# International Journal on Soft Computing

**ENRICHED
PUBLICATIONS**

# International Journal on Soft Computing

**Aims & Scope**

Soft computing is likely to play an important role in science and engineering in the future. The successful applications of soft computing and the rapid growth suggest that the impact of soft computing will be felt increasingly in coming years. Soft Computing encourages the integration of soft computing techniques and tools into both everyday and advanced applications. This Open access peer-reviewed journal serves as a platform that fosters new applications for all scientists and engineers engaged in research and development in this fast growing field.

# International Journal on Soft Computing

## Contents

# FACE SKETCH GENERATION USING EVOLUTIONARY COMPUTING

## N K Bansode 1 and P K Sinha 2

1Department of Computer Engineering, College of Engineering, Pune

## A B S T R A C T

*In this paper, an evolutionary genetic algorithm is used to generate face sketch from the face description. Face sketch generation without face image is extremely important for the law enforcement agencies. The genetic algorithm is used for generating face sketch through several iterations of the algorithm. The face image description is captured through graphical user interface just by clicking options for each face features. Face features are used to extract face images and generate initial population for the genetic algorithm. Genetic operators such as selection, crossover and mutation are used for next generation of the population. The Genetic algorithm cycle is repeated until the user is satisfied with face sketch generated. The novelty of the paper includes face sketch generation from face image description. The result shows that evolutionary based technique for sketch generation produces the desired face sketch.*

*KEYWORDS : Evolutionary Computing, Genetic Algorithm, Face Sketch Generation*

## 1. INTRODUCTION

Face plays an important role in person identification and conveys a verity of demographic information like age, gender, and emotions. We can recognize a familiar person and remember for several years. There are several applications of automatic face recognition such face authentication, face movement tracking, security, and surveillance. In Investigation, witness or victim of the crime provides the description of an attacker or any other source of information related to crime. Witness plays a very important role to give valuable information regarding the crime. Sometimes the attacker face image is not available in such cases, an artist help is taken to generate face sketch from the description given by the witness.

Employing artist for face sketch generation is time consuming and tedious tasks. The cognitive interview process is used to obtain information from a witness of the crime regarding the facial description of the suspect. To enable a computer to generate face sketch from the description involves an automation of conceptual sketching. Face composite systems were developed as alternate systems to the sketch artist. The face composite generation consists of a selection of the face features matching with target face and assembling together in the face frame. In the past several years, face composite systems were developed with the use of new technology for composite generation. The problem with face composite system is that, the limited number of the face features supported for face composite generation. The advanced face composite systems consist of large dataset of the facial features.

The recent face composites are generated using evolutionary genetic algorithm [1]. The evolutionary algorithm generates a variety of the faces by evolving face through several generations. The user (witness) selects the best matching face in the current population and new faces are generated through the process of face evolutions. The process stops when the generated faces are similar to the target face.

## 2. RELATED WORK

Face sketch generation systems developed in the past from sketch artist to the modern intelligent algorithm such as a genetic algorithm. The face composite systems developed using technological support available at the time of development of the composite systems. Due to technological advancement, a variety of the face composite systems developed over the period of the time [2]. The study of such systems presented in the following paragraphs. Xiaoou Tang and Xiaogang Wang [19] presented photo retrieval system using sketch drawing. Face photo recognized using sketches. Face features such as shape and texture calculated by eigen transformations. Hao Wang and Kangqiao Wang [6] used feature extraction and image based face drawing. Hough transforms used for face component detection such as eyes and intensity valley information to locate the pair of iris. Xiaoou Tang and Xiaogang Wang [20] described face recognition system using face sketches. A database of face photo and sketches of 188 people is used for photo retrieval. The image of face photo and sketch represents the different form of the image. Photo image represented by grayscale values and texture information, while sketches presented only by the grayscale values. In order to match sketch with the photo, the photo image converted into the sketch image. i.e. a database of photo image transformed into a database of sketch images. The eigenface recognition used to recognize face sketch in the database.

Hong Chen et al. [5,8,9] attempted to generate example based composite sketching of human portraits. The method for drawing face composite similar to the method used by the artist used for drawing of the picture. Fan Yang [3] presented non-parametric generation of example based human facial sketch. The conditional distribution of pixels in sketch image used to generate sketch. Hao Wang [7] attempted to draw a face using active shape models and parametric morphing.

Futoshi Sugimoto et al [4] presented drawing of a facial image in users mind using psychometric space model of the face. An image in user mind represented as psychological space model and image sketch considered into the different model of drawing. The genetic algorithm used to search image. Fuzzy reasoning is used to calculate the fitness of images generated by the genetic algorithm [15]. The process is repeated until the user gets satisfied or maximum number of generations are completed. Junji Nishino and Tomonori Kameyama [10] explained the process of caricature drawing using linguistic variables for the face features. Mayada F and Abdul Halim et al [11] described a system for facial composite generation using the genetic algorithm. The system consists of two step process. In step one, a database of facial part is created. In the second step, the genetic algorithm is used for reconstructing the facial composite image likeness to the facial composite image in mind of the witness. The recognition based strategy used to recognize the image rather than to recall the image. Additional tools for painting, smoothing, and sharpening used to enhance quality of the facial image.

Masashi Yamada et al [12] used the genetic algorithm to draw a logo. The picture of logo consists of one string and two images. Stuart Gibson et al [16] presented a facial composite system using an evolutionary algorithm. Global and local model for face features used for drawing face composite. Shape and texture for training images are derived and treated as the appearance model for the face. A witness is presented with virtual faces and allowed to determine likeness and ranking of each face by comparing with the target face. Three variants of evolutionary algorithm used and their performance measured using the virtual witness. Pong C, Yuen and C H Man [13] performed an experiment for human face searching using the face sketch images.

The literature survey for face composite generation indicates that several face composite systems were developed and the performance of these face composite systems was low. In this paper, we have implemented new approach for face sketch generation algorithm to generate face sketch from face image description.
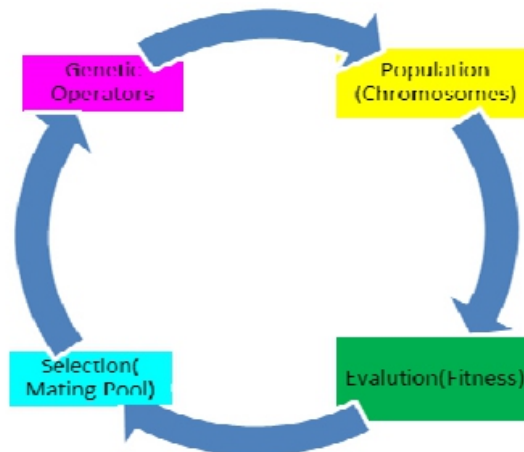
## 3. METHODOLOGY

The face sketch is generated using the process of the evolution. The genetic algorithm evolves the faces through several generations. The genetic algorithm is a search and optimization technique based on Darwin's principal of the "survival of the fittest". Face sketch generation system based on the genetic algorithm is an advanced process of the sketch generation which evolves the face using genetic operators. The Genetic algorithm is described in the following sections.

### 3.1 Genetic Algorithm

Genetic Algorithms are an evolutionary procedure that finds the solution to problems using the mechanics of natural selection. Genetic algorithms are used in the problem where the finding solution is difficult, but due to the probabilistic nature these algorithms gives optimal solutions. In the cycle of the genetic algorithm as shown in figure 1, it starts with an initial set of random solutions called population. Each individual in the population is called as chromosome, representing the solution to the problem to solve. During each generation, the chromosomes are evaluated using some measure of fitness [21]. This fitness of individual solution string is used to create the next generation, a new chromosome are formed by three essential operations: selection, crossover, mutation. The process of the genetic algorithm cycle is shown in Figure 1.

### 3.2 Genetic Operators

Selection is a process in which individual strings are copied according to their objective (fitness) function values. Copying strings according to their fitness value uses means that strings with a higher fitness value will have a higher probability contributing one or more offspring in the next generation. The crossover is a process of merging two chromosomes from current generation to from two similar offspring's. The mutation is a process of modifying a chromosome and occasionally one or more bits of a string are altered while the process is being performed. The flowchart of the genetic algorithm is shown in Figure 2



**Figure 1:** Genetic Algorithm Cycle

**Figure 2 :** Genetic Algorithm flowchart
**Table 1:** Face features description parameters

| Sr. No. | Features | | | |
|---|---|---|---|---|
| 1 | Gender | Male | Female | |
| 2 | Age Group | Child | Young | Old |
| 3 | Face Shape | Ellipse | Circle | Oval, Square, Triangle |
| 4 | Left Eye brow | Small | Normal | Large |
| 5 | Right Eyebrow | Small | Normal | Large |
| 6 | Left Eye | Thin | Medium | Large |
| 7 | Right Eye | Small | Normal | Large |
| 8 | Nose | Small | Normal | Large |
| 9 | Mouth | Thin | Medium | Large |

**Figure 3 :** Graphical user interface for face description

### 3.3 Face Sketch Generation Algorithm

The face sketch is generated using evolutionary genetic algorithm based on the facial feature description provided by the user. Table 1 shows the face features description and the possible parameter values. The graphical user interface is designed as shown in Figure 3 for capturing the facial parameters. The facial description entered through the GUI is used to extract faces which are resemble to the feature descriptions. The faces collected from the description used as initial population for the genetic algorithm. The genetic algorithm works on these faces is given below:

**Genetic algorithm for Sketch Generation**

// This algorithm generates the face composite from the face image
// Input    : Face Image
// Output : Face Sketch
//  Input parameters: Population Size, Crossover Rate, Mutation Rate, Max. Number of Generations

**Procedure**

1.  Begin
2.      Start Generation (t <0)
3.        Initialize Face Population (t)
4.          While (Not Termination Condition)
5.            begin
6.              t <-  t+1                        // Generation= Generation +1
7.              Select p(t) from p(t-1)   // Select the parent faces from the population
8.              Crossover p (t)              // Crossover the face to produce new faces

```
9.              Mutate p (t)          // Modify the face (genotype)
10.             Evaluate p (t)        // Find the fitness with face with the target face
11.        end                        // End of  Generation
12.     End                           // End of Maximum Number of Generations
13.  End                              // End Begin
```



**Figure 4 :**  Face Skecth Generation (Initial Young Population)



**Figure 5:** Face sketch Generation (Initial Old Populations)

## 4. RESULTS

The face sketch generation from the face features description is performed based on the evolutionary genetic algorithm. This system developed for automatic face sketch generation similar to the sketch artists. Face sketch generation for two types of population such as young and old population is implemented as shown in Figure 4 for young population and Figure 5 for old population. The genetic algorithm parameters such as population size, crossover rate and mutation rate are specified for each type of the population. The population after the 20 generations are converged and shown in Figure 6 for young population. The mean square error is measured for every generation and shown in Table 2 and 3. The two face datasets used for sketch generations [22, 23]. Figure 7 and 8 shows graph for the mean square error in each generation. The average mean square error is reduced in each generation.



**Figure 6: Population after 20 generations**

**Table 2: Final Population after 20 Generation (Young Population)**

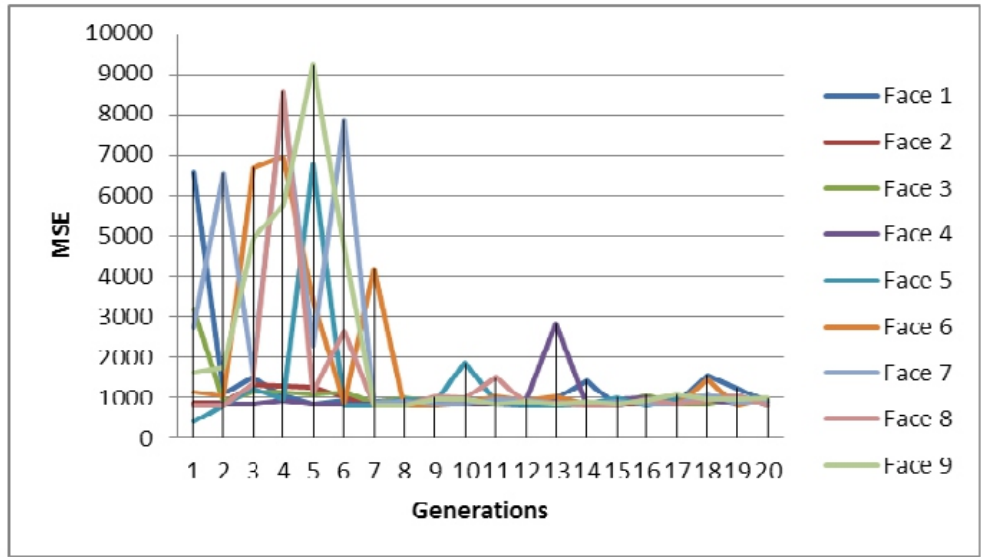| Face\Gen. | Face 1 | Face 2 | Face 3 | Face 4 | Face 5 | Face 6 | Face 7 | Face 8 | Face 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 6575.89 | 6865.27 | 3196.66 | 1816.162 | 1398.12 | 1137.18 | 2752.34 | 798.47 | 1635.24 |
| 2 | 1115.87 | 867.74 | 897.13 | 830.73 | 816.40 | 1060.80 | 6547.87 | 816.11 | 1758.06 |
| 3 | 1509.26 | 1322.38 | 1132.46 | 846.23 | 1227.64 | 6680.36 | 1605.06 | 1343.65 | 4933.15 |
| 4 | 1033.20 | 1295.46 | 1133.06 | 925.87 | 951.78 | 6942.66 | 8377.96 | 8553.80 | 5729.37 |
| 5 | 831.02 | 1266.25 | 1080.76 | 826.25 | 6773.00 | 3411.72 | 2296.55 | 1153.05 | 9212.70 |
| 6 | 915.20 | 1012.78 | 1126.04 | 824.05 | 820.64 | 852.79 | 7855.39 | 2654.95 | 4705.32 |
| 7 | 904.68 | 816.20 | 870.61 | 823.37 | 806.54 | 4184.95 | 870.52 | 813.93 | 814.66 |
| 8 | 971.35 | 855.15 | 976.11 | 817.09 | 950.79 | 818.30 | 902.94 | 818.20 | 819.93 |
| 9 | 901.76 | 831.25 | 889.93 | 900.24 | 900.59 | 819.94 | 823.06 | 1011.00 | 960.84 |
| 10 | 830.97 | 823.87 | 901.44 | 848.54 | 1851.98 | 864.03 | 835.12 | 973.43 | 948.34 |
| 11 | 832.28 | 869.25 | 913.56 | 866.51 | 889.02 | 1005.80 | 966.75 | 1499.73 | 847.55 |
| 12 | 819.63 | 860.48 | 965.04 | 913.38 | 818.14 | 919.31 | 937.65 | 870.19 | 884.35 |
| 13 | 1001.42 | 846.14 | 818.84 | 2821.5 | 817.88 | 1025.73 | 855.15 | 858.94 | 821.62 |
| 14 | 1424.40 | 864.74 | 853.94 | 816.37 | 850.78 | 832.70 | 849.02 | 816.61 | 899.64 |
| 15 | 818.18 | 823.99 | 821.93 | 944.27 | 1001.46 | 823.36 | 900.54 | 816.35 | 848.95 |
| 16 | 921.95 | 837.84 | 1058.30 | 972.45 | 819.52 | 901.47 | 825.33 | 968.29 | 915.40 |
| 17 | 899.33 | 1045.80 | 824.01 | 857.66 | 965.53 | 823.64 | 828.59 | 832.33 | 1063.82 |
| 18 | 1536.84 | 984.36 | 845.18 | 969.27 | 977.60 | 1440.4 | 1068.5 | 936.11 | 964.81 |
| 19 | 1234.40 | 979.99 | 964.17 | 825.51 | 1037.4 | 818.61 | 874.14 | 1047.76 | 966.43 |
| 20 | 890.07 | 929.83 | 919.91 | 968.92 | 817.83 | 964.79 | 917.45 | 819.08 | 982.88 |

**Figure 7: Graph of MSE (Young Population)**
**Table 3: Final Population after 20 Generation (Old Population)**
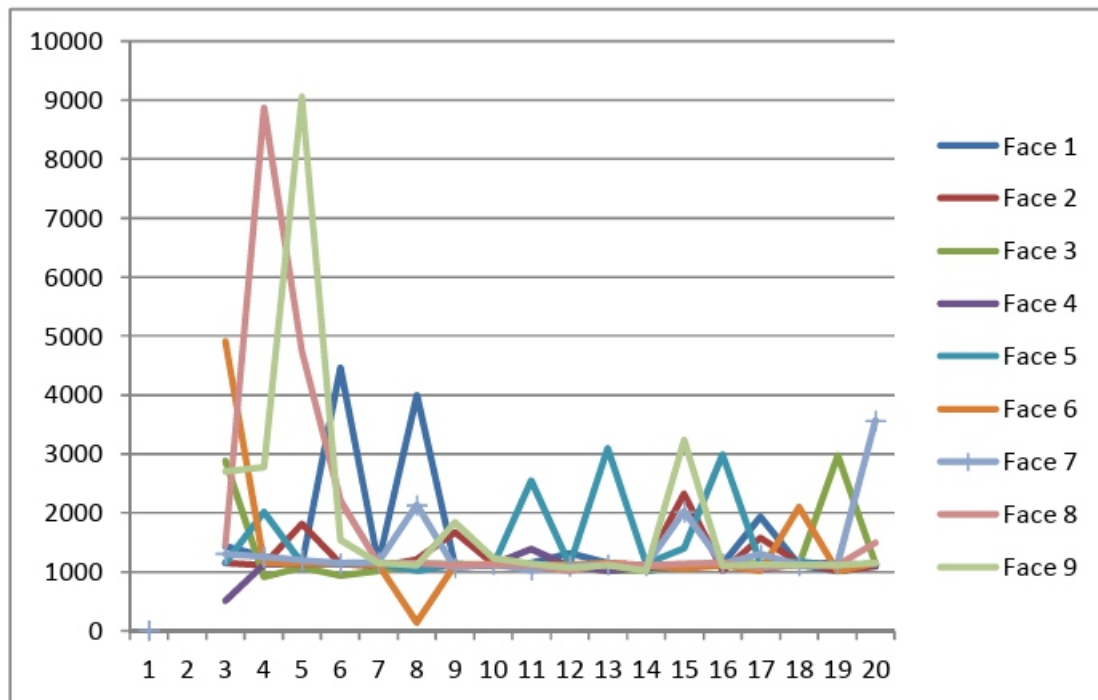
| Face\Gen | Face 1 | Face 2 | Face 3 | Face 4 | Face 5 | Face 6 | Face 7 | Face 8 | Face 9 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1431.23 | 1153.11 | 2878.64 | 511.15 | 1146.86 | 4914.74 | 1308.75 | 1461.32 | 2698.54 |
| 2 | 1257.42 | 1112.41 | 908.86 | 1130.82 | 2014.21 | 1153.73 | 1256.49 | 8873.04 | 2772.91 |
| 3 | 1120.96 | 1816.36 | 1076.51 | 1139.53 | 1145.32 | 1134.93 | 1194.01 | 4743.95 | 9060.10 |
| 4 | 4462.25 | 1133.22 | 940.66 | 1131.19 | 1143.41 | 1143.38 | 1143.11 | 2201.80 | 1538.10 |
| 5 | 1144.10 | 1087.46 | 1019.95 | 1127.60 | 1134.24 | 1126.31 | 1160.80 | 1147.81 | 1146.19 |
| 6 | 4000.30 | 1202.86 | 1138.55 | 1130.76 | 1015.02 | 143.60 | 2133.07 | 1143.55 | 1108.77 |
| 7 | 1125.39 | 1675.82 | 1142.28 | 1078.75 | 1123.78 | 1109.05 | 1070.18 | 1108.72 | 1836.99 |
| 8 | 1125.97 | 1142.91 | 1109.42 | 1123.56 | 1108.95 | 1140.99 | 1127.38 | 1122.18 | 1227.56 |
| 9 | 1143.67 | 1110.93 | 1129.70 | 1391.66 | 2553.60 | 1118.98 | 1036.19 | 1113.98 | 1141.56 |
| 10 | 1315.29 | 1043.99 | 1115.83 | 1083.91 | 1107.89 | 1124.31 | 1101.20 | 1013.83 | 1062.34 |
| 11 | 1145.15 | 1138.30 | 1018.18 | 1016.39 | 3091.19 | 1145.21 | 1124.62 | 1134.60 | 1109.93 |
| 12 | 1015.18 | 1114.12 | 1025.80 | 1064.26 | 1132.74 | 1112.51 | 1108.57 | 1113.12 | 1013.46 |
| 13 | 1112.33 | 2323.08 | 1108.53 | 1111.34 | 1399.08 | 1049.77 | 2020.83 | 1141.52 | 3234.54 |
| 14 | 1133.04 | 1015.27 | 1144.83 | 1117.88 | 2983.03 | 1108.11 | 1141.82 | 1145.14 | 1106.03 |
| 15 | 1926.40 | 1577.58 | 1140.37 | 1048.10 | 1135.61 | 1017.96 | 1297.80 | 1070.23 | 1129.32 |
| 16 | 1108.45 | 1144.62 | 1107.96 | 1141.11 | 1167.96 | 2102.49 | 1108.62 | 1117.61 | 1108.17 |
| 17 | 1022.07 | 1022.39 | 2972.77 | 1143.40 | 1099.93 | 1009.77 | 1141.56 | 1108.73 | 1107.82 |
| 18 | 1113.68 | 1095.35 | 1132.13 | 1106.48 | 1131.28 | 1144.20 | 3562.65 | 1490.80 | 1144.87 |
| 19 | 1108.46 | 1011.43 | 1122.13 | 1061.04 | 1118.53 | 1127.44 | 1125.32 | 1111.99 | 1121.25 |
| 20 | 1130.77 | 1086.23 | 1108.90 | 1137.04 | 1082.53 | 1119.67 | 1133.39 | 1129.49 | 1121.43 |

