research method, in broad outline, it discusses the flow of

the process of watershed segmentation combinations and the process flow segmentation watershed without a combination. In the implementation phase, it discusses the results of implementation of the watershed segmentation program. For datasets explain the types and examples of images used.

### 2.1 Design of Research Methods

In the research method phase, it discusses the design of the combination segmentation process flow between watershed algorithms, profile projection and filtering rather than objects with watershed algorithmic segmentation without combination. The design of the segmentation process flow is shown in Figure 1.

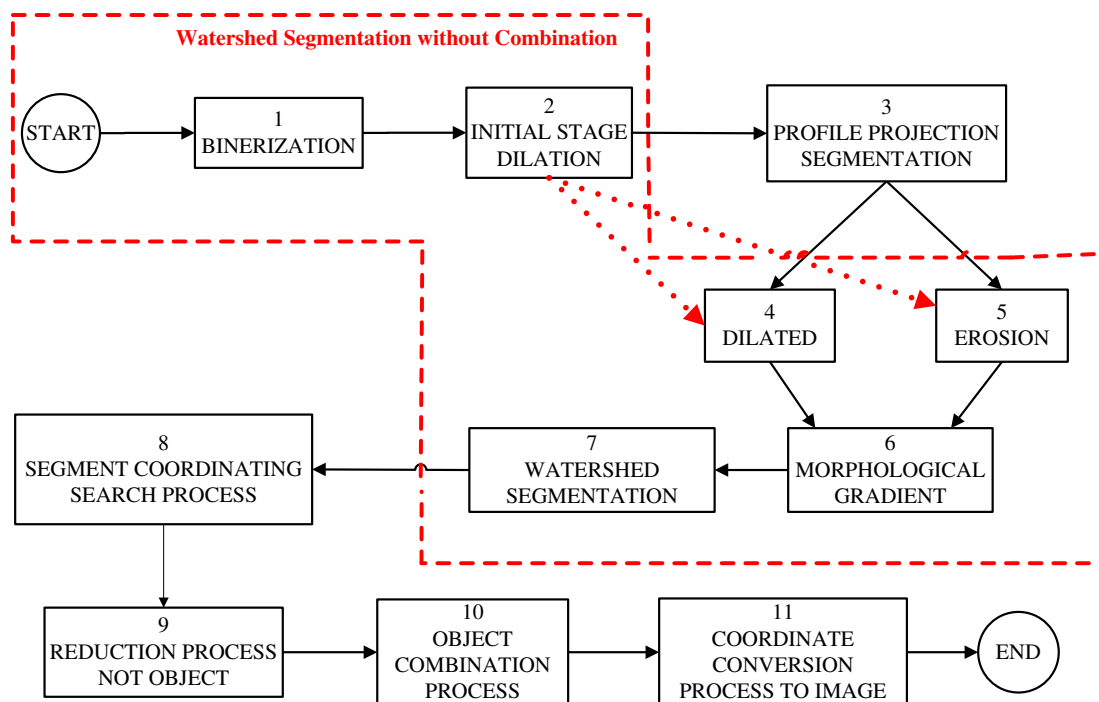


Figure 1. Designing the process flow of the image segmentation number notation

Figure 1 shows the flow of the combined segmentation process from the watershed algorithm, profile projection and filtering not objects that include binaryzation processes, dilated early stages until the process of converting coordinates to images, while the design of watershed algorithm segmentation without combination includes

binary processes, initial dilation, dilation, erosion, morphological gradient and watershed segmentation.

The binary process is the process of changing the gray scale format into a binary image (black and white). In the binary process, thresholding will be sought, then the point with a certain gray scale value range is changed to black and the rest to white. In this study using the otsu method [9]. The purpose of the otsu method is to divide the gray scale image histogram into two different regions automatically without requiring user assistance to enter the threshold value.

The initial dilation process is the process of enlarging the size of an object and repairing a separate object by adding layers around the object with 8-connected. The purpose of this process is to smooth the binary image so that separate pixels will join into a complete object.

Profile projection is the process of changing a binary image into a single dimension array (histogram) that is perpendicular to the  $x$  axis or  $y$  axis. Image profile projection is divided into 2, namely horizontal profile projections and vertical profile projections. Horizontal profile projection is the number of black pixels perpendicular to the  $x$  axis using equation (1):

$$P_h[j] = \sum_{i=1}^N S[i,j] \quad (1)$$

Whereas vertical profile projection is the number of black pixels perpendicular to the  $y$  axis using equation (2):

$$P_v[i] = \sum_{j=1}^M S[i,j] \quad (2)$$

Description of equations (1) and (2), namely  $S [i,j]$  are images in row  $i$  and column  $j$ , are vertical profile projections, are horizontal profile projections,  $M$  is column or image width,  $N =$  row or image height,  $i$  are indexes for rows (1,2,3, ...,  $N$ ) and  $j$  is index for columns (1,2,3, ...,  $M$ ).

Dilation process [9] is a process to increase the size of an object by adding layers around the object with 8-connected. The purpose of the dilation process is slightly different from the initial dilation process, which is to produce a morphological gradient

image. The process of erosion [9] is the process of reducing the size of an object by eroding the layer around the object with 8-connected. The purpose of this process is to produce a morphological gradient image. The morphological gradient process [9] is the process by which the new image produced is the result of the difference between the dilation process and the erosion process. The morphological gradient process aims to prevent excessive segmentation. The watershed segmentation process [9] is a regionally based segmentation process that is carried out by dam formation or watershed lines.

Segment coordinate search process is the process of finding the edge of each segment result represented in the form of a row index and segment image column index. The coordinate search process aims as the initial process for reduction rather than objects, combining objects and labeling sequence segments. The segment coordinate search process is carried out using profile projections horizontally and vertically without going through the cutting process. Segment coordinate search results generate a list in the form of a table which includes:

1. A segment number is the result of a numbered segment,
2. The segment line X1 represents the upper edge segment,
3. The segment line Y1 represents the left edge segment,
4. The segment line X2 represents the lower edge segment, and
5. The segment line Y2 represents the right edge segment

The reduction process is not an object is the process of removing the results of a segment not an object. The reduction process aims to improve the level of accuracy of segmentation. Because the results of segmentation are influenced by the morphological gradient process, the reduction process is divided into 2 stages, namely:

1. The first stage is to reduce not the object in the background segment, and
2. The second stage is to reduce segments not special objects in notations that have holes.

such as number notations 0 and 6, where each time a segmentation occurs, it will produce segments instead of objects in the form of holes that are considered as object segments.

