

**PROJECT REPORT ON
K-MEANS OUTLIER REMOVAL USING MAHALANOBIS DISTANCE
ALGORITHM**

Submitted in partial fulfilment of the requirements for
the award of the degree of

BACHELOR OF TECHNOLOGY

Submitted by

119003054

Hari priya R

Computer Science and Engineering



Under the Guidance of

**Smt Dr A Joy Christy
AP-II**

**SCHOOL OF COMPUTING
SHANMUGHA
ARTS, SCIENCE, TECHNOLOGY & RESEARCH ACADEMY
(SASTRA DEEMED TO BE UNIVERSITY)
(A University Established under section 3 of the UGC Act, 1956)
TIRUMALAISAMUDRAM
THANJAVUR – 613 401.
April(2019)**

SHANMUGHA
ARTS, SCIENCE, TECHNOLOGY & RESEARCH ACADEMY
(SASTRA DEEMED TO BE UNIVERSITY)
(A University Established under section 3 of the UGC Act, 1956)
TIRUMALAISAMUDRAM, THANJAVUR – 613401



BONAFIDE CERTIFICATE

Certified that this project work entitled **“K MEANS K-MEANS OUTLIER REMOVAL USING MAHALANOBIS DISTANCE ALGORITHM”** submitted to the Shanmugha Arts, Science, Technology & Research Academy (SASTRA Deemed to be University), Thirumalaisamudram - 613401 by Haripriya R (119003054), Computer Science and Engineering in partial fulfilment of the requirements for the award of the degree of **BACHELOR OF TECHNOLOGY** in their respective programme. This work is an original and independent work carried out under my guidance, during the period December 2018 –April 2019.

Smt Dr A Joy Christy

**ASSOCIATE DEAN
SCHOOL OF COMPUTING**

Submitted for Project Viva Voce held on_____

Examiner –I

Examiner-II

SHANMUGHA
ARTS, SCIENCE, TECHNOLOGY & RESEARCH ACADEMY
(SASTRA DEEMED TO BE UNIVERSITY)
(A University Established under section 3 of the UGC Act, 1956)
TIRUMALAISAMUDRAM, THANJAVUR – 613401



DECLARATION

I submit this project work entitled “**K MEANS K-MEANS OUTLIER REMOVAL USING MAHALANOBIS DISTANCE ALGORITHM**” to the Shanmugha Arts, Science, Technology & Research Academy (SASTRA) Deemed to be University, Tirumalaisamudram–613401, in partial fulfillment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY and declare that it is my original work carried out under the guidance of Prof. **Smt Dr A Joy Christy** School of Computing, SASTRA.

Name1: Hari Priya R

Signature1:

Reg. No.: 119003054

Date: May 5, 2019

Place: Thanjavur, India

ACKNOWLEDGEMENTS

The satisfaction that accompanies the successful completion of this project would be incomplete without the mention of the people who made it possible, without whose constant guidance and encouragement would have made efforts go in vain. We consider ourselves privileged to express my gratitude and respect towards all those guided us through the completion of this project. We cordially thank our Registrar **Dr. G. Bhalachandran** , Dean – Planning & Development **Dr. S. Vaidhyasubramaniam** and Dean – Sponsored Research – Dr. S. Swaminathan for providing us the infrastructure to carry out the project.

We would like to express my utmost gratitude to our Dean **Dr. A. Umamakeswari**, who gave us the opportunity to complete the project.

We convey thanks to our Project guide **Smt. Dr. Joy Christy A, Assistant Professor** for providing encouragement, constant support and guidance which was of great help to complete the project successfully.

We would also thank our institution and faculty members. We also extend my heartfelt thanks to my family, friends and well-wishers.

SYNOPSIS

K-MEANS OUTLIER REMOVAL USING MAHALANOBIS DISTANCE ALGORITHM

119003054 - Hari Priya R - B.Tech.- Computer Science and Engineering

Outlier detection is a part of data analytics that helps user to find discrepancies in working machine by applying outlier detection algorithm on the captured data for every fixed interval. An outlier is a data point that exhibits different properties from other points that are due to some external or internal forces. These outliers can be detected by clustering the data points. To detect outliers, optimal clustering of data points is important. Problem that arises quite frequently in statistics is identification of groups or clusters of data within a population or sample. The most widely used procedure to identify clusters in a set of observations is K-Means using Euclidean distance. However, Euclidean distance is not so efficient for finding anomaly in multivariate space. To remedy this shortfall in the K-Means algorithm the Mahalanobis distance metric was used to capture the variance structure of the clusters that is followed by the application of ESVA-LCPSD algorithm to detect the outliers. If this method serves as a significant improvement over its competitors, then it will provide a useful tool for detecting rare items, events or observations which raise suspicions by differing significantly from the majority of the data.

Keywords: Mahalanobis distance, Multivariate outlier detection, ESVA- LCPSD.

Ms Dr A Joy Christy, AP-II

Table of Contents

CHAPTER 1 : INTRODUCTON	vii
CHAPTER 2 : PROBLEM STATEMENT	x
CHAPTER 3	xi
LITERATURE REVIEW.....	xi
CHAPTER 4 : SOFTWARE AND HARDWARE SPECIFICATIONS	xiii
CHAPTER 5 : PROPOSED MODEL.....	xiv
CHAPTER 6 : METHODOLOGY AND APPROACH	xviii
CHAPTER 7 : RESULTS AND DISCUSSION	xxvii
CHAPTER 8 : CONCLUSION AND FUTURE WORK.....	xxxi
REFERENCES.....	xxxii

CHAPTER 1 : INTRODUCTION

In this Information era, it is believed that information leads to power and success. Future of many companies and government organizations relies on the information what they have with them. With the improvement in the storage techniques now it is possible to collect and store a tremendous volume of information. Organizations have been collecting an immeasurable data from simple text documents to more complex information such as Medical data, Satellite data, spatial data and multimedia data. Mining of these data, using sophisticated mathematical algorithms, provides much useful information regarding the probability of future events, unusual events that might be interesting or data errors that require further investigation.

Data mining is the process of uncovering patterns and finding anomalies and relationships in large datasets that can be used to make predictions about future trends. The main purpose of data mining is extracting valuable information from available data. It is also popularly known as Knowledge Discovery in Databases (KDD). Data Mining comprises of few steps starting from preliminary raw data collections to some form of identifying new knowledge. It is an iterative process and uses the following steps such as Data cleaning, Data integration, Data selection, Data transformation, Data mining, Pattern evaluation and Knowledge Representation. Once the extracted information is offered to the user, the assessment measures can be improved and further refined to get more fitting results.

One of the important applications of data mining is outlier detection. Outlier detection is the process of detecting and subsequently excluding inappropriate data from the given set of data. An outlier is a piece of data that deviates drastically from the standard norm or average of the data set. Outlier detection has two-steps viz., Clustering and detecting deviated data among the clustered sets. Therefore, the process of grouping observations into cluster is a foremost problem in analysing data sets. So far, the most widely used algorithm to identify clusters in a set of observations is K-Means. But, the main constraint of this algorithm is that it uses Euclidean distance metric, which is prone to noisy data and outliers, which in turn gives a non-spherical cluster. Also, this distance suite well only for univariate datasets. Hence, this paper introduces the technique of Mahalanobis distance (MD) to detect an observation having an unusual pattern. The MD measures the relative distance between two variables with respect to the mean of the multivariate data. These calculated distance values are used by Extreme Value Analysis algorithm ESVA-LCPSD to find outliers and thereby eliminating the need of deciding threshold value manually.