**Figure 8: Graph of MSE (Old Population)**

## 5. CONCLUSIONS

In this paper, we have performed the experiment for face sketch generation based on evolutionary genetic algorithm. The face is described using facial feature description such as gender, age and shape. The face other features such as size of left eye, right eye, left eyebrow, right eye brow, nose and mouth are captured through the graphical user interface as shown in Figure 3. The face sketch generation using face description based on genetic algorithm is novelty concept implemented in paper. Genetic algorithm works on the population of face and evolves through several generations. In each generation, the faces with higher fitness value retained and the faces with the lower fitness value are removed. Thus, genetic algorithm iterates through the several generation until the desired face is generated.

## REFERENCE

• *Charlie D Frowd, "Implement Holistic Dimensions for a Facial Composite system"*
• *Journal of Multimedia Vol. 1 no 3 June 2006.*
• *Douglas DeCarlo, Dimitris Metaxas, Mathew Stone, "An Anthropometric Face Model using Variational Techniques", Proceedings SIGGRAPH 98.*
• *Fan Yang "Non–parametric Generation Techniques of Example–based Human Facial sketch", Proceedings of the fifth world Congress on Intelligent control and Automation , June 15-19 2004 Hangzhou , P R China.*
• *Futoshi Sugimoto "A Human Interface to search and draw facial images in Mind by Using Psychometrical model of faces", IEEE International Fuzzy Systems Conference Proceedings august 22-25 1999 Seoul Korea.*
• *Hong Chen Ying-Qing Xu, Heung, "Example –based Facial Sketch Generation with nonparametric sampling", IEEE 2001.*
• *Hao Wang, Kangqiao Wang, "Facial Feature Extraction and Image-Based Face Drawing", International conference on Signal Processing", Proceedings 2002*

•*Hao Wang, "Image–Based Face Drawing Using Active Shape and Parametric Morphing" ,IEEE International conference on Neural Network and Signal Processing Nanjing, China, December 1417, 2003.*

# FUZZY-CLUSTERING BASED DATA GATHERING IN WIRELESS SENSOR NETWORK

**Arezoo Abasi  and Hedieh Sajedi**

Department of Mathematics, Statistics and Computer Science, College of Science, University of Tehran, Tehran, Iran

## A B S T R A C T

*Wireless Sensor Networks (WSN) is spatially distributed, collection of sensor nodes for the purpose of monitoring physical or environmental conditions, such as temperature, sound, pressure, etc. and to cooperatively pass their data through the network to a base station. The critical challenge is to minimize the energy consumption in data gathering and forwarding from sensor nodes to the sink. Cluster based data aggregation is one of the most popular communication protocols in this field. Clustering is an important procedure for extending the network lifetime in wireless sensor networks. Cluster Heads (CH) aggregate data from relevant cluster nodes and send it to the base station. A main challenge in WSNs is to select suitable CHs. Another communication protocol is based on a tree construction. In this protocol, energy consumption is low because there are short paths between the sensors. In this paper, Dynamic Fuzzy Clustering data aggregation is introduced. This approach is based on clustering and minimum spanning tree. The proposed method initially uses fuzzy decision making approach for the selection of Chs. Afterward a minimum spanning tree is constructed based on CHs. CHs are selected efficiently and accurately. The combining clustering and tree structure is reclaiming the advantages of the previous structures. Our method is compared to the well-known data aggregation methods, in terms of energy consumption and the amount of energy residuary in each sensor network lifetime. Our method decreases energy consumption of each node. When the best CHs selected and the minimum spanning tree is formed by the best CHs, the remaining energy of the nodes will be preserved. Node lifetime has an important role in WSN. Using our proposed data aggregation algorithm, survival of the network is improved.*

*KEYWORDS : Sensor networks; Energy efficiency; Data aggregation; Fuzzy decision making.*

## 1.INTRODUCTION

Wireless sensor network (WSN) is a collection of thousands of low-cost, low power electronically programmable devices, which are deployed in a monitored area in stochastic manner [2]. In this area, there is no opportunity for maintenance and battery replacement for the most of the applications, which use the sensor nodes to surveillance the remote field [1]. Potential applications of sensor networks include Industrial automation, Automated and smart homes, Video surveillance, Traffic monitoring, Medical device monitoring, Monitoring of weather conditions, Air traffic control, Robot control.

The most-distinguishing attributes of nodes used in WSNs are the limited power supply, storage capacity and communication bandwidth required. In WSN, bandwidth utilization and energy saving is a very important criterion for any existing and new applications. Normally, data collected from WSNs are large which makes it essential to eliminate redundant data, minimize the number of transmissions, and improve the energy consumption. The effort to reduce the number of data packet transmission with the in-network processing is called data aggregation [3].

Our sensor's battery is limited. The lifetime on each node depends on the power that has significantly affected the relationship between the nodes. One of the accurate requirements of these nodes is the efficient use of the saved energy. Multiple algorithms have been designed for impressive handling of nodes energy in WSNs using several clustering schemes [4, 5]. Optimal data aggregation can save nodes energy. In this sensor network, data are gathered by the sensor nodes from our study area. There is a data transmission method that merges data from several sensor nodes into one pack, which is data aggregation. Decreasing the disjointed communication at different levels and in turn to reduce the total energy consumption is the main aim of data aggregation.

There are dissipated different amounts of energy to process raw data. There are two popular protocols: Cluster based data aggregation [6] and Tree based data aggregation [7]. Some of WSNs consists of clusters, in which each cluster has a CH. CHs have a significant impress in network lifetime. An ideal CH is the one, which has the highest residual energy, maximum number of neighbour nodes around the CH and the shortest distance from the base station[8]. Whatever the selected CH is more similar to the ideal CH, network lifetime is increased.

We can use Multiple Attribute Decision Making (MADM) approach to select CHs with multi criteria [9]. This method quantitatively selects alternatives based on their multiple criteria. The main problem is the difficult estimation of the exact values of all the criteria. Synchronous consideration of all criteria in CHs selections can be used MADM approach. In case of multi criteria, fuzzy based MADM methodologies are efficient and impressive [10, 11].

In this paper, we proposed a hybrid approach called DFC data aggregation, which gathers and combines data and avoids redundant data transformations, therefore successively reduces power consumption and bandwidth.

Proposing DFC data aggregation, we preserve the advantages and minimize the disadvantages of the clustering and tree based approaches. We use DFC data aggregation to extend the lifetime of WSNs and energy consumption of sensor nodes. The optimized CHs are selected to spread energy efficiently using multi criteria. CHs are selected based on the residual energy, the number of neighbour nodes and distance from the base station. After cluster formation, CHs receive data from member nodes in clusters, aggregate data and send it to the base station. A spanning tree covers all the sides as vertices and consists no cycles. The tree is constructed in the procedure that the node with the smallest identifier is chosen as the root [10, 12]. All the nodes with the shortest path conjunct to the selected root. The protocol requires that each node exchange configuration message in a specific format, which contains its own identifier, its chosen root, and the distance to this selected root. Each node updates its configuration message upon identifying a root with a smaller identifier or the shortest-path neighbour. In addition, the neighbour for which the shortest route configuration message comes from is chosen as the parent of a node whenever it is detected.

In this paper, we employ multi criteria decision-making approach, Fuzzy Analytic Hierarchy Process (FAHP) and hierarchical fuzzy in clusters on WSNs [13, 14]. The Analytic Hierarchy Process (AHP) considers a set of assessment criteria, and a set of alternatives among which the best decision is to be made. The AHP generates a weight for each evaluation criteria according to the comparisons of the criteria. The superior the weight, the more significant the corresponding criterion. The AHP method could improve the network lifetime significantly.

In this research, we also analyse other methods, including Low Energy Adaptive Clustering Hierarchy (LEACH) [14], Tree Dara Aggregation (CTDA) [16], Modified Cluster based and Tree based Data Aggregation (MCTDA) [16], and Cluster based and Tree based Power Efficient Data

Collection and Aggregation (CTPEDCA)[17]. We compared our proposed method with these mentioned methods in terms of energy consumption and the amount of energy remaining in each sensor network lifetime.

Simulation conclusions illustrate that our proposed approach is more efficient than LEACH, CTDA MCTDA and CTPEDCA algorithms considering energy consumption. The remainder of this paper is organized as follows. Section 2 is about previous work. In Section 3, we describe our system model. The proposed algorithm is explained in Section 4. Section 5 describes the experimental results and discussions and finally Section 6 concludes the paper.

## 2.PREVIOUS WORK

Energy efficiency is one of the key design requirements in battery-powered wireless sensor networks. One important solution for it is to minimize the number of transmitted message in the network [18].

Several works have been done on data aggregation in wireless sensor network that reduce the power consumption. Clustering in WSNs is an effective procedure to decrease the energy consumption of sensor nodes. In cluster based routing algorithms for wireless networks, LEACH is famous because it is simple and efficient. In LEACH, CH nodes are selected randomly and all the non-CH nodes are formed based on the received signal power from the CHs. In LEACH each node can become a CH, there is no pattern in electing CHs and all nodes have the same chance to be a CH, thus LEACH is not efficient. CHs are selected randomly and the energy is divided between all the nodes equally. CHs aggregate all received data from all nodes in the clusters [15].

LEACH forms clusters based on the received signal strength and use the CHs as portals to the sink. All the data processing like data fusion and aggregation are locally accomplished into the cluster. CH is selected periodically among the nodes of the cluster. LEACH forms distributed clusters, where nodes make independent decisions without any concentrated control. In LEACH, each CH has a straight communicates with the base station no matter the distance is close or not.

When the network is massive, the communication between CHs and the base station consumes much energy for the long distance transmission. In LEACH, the size of clusters can be increased if the number of CHs is reduced. This makes induced excessive delays introduced by the number of nodes in the same cluster [19, 20].

The work in [21] presents the hybrid approach for cluster-based aggregation, which adaptively selects the appropriate data aggregation function. This paper shows an improvement in energy consumption with the velocity of the target. Dynamic clustering shows better performance when velocity of the target is high.

CTDA is a hybrid cluster and tree based algorithm and is proposed for data aggregation. It employs a data aggregation mechanism in the CH to lessen the amount of data transmitted. Therefore, CTDA decreases the energy dissipation in communication. CTDA decreases data transfer volume so it enhances energy efficiency and attains the purpose of saving energy of the sensor nodes. CTDA decreases the number of nodes, which directly send data to the base station.

In WSN with constrained energy, it is inefficient for sensors to select CHs randomly. CTDA method does not perform any calculation in choosing the CHs and select CHs randomly. It is nonoptimal to selected CHs by chance because it imposes an additional burden to the network. CTDA does not consider the amount of remaining energy in the nodes and it increases the wasted energy and decreases the lifetime of the network [16].

In MCTDA method, minimum spanning tree does not do data aggregation and only data of Chs by tree structure is sent to the base station [22]. CTPEDCA uses the full distribution in hierarchical WSNs. CTPEDCA is based on clustering and Minimum Spanning Tree routing strategy for CHs and the time complexity is small. CTPEDCA can balance the energy consumption of all the nodes, particularly the CH nodes in each round and extend the lifetime of the networks. In each round, CTPEDCA allows only one CH communicate directly to the base station. In CTPEDCA, a CH with the maximum remaining energy is selected as the base, CH0. CH0 constructs a minimum spanning tree between all CHs and broadcasts tree information for all the CHs. If the number of CH is K, K-1 CHs send data only to CH0 and Ch0 transmit data to the base station. The disadvantage of this method is the network is dependent on the CH0. CH0 is placed under pressure and needs a lot of energy. If CH0 is failed, the network also failed. When the base station is too far, this method is not useful [17].

In WSN, improving the energy performance and maximizing the networking lifetime are the main challenges. For this reason a hierarchical clustering scheme, called Location Energy Spectral Cluster Algorithm (LESCA) is proposed in [23]. LESCA specifies the number of clusters in a WSN automatically. It is based on spectral classification and considers the remaining energy and some properties of nodes. LESCA uses the K-way algorithm and proposes new features of the network nodes such as average energy, distance to the base station, and distance to cluster centers in order to determine the clusters and to elect the cluster's heads of a WSN. If the clusters are not constructed in an optimal way and/or the number of the clusters is greater or less than the optimal number of clusters, the total consumed energy of the sensor network per round is increased exponentially.

## 3. MODEL DESCRIPTION AND ASSUMPTIONS

In this work, we consider a multi-hop WSN consisting of n stationary and location-aware sensor nodes, denoted by {s1,s2,...,sn}, which are distributed randomly throughout an area. The network contains the sink node denoted by s0 in a preassigned location that collects data from all sensor nodes.

In our proposed study, we consider the following assumptions:
• All the nodes know their location and nodes are distributed randomly in the experimental area.
• The base station has no energy constraint and is located at the top of the area.
• The initial number of CHs is constant and does not change over time.

The superiority of protocols is changed because there are different presumptions about the radio features, such as energy dissipation in transmitter and receiver modes. In our plan, a simple model is used for the radio energy dissipation, which is the transmitter, power amplifier, and receiver dissipates energy to run the radio electronics [24]. The distance between the transmitter and the receiver is used for the free space (d2 power loss) and the multipath fading (d4 power loss) channel models.

In general, the free space (fs) model is used when the distance is less than a threshold d0 and if more than the threshold d0, the multipath (mp) model is used [24].

Therefore, when n bit data message is transmitted over a distance d to achieve an acceptable signal, the energy expended by the radio ETX can be expressed as Eq. (1).

$$E_{TX}(n, d) = \begin{cases} n\, E_{elec} + n\, \varepsilon_{fs}\, d^2 & d \leq d_0 \\ n\, E_{elec} + n\, \varepsilon_{mp}\, d^4 & d \geq d_0 \end{cases} \quad (1)$$

where, εfs is the energy consumed by the amplifier to transmit at a shorter distance. εmp is the energy consumed by the amplifier to transmit at a longer distance. Eelec is the energy dissipated in the electronic circuit to transmit or receive the signal, which relied on agents such as the digital coding, modulation, filtering and spreading of the signal. is the radio energy consumed to receive this message, which is calculated by Eq.(2).

$$ ERX\ (n) = n * E_{elec} \qquad (2) $$

## 4.DYNAMIC FUZZY CLUSTERING (DFC) DATA AGGREGATION

This paper proposes an algorithm for data aggregation called Dynamic Fuzzy Clustering (DFC) data aggregation. DFC data aggregation uses the concepts of cluster and tree based algorithms. The main idea of the cluster based routing is to lessen the amount of data transmission via engage the data aggregation mechanism in the CH. DFC data aggregation decreases the energy dissipation and saves the residual energy of the nodes. DFC data aggregation has three steps:

• CHs selection
• Cluster construction
• Tree formation of Chs

Our proposed method is inspired from two approaches named Pareto Optimal Solutions [25-26] and Fuzzy TOPSIS. At the beginning of the network, we select CHs based on Fuzzy TOPSIS [8]. The clusters are formed based on the distance between nodes and CHs. Then, the tree is organized due to CHs situation. This process continues until the first CH dies or the CH energy gets lower than a defined threshold. In this case, we determine CHs based on Fuzzy TOPSIS again. We determine the CHs based on maximizing the amount of energy efficiency. Although the initial number of CHs is assumed constant, it can decrease clustering the performing of the algorithm.

The number of CHs is related to the several parameters, i.e., network topology, residual energy of nodes, and the relative costs of calculation versus communication. The iteration of the mentioned steps creates rounds in DFC algorithm.  In the sequel, we describe the steps of DFC data aggregation.

### 4.1.CH selection

Multi Criteria Decision Making (MCDM) techniques have been applied to quantitative decision making problems [27]. MCDM can be divided into two main categories. Multi-attribute decisionmaking (MADM) approach [24] is one of the main categories of MCDM techniques. On the other hand, multi objective decision-making (MODM) [25] is another main category in MCDM techniques.

In this paper, we use MODM (Pareto optimal technique) and MADM (fuzzy TOPSIS) for selecting Chs.

### MODM (Pareto optimal technique)

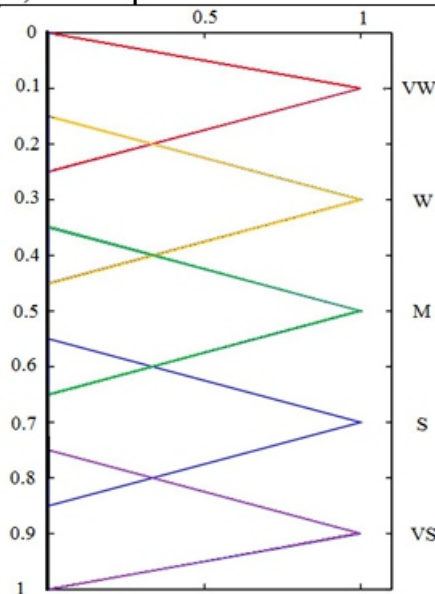The Pareto optimal solutions introduced by the economist Vilfredo Pareto [26]. Pareto technique determines the solution space which solutions are non-dominated. Pareto solution space specifies an area which comprising of all conceivable solutions in multi objective decision making problems. The solution space is classified into three groups, namely, completely dominated, neither dominant nor dominating and non-dominated.