So a non-object segment is a subset of an object segment so that it can be defined using equation (3):

$$\text{If segment } A \subset \text{segment } B, \text{ then delete segment } A \tag{3}$$

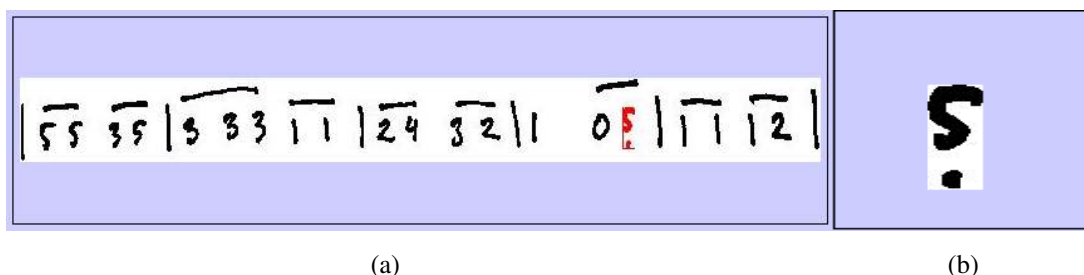
The process of combining objects is the process of uniting two objects resulting from a segment into a segment object. On the results of watershed segmentation in high notes or low tones 2 segments of objects will be formed so that the object merging process needs to be done. In the point segment and tested segment, a merger process can be carried out if vertically the point segment has 1 area or region with the tested segment, so that it can be defined using equation (4):

$$\text{If the point segment } \cap \text{ segment is tested, then join the two segments} \tag{4}$$

The process of converting coordinates to images is the final stage of the segmentation process. The process of coordinating to image is the process of changing the coordinate axis that has undergone a process of reduction rather than an object and the incorporation of objects into segment image files.

### 2.2 Implementation

In the results of the implementation of the combination watershed segmentation program is illustrated in Figure 2.



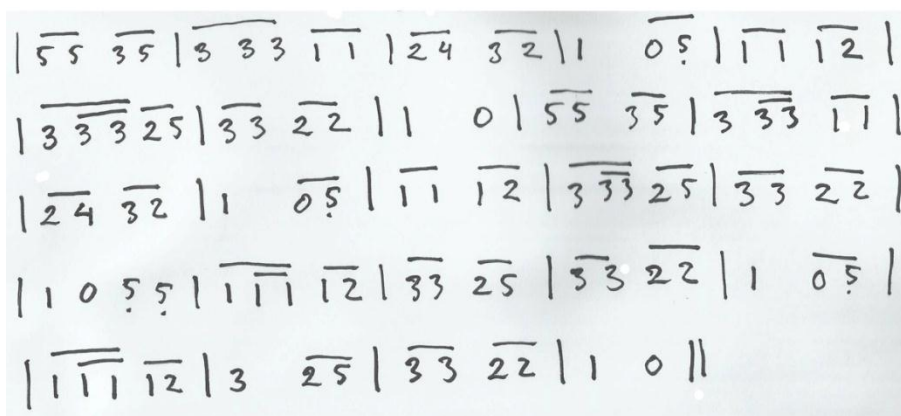
**Figure 2.** Illustration of segmentation results: (a) Line 1 sub segmentation results. (b) Results of segment image detail

**2.3 Dataset**

The dataset used is the image of handwritten number notation on paper (analog data), then the scan process and crop are carried out so that it produces the type of image file jpg (digital data) with a file size of 1800 × 900 pixels. Table 1 is a list of datasets that will be used as test data. The category column consists of 3 types of categories, namely simple, medium and complex. In the simple category is a score image that does not have a flat line segment or legato, the medium category is a simple category image and has a 1-level flat line segment or legato, and a complex category is a medium category score image and has a combination of 2-level flat line segments or combinations flat line with legato. Examples of original images of number notation is shown in Figure 3.

**Table 1 .** Feature image dataset number notation

No	Image name	Title	Origin	Category
1	Feature image 1	Cublak cublak suweng	Central Java	Simple
2	Feature image 2	Naik Naik Ke Puncak Gunung	Maluku	Simple
3	Feature image 3	Ayo Mama	Maluku	Medium
4	Feature image 4	Ampar Ampar Pisang	South Borneo	Medium
5	Feature image 5	Anak Kambing Saya	East Nusa Tenggara	Medium
6	Feature image 6	Gelang Sipaku Gelang	West Sumatra	Medium
7	Feature image 7	Suwe Ora Jamu	Yogyakarta	Complex
8	Feature image 8	Burung Tantina	Maluku	Complex
9	Feature image 9	Yamko Rambe Yamko	Papua	Complex
10	Feature image 10	Cik Cik Periok	West Borneo	Complex



**Figure 3.** Original images are handwritten number notations.

### 3 Results and Discussions

The results of segmentation testing and discussion were carried out by testing the level of accuracy of the combination watershed and watershed segmentation without combination. Testing the results of watershed segmentation accuracy without a combination only displays segmentation results. In testing the results of the accuracy of segmentation watershed combinations in addition to displaying the results of segmentation, it also displays the results of a non-object reduction process and object merging.

#### 3.1. Results of Segment Accuracy in Combined Watershed Segmentation

In Table 2 the results of the accuracy of the watershed segmentation combination with the average test results of the level of accuracy of segments in combination watershed segmentation is 99.74%. The factor that causes the level of accuracy of the segment not to reach 100% occurs in the 4th image which experiences errors of 6 objects.

**Table 2.** Results of Combined Watershed Segmentation Accuracy

Feature image	Number of Objects				True	False	Accuracy (%)
	Original image	Segmentation results image	Reduction Not Object	Object Merging			
1	79	105	18	8	79	0	100
2	124	141	14	3	124	0	100
3	169	197	13	15	169	0	100
4	159	195	28	8	155	6	97.4
5	174	203	17	12	174	0	100
6	108	119	11	0	108	0	100
7	112	122	6	4	112	0	100
8	110	139	20	9	110	0	100
9	215	253	33	5	215	0	100
10	167	190	18	5	167	0	100
Average							99.74

Figure 4 contains an example of the results of the image segmentation of the correct combination of watershed. Figure 4 point A is the result of watershed segmentation in the form of background segments, hole segments and score segments. Figure 4 point B

is an example of the results of a background segment and a hole segment is not an object so that it will experience a reduction process not an object. Figure 4 point C is an example of 2 segments which are low tones so that they will experience a process of combining objects. Figure 4 point D is an example of segmentation results that have taken the watershed combination process and the three segments are the correct segments.