MD incorporates multivariate approach to find distance using covariance between the features of the dataset. Formula that uses distances from each observation to the central mean (M):

$$D = [(X_i - M)^T C^{-1} (X_i - M)]^{0.5}$$

Where

D = Mahalanobis distance

X_i = an object vector

M = arithmetic mean vector

C^{-1} = Inverse Covariance matrix of independent variables

T = Indicates vector should be transposed

Some important characteristics of MD that helps to provide better clustering are:

- It accounts for the fact that the variances in each direction are different.
- It accounts for the covariance between variables.
- It reduces to the familiar Euclidean distance for uncorrelated variables with unit variance.

The Mahalanobis distance accounts for the variance of each variable and the covariance between variables. Geometrically, it does this by transforming the data into standardized uncorrelated data and computing the ordinary Euclidean distance for the transformed data. In this way, it provides a way to measure distances that takes into account the scale of the data. The distance obtained is analysed using ESVA-LCPD, an extreme value analysis based algorithm, to get the outliers with the help of algorithm itself instead of finding them with the help of threshold that is set manually.

CHAPTER 2 : PROBLEM STATEMENT

For detecting outliers in a set of observations, it is important to cluster the points accurately. Clusters are characterized by groups of data points which are in “close” proximity to one another. While it is much easier to visually detect clusters in univariate or bivariate data, the task becomes increasingly difficult as the dimensionality of the data increases. One of the largest used clustering algorithms is K-Means using Euclidean distance. Euclidean distance suffers from a scaling effect that describes a situation where the variability of one parameter masks the variability of another parameter and it happens when the measurement ranges or scales of two parameters are different, making it difficult to find the true outliers. Hence, to overcome this shortfall, the idea of using K-Means with Mahalanobis distance (MD) is put forth. The MD methodology distinguishes multivariable datagroups by a univariate distance measure, which is calculated from the correlated values of all the parameters that the dataset is dependent on. Thus, MD value is calculated using the normalized value of observations and their correlation coefficients thus making MD more efficient. After obtaining univariate distance, **an Extreme Value Analysis based algorithm**, ESVA-LCPSD is used to detect the outliers.

CHAPTER 3

LITERATURE REVIEW

Kadam and Pund [1] reviewed several approaches to detect outliers, including the cluster-based approach. Aparna and Nair [4] proposed the CHB-K-Means algorithm by using a weighted attribute matrix to detect outliers. Jiang et al. [5] proposed two initialization methods for the k-modes algorithm to choose initial cluster centers that are not outliers. Although much work has been done on outlier analysis, few of them perform clustering and detect outliers simultaneously. In this section, we focus on clustering methods with the built-in mechanism of outlier detection and give a review of those methods.

Jiang et al. [5] proposed a two-phase clustering algorithm for outlier detection. In the first phase, the k -means algorithm is modified to partition the data in such a way that a data point is assigned to be a new cluster centre if the data point is far away from all clusters. In the second phase, a minimum spanning tree is constructed based on the cluster centers obtained from the first phase. Clusters in small sub trees are considered as outliers. He et al. [6] introduced the concept of cluster-based local outlier and designed a measure, called cluster based local outlier factor (CBLOF), to identify such outliers.

Hautamäki et al. [7] proposed the ORC (Outlier Removal Clustering) algorithm to identify clusters and outliers from a dataset simultaneously. The ORC algorithm consists of two consecutive stages: the first stage is a purely k -means algorithm; the second stage iteratively removes the data points that are far away from their cluster centroids.

Zhou et al. [9] proposed a three-stage k -means algorithm to cluster data and detect outliers. In the first stage, the fuzzy c -means algorithm is applied to cluster the data. In the second stage, local outliers are identified and the cluster centers are recalculated. In the third stage, certain clusters are merged and global outliers are identified.

Ahmed and Naser [11] proposed the ODC (Outlier Detection and Clustering) algorithm to detect outliers. The ODC algorithm is a modified version of the k -means algorithm. In the ODC algorithm, a data point that is at least p times the average distance away from its centroid is considered as an outlier. Whang et al. [2] proposed the NEO- k -means (Non-exhaustive Overlapping k -means) algorithm, which is also able to identify outliers during the clustering process.

Guojun Gan and K. Michael [14] proposed k -Means Outlier Removal Algorithm that extends k -means algorithm to provide data clustering and outlier detection simultaneously by introducing an additional “cluster” to the algorithm to hold all outliers. Unlike most existing clustering algorithms with outlier detection, the KMOR algorithm assigns all outliers into a group naturally during the clustering process.

Jayakumar and Thomas [10] proposed an approach to detect outliers based on the Mahalanobis distance. Sachin Kumar [13] had employed a probabilistic approach for defining warning and fault threshold MD values to improve upon the traditional approaches where threshold MD values are decided by experts and demonstrates that the approach to define threshold MD value is a major improvement. Igor Melnykov [3] proposed a new algorithm to extend k -means algorithm with Mahalanobis distance to take into account of variance structure of the clusters.

CHAPTER 4 : SOFTWARE AND HARDWARE SPECIFICATIONS

This project has been implemented using Python 3.7.1 Compiler. Python provides many built-in packages for data analytics application and thereby making it easy to learn and code. Moreover, it has built-in packages for illustrating results in much clearer formats like graphs, charts, scatter-plots. Thus, python is considered for implementing this project.

Packages that are used in this project are Numpy, Sklearn, Matplotlib, mpl_toolkits. Numpy is used to handle very large datasets. Sklearn package contains cluster function that can be used to cluster the data. Matplotlib is used for graphical representation of datasets and how they are clustered using graphs, scatter-plots and others. Thus these packages should be pre-installed before implementing.

Main hardware requirement of the project is the type of Processor, Memory, RAM. 32-bit or 64-bit Operating System, having x32 or x64-based processor. Processor used is Intel(R) Core™ i3-6100U CPU @2.30GHz 2.30GHz.

CHAPTER 5 : PROPOSED MODEL

Aim of the proposed model is to find the outliers in the dataset more accurately. After clustering the points using k -Means with Mahalanobis distance algorithm, final step is to find the outliers based on the extreme value Analysis algorithm that takes the threshold for deciding outlier as a part of algorithm rather than fixing it manually. The steps to be followed to find the outliers are shown in the fig 5.1.

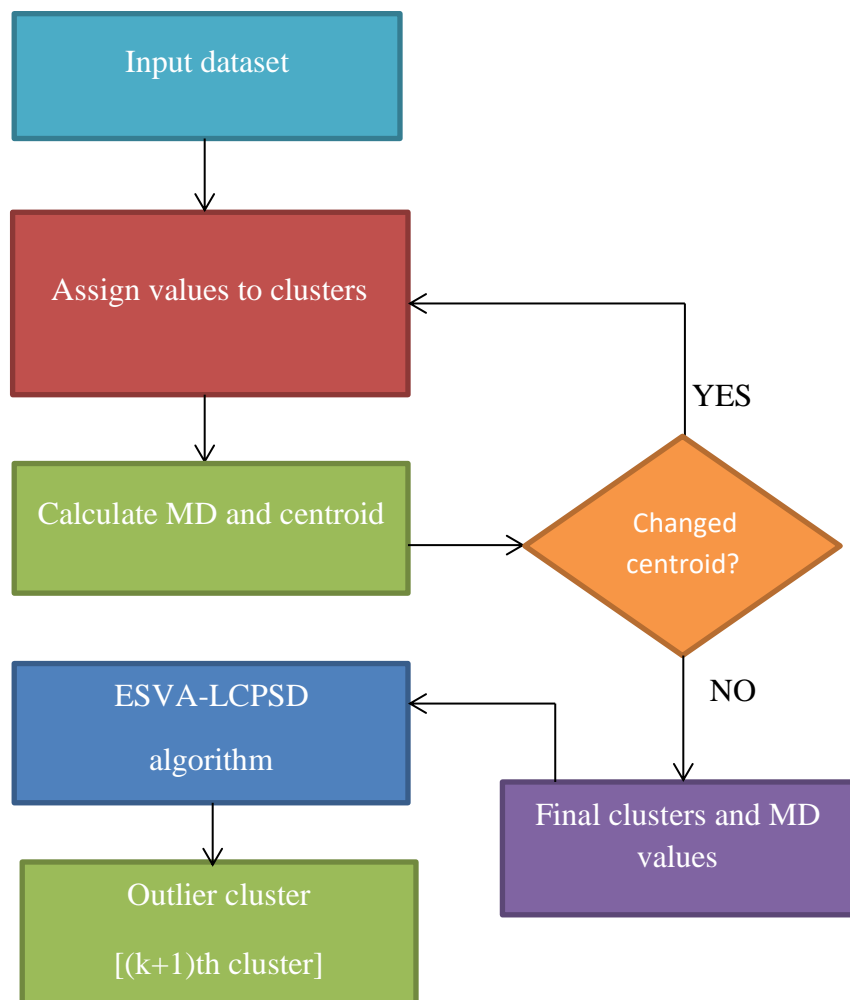


Fig 5.1 Steps to find outliers using k-means with Mahalanobis distance based clustering and ESVA-LCPSD algorithm