## MADM (fuzzy TOPSIS)

It is often difficult to determine the exact values of attributes of the sensor nodes [8]. Thus, we use a fuzzy approach to determine the comparative significance of criteria instead of exact values. In this algorithm, five fuzzy linguistic variables are considered. Figure 1 illustrates the fuzzy triangular functions. The triangular membership functions are determined in Table 1.

**Table1.** Transformation of fuzzy triangular membership function

| Rank | Triangular membership function |
|---|---|
| Very Weak (VW) | (0.00, 0.10, 0.25) |
| Weak (W) | (0.15, 0.30, 0.45) |
| Moderate (M) | (0.35, 0.50, 0.65) |
| Strong (S) | (0.55, 0.70, 0.85) |
| Very Strong (VS) | (0.75, 0.90, 1.00) |



**Figure 1.** Fuzzy triangular functions

TOPSIS approach has contributed to solving the decision-making problems. In fuzzy TOPSIS approach, decision matrix has "m" alternatives and "n" attributes that could be assumed as a problem of "n" dimensional hyperplane has "m" points whose location is given by the value of their attributes [10]. i, j, respectively are   and  . The decision matrix is as the following:

$$A = \begin{bmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & x_{ij} & \vdots \\ x_{m1} & \cdots & x_{mn} \end{bmatrix} \quad (3)$$

The weight of the th column of matrix A is shown as (4):

$$C = \left[ c_1, c_2, \cdots, c_j, \cdots, c_n \right] \qquad (4)$$

where and are fuzzy numbers. We have determined 0.5, 0.25, and 0.25 weights to the remaining energy, number of neighbours, and distance from the sink s0, respectively. P is a fuzzy decision matrix, which is normalized as the follow:

$$P = [p_{ij}]_{m \times n}$$

F is the weighted, normalized fuzzy decision matrix.

$$F = \begin{bmatrix} c_1 p_{11} & c_2 p_{12} & \cdots & c_n p_{1n} \\ c_1 p_{21} & c_2 p_{22} & \cdots & c_n p_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ c_1 p_{m1} & c_2 p_{m2} & \cdots & c_n p_{mn} \end{bmatrix} \qquad (5)$$

In order to simplify the above matrix we summarize it as follows:

$$F = \begin{bmatrix} f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \cdots & f_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ f_{m1} & f_{m2} & \cdots & f_{mn} \end{bmatrix} \qquad (6)$$

The best conceivable solution is the shortest distance from the ideal solution, and the worst conceivable solution is the furthest distance from the ideal solution.

The best and the worst solutions are obtained from the weighted, normalized fuzzy decision matrix given by (6). The Best Solutions are defined as and is an acronym for the Worst Solutions:

$$BS_j = \left\{ (\max f_{ij} | i = 1, 2, \ldots, m), j = 1, 2, \ldots, n \right\} \qquad (7)$$

The worst solutions are defined as:

$$WS_j = \left\{ (\min f_{ij} | i = 1, 2, \ldots, m), j = 1, 2, \ldots, n \right\} \qquad (8)$$

We select a solution which is the nearest from the best conceivable solution and the furthest from the worst ideal solution. The distances of each alternative from the best solution and the worst solution are the separation measures. Distance of Best Solutions (DBS) and Distance of Worst Solutions (DWS) are given as:

$$DBS_i = \sum_{j=1}^{n} d\left( f_{ij}, BS_j \right) \qquad i = 1, 2, \ldots, m \qquad (9)$$
$$DWS_i = \sum_{j=1}^{n} d\left( f_{ij}, WS_j \right) \qquad i = 1, 2, \ldots, m \qquad (10)$$

Rank indices of TOPSIS are estimated as:

$$Rank_i = \frac{DBS_i}{DWS_i + DBS_i} \qquad (11)$$

Superior TOPSIS rank nodes are selected as the CHs. Each selected CH gets a unique identifier (ID).

## 4.2. Cluster Construction

All the selected CHs disseminated identity message to non-CH nodes in the network. Each node calculates the distance from all the CHs then joins to the cluster, which has the minimum distance from its CH. K specifies the number of CHs. A distance matrix is used for re-clustering nodes based on the distance to the selected CHs. The distance metric used here is the Euclidean metric. The Euclidean distance between CH and a node is relying on their situations. Consider X and Y are two nodes, i and j demonstrates two node locations. Euclidean distance is calculated based on Eq. (12):

$$d(X, Y) = \sqrt{(X_i - Y_i)^2 + (X_j - Y_j)^2} \qquad (12)$$

Each element in the distance matrix represents the difference between the CH and the node. After cluster formation, each CH is accountable for gathering the data from all the nodes in the cluster. When a framework (of data) from all the nodes in the cluster is consummated and aggregation is performed, each CH dispatches the framework to the base station. The proceeds of reclustering and data transportation is continued for R rounds until all the nodes being dead. If the number of nodes in the cluster gets smaller than the predefined threshold, the cluster is merged with the neighbouring clusters.

## 4.3. Tree formation of CHs and Data transmission

After cluster formation, the CH sends message to all non-cluster nodes in wireless sensor network which includes the CH ID, location, cluster size (for example the number of nodes in each cluster), and remaining energy. CHs also send their data and location to the base station. Base station prepares a minimum spanning tree based on the position of CH nodes so the minimum spanning tree is between CH nodes and the base station. In this plan, CHs use free area channel model to send data to the base station. In each round, the minimum distance from a vertex to another vertex is chosen based on the location of CH nodes in the tree. Combining data from several sensors used for removing the redundant transmission. Non CH nodes send their data by the framework to the CH while they are in transmission mode, so data transmission is broken into frameworks. Nodes could dispatch their data without any collision in the network. In this research, we assumed that nodes are all the time synchronized by having the base station sent out synchronization pulses to each node. When the CH receives the data from all the non CH nodes, it performs data aggregation to produce a useful data message for sending to the base station. After aggregating data, CHs transmit their resultant data along the tree (by the minimum spanning tree between CH nodes). Finally, the base station receives the final resultant data. Non-CH nodes could leave clusters when its energy is finished. If any non-CH node leaves, the related cluster releases it. If CH node is dead or a new node is joined to the network, the CH selection algorithm should be re-run.

In this paper, we consider two versions of DFC, named DFC-1 and DFC-2. In DFC-1, a node consumes its finite energy budget during the algorithm. We consider a specific threshold in DFC2 for the CHs. When the amount of energy of a CH passes from the specified threshold, a new CH is selected. In DFC-
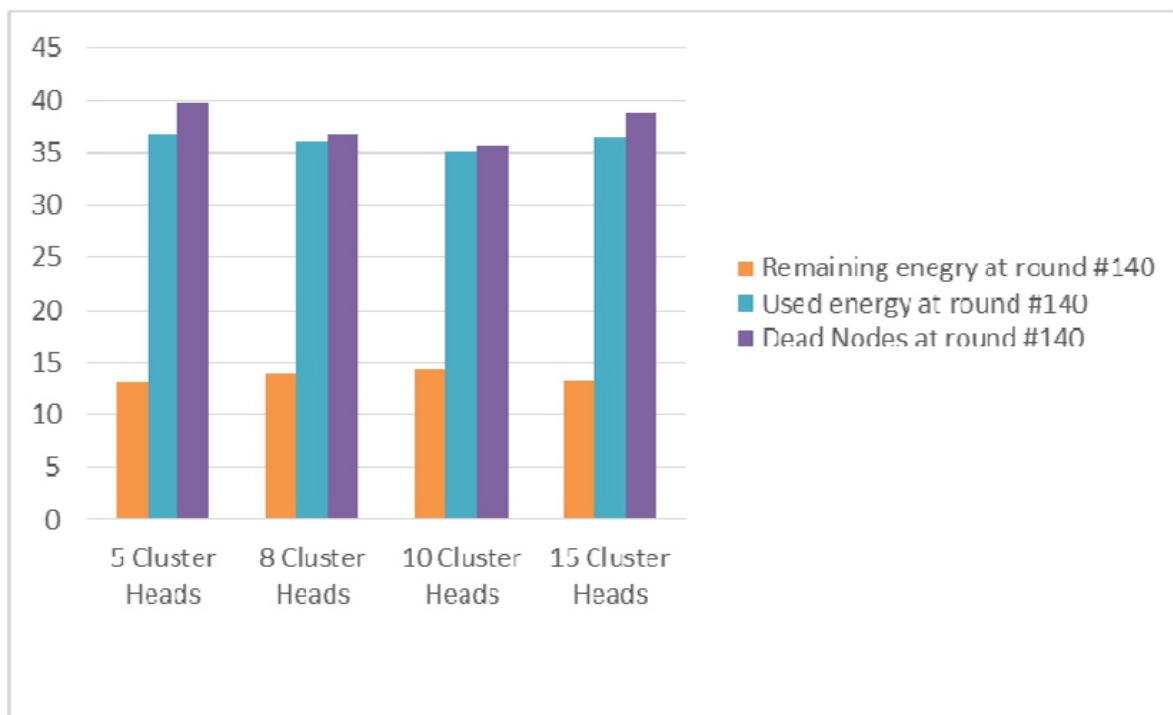
our threshold is achieved when the amount of energy of CH is reduced by half.

## 5.SIMULATION RESULTS AND ANALYSIS

In this section, we demonstrate detailed simulation experiments to evaluate the performance of the proposed algorithm. We adopt MATLAB as the platform tool, which is popularly used in the simulation experiments of wireless sensor networks [18]. In our proposed algorithm, the number of nodes is set to 100. The sink is situated far away from the area. In Cluster based approach, we consider ten CHs (K=10) in the network, which divide the nodes into ten clusters. We do an experiment in which different numbers of CHs are evaluated by the three criteria. The numbers of studied CHs are 5, 8, 10, and 15. We have compared them in remaining energy, energy consumption and the number of dead nodes. After R rounds, the most optimal CHs have more remaining energy, minimize the amount of consumed energy and the number of dead nodes. We set R=140. In Figure 2, experimental results show that K CHs are the most optimal conditions in comparison with other CHs. The selected optimal CHs have the lowest wasted energy and dead nodes, these CHs can keep more energy.



**Figure 2.** The effect of number of clusters in the DFC based on the number of dead nodes and used energy and remaining energy at round 140.

For selecting the best CHs, we have used Pareto optimal solution. Pareto optimal CHs are considered three criteria containing the remaining energy of the node, the minimum distance from the sink, and the number of adjacent nodes. Our properties of the criteria are normalized in the range [0, 1].

Membership function is applicable for converting the quantities into linguistic variables, afterward variables are converted into a fuzzy triangular membership function. We specify the fuzzy best solutions and fuzzy worst solutions. According to these quantities, we calculate the separation rate and rating indices for the selecting node. The lifetime of the network is extended in the period of the number of cycles until the first node in the network runs out of its complete energy. CHs are chosen for each node

until all the nodes expand their whole energy. In a Tree based data aggregation approach, an aggregated tree is constructed based on a minimum spanning tree which source nodes are thought out as leaves, so data are forwarded by the parent node for each node. The tree-based procedure has a low distance between each node and its parents, thereby wasted energy is diminished. Nevertheless, the depth of the tree is high. This hybrid method uses the advantages of the clustering and the tree structures while minimizing the disadvantages of them. A comparison of our proposed method against LEACH, CTDA, and CTPEDCA is represented that the present protocol is more effective than other mentioned methods in WSNs. We use four different routing protocols: LEACH protocol, CTDA algorithm, MCTDA algorithm, and our propounded method DFC algorithm for evaluating the performance of our proposed protocol with the rests. The simulation conclusions show the proposed approach has better efficiency than LEACH, CTDA, and CTPEDCA. We use a uniform simulation environment to facilitate comparison of energy savings and consume energy between protocols. Hundred sensor nodes are randomly spread in an area, which are shown in Figure 3 with coloured dots. The base station is placed far away from the area, at coordinates (50,200) which is shown in Figure 3 with an orange diamond.
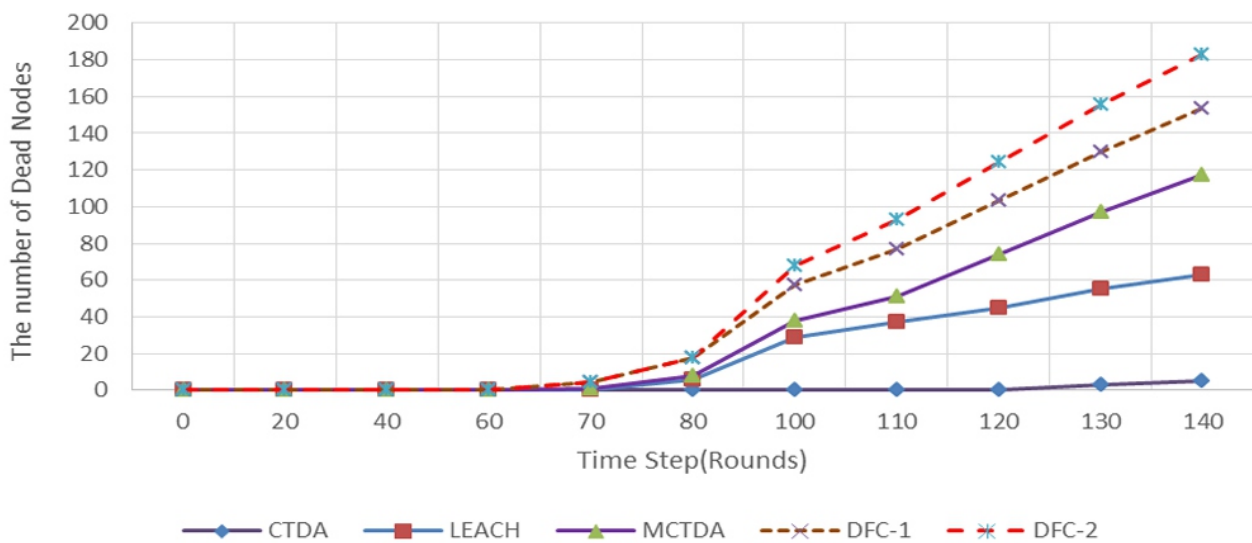


**Figure 3.** A simulated wireless sensor network with nodes and a base station at the top. In Table 2, the parameters used to develop the simulation in our experiments are listed.

**Table 2.** Simulation parameters

| Parameter | Value |
|---|---|
| Network Size | $100 \times 100$ m |
| Number of the nodes | 100 |
| $E_{elec}$ | 50nJ / bit |

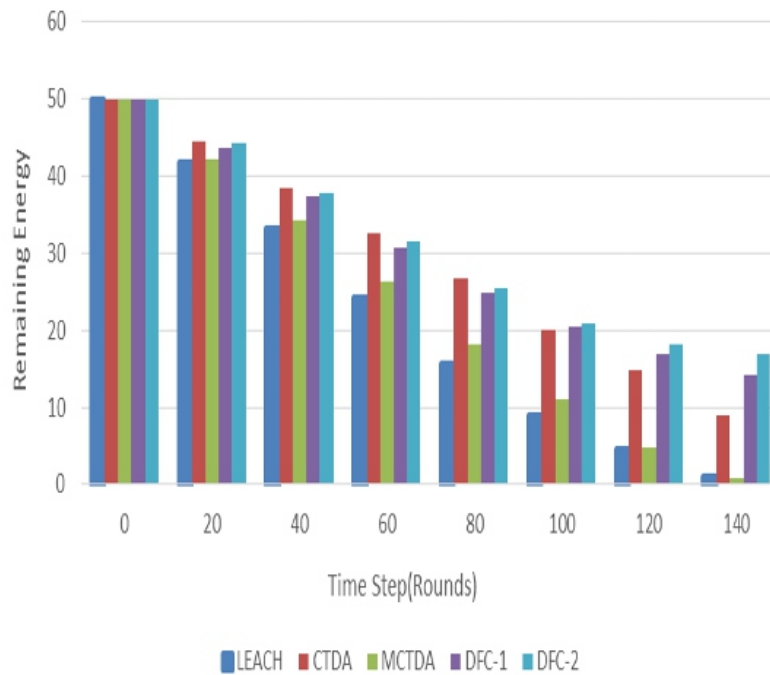| | |
|---|---|
| $\varepsilon_{fs}$ | 10 pJ/bit/m$^2$ |
| $\varepsilon_{mp}$ | 0.0013 pJ/bit/m$^4$ |
| BS location | (50,200) |
| EDA (data aggregation) | 5nJ / bit / signal |
| Control Packet size | 800 |
| Data Packet size | 4000 |
| R | 140 |
| K | 10 |

A node is considered "dead" when it spent all its energy in the transferring process and unable to send and receive the data. The simulation results of dead nodes are shown in Figure 4.



**Figure 4.** The number of dead nodes during the simulation.

Although the number of dead nodes in CTDA is low, but CTDA has many disadvantages. CTDA selects CHs randomly and it does not have any calculations to select the CHs. CTDA may select a CH with very low energy or choose a CH with the least number of neighbours. The number of dead nodes in DFC-1 and DFC-2 is less than LEACH and MCTDA. This pros is because of the CHs are calculated and elected based on three criteria the remaining energy, distance to the base station and the number of neighbours around.

According to the short distance between nodes in the proposed approach, network lifetime is increased. Furthermore, to decrease node solubility, DFC-1 and DFC-2 algorithms are more energy efficient all over the simulation. In DFC-2, we define a threshold for the amount of energy in CH, when the node's energy is less than the threshold, the new CH is replaced. The simulation results of residual energy are illustrated in Figure 5. Our findings show that the remaining energy is increased. Choosing the correct CH in the proposed method makes shorter distance between nodes. Nodes are selected as CHs, which have the largest number of neighbours. Thus, less energy are wasted so each node can hold more energy. Energy consumption of the nodes is reduced.

**Figure 5.** Remaining Energy of the nodes.

Figure 6 demonstrates the total dissipated energy by using LEACH, CTPEDCA, CTDA, and the proposed algorithm during the network simulation.

According to the results in Figure 6, CHs with the highest rank send data to the base station, which decrease the total energy consumption. In LEACH method, usually CHs attempt to send data to the base station due to the great distance through a multipath expunction channel model and consumption of energy is high. In MCTDA, accumulate data does not operate between the nodes in the tree and all data packets of CHs are sent to the base station, so energy consumption is high.

In DFC method, the amount of energy consumption is less because the CHs are selected intently. The minimum spanning tree constructed by elite CHs, for this reason DFC saves more energy than other mentioned algorithms.

**Figure 6.** Total Consumed Energy.

## 6.CONCLUSIONS AND FUTURE WORKS

A fundamental challenge in the design of WSNs is the proper utilization of resources that are scarce. In this paper, we employ Fuzzy TOPSIS method for finding the best CHs in WSNs. Three criteria contain remaining energy, distance of the nodes from the base station and the number of neighbour nodes. These criteria are discussed in order to optimize the number of CHs. The treebased method constructs a minimum spanning tree distance between CHs and the base station, which lead to decreasing energy dissipation. We proposed an energy effective algorithm in this paper, called DFC. DFC is a cluster and tree based data aggregation. Our proposed algorithm is compared with LEACH, CTDA and MCTDA protocols. The conclusions of this simulation demonstrate that the DFC is a considerable energy saving node which increase the network lifetime compared to the above- mentioned protocols.In the future, we will work on the extension of the DFC for mobility and heterogeneity of both nodes and sink.