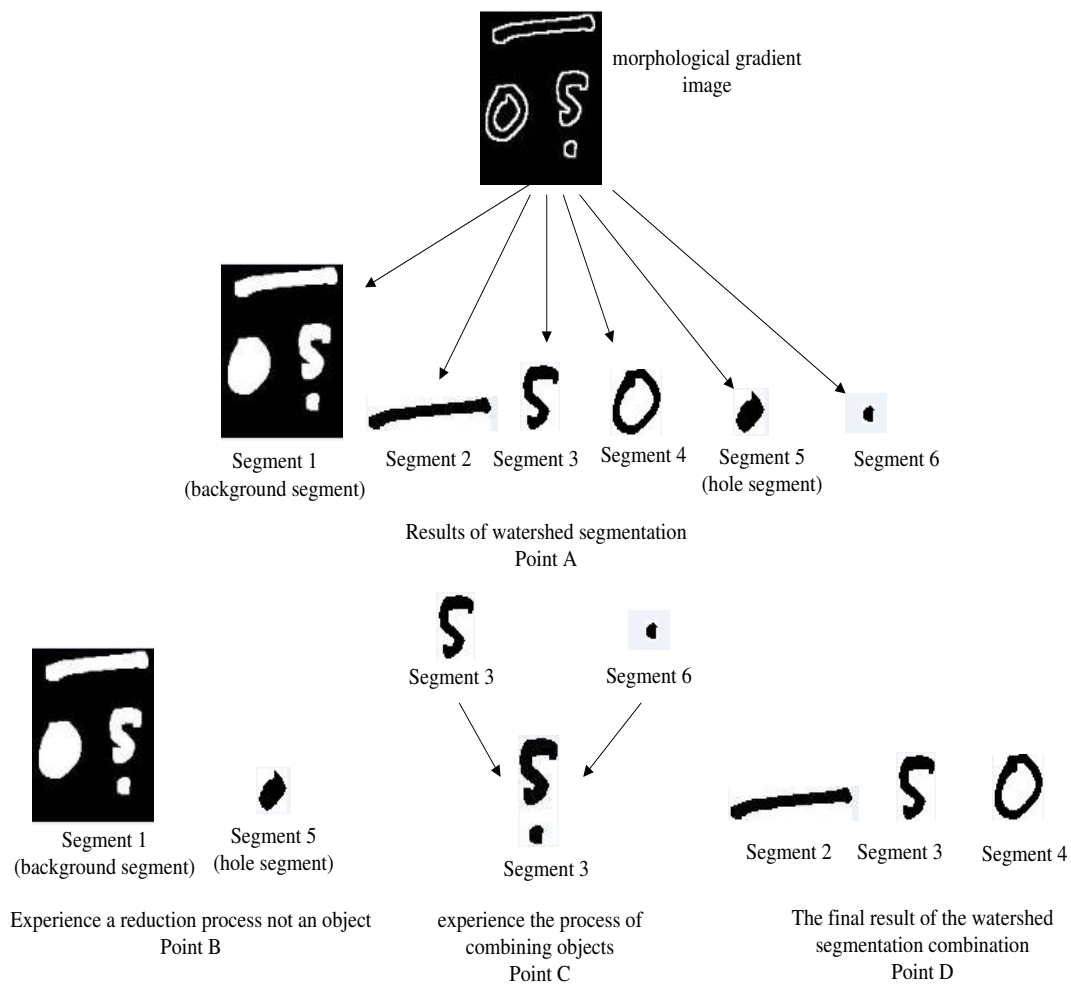


Figure 4. Examples of combined image results in combination watershed



**3.2. Results of Segment Accuracy in Watershed Segmentation without Combination**

In Table 3, the accuracy of 94.82% is obtained for the case of watershed segmentation without combination. The factor that causes decreased accuracy is the process like Figure 5. In Figure 5 point A is the result of watershed segmentation without combination. Figure 5 point B is the remainder of the non-object segment, namely the background segment and the hole segment does not undergo a reduction process. Figure 5 point C segment 3 and segment 5 do not experience the process of combining objects into low notes. Figure 5 point D is the final result of watershed segmentation without combination.

**Table 3.** Result of segment accuracy in watershed segmentation without combination

Feature image	Number of Objects			Accuracy (%)
	Original image	Segmentation results image	True False	
1	79	102	71 31	89.87
2	124	136	121 15	97.58
3	169	194	154 40	91.12
4	159	190	148 42	93.82
5	174	198	162 36	93.1
6	108	118	108 10	100
7	112	119	108 11	96.43
8	110	134	101 33	91.82
9	215	248	211 37	98.14
10	167	186	162 24	97.01
Average				94.82

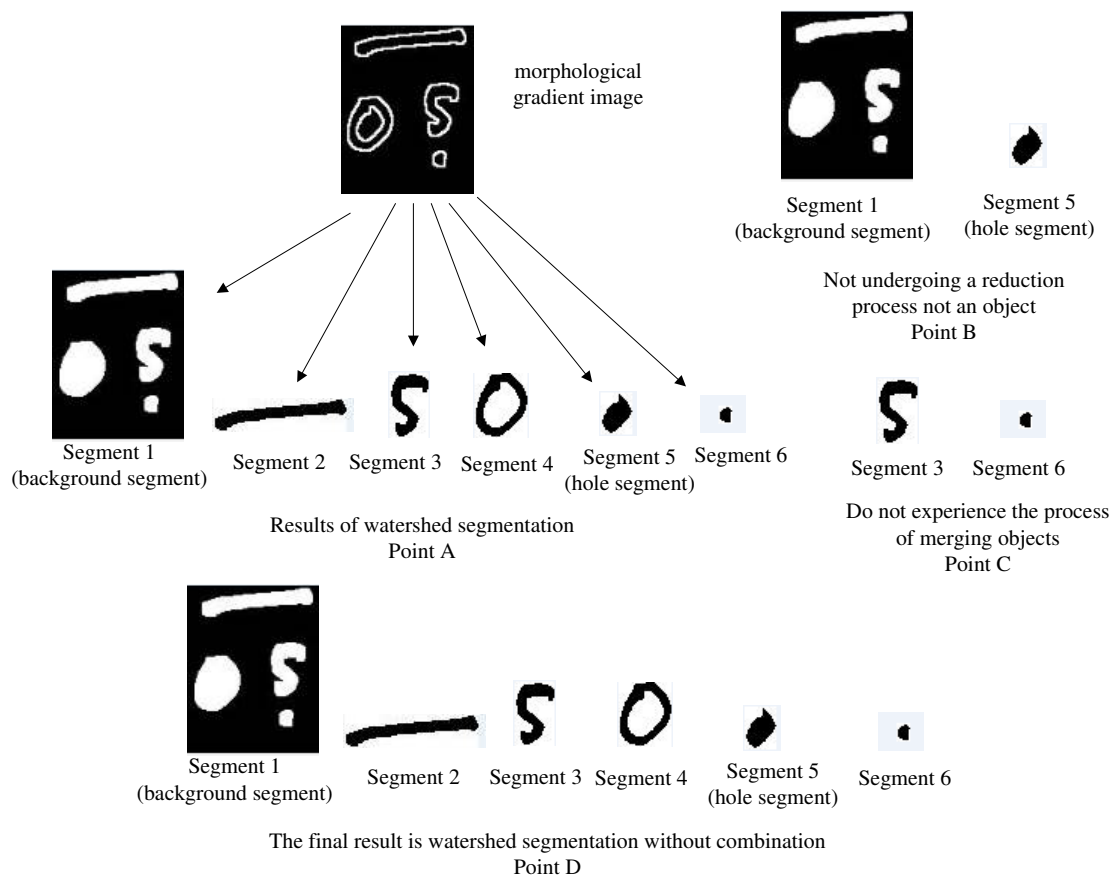


Figure 5. Examples of watershed image segmentation results without combination

## 4 Conclusions

Based on the results of research from system testing it can be concluded that the level of accuracy of segments in combination watershed segmentation is 99.74% higher than the level of accuracy of watershed segmentation without a combination of 94.82%. Increased accuracy of 4.92% indicates that the combination watershed segmentation algorithm is better when compared to the watershed segmentation algorithm without combination. So that the combination of watershed segmentation algorithms can be used in the object recognition process.