Final MD values are given as input for calculating threshold using Extreme value analysis algorithm, ESVA-LCPSD, to detect the outliers. Thus, this algorithm has three steps:

- K-means clustering
- Distance calculation using Mahalanobis Distance metric
- ESVA-LCPSD, Extreme value analysis to set threshold.

Brief of each step is depicted as follows. The K-Means clustering is portrayed in Fig 5.2.

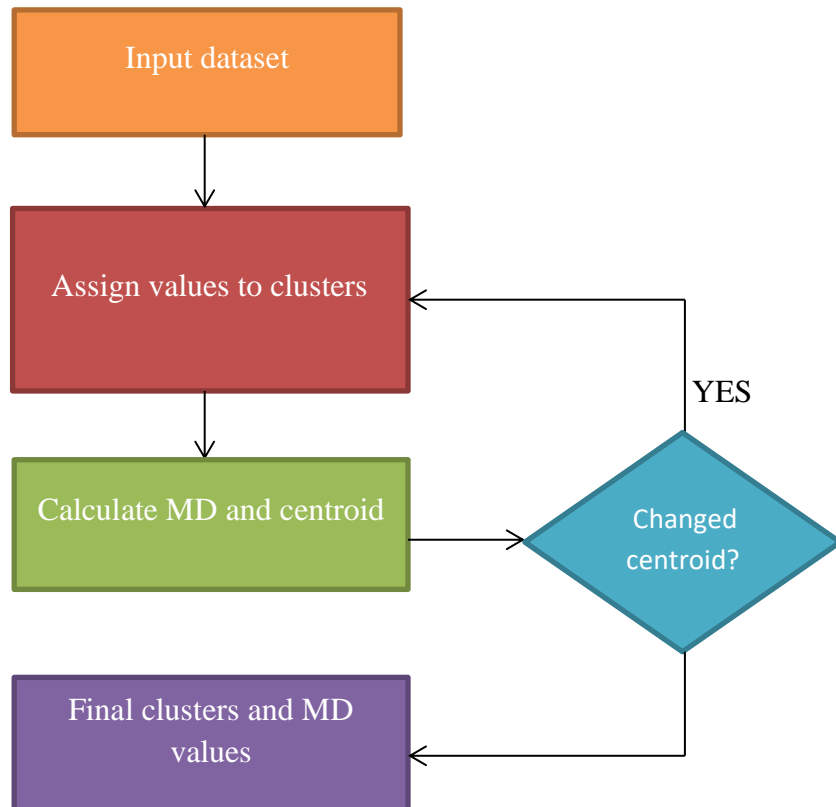


Fig 5.2 Algorithm of Modified k-means Clustering using Mahalanobis Distance

As shown in the figure K-Means algorithm takes ‘n’ points as input and clusters these ‘n’ points into each of the ‘k’ clusters as much as possible. If any of the point cannot be clustered, then it is taken as an outlier. As an extension of k-means algorithm, this algorithm maintains (k+1) clusters where k clusters are to hold the clustered points while (k+1)th cluster holds the outlier points. After, choosing the cluster centres at random, MD values and mean

for each cluster is recalculated at each step as a part of every iteration until the optimal cluster is obtained. Algorithm to calculate MD is shown in Fig 5.2.

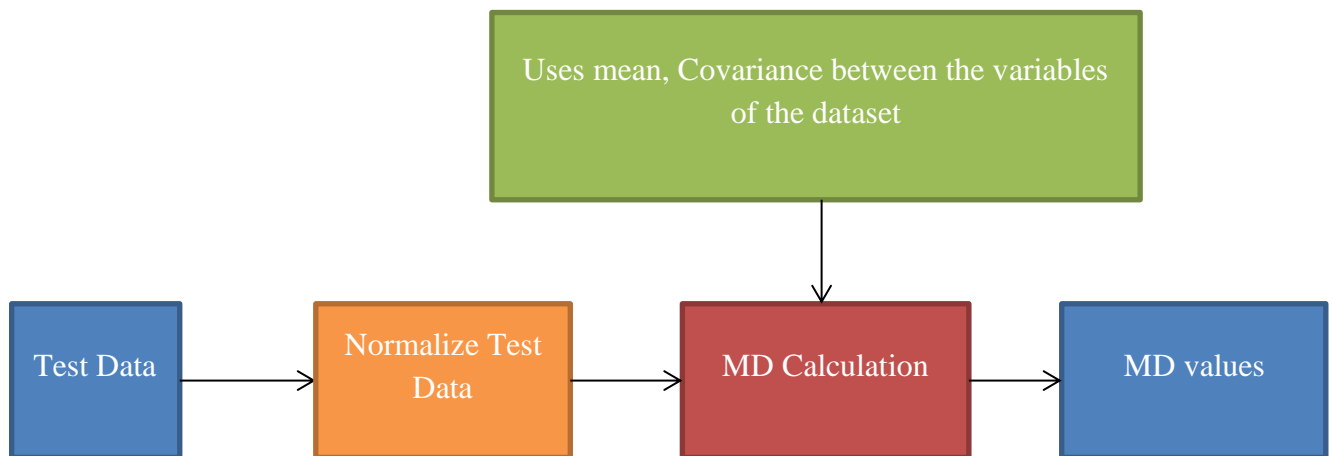


Fig 5.3 MD Calculation using test data

The Mahalanobis Distance measures the relative distance between two variables with respect to the mean of the multivariate data. Thus, the univariate distance of multivariate dataset is obtained with the Mahalanobis distance that can be used to cluster the test data using k-means clustering until the predefined number of iterations or till the difference between previous centroid and the current centroid is so marginal. Formula to calculate Mahalanobis distance is:

$$D = [(X_i - M)^T C^{-1} (X_i - M)]^{0.5}$$

Where

D = Mahalanobis distance

X_i = an object vector

M = arithmetic mean vector

C^{-1} = Inverse Covariance matrix of independent variables

T = Indicates vector should be transposed

From the MD values that is obtained, Extreme SR value analysis with single representation of data objects is implemented to identify the outliers. The ESVA-LCPSD method is used to find the threshold for the normal data and to detect the outliers. Normally, the threshold value is decided by the user manually which is a big shortfall since he/she cannot predict the threshold value always and that may in turn gives false outliers as output.