## REFERENCES

*[1] Akyildiz, I. F. Su, W. , (2002), "Wireless Sensor Networks: A Survey.", J. Computer Networks., Vol. 38, pp. 393-422.*

*[2] Chanak, P., Banerjee, I., (2016) , "Fuzzy rule-based faulty node classification and management scheme for large scale wireless sensor networks", Expert Systems with Applications, Vol. 45, pp. 307–321.*

*[3] S. Mantri, D., Rashmi Prasad, N., Prasad, R., (2015), "Bandwidth efficient cluster-based data aggregation for Wireless Sensor Network," Computers & Electrical Engineering, Vol. 41, pp. 256264.*

*[4] Azad P, Sharma V. (2013), "Cluster Head Selection in Wireless Sensor Networks under Fuzzy Environment", ISRN Sensor Networks, pp. 1-8.*

*[5] Abbasi, A. A, Younis M. (2007) "A survey on clustering algorithms for wireless sensor networks", Computer Communications, Vol. 30, No. 14-15, pp. 2826–2841.*

*[6] Asemani M, Esnaashari M., (2015) "Learning automata based energy efficient data aggregation in wireless sensor networks", Wireless Networks. pp. 256–264.*

*[7] Selvin S, Kumar S. (2012), "Tree Based Energy Efficient and High Accuracy Data Aggregation for Wireless Sensor Networks", Procedia Engineering, Vol. 38, pp. 3833-3839.*

*[8] Baykasoğlu A, Gölcük İ. , (2015), "Development of a novel multiple-attribute decision making model via fuzzy cognitive maps and hierarchical fuzzy TOPSIS", Information Sciences, Vol. 301, pp. 75-98.*

*[9] Rathod M, Kanzaria H.(2011) , "A methodological concept for phase change material selection based on multiple criteria decision analysis with and without fuzzy environment.", Materials & Design, Vol. 32, No. 6, pp.3578-3585.*

*[10] Yang, Taho, and Chih-Ching Hung, (2007), "Multiple-Attribute Decision Making Methods For Plant Layout Design Problem", Robotics and Computer-Integrated Manufacturing , Vol. 23, No.1, pp. 126137.*

*[11] Rajagopalan, R. Varshney, P. K.(2006) "Data Aggregation Techniques In Sensor Networks: A Survey", J. IEEE Communications Surveys & Tutorials, Vol. 8, No. 4, pp. 48-63.*

*[12] Abdullah L, Najib L.(2014), "A new type-2 fuzzy set of linguistic variables for the fuzzy analytic hierarchy process", Expert Systems with Applications, Vol. 41, No.7, pp. 3297-3305.*

*[14] Sun C. (2010), "A performance evaluation model by integrating fuzzy AHP and fuzzy TOPSIS methods", Expert Systems with Applications, Vol. 37, No. 12, pp. 7745-7754.*

*[15] Akkari W, Bouhdid B, Belghith A.,(2015), "LEACH: Low Energy Adaptive Tier Clustering*

*Hierarchy", Procedia Computer Science, Vol. 52, pp. 365-372.*

*[16] Sajedi H, Saadati Z., (2014), "A Hybrid Structure for Data Aggregation in Wireless Sensor Network", Journal of Computational Engineering, Vol. 2014, pp. 1-7.*

*[17] Wang W, Wang B, Liu Z, Guo L, Xiong W.,(2011), "A Cluster-based and Tree-based Power Efficient Data Collection and Aggregation Protocol for Wireless Sensor Networks", Information Technology J., Vol. 10, No. 3, pp. 557-564.*

*[18] Cheng,H., Su, Z., Xiong, N., Xiao, Y., (2016), "Energy-efficient node scheduling algorithms for wireless sensor networks using Markov Random Field model", Information Sciences, Vol. 329, pp. 461–477.*

*[19] Richard, W. G., (2009), Extending LEACH routing algorithm for Wireless Sensor Network. Data Communications Engineering, Makerere University.*

*[20] Batra, N., Jain, A., & Dhiman, S., (2011), "An optimized energy efficient routing algorithm for wireless sensor network", International Journal of Innovative Technology and Creative Engineering., Vol. 1.*

*[21] Woo-Sung, J., Keun-Woo,L., Young-Bae, K., Sang-Joon, P., (2011), "Efficient clustering-based data aggregation in wireless sensor networks", Journal of Wireless Networks, Vol. 17, No. 5. pp. 1387400.*

*[22] Ranjani, S.; Krishnan, S.; Thangaraj, C., (2012), "Energy-Efficient Cluster Based Data Aggregation for Wireless Sensor Networks", In Proc. of International Conference on Recent Advances in Computing and Software Systems, pp. 174-179*

*[23] Jorio A, El Fkihi S, Elbhiri B, Aboutajdine D., (2015), "An Energy-Efficient Clustering Routing Algorithm Based on Geographic Position and Residual Energy for Wireless Sensor Network", Journal of Computer Networks and Communications, pp.1-11.*

*[24] Heinzelman, W. Chandrakasan, A. Balakrishnan, H., (2000), "Energy-Efficient Communication Protocol for Wireless Microsensor Networks", In Proc. of IEEE International Conference on System Sciences, pp. 1-10.*

*[25] Chauhan A, Vaish R., (2013), "Pareto Optimal Microwave Dielectric Materials", Advanced Science, Engineering and Medicine, Vol. 5, No. 2, pp. 149-155.*

*[26] Kasprzak E, Lewis K. ,( 2001) "Pareto analysis in multiobjective optimization using the collinearity theorem and scaling method", Structural and Multidisciplinary Optimization, Vol. 22, No. 3, pp. 208218.*

*[27] Zanakis S, Solomon A, Wishart N, Dublish S., (1998) "Multi-attribute decision making: A simulation comparison of select methods", European Journal of Operational Research, Vol. 107, No. 3, pp. 507529.*

## AUTHORS

**Arezoo Abasi** was born in Tehran, Iran in 1994. She studies Computer Science at University of Tehran since 2012. She was awarded Khaje Nasir price in 2010. Her interests include artificial intelligence, soft computing and wireless sensor network.

**Hedieh Sajedi** received a B.Sc. degree in Computer Engineering from Poly Technique University of Technology in 2003, and M.Sc. and Ph.D degrees in Computer Engineering (Artificial Intelligence) from Sharif University of Technology, in 2006 and 2010, respectively. She is currently an Assistant Professor at the Department of Computer Science, University of Tehran, Iran. Her research interests include Computer Networks, Machine Learning, and Signal Processing.

# SPATIO-TEMPORAL CHARACTERIZATION WITH WAVELET COHERENCE : A NEXUS BETWEEN ENVIRONMENT AND PANDEMIC

**Iftikhar U. Sikder1 and James J. Ribero2**

1Department of Information Systems, Cleveland State University, USA

2 IBA, University of Dhaka, Bangladesh

## A B S T R A C T

*Identifying spatio-temporal synchrony in a complex, interacting and oscillatory coupled-system is a challenge. In particular, the characterization of statistical relationships between environmental or biophysical variables with the multivariate data of pandemic is a difficult process because of the intrinsic variability and non-stationary nature of the time-series in space and time. This paper presents a methodology to address these issues by examining the bivariate relationship between Covid-19 and temperature time-series in the time-localized frequency domain by using Singular Value Decomposition (SVD) and continuous cross-wavelet analysis. First, the dominant spatio-temporal trends are derived by using the eigen decomposition of SVD. The Covid-19 incidence data and the temperature data of the corresponding period are transformed into significant eigen-state vectors for each spatial unit. The Morlet Wavelet transformation is performed to analyse and compare the frequency structure of the dominant trends derived by the SVD. The result provides cross-wavelet transform and wavelet coherence measures in the ranges of time period for the corresponding spatial units. Additionally, wavelet power spectrum and paired wavelet coherence statistics and phase difference are estimated. The result suggests statistically significant coherency at various frequencies providing insight into spatio-temporal dynamics. Moreover, it provides information about the complex conjugate dynamic relationships in terms phases and phase differences.*

*KEYWORDS : Wavelet analysis, Cross-wavelet power, Wavelet coherence, Covid-19, Singular Value Decomposition*

## 1. INTRODUCTION

Transformation by localized wavelike function called 'wavelet' addresses the inefficiencies of Fourier transformation by using waveforms of shorter duration at higher frequencies and waveforms of longer duration at lower frequencies [1]. Fundamentally, wavelets analyze a signal or time series according to scale where high frequency is represented by low scale and low frequency by high scale resulting into better frequency resolution for low frequency events and better time resolution for high-frequency events. Additionally, wavelets capture features across a wide range of frequencies and enables one to analyze time series that contain non stationary dynamics at many different frequencies[2]. Wavelet transformation can be done in a smooth continuous way (continuous wavelet transform - CWT) or in discrete steps (discrete wavelet transform - DWT).

Due to their versatility in handling very irregular complex data series in absence of knowing the underlying functional structure, wavelet transform analysis can be applied to analyze diverse physical phenomena e.g. climate change, financial analysis, cardiac arrhythmias, seismic signal de-noising, video image compression and so forth [3] - [4].

In this paper, we have applied the wavelet transformation to elucidate the interconnection between the environment and the Covid-19 pandemic. The dynamics between the bio-physical or climatic variables specifically the temperature and the diffusion of Covid-19 cases is reported to be ambiguous [5], [6], [7]. There are various claims with regards to the dependency between the incidence or prevalence and environmental variables. It has often been argued that lower (cold) temperature act as a catalyst in significantly increasing the spread of Covid-19 [8], [9]. There also exist alternative claims that warm temperatures slow down the spread of Covid-19 [10]. In contrast to these claims, some scholars assert that temperature does not play any role in the spread of Covid-19 [11]. In this paper, we have examined some specific empirical relationships of such dependencies, namely wavelet coherence and its statistical significance, phases and phase differences using the dataset of the USA.
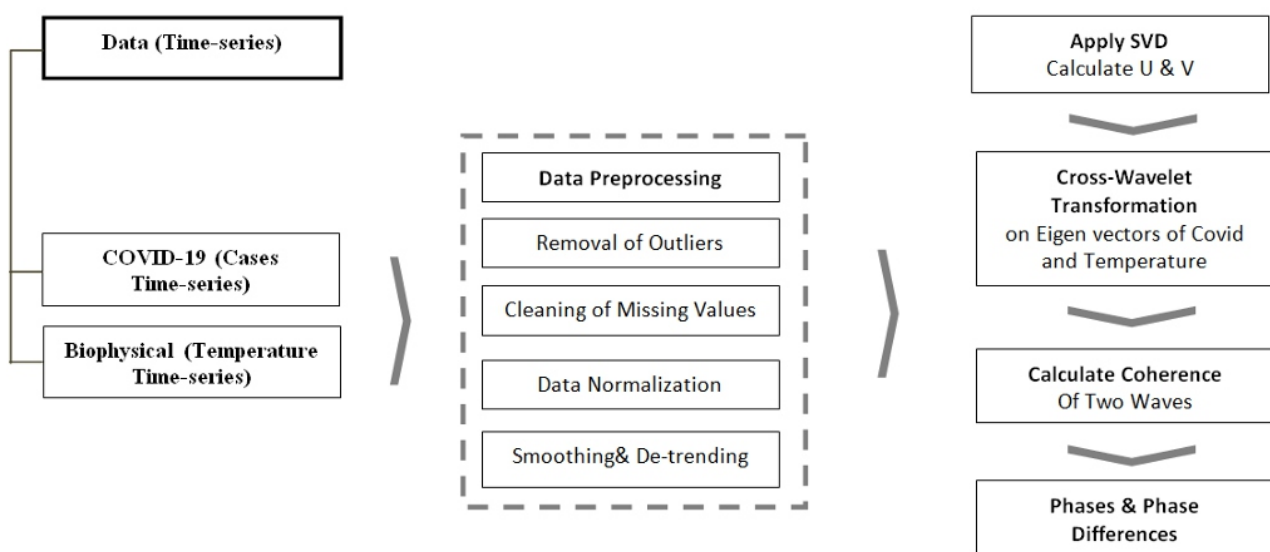
## 2. OBJECTIVES AND SCOPE

The primary objective of the paper is to characterize the dynamic relationship of Covid-19 and a time-variant bio-physical parameter namely the temperature. This paper aims to provide an empirical investigation that captures and analyzes the characteristic relationship of these variables. The study area was limited within the United States. The data for Covid-19 cases was collected from the fifty (50) states, and the corresponding data on temperature of the same period was collected from these states. The period covered was between Jan. 21, 2020, till date. Around 40,000 records (20000 Covid-19 data records and 20000 temporal temperature data records) have been collected and used for the research [12], [13].

## 3. METHODOLOGY

The variables used in the model are featured as time series data, and thus expected to fluctuate with an associated noise. Employing conventional smoothing technique involving amplitudebased statistical analysis would not be appropriate to achieve the research objective. Therefore, we adopt a Wavelet Transform algorithm not only to capture the periodicities of the variables over the time, but also to establish coherence among the variables in the frequency domain.

### 3.1. Modeling Framework

The Figure 1 below details the process flow of the research:

**Figure 1. Methodology flow chart**

### 3.1.1. Singular Value Decomposition and Wavelet Transformation

A widely adopted matrix factorization or dimension reduction technique namely Singular Value Decomposition (SVD) was used. SVD is one of the methods for matrix factorizations that generates eigen decomposition of high dimensional data. It enables low-rank approximation of a matrix. Given a data set $\mathbf{x} \in \mathbb{C}^{m \times n}$ columns $X_k \in \mathbb{C}^m$, the SVD generates a unique decomposition:

$$\mathbf{X} = U\Sigma V^* \quad (1)$$

where U is a m × m and V is n × n unitary matrices with orthogonal columns. $\Sigma$ is real, nonnegative matrix with unique diagonal entries which is known as singular values of $\mathbf{X}$ real, nonnegative matrix with unique diagonal entries which is known as singular values of X. The * denotes complex conjugate transpose. It is possible to represent the decomposition in terms of compact or economy SVD representation where $\Sigma$ compact or economy SVD representation where $\Sigma$ is square diagonal of size r × r , where rank r <= min (m, n). Thus SVD provides low-rank approximations. The approximated or truncated representation is given by the sum of rank-1 matrices:

$$X' = \sum_{k=1}^{r} \sigma_k u_k v_k^* = \sigma_1 u_1 v_1^* + \sigma_2 u_2 v_2^* + \sigma_3 u_3 v_3^* \quad (2)$$

In this study we selected $r = 3$ which accounts for significant variance of the original data matrix. The selected vectors ($v_1$, $v_2$, $v_3$) represent dominant time series signals which could be subject to further spectral analysis using wavelet analysis. Unlike Fourier analysis which permits to study the cyclical nature of time series in frequency domain without any time localization, wavelet analysis allows to characterize the periodicity of time series over time in terms of orthogonal basis by providing multi-resolution decompositions. The Wavelet function $\psi(t)$ is given by :

$$\psi_{a,b}(t) = \frac{1}{\sqrt{|a|}} \psi\left(\frac{t-b}{a}\right), a, b \in \mathbb{R}, a \neq 0 \quad (3)$$

Where $a$ is a scaling or dilation parameter and $b$ is a translation parameter. The scaling involves stretching (when $\lceil a \rceil > 1$) or compressing ($\lceil a \rceil < 1$), while translating means shifting position in time.

### 3.1.2. Data Pre-Processing

The data sets on Covid-19 cases and temperature across the 50 states accessed in a format that was not readily available for analysis. The pre-processing steps involved the basic data cleaning functions such as removal of irrelevant attributes and missing values, removal of outliers, de-meaning, linear detrending, and data normalization. The data processing was done in R environment using *WaveletComp* package [14]. The outputs of the pre-processing are data frames that were transformed into rectangular matrices, where the rows represent either Covid-19 cases by states or temperature, and the columns represent the date. The continuous wavelet transformation of a time series $x(t)$ with respect to wavelet function $\psi$ is :

$$W_{x,\psi}(a,b) = \int_{-\infty}^{\infty} x(t) \frac{1}{\sqrt{|a|}} \psi\left(\frac{t-b}{a}\right) dt \quad (4)$$

Thus, continuous wavelet transformation $W_{x,\psi}(a,b)$ is a convolution function that provides time localized frequency information. The wavelet power spectrum (scalogram) is simply given by :$|W_{x,\psi}(a,b)|^2$. For example, given a synthetic time series $y(t)$,

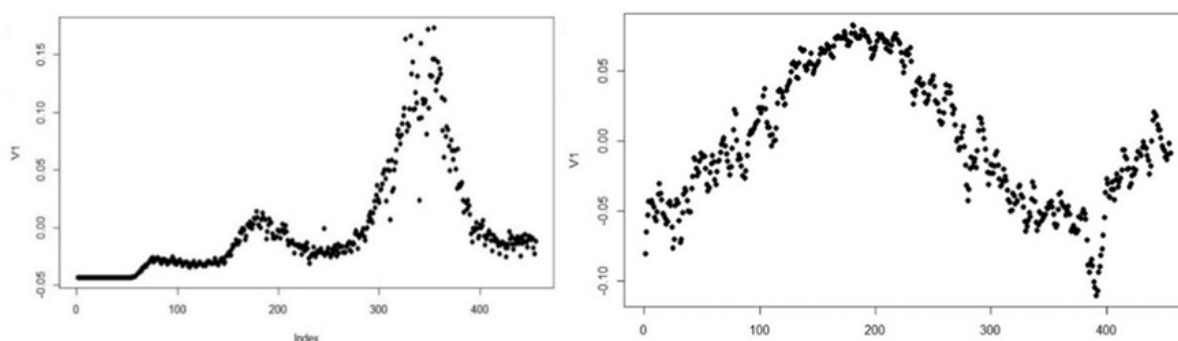$$y = \sin\left(\frac{2\pi}{p_1}\right) + \sin\left(\frac{2\pi}{p_2}\right) + \epsilon \quad (5)$$

which consists of composite of two sinusoids with period $p_1$ and $p_2$ and additive Gaussian noise $\epsilon$; for specified periods of $p_1 = 50$, $p_2 = 100$, the time series exhibit characteristic periodicities as shown in Figure 3a. A wavelet transformation of $y(t)$ and subsequent power spectrum clearly identifies two distinct regions of high wavelet power in red color in the scalogram (see figure 3c). The transformation result depends on the choice of wavelet function $\psi(t)$. In this case, we have selected *Morlet Wavelets* function, which is family of functions given by the equation:

$$\psi_{\omega_0}(t) = K e^{-i\omega_0 t} e^{\frac{-t^2}{2}} \quad (6)$$

Where, $K$ is the normalizing constant and with the value $K = \pi^{-1/4}$ ensures unit energy of $\psi_{\omega_0}(t)$. Essentially, Morlet Wavelet is the product of a complex time series with a Gaussian function normalized by a constant. Figure 3b shows the Morlet function in real axis as well as in imaginary axis viewed from different orientation.

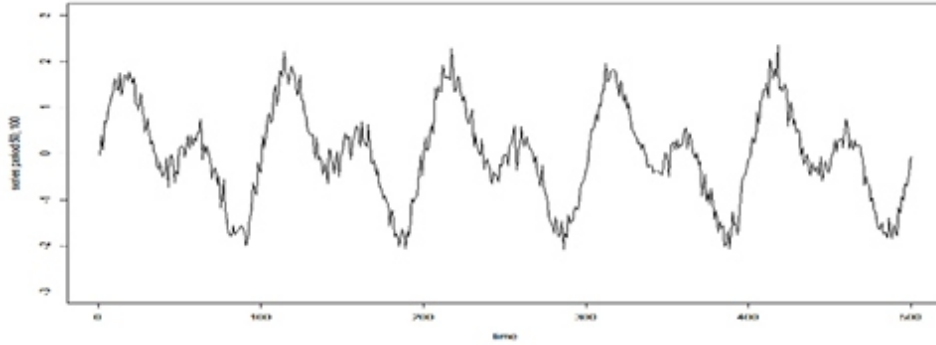### 3.1.3. Calculating Singular Value Decomposition (SVD)
The SVD technique was applied on both the Covid-19 and temperature data sets to compress the data into ortho normal eigen basis to rectangular matrices.Based on the top singular values, top three eigen states for both Covid-19 and temperature were selected, which in combination accounts for significant total variance of the original data. Plots of the transformed data set are displayed in Figure 2 below showing negative correlation (with correlation coefficientof-0.34) of cases and temperature.