## References

[1] S. Wijayanti, *Seni Budaya (Musik) Kelas X SMA Negeri 1 Pati*, Seni Budaya Kelas X SMA, p. 2006, (2006).

- [2] Murinto and A. Harjoko, “Segmentasi citra menggunakan watershed dan intensitas filtering sebagai pre processing,” *Seminar Nasional Informatika*, 43–47, Mei 2009.
- [3] D. Rahmawati, A. Harjoko, and R. Sumiharto, “Purwarupa sistem tracking sungai menggunakan unmanned aerial vehicle,” *Indonesian Journal of Electronics and Instrumentation Systems*, **3** (2), 157–164, 2013.
- [4] A. R. Himamunanto and A. R. Widiarti, “Javanese character image segmentation of document image of Hamong Tani,” *Digital Heritage International Congress*, **1**, 641–644, 2013.
- [5] A. R. Widiarti, A. Harjoko, S. Hartati, and Marsono, “Implementasi model segmentasi manuskrip beraksara Jawa pada manuskrip beraksara Batak,” *Proceeding Seminar Nasional Inovasi dan Teknologi Informasi*, 81–84, Oktober 2014.
- [6] T. Shreekanth and V. Udayashankara, “A two stage braille character segmentation approach for embossed double sided Hindi devanagari braille documents,” *International Conference on Contemporary Computing and Informatics*, 533–538, November 2014.
- [7] S. Biswas and A. K. Das, “Text segmentation from scanned land map images using radon transform based projection profile,” *Proceedings of the International Conference of Soft Computing and Pattern Recognition (SoCPaR)*, 413–418, 2011.
- [8] D. Min, “Research on numbered musical notation recognition and performance in a intelligent system,” *Institute of Information and Engineering Hunan University of Science and Engineering*, 2–5, 2011.
- [9] R. C. Gonzalez and R. E. Woods, “Digital Image Processing: 2nd,” *Publishing Company*, USA, 2002.
- [10] Y. Zhang and X. Cheng, “Medical image segmentation based on watershed and graph theory,” no. 1, 1419–1422, 2010.
- [11] K. Pinaryanto and A. R. Widiarti, “Implementasi segmentasi citra dokumen teks sastra Jawa menggunakan algoritma watershed,” *Undergraduate thesis*, Universitas Sanata Dharma, 2009.

This page intentionally left blank

# **Factors Influencing the Difficulty Level of the Subject: Machine Learning Technique Approaches**

Hari Suparwito

*Department of Informatics, Faculty of Science and Technology,*

*Sanata Dharma University, Yogyakarta, Indonesia*

*Corresponding Author: shirsj@jesuits.net*

(Received 07-05-2019; Revised 21-05-2019; Accepted 21-05-2019)

## **Abstract**

The difficulty level of a subject is needed either to understand the student acceptance of the subject and the highest level of student achievement in it. Some factors are considered, what kind of instructions, the readiness of the instructor and students in teaching and learning, evaluation and monitoring systems, and student expectations. Many factors are involved, and educators should know this. It is better if they can discern which are the prime factors and which the secondary factors. The purpose of the study is to find out the determinant factors in establishing the difficulty level of the subject from the students', teachers' and infrastructure point of view using three machine learning techniques. The MSE and the variable importance measurement were used to predict between some factors such as Attendance, Instructors, and other factors as independent variables and the difficulty level of the subject as a dependent variable. The study result showed that Gradient Boosting Machine obtained the MSE value result 1.14 and 1.30 for training and validation dataset. The model generated five variable importance as an independent factor, i.e. Attendance, Instructor, The course can give a new perspective to students, The quizzes, assignments, projects and exams

contributed to helping the learning, and The Instructor was committed to the course and was understandable. The Gradient Boosting Machine is superior to other methods with the lowest MSE and MAE values results. Two methods, Gradient Boosting Machine and Deep Learning, have produced the same five main factors that influenced the difficulty of the subject. It means these factors are significant and should get attention by the stakeholders

**Keywords:** machine learning, regression, deep learning, random forest, gradient boosting machine, data mining, education.

## **1 Introduction**

Education provides people with knowledge about life and the world. It helps build character and leads to illumination. Given the importance of education, researchers ask themselves what factors influence the process of teaching and the attitude of students so that the students can understand the subjects, and what factors help to measure the difficulty level of subjects. The difficulty level of subjects is needed both to understand either the student acceptance of their subject or to ascertain the highest level of the student achievement in them [1]

John D. et al. [2] have examined some aspects and conducted some reviews based on learning conditions, student characteristics, materials and criterion tasks for effective learning techniques. Another group of researchers [3] have found that the social context influenced effective teaching and learning. Some factors mentioned were direct instruction, frequent monitoring, sense of communities, and student expectations. There are many factors involve here.

Research on education using data mining are increasing and promising in the last years and mostly focusing the research on student's performance, the effectiveness of learning and students and teacher's perception of learning [4]. Romero et al. stated that the objective using data mining in education areas is to improve the learning itself and the actors are students and teachers with the subjects of learning and the way to deliver as a medium relates them. Vanthienen and De Witte [5] revealed that their study

showed the use of machine learning methods is advantageous especially when it faces a nonlinear interaction function such as the role of a school principal to accommodate the district size policies. Another research in education field using the machine learning technique was undertaken by Liao, Zingaro [6]. They stated that using machine learning techniques; they can identify students who are at risk of performing poorly in a course.

Moreover, the machine learning approach was also performed for evaluating and predicting the student's level of proficiency [7]. To successfully predict the quality of this type of educational process the authors use one of the machine learning techniques. They claimed that the proposed technique could be effectively used in the educational management when the online teaching strategy should be selected based on student's goals, individual features, needs and preferences. Finally, Cope and Kalantzis [8] claimed that the use of machine learning and big data analysis in research on education should be undertaken because these emerging sources of evidence of learning have significant implications for the relationships between assessment and instruction. Moreover, for educational researchers, these datasets are in some senses different from conventional evidentiary sources, and this raises a new approach and give a different point of view to the traditional research in education areas.