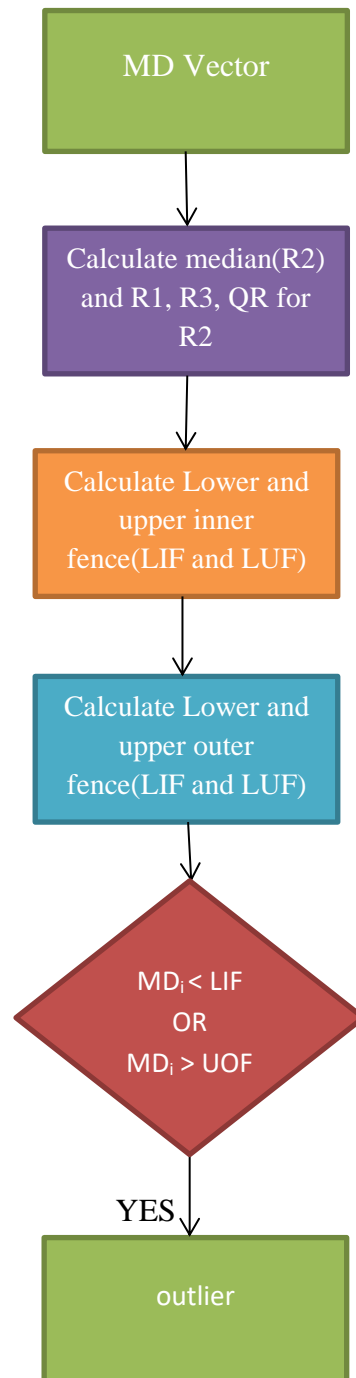


Fig 5.3 ESVA-LCPSD algorithm

CHAPTER 6 : METHODOLOGY AND APPROACH

In this chapter, algorithm of K-means outlier detection using Mahalanobis distance and ESVA-LCPSD algorithm is discussed. The dataset is obtained as the input for detecting the outliers. Modified k-means clustering using Mahalanobis distance takes dataset and k as input and outputs the deviated points. The pseudo code to follow is given below.

Algorithm 1

Input:

=> Dataset containing 'n' instances

=>k: the number of clusters

Output:

=> Outliers

Algorithm:

1. Initially, depending on the value of k, k random points are chosen as initial centroids.
The distances of each data point from the centroids chosen earlier are evaluated using the Mahalanobis distance using the formula $\mathbf{D} = [(\mathbf{X}_i - \mathbf{X}_j)^T \mathbf{C}^{-1}(\mathbf{X}_i - \mathbf{X}_j)]^{0.5}$.
 2. The distance values are compared and the data point is assigned to the centroid which has shortest Mahalanobis distance (MD) value.
 3. The previous steps are repeated. The process is stopped if the clusters obtained are same as that of the previous step.
 4. Use Extreme Value Analysis algorithm ESVA-LCPSD [**Algorithm 2**] to find the outliers that are put into (k+1)th cluster.
-

To cluster the dataset, Mahalanobis distance is metric is used that incorporates covariance between the features of the dataset to achieve much efficient clustering in multivariate space.

The formula to get the distance between the point and the cluster centre is,

$$D = [(X_i - M)^T C^{-1} (X_i - M)]^{0.5}$$

Where

D = Mahalanobis distance

X_i = an object vector

M = arithmetic mean vector

C^{-1} = Inverse Covariance matrix of independent variables

T = Indicates vector should be transposed

Final cluster's MD values are given as input to ESVA-LCPSD algorithm that uses Extreme value Analysis to find deviated data objects in the dataset.

Algorithm 2

Input:

=> MD values

Output:

=> outliers

Algorithm:

1. Calculate median $R2$ from MD values.
2. Compute lower and upper quartile range $R1$ and $R3$ respectively from $R2$.
3. Computer the range of inter-quartile RQ by subtracting $R3 - R1$.
4. Compute Lower Inner Fence(LIF) as $R1 - 1.5 * RQ$.
5. Compute Upper Inner Fence(UIF) as $R3 + 1.5 * RQ$.

6. Compute Lower Outer Fence(LOF) as $R1-3*RQ$.
 7. Compute Upper Outer Fence(UOF) as $R+3*RQ$.
 8. Outliers are the data objects goes beyond LIF and UOF of distribution.
-

The algorithm is approached using python to get better results. Python has built-in packages for data analytics application and thereby making it easy to learn and code. Moreover, it has built-in packages for illustrating results in much clearer formats like graphs, charts, scatter-plots. Thus, python is considered for implementing this project. Source code of the project is:

```
import numpy as np

import random

import sys

import pylab as pl

from matplotlib import pyplot as plt

from sklearn.cluster import KMeans

from mpl_toolkits.mplot3d import axes3d, Axes3D

def kMeans(X, K, maxIters = 10, plot_progress = None):

    centroids = X[np.random.choice(np.arange(len(X)), K), :]

    for i in range(maxIters):

        # Cluster Assignment step

        C = np.array([np.argmin([np.dot(x_i-y_k, x_i-y_k) for y_k in centroids]) for x_i in X])
```

```

# Move centroids step

centroids = [X[C == k].mean(axis = 0) for k in range(K)]

if plot_progress != None: plot_progress(X, C, np.array(centroids))

return np.array(centroids) , C

```

```

def normalize(original):

    dsz1=14500

    dsz2=10

    for j in range(0,dsz2):

        mx=0

        mn=10000000000

        for i in range(0,dsz1):

            if(original[i][j]>mx):

                mx=original[i][j]

            if(original[i][j]<mn):

                mn=original[i][j]

        for i in range(0,dsz1):

            original[i][j]=(original[i][j]-mn)/(mx-mn)

    return original

```

```

def MahalanobisDist(data, verbose=False):

    covariance_matrix = np.cov(data, rowvar=False)

    if is_pos_def(covariance_matrix):

        inv_covariance_matrix = np.linalg.inv(covariance_matrix)

```

```

if is_pos_def(inv_covariance_matrix):

    vars_mean = []

    for i in range(data.shape[0]):

        vars_mean.append(list(data.mean(axis=0)))

    diff = data - vars_mean

    md = []

    for i in range(len(diff)):

        md.append(np.sqrt(diff[i].dot(inv_covariance_matrix).dot(diff[i])))

    if verbose:

        print("Covariance Matrix:\n { }\n".format(covariance_matrix))

        print("Inverse of Covariance Matrix:\n { }\n".format(inv_covariance_matrix))

        print("Variables Mean Vector:\n { }\n".format(vars_mean))

        print("Variables - Variables Mean Vector:\n { }\n".format(diff))

        print("Mahalanobis Distance:\n { }\n".format(md))

    return md

else:

    print("Error: Inverse of Covariance Matrix is not positive definite!")

else:

    print("Error: Covariance Matrix is not positive definite!")

```

```

def is_pos_def(A):

    if np.allclose(A, A.T):

        try:

```

```

        np.linalg.cholesky(A)

    return True

except np.linalg.LinAlgError:

    return False

else:

    return False


def update_U(res,Z,dsz1,dsz2,k,d_avg,U):

    outlier=0

    cur_dist=float(0)

    for i in range(0,k):

        U[i][1]=1000000000

        extreme=False

        #for j in range(0,10):

        MD = MahalanobisDist(data, verbose=True)

        om=np.sort(MD)

        print(om)

        q2=np.median(om)

        print("median",q2)

        q75, q25 = np.percentile(om, [75,25])

        print(q25,q75)

        iqr = q75 - q25

        print("iqr:" ,iqr)

        lif=q25-1*iqr

        print(lif)

```

```

uif=q75+1*iqr
print(uif)
lof=q25-2*iqr
print(lof)
uof=q75+2*iqr
print(uof)
outliers = []
for i in range(len(MD)):
    if ((MD[i] <=lof) or (MD[i]>=uof)):
        outliers.append(i)
        U[i][0]=1
    else:
        U[i][0]=0
        U[i][1]=MD[i]
#print(U)

print("Outliers Indices: { }\n".format(outliers))

for ii in outliers:
    print(ii,data[ii])

print("Total number of dataset:",len(data))
print("Count of outliers:",len(outliers))

def kmor(res,k,gama,no,sigma,nmax,dsz1,dsz2):
    c=0