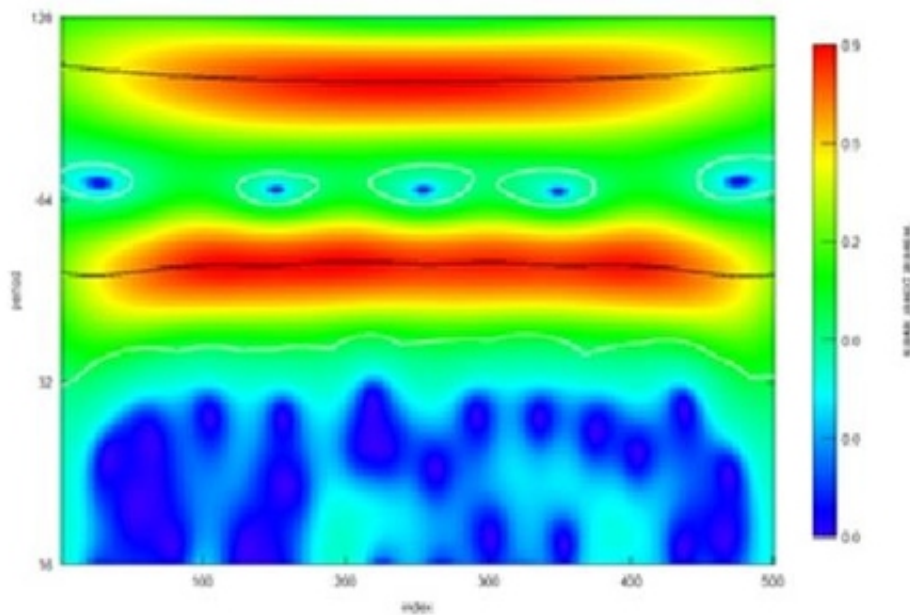**Figure 2.** Significant eigen vectors of temperature and Covid-19 cases

### 3.1.4. Calculating Wavelet Transformation
A Wavelet transform decomposes a time series into a set of wavelets localized in time. The Wavelet transformation was performed on both Covid-19 and temperature time – series of eigen vectors. We applied Morlet continuous Wavelet transforms to the transformed data using the WaveletComp R package (Figure 3).

**Figure 3a.** Synthetic time series with period p1= 50 and p2 =100 with additive Gaussian noise



**Figure 3b.** Wavelet power spectrum (Scalogram) of Synthetic time series with period p1= 50 and p2 =100 with additive Gaussian noise

Wavelet transformation leads to a continuous, complex-valued output of the time series that preserves both time and frequency resolution parameters. The transform is separable into its real part and imaginary part providing information on both local amplitude and instantaneous phase. The separation allows for the investigation of coherency between the two time series. Given two time series $X(t)$ and $Y(t)$, and corresponding *wavelet spectrums* $W_x(a,b)$ and $W_y(a,b)$ which could be considered as localized energy spectrum varying with scale $a$, and translation $b$, and associated frequency $\omega$ and time $t$. The cross-wavelet transformation $W_{xy}(a,b)$ is associated with complex-valued wavelet coherency [15], [16]:

$$\Upsilon(a,b) = \frac{\langle W_{xy}(a,b) \rangle}{\sqrt{\langle W_x(a,b) \rangle * \langle W_y(a,b) \rangle}} \qquad (7)$$

and the normalizing wavelet power spectra coherence is:

$$\Upsilon^2(a, b) = \frac{[\langle Re(W_{xy}(a,b)) \rangle]^2 + [\langle Im(W_{xy}(a,b)) \rangle]^2}{\langle W_x(a,b) \rangle * \langle W_x(a,b) \rangle} \qquad (8)$$

The angle brackets $\langle \; \rangle$ denotes the smoothing operation over time. The squaring of the amplitude component gives us the wavelet power spectrum $0 \leq \Upsilon^2(a, b) \leq 1$, which is somewhat analogous to conventional correlation coefficient.
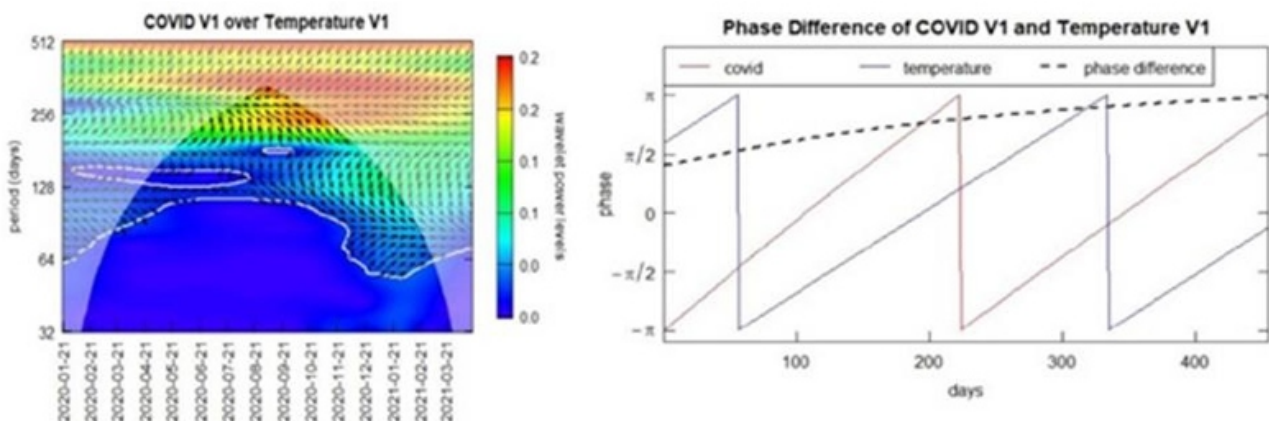
The corresponding wavelet phase difference is given by:

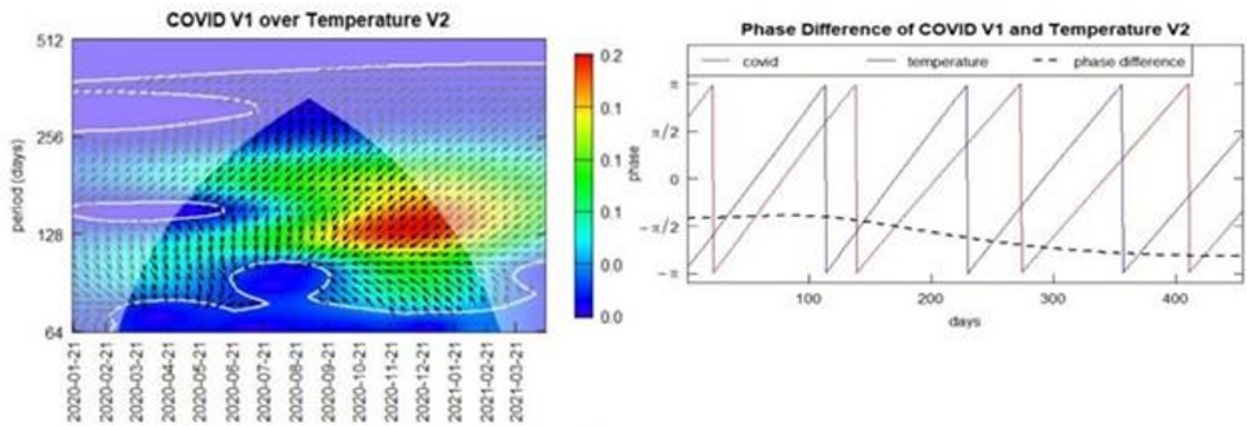$$\varphi_{xy} = \tan^{-1} \frac{Im \langle W_{xy} \rangle}{Re \langle W_{xy} \rangle} \qquad (9)$$

After computing the wavelet power spectrum for each of eigenvector, we analyze the coherence of the paired waves of Covid-19 and temperature using the coherence function. The phase lags between the variables were also computed. The cross-wavelet transformation provided crossmagnitude, phase differences, non-stationarity, and coherency between signals. Using these results of the cross-wavelet transformation, a series 'synchronicity at certain periods and across certain ranges of time was analyzed.

## 4. RESULTS AND MODEL INTERPRETATION

The cross-wavelet analysis generated coherence plot that shows that that there is a coherence (correlation) between Covid-19 and temperature, and these relationships are statistically significant (the region enveloped or bounded by the white line). The phases and phase differences show varied results. Figures 4a, 4b, 4c shows Covid-19 and temperature are out of phase with varying phase lags while Figure 4d and 4e shows that are in phase. Comparing the result of plots 4a and 4e, temperature were both out of phase in 4a, with temperature leading and Covid-19 lagging by 96 days, while from 4e both time series were in phase. Though temperature was leading, the lag period was much narrow (around 5 days) compared to 4a.



**Figure 4a.** Cross-wavelet analysis generated coherence plots (4a. Covid-19 V1 & Temperature V1)

**Figure 4b.** Cross-wavelet analysis generated coherence plots (Covid-19 V1 & Temperature V2)



**Figure 4c.** Cross-wavelet analysis generated coherence plots (Covid-19 V1 & Temperature V3)



**Figure 4d.** Cross-wavelet analysis generated coherence plots (Covid-19 V2 & Temperature V2)

**Figure 4e.** Cross-wavelet analysis generated coherence plots (Covid-19 V2 & Temperature V3)

## 5. CONCLUSION

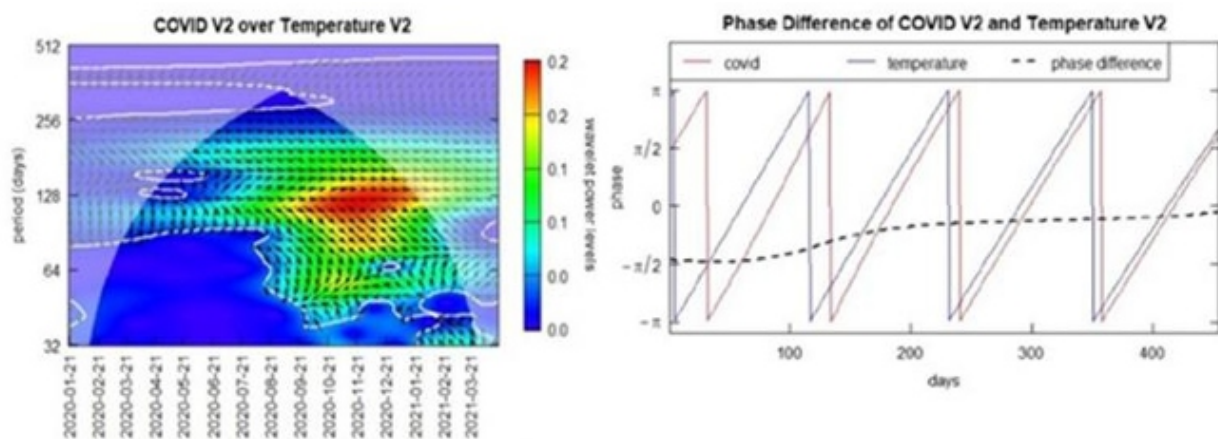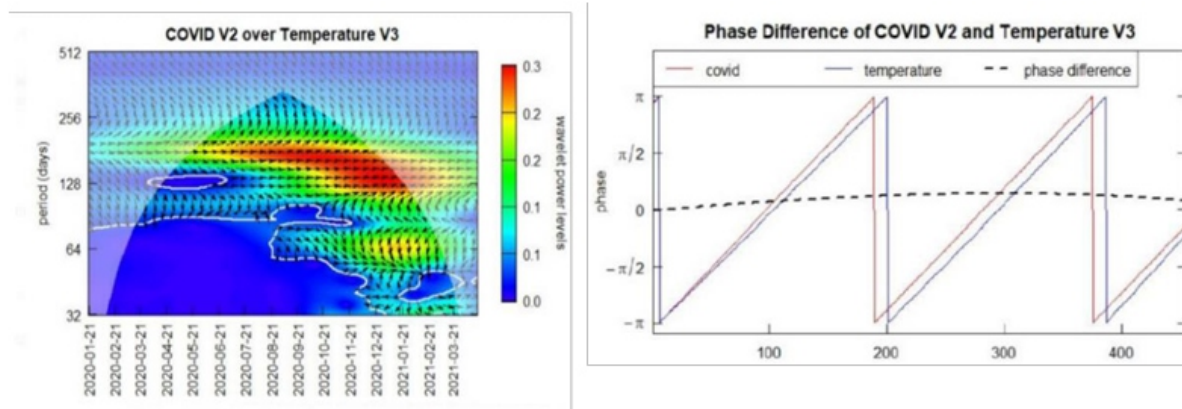The complex oscillatory interactions of the environmental variables and the incidence of pandemic make it difficult to characterize the subtle synchronization of the coupled system. In contrast with standard method, we used wavelet coherence because it is appropriate for analyzing non stationary signals. The advantage of the proposed methodology is that it does not require the stationarity assumption of time-series which is often very difficult to fulfil. In this research, the space-time dynamics between the two time-series, namely the Covid-19 and the corresponding temperature is investigated in the time-localized frequency domain using SVD and analytic Morlet wavelet transforms. The results from the continuous cross-wavelet transform shows power spectrum strengths and coherence corresponding at various frequencies (periods). The coherence statistics suggest statistically significant relationship. The results also show varying phases and phase lags with leading and lagging behavior showing complex conjugate dynamics. Future studies focusing on spatially explicit mapping of coherence and other signal processing techniques e.g. Singular Spectrum Analysis (SSA), Empirical Mode Decomposition (EMD) etc. could provide additional explanatory schemes and better understanding of the spatio-temporal dynamics of the disease.

## REFERENCES

• *A. Grossmann and J. Morlet, "Decomposition of Hardy functions into square integrable wavelets of constant shape," SIAM journal on mathematical analysis, vol. 15, pp. 723-736, 1984.*

• *I. Daubechies, "Ten lectures on wavelets (SIAM, Philadelphia, 1992)," ed: Crossref, 1992.*

• *D. L. Donoho and J. M. Johnstone, "Ideal spatial adaptation by wavelet shrinkage," biometrika, vol. 81, pp. 425-455, 1994.*

• *J. B. Ramsey, "The contribution of wavelets to the analysis of economic and financial data," Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences, vol. 357, pp. 2593-2606, 1999.*

• *P. K. Sahoo, M. A. Powell, S. Mittal, and V. Garg, "Is the transmission of novel coronavirus disease (COVID-19) weather dependent?," Journal of the Air & Waste Management Association, vol. 70, pp. 1061-1064, 2020.*

• *G. Kroumpouzos, M. Gupta, M. Jafferany, T. Lotti, R. Sadoughifar, Z. Sitkowska, et al., "COVID19: A relationship to climate and environmental conditions?," Dermatologic therapy, vol. 33, pp. e13399-e13399, 2020.*

• *D. K. Rosario, Y. S. Mutz, P. C. Bernardes, and C. A. Conte-Junior, "Relationship between COVID-19*

and weather: Case study in a tropical country," *International journal of hygiene and environmental health, vol. 229, p. 113587, 2020.*

• *M. Wang, A. Jiang, L. Gong, L. Luo, W. Guo, C. Li, et al., "Temperature significant change COVID-19 transmission in 429 cities. medRxiv 2020.02. 22.20025791," ed, 2020.*

• *D. N. Prata, W. Rodrigues, and P. H. Bermejo, "Temperature significantly changes COVID-19 transmission in (sub) tropical cities of Brazil," Science of the Total Environment, vol. 729, p. 138862, 2020.*

• *B. Oliveiros, L. Caramelo, N. C. Ferreira, and F. Caramelo, "Role of temperature and humidity in the modulation of the doubling time of COVID-19 cases," MedRxiv, 2020.*

• *T. Jamil, I. Alam, T. Gojobori, and C. M. Duarte, "No evidence for temperature-dependence of the COVID-19 epidemic," Frontiers in public health, p. 436, 2020. "COVID-19 Dataset," ed. "Temporal Temperature Dataset," ed.*

• *A. Roesch, H. Schmidbauer, and M. A. Roesch, "Package 'WaveletComp'," The Comprehensive R Archive Network2014, 2014.*

• *D. Maraun, J. Kurths, and M. Holschneider, "Nonstationary Gaussian processes in wavelet domain: Synthesis, estimation, and significance testing," Physical Review E, vol. 75, p. 016707, 2007.*

• *M. Chavez and B. Cazelles, "Detecting dynamic spatial correlation patterns with generalized wavelet coherence and non-stationary surrogate data," Scientific reports, vol. 9, pp. 1-9, 2019.*

## AUTHORS

**Iftikhar U. Sikder** is an Associate Professor jointly appointed in the department of Information Systems and the Electrical Engineering and Computer Science at Cleveland State University, USA. He holds a PhD in Computer Information Systems from the University of Maryland, Baltimore. His research interests include soft computing, granular computing, spatial databases, spatio-temporal data mining and adaptive complex systems.

His papers appeared in the journal of Knowledge-Based Systems, Risk Analysis, Expert Systems with Applications, Intl. journal of Digital Earth, International Journal of Mobile Communications, and Information Resources Management Journal, Intl. Journal of Management & Decision Making, and Intl. Journal of Aerospace Survey and Earth Sciences and many others journals. He has authored many book chapters and presented papers in many national and international conferences. Dr. Sikder is also currently serving in the editorial board of International Journal of Computational Models and Algorithms in Medicine and Intl. Journal of Computers in Clinical Practice.

**James J. Ribero** is an Adjunct Faculty at IBA, Dhaka University, Bangladesh. He holds MBBS, MS(Microbiology) and MBA from Dhaka University, Bangladesh. His research interests include applications of Machine Learning and Big Data technologies into medical and life sciences. He has authored book chapters, monographs, and presented papers in many national and international conferences. He was the former executive editor of The Orion (ISSN 1606-9722) medical journal.

# AI TESTING: ENSURING A GOOD DATA SPLIT BETWEEN DATA SETS (TRAINING AND TEST) USING K-MEANS CLUSTERING AND DECISION TREE ANALYSIS

**Kishore Sugali, Chris Sprunger and Venkata N Inukollu**

Department of Computer Science and, Purdue University, Fort Wayne, USA

## A B S T R A C T

*Artificial Intelligence and Machine Learning have been around for a long time. In recent years, there has been a surge in popularity for applications integrating AI and ML technology. As with traditional development, software testing is a critical component of a successful AI/ML application. The development methodology used in AI/ML contrasts significantly from traditional development. In light of these distinctions, various software testing challenges arise. The emphasis of this paper is on the challenge of effectively splitting the data into training and testing data sets. By applying a k-Means clustering strategy to the data set followed by a decision tree, we can significantly increase the likelihood of the training data set to represent the domain of the full dataset and thus avoid training a model that is likely to fail because it has only learned a subset of the full data domain.*

*KEYWORDS : Artificial Intelligence (AI), Machine Learning (ML), Software Testing*

## 1. INTRODUCTION

### 1.1. Overview of Artificial Intelligence and Machine Learning.

Artificial Intelligence (AI), a rapidly emerging branch of Computer Science, emphasizes on the modelling and programming of human intelligence in machines, and, to enable them to think and function like rational intelligent systems. AI can be defined as the capability of a machine to imitate intelligent human behavior. Think about this - a machine than can easily execute simple to complex tasks on a daily basis without much of human intervention. AI has made several breakthroughs in the recent years and is gaining traction for using computers to decipher otherwise complex problems, and, thus surpassing the quality of current computer systems [5]. In [6], Derek Partridge demonstrates various major classes of association that exist between artificial intelligence (AI) and software engineering (SE). These areas of communication are software support environments; AI tools and techniques in standard software; and the use of standard software technology in AI systems. Mark Kandel and Bunke, H, in [7], have also tried to correlate AI and software engineering at certain levels and discussed whether AI can be directly applied to SE problems, and if SE Processes are equipped for taking advantage of AI techniques.