The objective of this research is to find out the determinant factors that affect the student's acceptance focusing on the difficulty level of students understanding of the subjects. Instead of using a statistical approach in this present study we performed three machine learning techniques, i.e. Deep Learning, Random Forest, and Gradient Boosting Machine. Another purpose of this research is to introduce and compare the results of three machine learning methods in education areas. As the data set, we collected the dataset from the student evaluation at Gazi University Ankara [9] and was taken from the UCI repository dataset. This data set will be examined by three machine learning techniques using H2O platforms.

This paper is organised as follows. In section 2, we describe the research methodology with the following process in data mining approaches and then the results based on the H2O data mining tools calculation are presented and discussed in section 3. In chapter 4, we provide the conclusion and the subsequent work research outcome.

## **2 Research Methodology**

In general, the steps in this study follows the model of data mining techniques [10]:

### ***2.1 Objective determination***

The first step was to discover the real-world problems. This study will attempt to answer the educational question of how to understand and measure the difficulty level of the subject from the students', teachers' and infrastructures' point of view. To be more precise, the following research question was raised: What is the determinant factors which make students think and establish that this subject is difficult or easy?

A hypothesis was created to test which attributes in the data set gives a significant contribution toward the research question: Students think that the level of the subject difficulty is more likely to be influenced by the subject syllabus, activities and interactions between students and instructors and the readiness of students and teachers to engage in the learning process.

By analysing and testing this hypothesis, it shall know the determinant factors to answer the question of why do the students think that the subject is difficult to understand? Moreover, what should be done by the teachers so that the students can accept and understand the subject materials more easily?

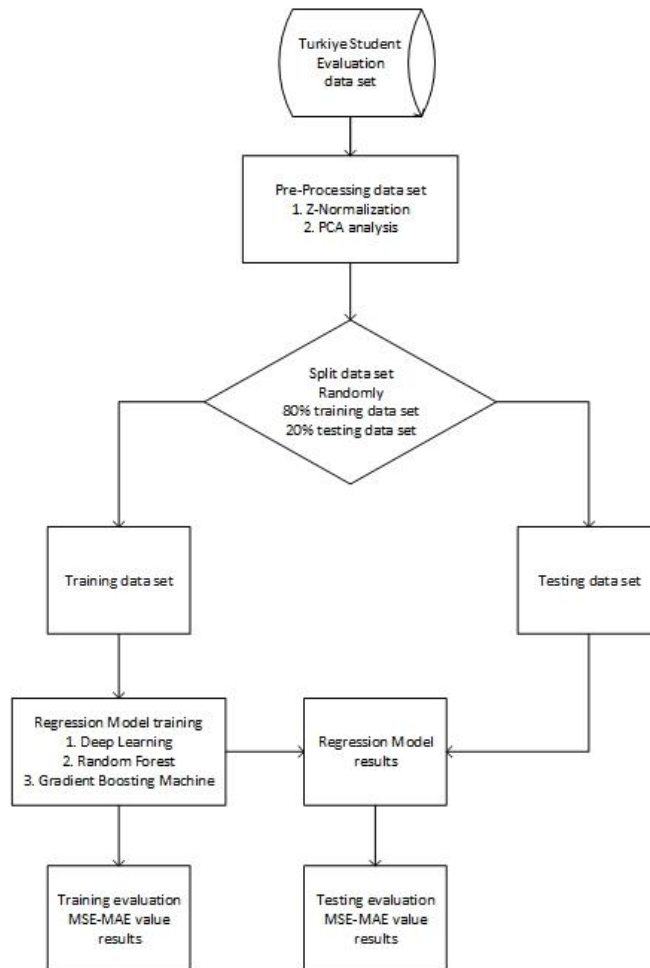
### ***2.2 The proposed work***

To examine three machine learning models we selected the dataset from the UCI machine learning repository about Turkiye Students Evaluation data set [9]. Furthermore, the dataset was analysed for reducing its dimensional features by using Principle Component Analysis (PCA) and then followed by performing a data normalisation using z-normalization. Moreover, the dataset was randomly split into ratio 80% : 20% from the data population as training and validation dataset.

Three machine learning techniques then were applied to training dataset obtaining the regression model, the MSE and MAE value results and the variable importance of each method. Using the model, we observed validation dataset to find out the MSE and



MAE values results and the variable significance for the testing dataset. All processes can be seen in the diagram below.



**Figure 1.** The proposed work. It was started by selecting the dataset to deliver the MSE and MAE value results and the variable importance rank

### 2.3 Data pre-processing

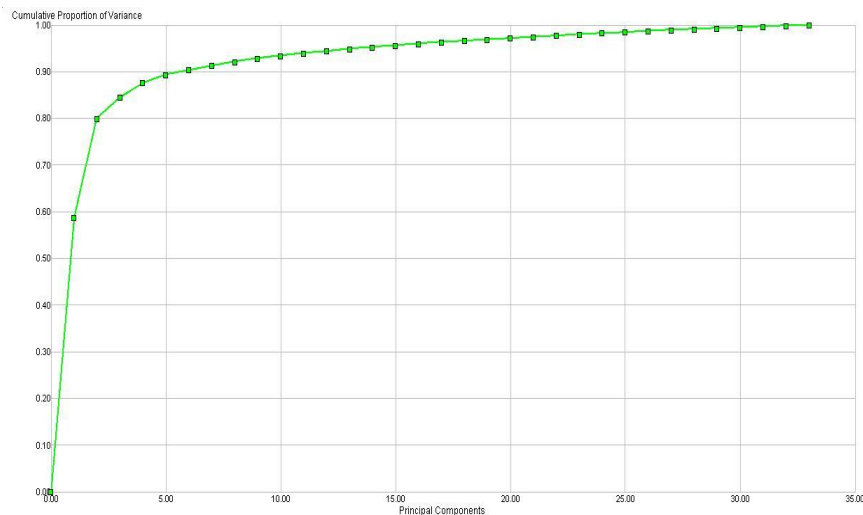
The research used the data result of the student questionnaire at Gazi University Ankara Turkey [9]. The dataset was obtained from the UCI machine learning repository dataset (<https://archive.ics.uci.edu/ml/index.php>). There are 5820 instances in the data set with 33 attributes where 28 attributes are formed in a Likert-type scale with the value from 1 to 5. The Likert-type scale values 1 equals to a strongly disagree value, and the value 5 equals to a strongly agree value. The five other attributes are questions with the answers in the natural numbers data format. The questions can be grouped into

three substantial group questions based on students', teachers' and infrastructures' point of view.

Next, we undertook a PCA analysis for features reduction. Matrix correlation from the PCA analysis showed each eigenvalue of the features. A new variable (principal component) was calculated based on eigenvalues with the values bigger than one. The PCA analysis result for a new variable is five principal components. We analysed and found that five principles components can be grouped into Attendance, Instructor, subject preparation, quizzes or exams, and the relationship between students and instructors.