```



```

r=0

prev_p=float(0)

cur_p=float(0)

d_avg=float(0)

s=[]

U=[[0]*2]*dsz1

Z=[[0.0]*10]*dsz1

while(c<k):

    r=random.randint(50,100)%(dsz1+1)

    if(r not in s):

        s.append(r)

        c=c+1

    for i in range(r):

        Z.append(res[i])

    update_U(res,Z,dsz1,dsz2,k,1000000000,U)

```

```

def show(X, C, centroids, keep = False):

    import time

    time.sleep(0.5)

    plt.cla()

    plt.plot(X[C == 0, 0], X[C == 0, 1], '*b',

             X[C == 1, 0], X[C == 1, 1], '*r')

    plt.plot(centroids[:,0],centroids[:,1],'*m',markersize=10)

    plt.draw()

    if keep :

```

```
plt.ioff()
```

```
plt.show()
```

```
data = np.genfromtxt('C:\TDM-GCC-64\hp.csv',delimiter=' ')
```

```
print("data:\n { }\n".format(data))
```

```
U=[[0]*2]*14500
```

```
k=1
```

```
gama=3
```

```
nmax=100
```

```
sigma=1e-6
```

```
dsz1 =14500
```

```
dsz2=10
```

```
n0=int(dsz1*0.1)
```

```
kmor(data,k,gama,n0,sigma,nmax,dsz1,dsz2)
```

```
centroids, C = kMeans(data, K = 1, plot_progress = show)
```

```
show(data, C, centroids, True)
```

CHAPTER 7 : RESULTS AND DISCUSSION

Implementation of outlier detection using Mahalanobis distance produces much more accurate results, than K-Means Outlier Removal algorithm proposed by Guojun Gan and K.Michael [14]. Snapshots of the output got by both algorithms are embedded here for clarification.

To test the working of algorithm, shuttle dataset is used as the input and it contains 14500 test data that were captured by the sensors at regular time intervals.

FIRST PHASE:

```
C:\Users\user\AppData\Local\Programs\Python\Python37>python hel12402.py
data:
[[ 55.  0. 81. ... 88. 64.  4.]
 [ 56.  0. 96. ... 44.  4.  4.]
 [ 50. -1. 89. ... 40.  2.  1.]
 ...
 [ 55.  0. 77. ... 65. 42.  4.]
 [ 37.  0. 103. ... 85. 20.  1.]
 [ 56.  2. 98. ... 46.  4.  4.]]
```

First phase of output displays all the 14500 sets of original dataset.

SECOND PHASE:

Command Prompt - python hell2402.py

```
2.175638001635369, 2.5226366136834093, 3.3156424745806086, 3.521078157914489, 0.8733478012205542, 2.9617195759121366, 2.10715594877724, 1.980616016644526, 2.865166937421441, 3.8039747953718774, 1.6412619420955703, 2.407227930387155, 2.3118017250692553, 2.257856129425348, 2.044132477288264, 3.684313738041143, 0.8588170460281077, 2.255369546456509, 1.72851362733459, 4.03137607051997, 1.6328154269847727, 2.0524745676662874, 1.8869792268034433, 2.6943989722774857, 3.1079375081305556, 2.4057993406336173, 2.089657608440936, 1.8597060490465866, 3.174977907979643, 3.080840936322449, 2.7875916124053512, 1.566288586604041, 1.9516219465581879, 1.9373248063523325, 2.48653627082613, 1.8873725206963465, 1.7502796031412367, 2.9916118378247187, 2.7958681448841287, 3.024127553872125, 2.09119125540545, 3.6406722620575693, 1.9066474734580117, 2.301306390292733, 3.562870739166987, 2.5462266026169464, 2.666303520329139, 1.610024024086816, 2.3413843599887394, 3.4233953518738764, 2.0837561648409637, 2.424042108669978, 2.48139324896679, 1.949515321040017, 2.347319094978832, 1.4422653824050207, 2.0189985493226876, 2.891658261167018, 2.0235902101247873, 3.243849213816416, 3.1483079370461784, 3.90301650346057557, 2.712174939605234, 2.4107613756082205, 0.9058128312050764, 1.5005462128719784, 0.7623163016173509, 2.7431406665683203, 3.356117959372517, 3.484749258520835, 2.0881675789594247, 3.094228164267129, 3.5567702919596047, 0.7149862145982218, 2.940454107055712, 5.0688212573131946, 2.8943284346643425, 1.7795727673140314, 2.3860501749908036, 2.1669599543769977, 1.7028821025939862, 1.760818976816508, 5.28668721125024, 7.64106102303756, 2.2059985975079353, 1.775937085830499, 0.8456800289402449, 3.9278797325337065, 4.254646117187118, 3.8538679247721075, 1.0402763238967971, 2.833316954968106, 1.9415358303058925, 0.8545292606975529, 2.7401985859543534, 2.070820956383608, 2.2057001694699485, 46.540917067991025, 1.9610837839949928, 3.622833404982313, 3.530768304759665, 2.2340406072363213, 3.073959492664295, 3.7755368364952466, 1.678393345182828, 2.9176322425146184, 1.8698343318173039, 3.3777480283384853, 1.4443311112066182, 1.879643566339281, 5.236230016739976, 2.978523972237599, 1.7213968551508851, 2.2084646071305323, 1.7203162056546166, 1.907391841164367251, 0.9416633937370266, 0.8873787122040283, 3.152174320153582, 3.3992592538780824, 1.3478868560766863, 2.40813219486337, 1.8628966217791432, 0.8544016077630662, 3.2932665235194167, 1.7386485607214759, 0.852208295218506, 1.6393106491438492, 1.75979397878748825, 2.0269899757423593, 2.54099248477471, 5.735730004092261, 3.2565145254986674, 0.8754633834332034, 3.049110184828547, 1.2903483570518, 2.9617224762213405, 3.0880401360559717, 3.635447015908758, 2.3330824041909524, 2.187760830265934, 3.0485263882632894, 1.90205126248157, 1.8105486542429075, 1.7235694579199416, 2.5668568071134157, 2.823370152858219, 1.7127063100302804, 3.0210651671321913, 3.392897041313347, 3.1199361381268496, 1.9386537164268682, 3.810827008446047, 2.01001453486684, 3.659062615569902, 1.8024043553483
```

This phase gives the distance value for all 14500 datasets.