### 1.2. How AI Impacts Software Testing?

Research illustrates that software testing utilizes enterprise resources and adds no functionality to the application. If regression testing discloses a new error introduced by a revision code, a new cycle of regression begins. Additionally, most software applications require engineers to write testing scripts, and their skills must be on par with the developers who initially code the app. This extra overhead expense in the quality assurance process is consistent with the growing complexity of software products.[6]

To minimize the costs, the automated testing focuses more on AI capacity, efficiency, and speed. New applications progressively provide AI functionality, sparing the challenge for human testers to comprehensively evaluate the entire product. Either data or market trends, AI will be increasingly needed to certify intelligence -containing systems, partly because the spectrum of input and output possibilities is so wide.

The intent of this paper is to focus on one of the issues and challenges that testers face in effectively testing an AI application. Currently there are various AI methods such as classification and clustering algorithms that rely primarily on monotonous data to train models to predict accurate results [8].

## 2. ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

### 2.1. Types of Machine Learning Algorithms

Machine Learning Algorithms have been divided into different categories based on their purpose.
The following are the key categories: [3]

• **Supervised Learning -** ML tries to model relationships and dependencies between the target prediction output and the input features. Input data is called training data and has a known label or result.
• Algorithms include - Nearest Neighbor, Naive Bayes, Decision Trees, Linear Regression, Support Vector Machines (SVM), Neural Networks.
• **Unsupervised Learning -** Input data are not labeled and do not have a known result. Mainly used in pattern detection and descriptive modeling. Algorithms includes - k-means clustering and association rules.
• **Semi-supervised Learning -** In the previous two categories, either there are no labels for all the observations in the dataset or labels are present for all the observations. Semisupervised learning falls in between the two. Input data is a mix of labeled and unlabeled.
• **Reinforcement Learning -** It allows machines and software agents to automatically determine the ideal behavior within a specific context, in order to maximize its performance. Algorithms includes Q-Learning, Temporal Difference (TD), and Deep Adversarial Networks.

## 3. CHALLENGE OF SPLITTING DATA INTO TRAINING AND TEST SETS

It can be a major challenge to effectively split a data set into a training data set and a testing data set. A good split leads to a productive training of the model while a poor split is more likely to lead to an inefficient model.

### 3.1. Splitting data into Training and Testing sets overview

Typically, a data scientist would use a framework for automatically splitting the available data into mutually exclusive datasets for training and testing. According to [1] one popular framework used to do this is SciKit-Learn, which allows developers to split the size of dataset by random selection. When assessing different models, or retraining models, it is important to override the random seed used to split the data. The outcomes would not be consistent, equivalent, or reproducible if not performed precisely. Typically, 70% - 80% of the data is used for training the model, with the remainder reserved for evaluation. There are numerous advanced methods available to ensure that the division of training and testing has been conducted in a representative way. When considering the coverage of model testing, it should be measured in terms of data, rather than lines of code.

The design of regression testing activities demands attention. In traditional software development the risk of functional regression is typically low unless significant changes are made. In the case of AI almost any change to the algorithm, model parameters, or training data usually needs the model to be rebuilt from scratch, and the risk of regression is very high for previously tested functionality. This is

because 100% of the model could potentially change instead of a small percentage of the model based on the necessary modifications. In synopsis:

• The way that the initial data has been gathered is important to understand whether it is representative.

• It is important that the training data is not used to test the model otherwise testing will only appear to pass.

• The data split for training and testing may need to be made reproducible.
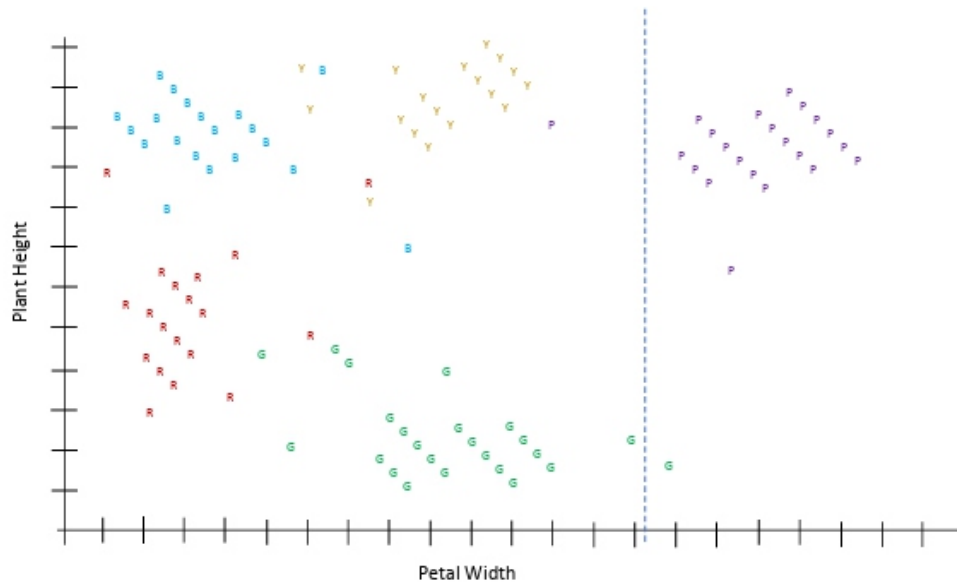
Improper splitting of datasets into training, validation, and testing sets will contribute to overfitting and underperformance in production. For instance, if a trained model expects a certain input feature, but at inference time if that feature is not passed to the, the model will fail to render a prediction. Other times however, the model will quietly crash. One of the most critical challenges for a data scientist to consider is data leakage. If one does not how to prevent data leakage, the leakage will come up often, and will destroy the models in the most subtle and dangerous ways. Particularly, leakage causes a model to look accurate and precise until one starts making decisions with the model, and then the model becomes very imprecise.

In this paper, we attempt to understand and propose a means to ensure that the data that is split between the training data set and testing data set robustly represents the domain of the full dataset. In our examples, we are assuming that there is not data leakage where the same data points appear in both the training and the testing data sets. Random splitting of the data into training and test data sets while not leaving out key parts of the total data domain is a common challenge in AI applications. In [2], we see an example in source code summarization. In attempting to generate natural language descriptions in source code, LeClair and McMillan explore the efficacy of splitting the data by method or by project. This is a very specific and narrow use case. We are seeking a methodology that is more generic and has the potential to be of use in a variety of use cases.
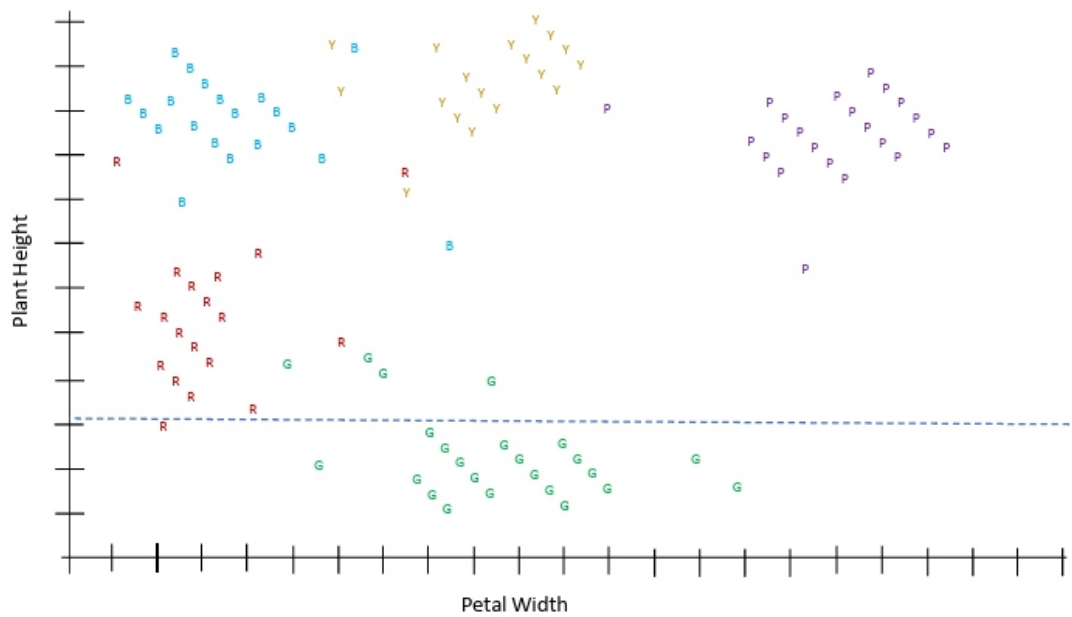
## 3.2. Example of a Poor Split

To help explain the challenge, let us consider a dataset representing flowers. To help visualize the problem, Figures 1 and 2 display a constructed collection of data. In both figures, the letter of the datapoint represents a particular type of flower. The y-axis represents plant height and the xaxis represents petal width. There may be several other attributes included in the dataset, such as leaf shape, flower color, and thorn presence. It is popular to use an 80/20 split when dividing data into training and test data sets, where we use 80 percent of the data to train the model and reserve 20 percent of the data to test the model. Both Figure 1 and Figure 2 show a splitting of the data by brute force that results in a poor split. In figure 1, very little data for the P flower variety is included in the test data set while in Figure 2, very little data for G variety of flower is seen. This poses a serious problem for training the model. There is one flower variety in either case that is highly under-represented in the test data. In training, we can expect the model to learn effectively on 4 of the 5 flower varieties, but it won't learn much about the fifth variety.

Almost all the test results, however, will be for a flower variety that the model did not learn. Thus, we would expect a prominent loss during the testing of the model and ultimately an ineffective model.

**Figure 1** Data split example



**Figure 2** Data split example

Of course, in the real world, none would split data in the above presented brute force manner. We would really like to pick data points for the training and for the test data sets at random. A straight random data split could have a complete data domain coverage such that there are no holes that will fail to learn about an important subset of the data from the model. However, it is not difficult to imagine that there is a chance of under representing 1 or 2 of the flower varieties by a strictly random split of the data. It would be ideal if we could apply a bit of intelligence to the splitting of the data while retaining an approach that still arbitrarily selects individual data points randomly.

## 3.3. k-Means Clustering

Clustering is a technique of Machine Learning technique that involves data point grouping. Theoretically, data points that are in the same group should have similar properties or characteristics, while data points from different clusters should have dissimilar properties or characteristics. In clustering, we do not have a target to predict; rather the model will understand the data and try to club similar observations and form different clusters. Hence, it is an unsupervised learning model. There are different clustering algorithms in AI that are available, and each algorithm has its own purpose.

We will be concentrating on K-Means clustering in our paper. K-means is one of the simplest algorithms for unsupervised learning that follows a simple and uncomplicated way of classifying a data set through a variety of clusters. The main objective is to define k centers, one for each cluster. As per [14], the algorithm is comprised of the following steps:

1) Let $X = \{x1, x2, x3,…,xn\}$ represent the set of data points, and ,$V=\{v1,v2,v3…,vn\}$ the set of centers.
2) Randomly select 'c' cluster centers.
3) Calculate the distance between each data point and cluster centers.
4) Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers.
5) Recalculate the new cluster center using:

$$v_i = (1/c_i) \sum_{j=1}^{c_i} x_i$$

Where, 'ci' represents the number of data points in ith cluster.
6) Recalculate the distance between each data point and newly obtained cluster centers.
7) If no data point was reassigned then stop, else, repeat from step 3.

The algorithm is significantly sensitive to the initially selected random cluster centers. The kmeans algorithm can be run multiple times to reduce this effect. The results depend on the value of k and there is no optimal way to describe a best "k".

## 3.4. Decision Tree

Decision tree is a machine learning prediction technique. Decision tree builds by repeatedly splitting data into smaller and smaller samples, Decision trees are trained by passing data down from a root node to leaves. The data is then repeatedly split according to predictor variables so that the child nodes are more "pure" or identical in terms of the outcome variables. One of the predictor variables will be chosen to make the root split. This creates a leaf node, which will further split into child nodes. All the leaves either contain only one class of outcome or they are too small to split further. At every node, a set of possible split points is identified for every predictor variable. The algorithm calculates the improvement in purity of the data that would be created by each split point of each variable. The split with the greatest improvement is chosen to partition the data and create child nodes. Calculating the improvement for a split [15], when the outcome is numeric, the relevant improvement is the difference in the sum of squared errors between the node and its child nodes after the split. For any node, the squared error is:

$$\sum_{i=1}^{n} (y_i - c)^2$$

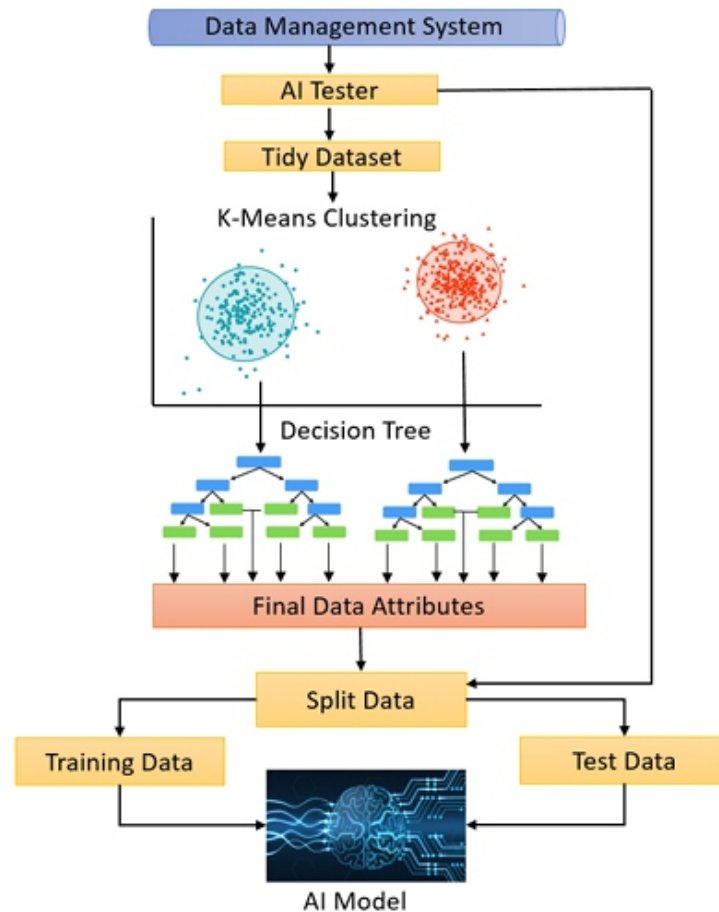$$Gini\ impurity = \sum_{i=1}^{k} p_i(1 - p_i) \qquad cross - entropy = -\sum_{i=1}^{k} p_i log(p_i)$$

Where n is the number of cases at that node, c is the average outcome of all cases at that node, and yi is the outcome value of the ith case. If all the yi are close to c, then the error is low. A good clean split will create two nodes, both which have all case outcomes close to the average outcome of all cases at that node. When the outcome is categorical, the split may be based on either the improvement of Gini impurity or cross-entropy:

where k is the number of classes and pi is the proportion of cases belonging to class i. These two measures give similar results and are minimal when the probability of class membership is close to zero or one. Let us consider the class's red and blue with sample data points from the example above and calculate the Gini impurity as shown below:

| Node | Class: red | Class: blue | Gini |
|------|------------|-------------|------|
| count | 10 | 5 | |
| proportion | 67% | 33% | |
| p (1 - p) | 0.222 | 0.222 | 0.444 |

| Child 1 | Class: red | Class: blue | Gini |
|---------|------------|-------------|------|
| count | 7 | 1 | |
| proportion | 88% | 13% | |
| p (1 - p) | 0.109 | 0.109 | 0.219 |

| Child 2 | Class: red | Class: blue | Gini |
|---------|------------|-------------|------|
| count | 3 | 4 | |
| proportion | 43% | 57% | |
| p (1 - p) | 0.245 | 0.245 | 0.490 |

The initial node contains 10 red and 5 blue cases and has a Gini impurity of 0.444. The child nodes have Gini impurities of 0.219 and 0.490. Their weighted sum is (0.219 * 8 + 0.490 * 7) / 15 = 0.345. Because this is lower than 0.444, the split is an improvement. Similarly, at every node, the purity will be determined and the model will decide whether further splits are needed.

### 3.5. Application of k-Means Clustering and Decision Tree for Accurate Data Split

The above architecture demonstrates the schematic overview of the main phases of our data split process. The architecture consists of five phases: Tidying phase, Clustering (k-Means clustering) phase, Decision tree phase, Data split phase and Training phase.

Data to be split for training and testing phases are stored in a data management system, which is confined to an unsupervised specific type of data. In unsupervised data, there are no output variables to predict. Input data are not labeled and do not have a known result. The first phase of our architecture is tidying the dataset, which is a crucial part of our proposal. Tidy data is a standard way of mapping the meaning of a dataset to its structure. A dataset is messy or tidy depending on how rows, columns, and tables line up with observations or data points, variables, and types. Around 80% of their time is spent by Data Scientists in cleaning, structuring and organizing the data. Tidy data is a way of structuring datasets to simplify analysis. In tidy data: Each variable must have its own column, each observation must have its own row, and each type of observational unit forms a table. [16] Messy data is any other arrangement of the data and it can be of these types:

• Column headers are values, not variable names.
• Multiple variables are stored in one column.
• Variables are stored in both rows and columns.
• Multiple types of observational units are stored in the same table.
• A single observational unit is stored in multiple tables.

There are more types of messy data not mentioned here, but they can be tidied in a similar way. In our approach, the AI tester is responsible for tidying the dataset, to perform this activity. The AI tester must have a sound knowledge of Artificial Intelligence and Data Mining. It is also important to train the tester

on how to understand the data, how to clean the data, and how to restructure messy data into a comprehensible format. There are several dataset tidying tools available, XLMiner being the robust data mining add-in for Excel that could be used to tidy the dataset. Once the dataset is tidy in nature, the AI tester can pass the dataset to the Clustering algorithm, which is the next phase of our architecture.

The K-Means clustering phase involves the grouping of observations that have similar properties or features, while data points from different clusters should have different properties or features. In clustering, we do not have a target variable to predict, rather the algorithm will understand the data and will group similar observations or characteristics to form different clusters. K-means is one of the simplest unsupervised learning algorithms that follows simple and uncomplicated methods to classify a data set into a variety of clusters. The main objective is to define k centers, one for each cluster. K value can be passed as an input parameter to the clustering algorithm to determine how many clusters we need, but if we do not pass k value, after going through many iterations, the algorithm itself will attempt to group the observations into different clusters. In our architecture, we have considered two clusters as an example of grouping all the similar observations in to two distinct clusters. After the clusters are created, the output data from each cluster will be passed as an input to the affiliated decision tree for next phase.