**Table 1.** Table of principle component

Component	Standard Deviation	Proportion of Variance	Cumulative of Variance
PC1	6.140	0.588	0.588
PC2	3.686	0.212	0.800
PC3	1.701	0.045	0.845
PC4	1.411	0.031	0.876
PC5	1.059	0.017	0.894



**Figure 2.** The cumulative proportion of variance versus principle component.

From five principal components, we selected which features have a high rank based on the eigenvector values of each feature. Finally, we found 15 features that can be used

in this study. Therefore, the number of features was reduced from 33 features to 15 features only. A new reduced feature is shown in the following table.

**Table 2.** PCA analysis results

Features	Name of features	
	Difficulty (target label)	
	Attendance	
	Instructors	
Q1	The semester course content, teaching method and evaluation system were provided at the start	
Q4	The course was taught according to the syllabus announced on the first day of class.	
Q5	The class discussions, homework assignments, applications and studies were satisfactory.	
Q7	The course allowed fieldwork, applications, laboratory, discussion and other studies.	
Q8	The quizzes, assignments, projects and exams contributed to help the learning.	
Q12	The course helped me look at life and the world with a new perspective.	
Q16	The Instructor was committed to the course and was understandable.	
Q21	The Instructor demonstrated a positive approach to students.	
Q22	The Instructor was open and respectful of the views of students about the course	
Q24	The Instructor gave relevant homework assignments/projects and helped/guided students.	
Q25	The Instructor responded to questions about the course inside and outside of the course.	
Q27	The Instructor provided solutions to exams and discussed them with students.	
Q28	The Instructor treated all students in a right and objective manner.	

#### **2.4 Data mining**

The next process after the data pre-processing was to decide the kind of evaluation to be applied to the data set. The regression task is chosen because the data set is already classified in attributes and the questionnaire’s answer is on a Likert-type scale from 1 to 5 means already classified too. Another reason is that this study’ goal is directed to discover which attributes are the determinant factors of the difficulty level of the subject.

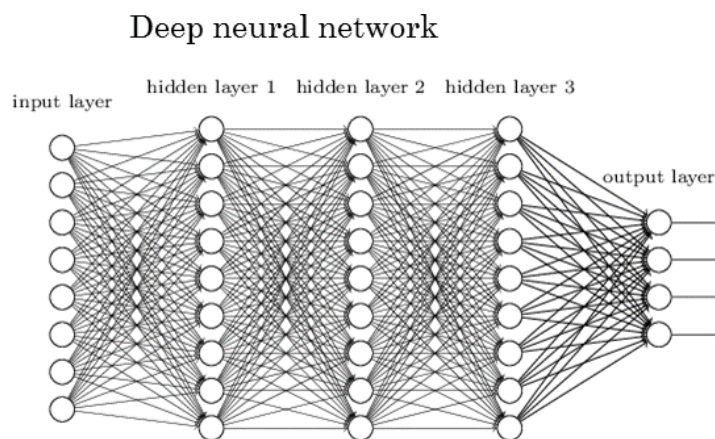
Three machine learning techniques that are Deep Learning (DL), Random Forest (RF) and Gradient Boosting Machine (GBM) were used to examine the data set focusing on the regression analysis between 15 attributes as an independent variable and the difficulty level of the subject as a target or dependent variable.

##### **2.4.1. Deep Learning**

Introduced the first time by Hinton et al. DL becomes more and more popular as one method to solve the problems in machine learning areas [11]. Deep learning is a part of machine learning techniques that aim to imitate the work of the human brain using an

artificial neural network. Different from other machine learning programs, the deep learning algorithm is made by a complex and high capability to learn, work and classify data.

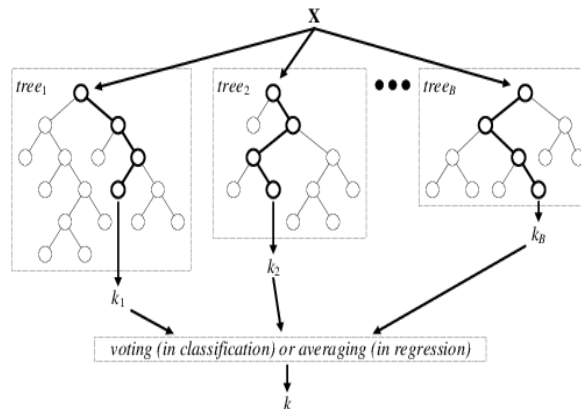
In general, DL consist of 3 main layers: input-hidden-output. Input layers work for containing raw data as input data. Hidden layers are applied for observing, learning and classifying data based on the references, in case of DL hidden layers usually consist of more than three layers. Output layers present the results.



**Figure 3.** Deep Learning diagram (the picture was taken from <https://www.kdnuggets.com/2017/05/deep-learning-big-deal.html>).

#### **2.4.2. Random Forest**

Random Forest is an ensemble learning technique for classification [12]. RF works by constructing a collection of decision tree at training time and returning the class that is the mode of all of the classes of the individual trees. Like DL, the RF algorithm has a significant advantage when analysing many of the datasets. It can address high-dimensional data with an excellent ability to learn from a large amount of data, and it can realise learning regression and classification for nonlinear sample data.



**Figure 4.** Random Forest architecture for classification and regression analysis (picture was taken from [https://www.researchgate.net/figure/Architecture-of-the-random-forest-model\\_fig1\\_301638643](https://www.researchgate.net/figure/Architecture-of-the-random-forest-model_fig1_301638643))

#### 2.4.3. Gradient Boosting Machine

Gradient boosting is a form of machine learning boosting. Boosting means target outcomes for each case are set based on the gradient of the error to the prediction. The idea behind GBM is to set the target outcomes for the next model in order to minimise the error. Each new model performs in the direction that minimises prediction error [13]. Even though RF and GBM are an ensemble learning method, GBM and RF differ in the way the trees are created: the order and the way the results are combined. GBM tries to add new trees that compliment the already built ones. This usually gives a better accuracy with fewer trees. Therefore, GBM performs better than RF if parameters tuned carefully [14].

#### 2.4.4 Cross-Validation

The goal of cross-validation is to test the model's ability to predict new data and to give an insight into how the model will generalise to an independent dataset. In each machine learning model was undertaken the K-fold Cross-Validation (CV) method and it was applied to training and testing data set. The K-fold CV method was selected for the data sampling method because data instances should be evaluated in training and testing data set. The number of instances is quite large so when the K-fold CV does the data sampling to the training and testing data set K-fold CV can do quite well. This

experiment was repeated many times, in this case, the repeating times was expressed by the K values. Even for some scientists argued that K=10 is the best value but in this research, the selection of the best K value in K-fold CV done by repeating many times experiment using various K values [15]. In this study, K-fold CV equal to 10 was applied.