THIRD PHASE:

```
[ 0.61076963  0.61274927  0.61688886 ... 58.21478679 93.77598032
 103.68157268]
median 2.352430296995397
R1: 1.8614581705548174
R3: 3.033105451772817
IQR: 1.1716472812179994
LIF: 0.689810889336818
UIF: 4.2047527329908165
LOF: -0.48183639188118144
UOF: 5.376400014208816
```

```
Outliers Indices: [56, 387, 529, 564, 610, 614, 695, 787, 1930, 2054, 2062, 2070, 2121, 2140, 2498, 2540, 2646, 2655, 4434, 5119, 5164, 5217, 5236, 5315, 5579, 5691, 5715, 5724, 5761, 5945, 6024, 6311, 6358, 6550, 6594, 6615, 6616, 6617, 6618, 6619, 6620, 6621, 6622, 6623, 6624, 6625, 6626, 6627, 6628, 6629, 6630, 6631, 6632, 6633, 6634, 6635, 6636, 6637, 6638, 6639, 6640, 6641, 6642, 6643, 6644, 6645, 6646, 6647, 6648, 6649, 6650, 6651, 6652, 6653, 6654, 6655, 6656, 6657, 6658, 6659, 6660, 6661, 6662, 6663, 6664, 6665, 6666, 6667, 6668, 6669, 6670, 6671, 6672, 6673, 6674, 6675, 6676, 6677, 6678, 6679, 6680, 6681, 6682, 6683, 6684, 6685, 6686, 6687, 6688, 6689, 6690, 6691, 6692, 6693, 6694, 6695, 6696, 6697, 6698, 6699, 6700, 6701, 6702, 6703, 6704, 6705, 6706, 6707, 6708, 6709, 6710, 6711, 6712, 6713, 6714, 6715, 6716, 6717, 6718, 6719, 6720, 6721, 6722, 6723, 6724, 6725, 6726, 6727, 6728, 6729, 6730, 6731, 6732, 6733, 6734, 6735, 6736, 6737, 6738, 6739, 6740, 6741, 6742, 6743, 6744, 6745, 6746, 6747, 6748, 6749, 6750, 6751, 6752, 6753, 6754, 6755, 6756, 6757, 6758, 6759, 6760, 6761, 6762, 6763, 6764, 6765, 6766, 6767, 6768, 6769, 6770, 6771, 6772, 6773, 6774, 6775, 6776, 6777, 6778, 6779, 6780, 6781, 6782, 6783, 6784, 6785, 6786, 6787, 6788, 6789, 6790, 6791, 6792, 6793, 6794, 6795, 6796, 6797, 6798, 6799, 6800, 6801, 6802, 6803, 6804, 6805, 6806, 6807, 6808, 6809, 6810, 6811, 6812, 6813, 6814, 6815, 6816, 6817, 6818, 6819, 6820, 6821, 6822, 6823, 6824, 6825, 6826, 6827, 6828, 6829, 6830, 6831, 6832, 6833, 6834, 6835, 6836, 6837, 6838, 6839, 6840, 6841, 6842, 6843, 6844, 6845, 6846, 6847, 6848, 6849, 6850, 6851, 6852, 6853, 6854, 6855, 6856, 6857, 6858, 6859, 6860, 6861, 6862, 6863, 6864, 6865, 6866, 6867, 6868, 6869, 6870, 6871, 6872, 6873, 6874, 6875, 6876, 6877, 6878, 6879, 6880, 6881, 6882, 6883, 6884, 6885, 6886, 6887, 6888, 6889, 6890, 6891, 6892, 6893, 6894, 6895, 6896, 6897, 6898, 6899, 6900, 6901, 6902, 6903, 6904, 6905, 6906, 6907, 6908, 6909, 6910, 6911, 6912, 6913, 6914, 6915, 6916, 6917, 6918, 6919, 6920, 6921, 6922, 6923, 6924, 6925, 6926, 6927, 6928, 6929, 6930, 6931, 6932, 6933, 6934, 6935, 6936, 6937, 6938, 6939, 6940, 6941, 6942, 6943, 6944, 6945, 6946, 6947, 6948, 6949, 6950, 6951, 6952, 6953, 6954, 6955, 6956, 6957, 6958, 6959, 6960, 6961, 6962, 6963, 6964, 6965, 6966, 6967, 6968, 6969, 6970, 6971, 6972, 6973, 6974, 6975, 6976, 6977, 6978, 6979, 6980, 6981, 6982, 6983, 6984, 6985, 6986, 6987, 6988, 6989, 6990, 6991, 6992, 6993, 6994, 6995, 6996, 6997, 6998, 6999, 7000, 7001, 7002, 7003, 7004, 7005, 7006, 7007, 7008, 7009, 7010, 7011, 7012, 7013, 7014, 7015, 7016, 7017, 7018, 7019, 7020, 7021, 7022, 7023, 7024, 7025, 7026, 7027, 7028, 7029, 7030, 7031, 7032, 7033, 7034, 7035, 7036, 7037, 7038, 7039, 7040, 7041, 7042, 7043, 7044, 7045, 7046, 7047, 7048, 7049, 7050, 7051, 7052, 7053, 7054, 7055, 7056, 7057, 7058, 7059, 7060, 7061, 7062, 7063, 7064, 7065, 7066, 7067, 7068, 7069, 7070, 7071, 7072, 7073, 7074, 7075, 7076, 7077, 7078, 7079, 7080, 7081, 7082, 7083, 7084, 7085, 7086, 7087, 7088, 7089, 7090, 7091, 7092, 7093, 7094, 7095, 7096, 7097, 7098, 7099, 7100, 7101, 7102, 7103, 7104, 7105, 7106, 7107, 7108, 7109, 7110, 7111, 7112, 7113, 7114, 7115, 7116, 7117, 7118, 7119, 7120, 7121, 7122, 7123, 7124, 7125, 7126, 7127, 7128, 7129, 7130, 7131, 7132, 7133, 7134, 7135, 7136, 7137, 7138, 7139, 7140, 7141, 7142, 7143, 7144, 7145, 7146, 7147, 7148, 7149, 7150, 7151, 7152, 7153, 7154, 7155, 7156, 7157, 7158, 7159, 7160, 7161, 7162, 7163, 7164, 7165, 7166, 7167, 7168, 7169, 7170, 7171, 7172, 7173, 7174, 7175, 7176, 7177, 7178, 7179, 7180, 7181, 7182, 7183, 7184, 7185, 7186, 7187, 7188, 7189, 7190, 7191, 7192, 7193, 7194, 7195, 7196, 7197, 7198, 7199, 7200, 7201, 7202, 7203, 7204, 7205, 7206, 7207, 7208, 7209, 7210, 7211, 7212, 7213, 7214, 7215, 7216, 7217, 7218, 7219, 7220, 7221, 7222, 7223, 7224, 7225, 7226, 7227, 7228, 7229, 7230, 7231, 7232, 7233, 7234, 7235, 7236, 7237, 7238, 7239, 7240, 7241, 7242, 7243, 7244, 7245, 7246, 7247, 7248, 7249, 7250, 7251, 7252, 7253, 7254, 7255, 7256, 7257, 7258, 7259, 7260, 7261, 7262, 7263, 7264, 7265, 7266, 7267, 7268, 7269, 7270, 7271, 7272, 7273, 7274, 7275, 7276, 7277, 7278, 7279, 7280, 7281, 7282, 7283, 7284, 7285, 7286, 7287, 7288, 7289, 7290, 7291, 7292, 7293, 7294, 7295, 7296, 7297, 7298, 7299, 7300, 7301, 7302, 7303, 7304, 7305, 7306, 7307, 7308, 7309, 7310, 7311, 7312, 7313, 7314, 7315, 7316, 7317, 7318, 7319, 7320, 7321, 7322, 7323, 7324, 7325, 7326, 7327, 7328, 7329, 7330, 7331, 7332, 7333, 7334, 7335, 7336, 7337, 7338, 7339, 7340, 7341, 7342, 7343, 7344, 7345, 7346, 7347, 7348, 7349, 7350, 7351, 7352, 7353, 7354, 7355, 7356, 7357, 7358, 7359, 7360, 7361, 7362, 7363, 7364, 7365, 7366, 7367, 7368, 7369, 7370, 7371, 7372, 7373, 7374, 7375, 7376, 7377, 7378, 7379, 7380, 7381, 7382, 7383, 7384, 7385, 7386, 7387, 7388, 7389, 7390, 7391, 7392, 7393, 7394, 7395, 7396, 7397, 7398, 7399, 7400, 7401, 7402, 7403, 7404, 7405, 7406, 7407, 7408, 7409, 7410, 7411, 7412, 7413, 7414, 7415, 7416, 7417, 7418, 7419, 7420, 7421, 7422, 7423, 7424, 7425, 7426, 7427, 7428, 7429, 7430, 7431, 7432, 7433, 7434, 7435, 7436, 7437, 7438, 7439, 7440, 7441, 7442, 7443, 7444, 7445, 7446, 7447, 7448, 7449, 7450, 7451, 7452, 7453, 7454, 7455, 7456, 7457, 