During the decision tree phase, each decision tree builds by repeatedly splitting the input data passed from the respective cluster into smaller and smaller samples. Decision trees are typically trained by passing data down to leaf nodes from a root node. The data splits repeatedly according to predictor variables so that child nodes are more "pure" or identical in terms of the outcome variables. All the leaves either contain only one class of outcome or are too small to split further. As mentioned in the previous section regarding the Gini impurity, the decision tree will split the nodes based on the Gini value and determine whether or not the split is necessary. Eventually, the decision tree produces different output attributes or variables which contain only one class of values per attribute with good purity.
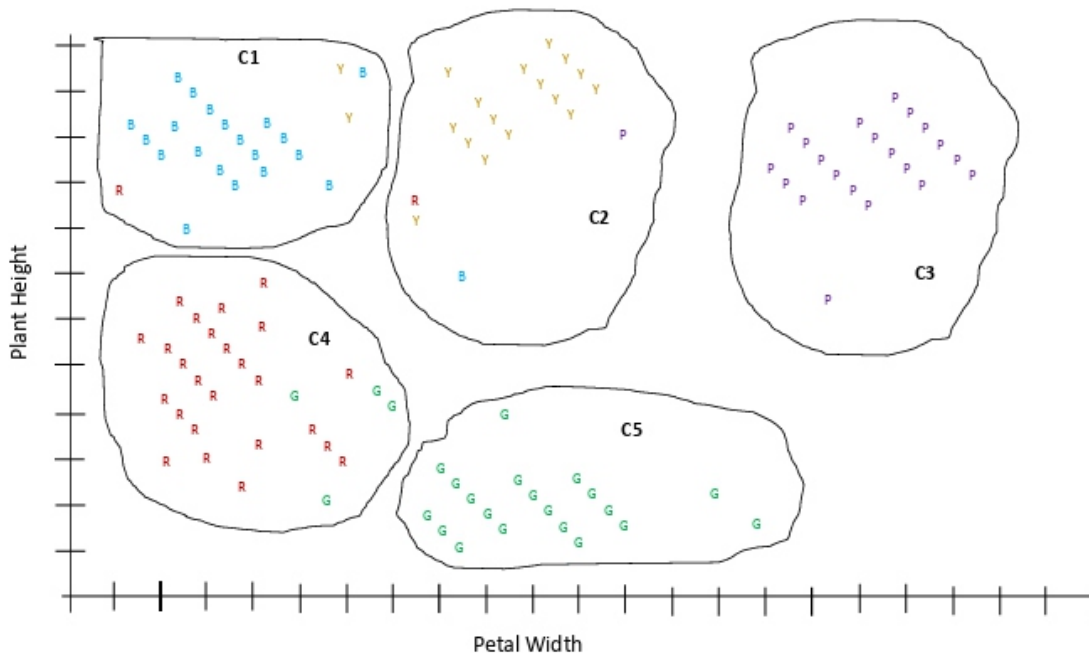
In the next data split and training phase, the AI tester collects the data attributes from each leaf of the decision tree and splits the data into training and testing sets. There is no standard percentage to decide how much to select for training and testing sets. It all depends on how much data the AI tester has available. If the AI tester has a large data set, the 75:25 training and testing ratio is okay to consider. Even if you get a very good precision after training the model with the training dataset, there is no guarantee that your trained model is a generalized one. One explanation for a non-generalized model is having a small training set, and the models appear to over-fit for small

training sets. Testing your model with distinct input combinations will therefore allow you to conclude whether your model is generalized. If you have a tiny data set, however, it is easier to go with 90:10. Another point to consider is Cross-validation, sometimes called rotation estimation, or out-of-sample testing. [17] It helps to assess how the findings of a statistical analysis will generalize to an independent data set. In the next section we will go through an example of exactly how our proposed architecture generates the accurate output.

## 3.6. Example Results

Let us now return to the example of flower data and apply the architecture to use k-means clustering and then the decision tree. To begin with as shown in Figure 4 we want the dataset to be tidy. Each column is a variable whereas each row is an observation of one flower. The AI Tester must determine how much data to use for training and how much to reserve for testing. A common selection is 80% for training and 20% for testing, but this is really at the discretion of the tester. The AI tester must also have data domain knowledge and be sufficiently familiar enough with the k-means clustering principle to select a

reasonable value for k.  In this example, we have chosen k = 5.  Figure 3 then shows a likely result of clustering the flower data into 5 clusters.
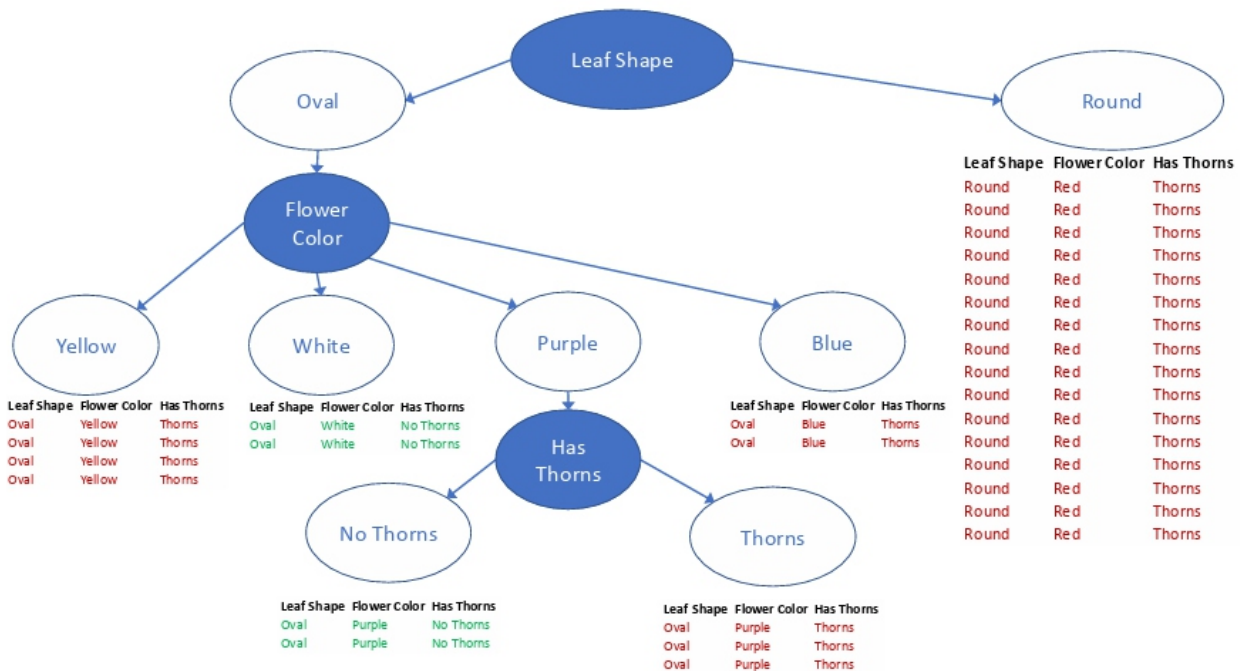


**Figure 3** Data Clustered

Next, we move on to the decision tree.  Here we begin to consider the flower data's additional key attributes.  For this example, we will look at the leaf shape, flower color, and the presence of thorns as three additional key attributes to use in the decision tree.  We must apply a decision tree to each of the five clusters.  For the sake of brevity, let us concentrate on only one of the clusters.  The focus will be Cluster 4.

| Leaf Shape | Flower Color | Has Thorns |
|------------|--------------|------------|
| Round | Red | Thorns |
| Round | Red | Thorns |
| Round | Red | Thorns |
| Oval | Yellow | Thorns |
| Round | Red | Thorns |
| Oval | Purple | Thorns |
| Round | Red | Thorns |
| Oval | Yellow | Thorns |
| Oval | Purple | No Thorns |
| Round | Red | Thorns |
| Oval | Purple | Thorns |
| Round | Red | Thorns |
| Oval | Blue | Thorns |
| Round | Red | Thorns |
| Round | Red | Thorns |
| Round | Red | Thorns |

| | | |
|---|---|---|
| Oval | Purple | No Thorns |
| Oval | Yellow | Thorns |
| Oval | White | No Thorns |
| Round | Red | Thorns |
| Round | Red | Thorns |
| Round | Red | Thorns |
| Oval | Yellow | Thorns |
| Oval | Purple | Thorns |
| Round | Red | Thorns |
| Oval | White | No Thorns |
| Round | Red | Thorns |
| Oval | Blue | Thorns |
| Round | Red | Thorns |

**Figure 4** Cluster 4 Attributes

In cluster 4, there are two varieties of flower represented. The R variety is the bulk of the data points, but there are also some data points of the G variety. Application of a decision tree to this data helps to identify key subgroups within the cluster. The end result of implementing a decision tree is shown in Figure 5.



**Figure 5** Decision Tree Applied to Cluster 4

Six key data subgroups within cluster 4 were revealed by the decision tree. At this point, the tester applies the chosen percentage split to each key data subgroup from each cluster. The effect is a test data set that covers the entire data set domain robustly and, essentially, an efficient model where the loss seen when testing the model is very similar to the loss seen when training the model.

## 3.7. Related Work

Moderate research exists in terms of software testing in AI. In [14] the key challenges of validating the quality of AI software have been summarized, which includes how to define quality assurance standard systems and develop adequate quality test coverage [9], how to connect quality assurance requirements and testing coverage criteria for AI systems based on big data, and how to use systematic methods to develop quality test models. [11] presents several new challenges that AI systems are facing in predicting system behavior such as determining the exact pre-conditions and inputs to obtain an expected result, defining expected results and verifying the accuracy of test outputs, and measuring the test coverage. Also [11] has interpreted the numerous challenges in test generation on the code, unit, module, or component levels.

Generating tests from code makes it very challenging for AI to understand the state of the software and its data and necessary dependencies. Parameters can be complex and goals and output optimizations may be unclear. The challenges that the model has in adapting itself to function more accurately in identification of gender images [12] have been clearly described. [1] Focused on the challenges faced in terms of testing facial recognition in AI systems. Data privacy is another big challenge in AI methodologies because of predictive analysis. Organizations are concerned with the transparency of their personal data [13]. Due to the disparity in training data and test data collection, overfitting is another problem facing AI models today[4]. [11] defines issues with the integration testing of the AI system in terms of transformation, cleaning, extraction, and normalization of data.

## 4. CONCLUSION

Software testing is just as critical in AI/ML applications as it is in any other type of software development. Due to the nature of how AI systems work and are developed, there are many difficulties that the software testers face with AI applications This paper attempted to lay out one potential solution to the problem of splitting data into training and testing data sets to ensure that the training data set is selected in such a way that it effectively covers the entire dataset domain. The aim is to train a successful model. The architecture in this paper sets out a methodology that increases the odds of a quality data split. It begins with a tidy data set that is the input to the kmeans clustering phase to identify natural groupings of data points with similar characteristics. Next, each cluster of data becomes an input to the decision tree phase to further break each cluster into smaller samples of related data points. Finally, on each leaf node of each decision tree, we perform the actual splitting of the data. We have provided an example of how this could work. We look forward to applying the architecture to a number of applications and working to perfect it for further study.

## REFERENCES

*[1] H.Zhuy et al., "Datamorphic Testing: A Methodology for Testing AI Applications", Technical Report*
*OBU-ECM-AFM-2018-02.*
*[2] Chen, TsongYueh, et al. "Metamorphic testing: A review of challenges and opportunities." ACM Computing Surveys (CSUR) 51.1 (2018): 1-27.*
*[3] David Fumo, https://towardsdatascience.com/types-of-machine-learning-algorithms-you-shouldknow-953a08248861.*
*[4] Salman, Shaeke, and Xiuwen Liu. "Overfitting mechanism and avoidance in deep neural networks." arXiv preprint arXiv:1901.06566 (2019).*
*[5] Partridge D. "To add AI, or not to add AI? In Proceedings of Expert Systems, the 8th Annual BCS SGES Technical conference", pages 3-13, 1988.*

[6] Partridge, D.” The relationships of AI to software engineering”, Software Engineering and AI (Artificial Intelligence), IEE Colloquium on (Digest No.087), Apr 1992.

[7] Last, M., Kandel, A., and Bunke, H. 2004 Artificial Intelligence Methods in Software Testing (Series in Machine Perception & Artifical Intelligence “Vol. 56). World Scientific Publishing Co., Inc.

[8] Manjunatha Kukkuru, “Testing Imperatives for AI Systems”, https://www.infosys.com/insights/aiautomation/testing-imperative-for-ai-systems.html.

[9] Du Bousquet, Lydie, and Masahide Nakamura. "Improving Testability of Software Systems that Include a Learning Feature." Learning 12 (2018): 13.

[10] J. Gao, et. al., “What is AI testing? and why”, 2019 IEEE International Conference on ServiceOriented System Engineering (SOSE), DOI: 10.1109/SOSE.2019.00015.

[11] Alliance For Qualification, “AI and Software Testing Foundation Syllabus”, 17 September 2019,

[12] S. Wojcik, E. Remy, “The challenges of using machine learning to identify gender in images”, 5 September 2019,

[13] Tom Simonit, "Microsoft and Google Want to Let AI Loose on Our Most Private Data", "MIT Technology Review", April 19, 2016,

[14] Sanatan Mishra, Technical Article “Unsupervised Learning and Data Clustering” 2019

[15] Jake Hoare, Data scientist https://www.displayr.com/how-is-splitting-decided-for-decision-trees/

[16] Rodrigo Mariono, https://towardsdatascience.com/whats-tidy-data-how-to-organize-messy-datasetsin-python-with-melt-and-pivotable-functions-5d52daa996c9

[17] Research Scholar Analysis review “ AI Data Split” Research Gate, 2019.

# AN INVENTORY MANAGEMENT SYSTEM FOR DETERIORATING ITEMS WITH RAMP TYPE AND QUADRATIC DEMAND: A STRUCTURAL COMPARATIVE STUDY

**Biswaranjan Mandal**

Associate Professor of Mathematics Acharya Jagadish Chandra Bose College, Kolkata

## A B S T R A C T

*The present paper deals with an inventory management system with ramp type and quadratic demand rates. A constant deterioration rate is considered into the model. In the two types models, the optimum time and total cost are derived when demand is ramp type and quadratic. A structural comparative study is demonstrated here by illustrating the model with sensitivity analysis.*

*KEYWORDS : Inventory, ramp type, quadratic, deterioration.*

## 1. INTRODUCTION

The economic order quantity model is the oldest inventory management model. An extensive research work has already been done by many researchers like Donaldson W.A.[1], Silver E.A.[2], Mandal and Pal[3]etc on inventory model assuming mostly on two types of time dependent demands – linear and exponential. Linear time dependent demand implies a steady increase or decrease in demand which is not realistic in real market. Again an exponential time varying demand is not realistic because in the real market situations, demand is unlikely to vary with a rate which is so high as exponential. Therefore we developed here two types of inventory models. First, effort has been made to analyze an inventory model for anitem that deteriorates at a constant rate, assuming the demand rate a ramp typefunction of time. Such type of demand pattern is generally seen in the case ofany new brand of consumer goods coming to the market.

The demand rate forsuch items increases with time upto a certain time and then ultimately stabilizes and becomes constant. It is believed that such type of demand rate is quiterealistic, vide, Hill [61]. Second, the quadratic time dependent demand seems to be the better representation of time-varying market demand.In this context few researchers like BiswaranjanMandal[4], M Cheng et al [ 5] and Ghosh et al [6] are noteworthy.

Finally numerical examples are done to illustrate the theory. The sensitivity of the optimal solution to change in the parameter values is examined and discussed. Also the comparative study between two types of inventory models are done along with finding its conclusion.

## 2. ASSUMPTIONS AND NOTATIONS

The mathematical models are developed under the following assumptions and notations:

(i) Replenishment size is constant and replenishment rate is infinite.

(ii) Lead time is zero.

(iii) T is the fixed length of each production cycle.

(iv) Ch is the inventory holding cost per unit per unit time.

(v) C0 is the ordering cost/order.

(vi) Cd is the cost of each deteriorated unit.

(vii) T is the cycle time.

(viii) TC is the average total cost per unit time.

(ix) A constant fraction of the on-hand inventory deteriorates per unit time. A deteriorated item is lost.

(x) Shortages are not allowed.

(xi) The demand rate R(t) is assumed in the model

(xii) There is no repair or replacement of the deteriorated items.

## 3. MATHEMATICAL MODELS

### Model 1: An inventory model with ramp type demand rate.

In this mode, the demand rate R(t) is assumed to be a ramp type function of time :

$$R(t) = D_0 [t - (t - \mu)H(t - \mu)], \; D_0 > 0$$

where $H(t - \mu)$ is the well-known Heaviside's function defined as follows :

$$H(t - \mu) = 1, \; t \geq \mu$$
$$= 0, \; t < \mu$$

Let I(t) be the on-hand inventory at any time t. The differential equations which the on-hand inventory I(t) must satisfy during the cycle time T is the following

$$\frac{dI(t)}{dt} + \theta I(t) = -R(t), 0 \leq t \leq T \tag{1}$$

In this model, we assume $\mu < t1$ and therefore the above governing equation becomes

$$\frac{dI(t)}{dt} + \theta I(t) = -D_0 t, 0 \leq t \leq \mu \tag{2}$$

and $$\frac{dI(t)}{dt} + \theta I(t) = -D_0 \mu, \mu \leq t \leq T \tag{3}$$

Solutions of the equations (2) and (3) are the following:

$$I(t) = -D_0 \left(\frac{t}{\theta} - \frac{1}{\theta^2}\right) + e^{-\theta t}\left(Q - \frac{D_0}{\theta^2}\right), 0 \leq t \leq \mu \tag{4}$$

And $$I(t) = \frac{D_0 \mu}{\theta}(e^{\theta(T-t)} - 1), \mu \leq t \leq T \tag{5}$$

From (4) and (5), we get

$$Q = \frac{D_0}{\theta^2}(1 - e^{\theta \mu}) + \frac{D_0 \mu}{\theta}(1 + e^{\theta T} - e^{\theta \mu})$$

Or $\quad Q = D_0 \mu (\dfrac{\theta}{2} T^2 + T - \dfrac{\theta}{2} \mu^2 - \dfrac{3\mu}{2})$ ( by using $e^{\theta} = 1 + \theta + \dfrac{\theta^2}{2}$ as O($\theta^3$)<<1) (6)

**Total Cost:** The total cost comprises of the sum of the setup cost, holding cost and deteriorating cost.

1. Setup Cost $\quad = \dfrac{C_o}{T}$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ (7)

2. Holding cost $= \dfrac{C_h}{T} \displaystyle\int_0^T I(t)\, dt = \dfrac{C_h}{T} [\displaystyle\int_0^{\mu} I(t)\, dt + \int_{\mu}^T I(t)\, dt]$

$$= \dfrac{C_h}{T} [\dfrac{D_0 \mu^2}{2\theta} - \dfrac{D_0 \mu T}{\theta} + \dfrac{D_0 \mu}{\theta^2} e^{\theta(T-\mu)} - (Q - \dfrac{D_0}{\theta^2}) \dfrac{1}{\theta} (e^{-\theta\mu} - 1)]$$

$$= \dfrac{C_h D_0 \mu}{T} (-\mu^2 + \dfrac{\mu^3 \theta}{4} + \dfrac{\mu^4 \theta^2}{4} - \dfrac{\mu^2 \theta T}{2} + \dfrac{T^2}{2} + \dfrac{\mu\theta T^2}{2} - \dfrac{\mu^2 \theta^2 T^2}{4}]$$

(by using $e^{\theta} = 1 + \theta + \dfrac{\theta^2}{2}$, as O($\theta^3$)<<1 ) $\qquad\qquad$ (8)

3. Deteriorating cost $= \dfrac{\theta C_d}{T} \displaystyle\int_0^T I(t)\, dt$

$$= \dfrac{C_d D_0 \mu \theta}{T} (-\mu^2 + \dfrac{\mu^3 \theta}{4} + \dfrac{\mu^4 \theta^2}{4} - \dfrac{\mu^2 \theta T}{2} + \dfrac{T^2}{2} + \dfrac{\mu\theta T^2}{2} - \dfrac{\mu^2 \theta^2 T^2}{4}]$$ $\quad$ (9)

Therefore the average total cost per unit time is given by

TC(T) = Setup cost + holding cost + deteriorating cost

$$= \dfrac{C_o}{T} + \dfrac{(C_h + \theta C_d) D_0 \mu}{T} (-\mu^2 + \dfrac{\mu^3 \theta}{4} + \dfrac{\mu^4 \theta^2}{4} - \dfrac{\mu^2 \theta T}{2} + \dfrac{T^2}{2} + \dfrac{\mu\theta T^2}{2} - \dfrac{\mu^2 \theta^2 T^2}{4}] \quad (10)$$

For minimum, the necessary condition is $\dfrac{dTC}{dT} = 0$

$$\text{Or,} \quad \dfrac{D_0 \mu}{2} (1 + \mu\theta - \dfrac{\mu^2 \theta^2}{2}) T^2 + D_0 \mu^3 (1 - \dfrac{\mu\theta}{4} - \dfrac{\mu^2 \theta^2}{4}) - \dfrac{C_0}{C_h + \theta C_d} = 0 \quad (11)$$

Which is the equation for optimum solution.