Machine learning methods worked by using some parameters and finding the best result, each machine learning method has specific parameters to adjust. We used data grid analysis to find the best parameters to provide the optimum results. The following table shows the grid search parameters applied for

**Table 3.** Grid parameters values model

<b>Model</b>	<b>Grid Parameter values</b>
DL	Function – Rectifier; Tanh Hidden layers – 200, 200, 100, 50; 100,100,50; 50, 100, 100, 50 Epochs – 50; 100; 200 CV – 5; 10
RF	nTrees – 50; 100; 200 Epochs – 50; 100; 200 CV – 5; 10
GBM	nTrees – 50; 100; 200 Epochs – 50; 100; 200 CV – 5; 10

The best performance from each model showed by the following parameters

**Table 4.** Parameters values model

<b>Model</b>	<b>Parameter values</b>
DL	Function – Rectifier Hidden layers – 200, 200, 100, 50 Epochs – 200 CV – 10
RF	Input dropout – 0.2 nTrees – 200 Epochs – 100 CV – 10
GBM	nTrees – 50 Epochs – 50 CV – 10

Tabel 4 shows the best parameters gave by the grid search analysis.

### 3 Results and Discussions

Three machine learning methods were used to examine the dataset. The results obtained were the MSE and MAE values of each method and the variable importance. The Mean Squared Error (MSE) value was used to find the difference between the estimator and what is estimated. The MSE is achieved by applying the following formula:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \tag{1}$$

Where  $\hat{Y}$  is a vector of  $n$  prediction and  $Y$  is the vector of observed values corresponding to the input to the function that created the predictions.  $Y_i$  is the  $i$ -th value of the vector.

In this study, the training dataset was the data obtained from 80% number of data population, while the dataset from the rest of the number of populations (20%) was used as a testing dataset. H2O machine learning tools were performed for training and testing dataset, and the MSE value results are presented in the following table.

**Table 5.** MSE and MAE values of three machine learning models

Models	Training data set		Validation data set	
	MSE	MAE	MSE	MAE
DL	1.25	0.89	1.33	0.92
RF	1.31	0.92	1.38	0.91
GBM	1.14	0.84	1.30	0.90

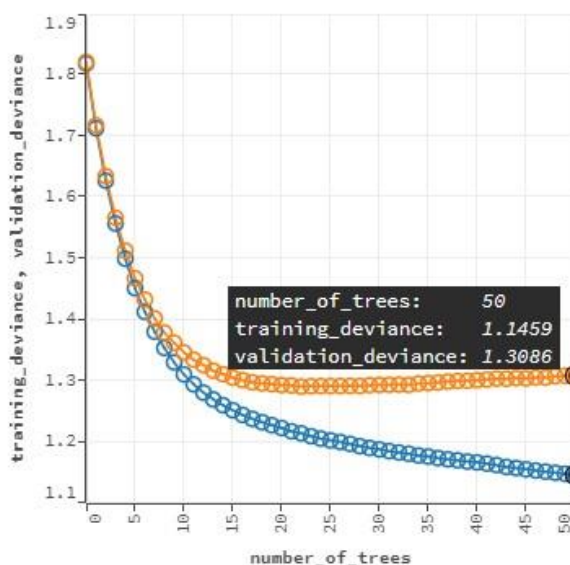
The lowest MSE values are the best result because it describes the similarity between the real values and the prediction values. In other words, the lower the MSE, the higher the accuracy of prediction as there would be an excellent match between the actual and predicted data set. In this study, the lowest MSE value is obtained by GBM models.

Like the MSE value, the MAE value obtained by the formula

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - x_i|}{n} = \frac{\sum_{i=1}^n |e_i|}{n} \tag{2}$$

Where  $x$  and  $y$  values are observed and predicted values. The lower MAE value also indicates better performance of the models.

Understanding the best model for the prediction can be performed by using deviance of training and testing dataset [16]. Deviance measurement is used for measure how well the model to predict It attempt is a generalisation of the idea of using the sum of squares of residuals in ordinary least square to cases where model-fitting is obtained by maximum likelihood. The following picture shows the deviance score for each number of trees in GBM.



**Figure 5.** GBM deviance score for each number of trees. We show the GBM model result only because GBM method obtained the best result

### 3.1. Variable importance

Wei, Lu [17] stated that it is essential to know which the more significant factor or variable in the regression or prediction analysis. Whereas Grömping [18] argued that predictive analysis would be more convincing when the most influential predictor variable obtained, though the way to find variable importance is challenging and some regression models are not directly planned to find the variable importance. Therefore, another method needs to be used to find the variable importance. Some techniques in machine learning could be used as an alternative way to find the variable importance, especially when dealing with high-dimensional input data and the categorical output.

Which variables are more significant in predicting the difficulty of the subject? Three ML methods were applied in this study. The percentage of Mean Square Error (MSE)



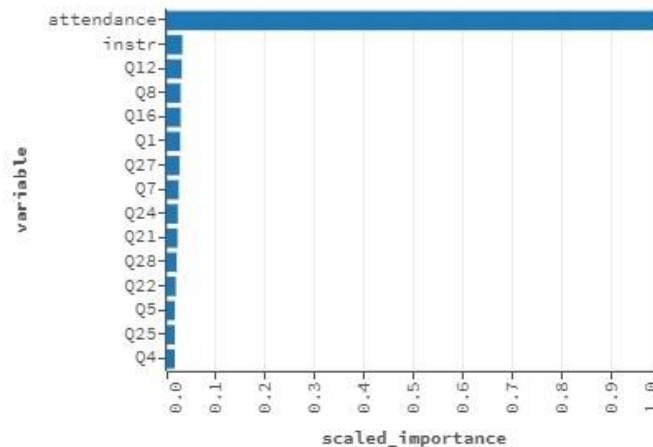
and Mean Absolute Error (MAE) was measured, which indicates which variable has a more significant influence compared with other variables in predicting the difficulty of the subject values. Table 6 shows the rank of the variable importance results and it also is given for example the graph of the variable importance from the GBM result in fig. 6

**Table 6.** variable importance results of each models

<b>Models</b>	<b>Variable importance</b>
<b>DL</b>	<ol style="list-style-type: none"><li>1. Attendance</li><li>2. Instructure</li><li>3. Q12 - The course helped me look at life and the world with a new perspective.</li><li>4. Q16 - The Instructor was committed to the course and was understandable</li><li>5. Q8 - The quizzes, assignments, projects and exams contributed to help the learning.</li></ol>
<b>RF</b>	<ol style="list-style-type: none"><li>1. Attendance</li><li>2. Q22 - The Instructor was open and respectful of the views of students about the course.</li><li>3. Q25 - The Instructor responded to questions about the course inside and outside of the course.</li><li>4. Q21 - The Instructor demonstrated a positive approach to students.</li><li>5. Instructure</li></ol>
<b>GBM</b>	<ol style="list-style-type: none"><li>1. Attendance</li><li>2. Instructure</li><li>3. Q12 - The course helped me look at life and the world with a new perspective.</li><li>4. Q8 - The quizzes, assignments, projects and exams contributed to help the learning.</li><li>5. Q16 - The Instructor was committed to the course and was understandable.</li></ol>

DL and GBM models have the same variable importance even though for Q8, Q12 and Q16 have a different rank. However, the main five factors are the same that was produced by DL and GBM analysis. For three machine learning models, two main factors are attendance and instructors have a significant influence in determining the difficulty level of the subject. It means these two factors are the most important predictor for the difficulty of the subject variable.