7458, 7459, 7460, 7461, 7462, 7463, 7464, 7465, 7466, 7467, 7468, 7469, 7470, 7471, 7472, 7473, 7474, 7475, 7476, 7477, 7478, 7479, 7480, 7481, 7482, 7483, 7484, 7485, 7486, 7487, 7488, 7489, 7490, 7491, 7492, 7493, 7494, 7495, 7496, 7497, 7498, 7499, 7500, 7501, 7502, 7503, 7504, 7505, 7506, 7507, 7508, 7509, 7510, 7511, 7512, 7513, 7514, 7515, 7516, 7517, 7518, 7519, 7520, 7521, 7522, 7523, 7524, 7525, 7526, 7527, 7528, 7529, 7530, 7531, 7532, 7533, 7534, 7535, 7536, 7537, 7538, 7539, 7540, 7541, 7542, 7543, 7544, 7545, 7546, 7547, 7548, 7549, 7550, 7551, 7552, 7553, 7554, 7555, 7556, 7557, 7558, 7559, 7560, 7561, 7562, 7563, 7564, 7565, 7566, 7567, 7568, 7569, 7570, 7571, 7572, 7573, 7574, 7575, 7576, 7577, 7578, 7579, 7580, 7581, 7582, 7583, 7584, 7585, 7586, 7587, 7588, 7589, 7590, 7591, 7592, 7593, 7594, 7595, 7596, 7597, 7598, 7599, 7600, 7601, 7602, 7603, 7604, 7605, 7606, 7607, 7608, 7609, 7610, 7611, 7612, 7613, 7614, 7615, 7616, 7617, 7618, 7619, 7620, 7621, 7622, 7623, 7624, 7625, 7626, 7627, 7628, 7629, 7630, 7631, 7632, 7633, 7634, 7635, 7636, 7637, 7638, 7639, 7640, 7641, 7642, 7643, 7644, 7645, 7646, 7647, 7648, 7649, 7650, 7651, 7652, 7653, 7654, 7655, 7656, 7657, 7658, 7659, 7660, 7661, 7662, 7663, 7664, 7665, 7666, 7667, 7668, 7669, 7670, 7671, 7672, 7673, 7674, 7675, 7676, 7677, 7678, 7679, 7680, 7681, 7682, 7683, 7684, 7685, 7686, 7687, 7688, 7689, 7690, 7691, 7692, 7693, 7694, 7695, 7696, 7697, 7698, 7699, 7700, 7701, 7702, 7703, 7704, 7705, 7706, 7707, 7708, 7709, 7710, 7711, 7712, 7713, 7714, 7715, 7716, 7717, 7718, 7719, 7720, 7721, 7722, 7723, 7724, 7725, 7726, 7727, 7728, 7729, 7730, 7731, 7732, 7733, 7734, 7735, 7736, 7737, 7738, 7739, 7740, 7741, 7742, 7743, 7744, 7745, 7746, 7747, 7748, 7749, 7750, 7751, 7752, 7753, 7754, 7755, 7756, 7757, 7758, 7759, 7760, 7761, 7762, 7763, 7764, 7765, 7766, 7767, 7768, 7769, 7770, 7771, 7772, 7773, 7774, 7775, 7776, 7777, 7778, 7779, 7780, 7781, 7782, 7783, 7784, 7785, 7786, 7787, 7788, 7789, 7790, 7791, 7792, 7793, 7794, 7795, 7796, 7797, 7798, 7799, 7800, 7801, 7802, 7803, 7804, 7805, 7806, 7807, 7808, 7809, 7810, 7811, 7812, 7813, 7814, 7815, 7816, 7817, 7818, 7819, 7820, 7821, 7822, 7823, 7824, 7825, 7826, 7827, 7828, 7829, 7830, 7831, 7832, 7833, 7834, 7835, 7836, 7837, 7838, 7839, 7840, 7841, 7842, 7843, 7844, 7845, 7846, 7847, 7848, 7849, 7850, 7851, 7852, 7853, 7854, 7855, 7856, 7857, 7858, 7859, 7860, 7861, 7862, 7863, 7864, 7865, 7866, 7867, 7868, 7869, 7870, 7871, 7872, 7873, 7874, 7875, 7876, 7877, 7878, 7879, 7880, 7881, 7882, 7883, 7884, 7885, 7886, 7887, 7888, 7889, 7890, 7891, 7892, 7893, 7894, 7895, 7896, 7897, 7898, 7899, 7900, 7901, 7902, 7903, 7904, 7905, 7906, 7907, 7908, 7909, 7910, 7911, 7912, 7913, 7914, 7915, 7916, 7917, 7918, 7919, 7920, 7921, 7922, 7923, 7924, 7925, 7926, 7927, 7928, 7929, 7930, 7931, 7932, 7933, 7934, 7935, 7936, 7937, 7938, 7939, 7940, 7941, 7942, 7943, 7944, 7945, 7946, 7947, 7948, 7949, 7950, 7951, 7952, 7953, 7954, 7955, 7956, 7957, 7958, 7959, 7960, 7961, 7962, 7963, 7964, 7965, 7966, 7967, 7968, 7969, 7970, 7971, 7972, 7973, 7974, 7975, 7976, 7977, 7978, 7979, 7980, 7981, 7982, 7983, 7984, 7985, 7986, 7987, 7988, 7989, 7990, 7991, 7992, 7993, 7994, 7995, 7996, 7997, 7998, 7999, 8000, 8001, 8002, 8003, 8004, 8005, 8006, 8007, 8008, 8009, 8010, 8011, 8012, 8013, 8014, 8015, 8016, 8017, 8018, 8019, 8020, 8021, 8022, 8023, 8024, 8025, 8026, 8027, 8028, 8029, 8030, 8031, 8032, 8033, 8034, 8035, 8036, 8037, 8038, 8039, 8040, 8041, 8042, 8043, 8044, 8045, 8046, 8047, 8048, 8049, 8050, 8051, 8052, 8053, 8054, 8055, 8056, 8057, 
```

```

Select Command Prompt - python hell2402.py
11087 [105. -4. 107. -2. 70. 0. 2. 37. 34. 5.]
11099 [105. 5. 107. 9. 72. 5. 1. 35. 34. 5.]
11160 [ 44. -536. 88. 0. 44. 0. 44. 44. 44. 0. 1.]
11264 [105. -5. 106. 1. 72. 3. 2. 34. 32. 5.]
11313 [104. 5. 106. 6. 70. -29. 1. 36. 34. 5.]
11394 [ 82. 2561. 106. -4. 34. -20. 24. 72. 48. 6.]
11397 [102. 4. 102. 1. 72. 18. 1. 30. 30. 5.]
11560 [123. 656. 105. 1. 36. -4. -18. 69. 86. 5.]
11700 [105. -5. 107. 0. 72. 0. 1. 35. 34. 5.]
11750 [ 4.5000e+01 0.0000e+00 1.1800e+02 1.7510e+03 3.1000e+02 1.5164e+04
 7.3000e+01 -1.9100e+02 -2.6400e+02 1.0000e+00]
11770 [105. -2. 106. -3. 72. 1. 1. 33. 32. 5.]
11802 [107. 1. 108. 0. 70. 0. 1. 38. 36. 5.]
11804 [105. 5. 106. 5. 72. 6. 1. 34. 32. 5.]
11817 [104. 0. 106. 0. 70. 0. 1. 36. 34. 5.]
11900 [ 37. -1. 106. 258. 34. -9. 69. 72. 4. 1.]
12080 [ 7.900e+01 0.000e+00 8.300e+01 0.000e+00 9.800e+01 8.098e+03
 4.000e+00 -1.400e+01 -1.800e+01 5.000e+00]
12110 [102. 3. 102. 0. 72. 12. 1. 30. 30. 5.]
12423 [108. 1. 109. 0. 72. 1. 1. 36. 36. 5.]
12442 [107. 0. 108. 0. 72. 8. 1. 36. 34. 5.]
12471 [104. 0. 104. -4. 70. 0. 1. 35. 34. 5.]
12629 [105. 0. 106. 0. 70. 0. 1. 36. 34. 5.]
12838 [ 85. 0. 88. -273. 0. 0. 3. 88. 84. 5.]
12951 [106. 0. 108. 0. 70. 0. 1. 38. 36. 5.]
13196 [106. -5. 108. 2. 72. 0. 2. 35. 34. 5.]
13272 [102. 0. 102. -6. 72. 25. 1. 30. 30. 5.]
13308 [106. 0. 108. 1. 72. 2. 1. 35. 34. 5.]
13397 [ 82. 0. 86. 0. -42. -29. 5. 130. 126. 5.]
13530 [105. 4. 106. 2. 72. 6. 1. 33. 32. 5.]
13554 [106. 0. 108. 0. 70. -1. 1. 38. 36. 5.]
13650 [104. -1. 106. -3. 70. -4. 1. 36. 34. 5.]
13839 [108. 2. 109. 0. 72. 4. 1. 36. 36. 5.]
14009 [105. -5. 106. 2. 70. 0. 2. 36. 34. 5.]
14036 [123. 170. 105. 0. 36. -1. -18. 69. 86. 5.]
14058 [ 37. 0. 76. 737. 20. 0. 40. 55. 16. 1.]
14112 [105. -4. 106. 0. 70. 0. 2. 36. 34. 5.]
14331 [105. 4. 106. 4. 72. 5. 1. 34. 32. 5.]
Total number of dataset: 14500
Count of outliers: 124

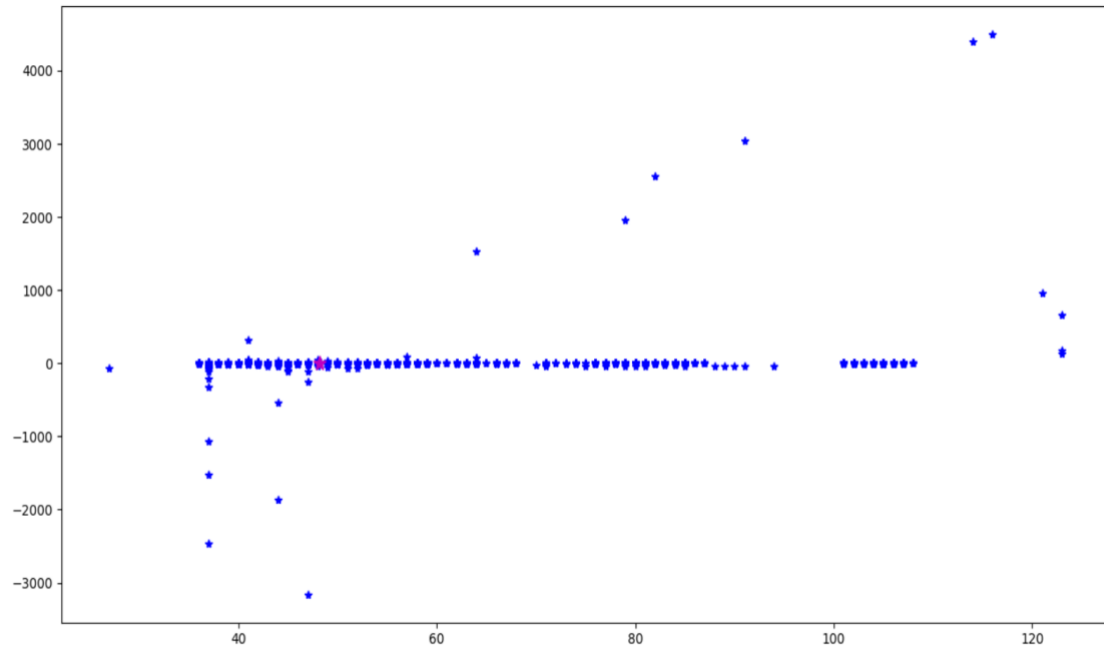
```