Let $T^*$ be the positive real root of the above equation (11), then $T^*$ is the optimum cycle time.

The optimum average total cost of TC(T) is TC($T^*$).

Note: If there be no deterioration i.e. $\theta = 0$, the equation (10) becomes

$$\frac{D_0\mu}{2} T^2 + D_0\mu^3 - \frac{C_0}{C_h} = 0 \text{ or } T = \sqrt{(\frac{2C_o}{D_0\mu C_h} - 2\mu^2)} \quad (12)$$

## 4. NUMERICAL EXAMPLE

For an inventory system, let the values of parameters be as follows:

$C_o$ = \$100; $C_h$ = \$10; $C_d$ = \$100; $D_0$ =100 units; $\mu$ = 0.12 year; $\theta$ = 0.01

Solving the quadratic equation (11) with the above numerical values, we find the optimum values of T as $T^*$ = 1.22 year.

The optimum values of Q, setup cost, holding cost and deteriorating cost are $Q^*$ = 12.57 units, Setup cost* = Rs 81.97, Holding cost*= Rs 71.86 and Deteriorating cost* = Rs 7.19.

The minimum average total cost per unit time is found to be $TC(T^*)$ = Rs 161.02.

## 5. SENSITIVITY ANALYSIS

We now study the effects of increasing rate of deteriorating items on the optimal average costs and cycle time. The results of this analysis are shown in the following table.

| $\theta$ | Optimum values of | | | | | |
|---|---|---|---|---|---|---|
| | T | Q | Setup cost | Holding cost | Deteriorating cost | Total cost |
| 0.01 | 1.22 | 12.57 | 81.97 | 71.86 | 7.19 | 161.02 |
| 0.02 | 1.16 | 11.92 | 86.21 | 68.26 | 13.65 | 168.12 |
| 0.03 | 1.12 | 11.50 | 89.29 | 65.87 | 19.76 | 174.92 |
| 0.04 | 1.08 | 11.08 | 92.59 | 63.48 | 25.39 | 181.46 |
| 0.05 | 1.04 | 10.96 | 96.15 | 61.07 | 30.54 | 187.76 |
| 0.06 | 1.00 | 10.19 | 100.00 | 58.65 | 35.19 | 193.84 |
| 0.07 | 0.97 | 9.87 | 103.09 | 56.85 | 39.79 | 199.73 |
| 0.08 | 0.94 | 9.54 | 106.38 | 44.05 | 44.03 | 205.45 |
| 0.09 | 0.92 | 9.33 | 108.70 | 53.84 | 48.46 | 211.00 |
| 0.10 | 0.89 | 8.99 | 112.36 | 52.01 | 52.01 | 216.38 |

(i) When the rate of deterioration θincreases, the optimum values Setup cost, Deteriorating cost and Total cost increase.

(ii) The optimum values of T,Q and Holding cost decrease with increases in the values of rate of deterioration θ.

**Model 2: An inventory model with quadratic demand rate.**

In the present model, we discussed a deterministic inventory model having a quadratic demand function with a constant deteriorating items. The demand rate function is considered as

$$R(t) = a + bt + ct^2 \text{, where a>0, } b \neq 0 \text{ , } c \neq 0 \text{ at time t and "a" is initial demand.}$$

The differential equation which the on-hand inventory I(t) must satisfies in the cycle time T is the following:

$$\frac{dI(t)}{dt} + \theta I(t) = -(a + bt + ct^2), 0 \leq t \leq T \tag{13}$$

The boundary conditions are I(0) = Q , I(T) = 0.

The solution of the equation (13) is

$$I(t) = e^{-\theta t}[\int_t^T (a + bt + ct^2) e^{\theta t} dt]$$

$$= (\frac{a + bT + cT^2}{\theta} - \frac{b + 2ct}{\theta^2} + \frac{2c}{\theta^3}) e^{\theta(T-t)} - \frac{a + bt + ct^2}{\theta} + \frac{b + 2ct}{\theta^2} - \frac{2c}{\theta^3} \tag{14}$$

Using I(0) = Q , we get

$$Q = (\frac{a + bT + cT^2}{\theta} - \frac{b}{\theta^2} + \frac{2c}{\theta^3}) e^{\theta T} - \frac{a}{\theta} + \frac{b}{\theta^2} - \frac{2c}{\theta^3}$$

$$= (a + \frac{2c}{\theta^2})T + (\frac{a\theta}{2} + \frac{b}{2} + \frac{2c}{\theta})T^2 + (\frac{b\theta}{2} + c)T^3 + \frac{c\theta}{2}T^4 \tag{15}$$

$$= (a + \frac{2c}{\theta^2})T + (\frac{a\theta}{2} + \frac{b}{2} + \frac{2c}{\theta})T^2 + (\frac{b\theta}{2} + c)T^3 + \frac{c\theta}{2}T^4 \tag{15}$$

$$\text{(by using } e^\theta = 1 + \theta + \frac{\theta^2}{2} \text{ as O}(\theta^3)<<1 )$$

**Total Cost:** The total cost comprises of the sum of the setup cost, holding cost and deteriorating cost.

1. Setup cost = $\dfrac{C_o}{T}$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ (16)

2. Holding cost = $\dfrac{C_h}{T}\int_0^T I(t)\,dt$

$$= \frac{C_h}{T} \int_0^T [(\frac{a+bT+cT^2}{\theta} - \frac{b+2cT}{\theta^2} + \frac{2c}{\theta^3})e^{\theta(T-t)} - \frac{a+bt+ct^2}{\theta} + \frac{b+2ct}{\theta^2} - \frac{2c}{\theta^3}]dt$$

$$= \frac{C_h}{T}[\frac{aT^2}{2} + \frac{bT^3}{2} + \frac{cT^4}{2}] = \frac{C_h}{2}[aT + bT^2 + cT^3] \text{ (by using } e^\theta = 1+\theta+\frac{\theta^2}{2}, \text{ as O}(\theta^3) << 1 \text{)}$$

$$\tag{17}$$

3.  Deteriorating cost $= \frac{\theta C_d}{T} \int_0^T I(t)\,dt = \frac{\theta C_d}{2}[aT + bT^2 + cT^3]$  (18)

Therefore the average total cost per unit time is given by

TC(T) = Setup cost + holding cost + deteriorating cost

$$= \frac{C_o}{T} + \frac{(C_h + \theta C_d)}{2}[aT + bT^2 + cT^3]$$  (19)

For minimum, the necessary condition is $\frac{dTC}{dT} = 0$

Or, $3cT^4 + 2bT^3 + aT^2 - \frac{2C_0}{C_h + \theta C_d} = 0$  (20)

Which is the equation for optimum solution.

## 6. Numerical Example

For an inventory system, let the values of parameters be as follows:
$C_o$ = \$100; $C_h$ = \$10; $C_d$ = \$100; $D_0$ = 100 units; $\theta$ = 0.01; a = 5; b = 3; c = 2.
Solving the fourth order equation (20) with the above numerical values, we find the optimum values of T as $T^*$ = 1.02 year.
The optimum values of Q, setup cost, holding cost and deteriorating cost are $Q^*$ = 41225 units, Setup cost* = Rs 98.04, Holding cost*= Rs 51.72 and Deteriorating cost* = Rs5.17.
The minimum average total cost per unit time is found to be $TC(T^*)$ = Rs 154.93.

## 7. SENSITIVITY ANALYSIS
We now study the effects of increasing rate of deteriorating items on the optimal average costs and cycle time. The results of this analysis are shown in the following table.

| $\theta$ | Optimum values of | | | | | |
|---|---|---|---|---|---|---|
| | T | Q | Setup cost | Holding cost | Deteriorating cost | Total cost |
| 0.01 | 1.02 | 41225.00 | 98.04 | 51.72 | 5.17 | 154.93 |
| 0.02 | 0.99 | 10104.48 | 101.01 | 49.15 | 9.83 | 159.99 |

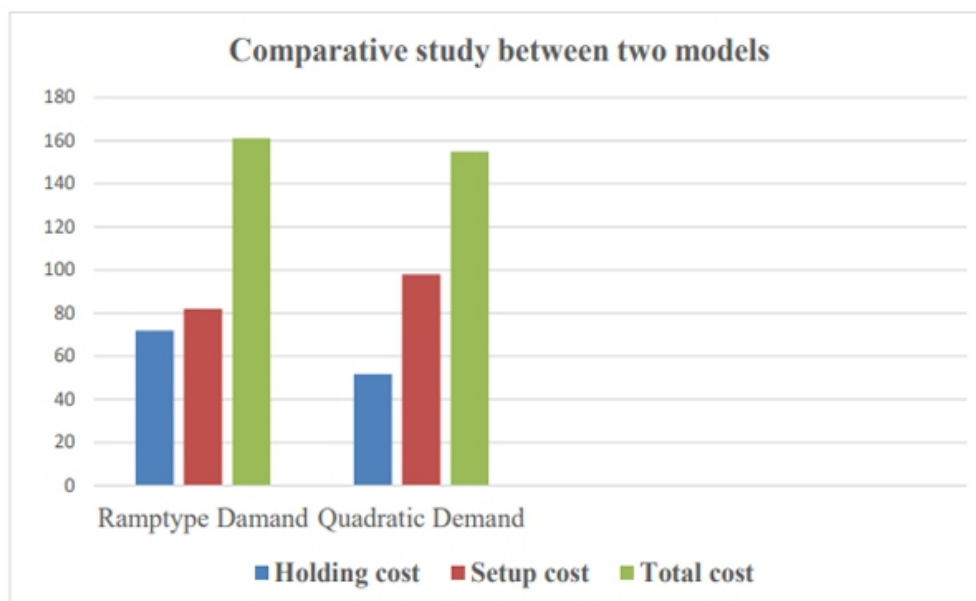| | | | | | | |
|---|---|---|---|---|---|---|
| 0.03 | 0.97 | 4444.80 | 103.10 | 47.49 | 14.25 | 164.84 |
| 0.04 | 0.94 | 2446.21 | 106.38 | 45.06 | 18.02 | 169.46 |
| 0.05 | 0.92 | 1559.88 | 108.70 | 43.48 | 21.74 | 173.92 |
| 0.06 | 0.90 | 1061.40 | 111.11 | 41.94 | 25.16 | 178.21 |
| 0.07 | 0.88 | 769.79 | 113.64 | 40.43 | 28.30 | 182.37 |
| 0.08 | 0.87 | 588.68 | 114.94 | 39.69 | 31.75 | 186.38 |
| 0.09 | 0.85 | 458.72 | 117.65 | 38.23 | 34.41 | 190.29 |
| 0.10 | 0.84 | 370.98 | 119.05 | 37.51 | 37.51 | 194.07 |

Analyzing the results given in the above table, the following observations are made:

(i) When the rate of deterioration $\theta$ increases, the optimum values Setup cost, Deteriorating cost and Total cost increase.

(ii) The optimum values of T, Q and Holding cost decrease with increases in the values of rate of deterioration $\theta$.

(iii) The similar nature is observed for both the model 1 and model 2.

## 8. COMPARATIVE STUDY BETWEEN TWO MODELS
The comparative study in numerical and graphical is performed here on the basis of the above data.

| Items | Ramp type demand | Quadratic demand |
|---|---|---|
| Holding cost | 71.86 | 51.72 |
| Setup cost | 81.97 | 98.04 |
| Total cost | 161.02 | 154.93 |

## CONCLUDING REMARKS

In this study, we have carried out two types of inventory models for deteriorating items with ramptype demand and quadratic demand in nature. The models are developed analytically as well as computationally. Numerical examples and sensitivity analysis of the solutions have been performed separately.

We have also done comparative study on holding cost, setup cost and total cost of the two types of models. It is observed from the analytical and graphical presentation that holding cost and total cost for model having quadratic demand rate are less than that of model having ramptype demand rate. On the other hand setup cost behaviour is opposite. So, holding cost and total cost in quadratic function demand are better to compare of ramptype demand and special attention is made on the inventory model having quadratic function of demand rate.

## REFERENCES

*(1) Donaldson W.A.,(1977) "Inventory replenishment policy for a linear trend in demand – An analytical solution", Operations Research Quarterly,Vol. 28, pp663-670.*

*(2) Silver E.A.,(1979) "A simple inventory decision rule for a linear trend in Demand",J. Operational Research Society, Vol. 30, pp71-75.*

*(3) Mandal B.& Pal A.K., (1988) "Order Level inventory system with ramp type demand rate", J. International Mathematics, Vol 1, No. 1, pp49-66.*

*(4) BiswaranjanMandal, (2010) "An EOQ inventory model for weibull distributed deteriorated items under ramp type demand and shortages", OPSEARCH, Indian J. of Operations Research, Vol 47, No. 2, pp158-165.*

*(5) Minghao Cheng, Bixi Zhang &Guoquing Wang, (2011) "Optimal policy for deteriorating items with trapezoidal type demand and partial backlogging", Applied Mathematical Modelling,Vol. 35, pp35523560.*

*(6) Ghosh S. K., Sarkar T. &Chaudhuri K.S., (2012) " An optimum inventory replenishment policy for a deteriorating item with time-quadratic demand and time dependent partial backlogging with shortages in all cycles", Applied Mathematics and Computation, Vol. 218, pp 9147-9155.*

# Instructions for Authors

**Essentials for Publishing in this Journal**

1  Submitted articles should not have been previously published or be currently under consideration for publication elsewhere.

2  Conference papers may only be submitted if the paper has been completely re-written (taken to mean more than 50%) and the author has cleared any necessary permission with the copyright owner if it has been previously copyrighted.

3  All our articles are refereed through a double-blind process.

4  All authors must declare they have read and agreed to the content of the submitted article and must sign a declaration correspond to the originality of the article.

**Submission Process**

All articles for this journal must be submitted using our online submissions system. http://enrichedpub.com/ . Please use the Submit Your Article link in the Author Service area.

---

**Manuscript Guidelines**

The instructions to authors about the article preparation for publication in the Manuscripts are submitted online, through the e-Ur (Electronic editing) system, developed by **Enriched Publications Pvt. Ltd**. The article should contain the abstract with keywords, introduction, body, conclusion, references and the summary in English language (without heading and subheading enumeration). The article length should not exceed 16 pages of A4 paper format.

**Title**

The title should be informative. It is in both Journal's and author's best interest to use terms suitable. For indexing and word search. If there are no such terms in the title, the author is strongly advised to add a subtitle. The title should be given in English as well. The titles precede the abstract and the summary in an appropriate language.

**Letterhead Title**

The letterhead title is given at a top of each page for easier identification of article copies in an Electronic form in particular. It contains the author's surname and first name initial .article title, journal title and collation (year, volume, and issue, first and last page). The journal and article titles can be given in a shortened form.

**Author's Name**

Full name(s) of author(s) should be used. It is advisable to give the middle initial. Names are given in their original form.

**Contact Details**

The postal address or the e-mail address of the author (usually of the first one if there are more Authors) is given in the footnote at the bottom of the first page.

**Type of Articles**

Classification of articles is a duty of the editorial staff and is of special importance. Referees and the members of the editorial staff, or section editors, can propose a category, but the editor-in-chief has the sole responsibility for their classification. Journal articles are classified as follows:

**Scientific articles:**

1. Original scientific paper (giving the previously unpublished results of the author's own research based on management methods).

2. Survey paper (giving an original, detailed and critical view of a research problem or an area to which the author has made a contribution visible through his self-citation);

3. Short or preliminary communication (original management paper of full format but of a smaller extent or of a preliminary character);

4. Scientific critique or forum (discussion on a particular scientific topic, based exclusively on management argumentation) and commentaries. Exceptionally, in particular areas, a scientific paper in the Journal can be in a form of a monograph or a critical edition of scientific data (historical, archival, lexicographic, bibliographic, data survey, etc.) which were unknown or hardly accessible for scientific research.

**Professional articles:**

1. Professional paper (contribution offering experience useful for improvement of professional practice but not necessarily based on scientific methods);

2. Informative contribution (editorial, commentary, etc.);

3. Review (of a book, software, case study, scientific event, etc.)

## Language

The article should be in English. The grammar and style of the article should be of good quality. The systematized text should be without abbreviations (except standard ones). All measurements must be in SI units. The sequence of formulae is denoted in Arabic numerals in parentheses on the right-hand side.

## Abstract and Summary

An abstract is a concise informative presentation of the article content for fast and accurate Evaluation of its relevance. It is both in the Editorial Office's and the author's best interest for an abstract to contain terms often used for indexing and article search. The abstract describes the purpose of the study and the methods, outlines the findings and state the conclusions. A 100- to 250-Word abstract should be placed between the title and the keywords with the body text to follow. Besides an abstract are advised to have a summary in English, at the end of the article, after the Reference list. The summary should be structured and long up to 1/10 of the article length (it is more extensive than the abstract).

## Keywords

Keywords are terms or phrases showing adequately the article content for indexing and search purposes. They should be allocated heaving in mind widely accepted international sources (index, dictionary or thesaurus), such as the Web of Science keyword list for science in general. The higher their usage frequency is the better. Up to 10 keywords immediately follow the abstract and the summary, in respective languages.

## Acknowledgements

The name and the number of the project or programmed within which the article was realized is given in a separate note at the bottom of the first page together with the name of the institution which financially supported the project or programmed.

## Tables and Illustrations

All the captions should be in the original language as well as in English, together with the texts in illustrations if possible. Tables are typed in the same style as the text and are denoted by numerals at the top. Photographs and drawings, placed appropriately in the text, should be clear, precise and suitable for reproduction. Drawings should be created in Word or Corel.

## Citation in the Text

Citation in the text must be uniform. When citing references in the text, use the reference number set in square brackets from the Reference list at the end of the article.

## Footnotes

Footnotes are given at the bottom of the page with the text they refer to. They can contain less relevant details, additional explanations or used sources (e.g. scientific material, manuals). They cannot replace the cited literature.
The article should be accompanied with a cover letter with the information about the author(s): surname, middle initial, first name, and citizen personal number, rank, title, e-mail address, and affiliation address, home address including municipality, phone number in the office and at home (or a mobile phone number). The cover letter should state the type of the article and tell which illustrations are original and which are not.