The previous study also revealed that student’s performance was not only dependent on their academic effort but also some other aspect that has a similar influence as well [19].



**Figure 6.** GBM variable importance

To answer the main question in the first section, now we can see the rank of the variable importance, especially from DL and GBM results. Moreover, if we observe which features have a significant influence, we can draw some points here,

- a) Attendance has the most significant impact. The respondent thought that attendance whether by students or by instructors have an important role and it can make their presumption about the subjects. Attendance means participation and involvement between students and instructors.
- b) Instructors and their attitudes or approach to the students are related to the subjects. The students are convinced that the instructors have a significant impact on delivering the subjects to them whether it was easy or difficult to be understood by them. This aspect is also related to the instructors' attitude such as how the instructor was committed to the course, how they respond if students are asking the subject in or out classes, how they can encourage the students to do the best with the selected subjects. The previous study by Martin, Wang [20] stated that instructors become an essential factor to make the subjects were easy or difficult in front of their students.
- c) The course can give a new perspective to students. A new perspective could be driven by the students. Therefore, they would focus on learning the subject and the next it will make the subject was easy to learn. In other words, giving a new

perspective for life become a stimulus to the students to learn and love the subjects.

- d) The quizzes, assignments, projects and exams contributed to help the learning. The students need the way to express their ability in understanding the subjects. The students felt that reading some theories were not enough, they needed some exercises, and by doing the exercises, they can understand the subject more. These aspects were also mentioned by Henderson and Harper [21] in their research. They revealed that some correction, assessment, and teacher's feedback on student's quizzes could help the students to prepare their exams better.

## 4 Conclusions

Three machine learning algorithms, i.e. Deep Learning, Random Forest, and Gradient Boosting Machine with K-folds CV data sampling methods have been applied to analyse the difficulty level of the subject based on students', teachers' and infrastructures' point of view. The data set is collected from the student questionnaire result at Gazi University Ankara. The result revealed that there are five determinant factors, i.e. Attendance, Instructors, the course helped me look at life and the world with a new perspective, the quizzes, assignments, projects and exams contributed to helping the learning, and the Instructor was committed to the course and was understandable. These five determinant factors can affect student's and instructor's perspective on the difficulty level of the subject. The two main factors are Attendance and Instructors. This study also demonstrated that data mining methods could be employed in the education field. However, the ability to understand data and how to work with them is very crucial. Data mining processes are important especially step by step at the stage model of data mining can be used as guidance on how to work with the data mining to solve the real-world problems.

In the subsequent study, it is possible to discover and compare these techniques with another algorithm in classification and regression tasks. Another possibility is also to compare some other tools such as Orange and Rapidminer tools where these two tools work on machine learning algorithm for solving the same problem.

## **Acknowledgements**

This research was supported by Department of Informatics Engineering, Sanata Dharma University. We would also like to thank the anonymous reviewers; whose comments greatly improved the manuscript.

## **References**

- [1] M. T. Tillery and A. Fishbach, “How to measure motivation: a guide for experimental social psychologist,” *Social and Personality Psychology Compass*, **8** (7), 328–341, 2014.
- [2] J. Dunlosky, K. A. Rawson, E. J. Marsh, M. J. Nathan, and D. T. Willingham, “Improving students' learning with effective learning techniques: promising directions from cognitive and educational psychology,” *Psychological Science in the Public Interest*, **14** (1), 4–5, 2013.
- [3] P. Hallinger and J. F. Murphy, “The social context of effective schools,” *American Journal of Education*, **94** (3), 328–355.
- [4] C. Romero and S. Ventura, “Data mining in education,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, **3** (1), 12–27, 2013.
- [5] J. Vanthienen and K. D. Witte, “Data analytics applications in education,” *CRC Press Taylor & Francis Group*, 2017.
- [6] Liao, S.N., et al., “A robust machine learning technique to predict low-performing students,” *ACM Transactions on Computing Education (TOCE)*, **19** (3), 18, 2019.
- [7] N. Kushik, N. Yevtushenko, and T. Evtushenko, “Novel machine learning technique for predicting teaching strategy effectiveness,” *International Journal of Information Management*, (2016). <https://doi.org/10.1016/j.ijinfomgt.2016.02.006>
- [8] B. Cope and M. Kalantzis, “Big data comes to school: implications for learning, assessment, and research,” *AERA Open*, **2** (2), 1–19, 2016.
- [9] G. Gunduza and E. Fokoue, *Turkiye student evaluation I*, University of California, School of Information and Computer Sciences, 2013.

- [10] P. Cabena, P. Hadjinian, R. Stadler, J. Verhees, and A. Zanasi, “Discovering data mining: from concept to implementation,” *Englewood Cliffs, N. J. Prentice Hall*, 1998.
- [11] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, **521** (7553), 436–444, 2015.
- [12] A. Liaw and M. Wiener, “Classification and regression by randomforest,” *R news*, **2** (3), 18–22, 2002.
- [13] J. H. Friedman, “Greedy function approximation: a gradient boosting machine,” *Annals of statistics*, **29** (5), 1189–1232, 2001.
- [14] R. E. Schapire, “The boosting approach to machine learning: an overview, in nonlinear estimation and classification,” *Springer*, 149–171, 2003.
- [15] R. Kohavi and F. Provost, “Confusion matrix,” *Machine Learning*, **30** (2-3), 271–274, 1998.
- [16] G. Ritschard, “Computing and using the deviance with classification trees,” *COMPSTAT 2006-Proceedings in Computational Statistics*, 55–66, August 2006.
- [17] P. Wei, Z. Lu, and J. Song, “Variable importance analysis: a comprehensive review,” *Reliability Engineering & System Safety*, **142**, 399–432, 2015.
- [18] U. Grömping, “Variable importance in regression models,” *Wiley Interdisciplinary Reviews: Computational Statistics*, **7** (2), 137–152, 2015.
- [19] A. A. Saa, “Educational data mining & students’ performance prediction,” *International Journal of Advanced Computer Science and Applications*, **7** (5), 212–220, 2016.
- [20] F. Martin, C. Wang, and A. Sadaf, “Student perception of helpfulness of facilitation strategies that enhance instructor presence, connectedness, engagement and learning in online courses,” *The Internet and Higher Education*, **37**, 52–65, 2018.
- [21] C. Henderson and K. A. Harper, “Quiz corrections: improving learning by encouraging students to reflect on their mistakes,” *The Physics Teacher*, **47** (9), 581–586, 2009.

This page intentionally left blank

## **AUTHOR GUIDELINES**

Author guidelines are available at the journal website:

<http://e-journal.usd.ac.id/index.php/IJASST/about/submissions#authorGuidelines>