The output of Third phase displays the number of outliers present in the given dataset that was calculated using ESVA-LCPSD algorithm.

FOURTH PHASE:

This phase prints the output in graphical form with the help of built-in matplotlib package.

The dataset centre is marked using violet color and the data points are in blue. Outliers can be seen easily as they are scattered away from other data points.



In KMOR method, the number of outliers obtained was 54. While in this case, the algorithm has encountered 124 outliers. This proves that the proposed algorithm performs better than all other existing algorithms.

ALGORITHM	OUTLIERS
K-Means Outlier Removal Algorithm	54
Modified K-Means Using Mahalanobis Distance Outlier Detection Algorithm	124

CHAPTER 8 : CONCLUSION AND FUTURE WORK

Outlier contains noisy data which is researched in various domains like marketing, web mining, trend analysis and more. Various techniques are already being researched that is more generic. We surveyed on various algorithms which is more useful for the beginners. This paper compares KMOR and Modified K-means Clustering using Mahalanobis distance outlier detection algorithm to prove the efficiency of Mahalanobis distance based clustering outlier detection over other algorithms by implementing these algorithms on the real time shuttle dataset that contains 43,000 datasets roughly.

Modified K-means Clustering using Mahalanobis distance outlier detection algorithm is used to cluster and detect outliers using Mahalanobis distance. Mahalanobis distance could be used as an alternative for Euclidean distance as it employs multivariate approach. One of the major advantages is ESVA-LCPSD algorithm to find the outlier points. This is because threshold is set by the user manually in most of the cases that may not be optimal all the time. Hence, this algorithm decides the threshold using extreme value analysis automatically relieving the user from the burden of setting the threshold. Hence, this paper proves that Modified K-means clustering using Mahalanobis distance outlier detection algorithm is much efficient in case of both accuracy and executing time. It provides clusters of spherical shape and find true outliers that that KMOR algorithm has failed to provide .

One of the issues with using Mahalanobis distances is that accuracy of the distance is sensitive to initialization. Thus, for further enhancements, initialization procedure should be altered to select points close to cluster centres. Research has been going to identify points that will help finding the point close to clusters which have more neighbours to be favoured in the initialization step.

REFERENCES

1. **N.V. Kadam , M.A. Pund** , Joint approach for outlier detection, *Int. J. Computer Science Applications*, 6 (2) (2013) 445–448.
2. **J. Whang , I.S. Dhillon and D. Gleich** , Non-exhaustive overlapping k -means, in: “*SIAM International Conference on Data Mining (SDM)*”, 2015 .
3. **Igor Melnykov, Volodymyr Melnykov**, On k-means algorithm with the use of Mahalanobis distances, *Statistics & Probability Letters*. 84 (2014) 88-95.
4. **K. Aparna, M.K. Nair**, *Computational Intelligence in Data Mining*, vol. 2, Springer, pp. 25–35.
5. **M. Jiang, S. Tseng and C. Su** , Two-phase clustering process for outliers detection, *Pattern Recognition Letter*. 22 (6–7) (2001) 691–700.
6. **Z. He, X. Xu and S. Deng** , Discovering cluster-based local outliers, *Pattern Recognition. Letter*, 24 (9–10) (2003) 1641–1650.
7. **V. Hautamaki, S. Cherednichenko , I. Kärkkäinen , T. Kinnunen , P. Fränti** , Improving k-means by outlier removal, in: *Proceedings of the “14th Scandinavian Conference on Image Analysis, SCIA’05”*, 2005, pp. 978–987 .
8. **G. Gan** , *Data Clustering in C ++ : An Object-Oriented Approach*, *Data Mining and Knowledge Discovery Series*, Chapman & Hall/CRC Press, Boca Raton, FL, USA, 2011 .
9. **rednichenko , I. Kärkkäinen , T. Kinnunen , P. Fränti** , Improving k-means by outlier removal, in: *Proceedings of the “14th Scandinavian Conference on Image Analysis”* , 2005 .
10. **S.-Y. Jiang, Q. An**, Clustering-based outlier detection method, in: “*Fifth International Conference on Fuzzy Systems and Knowledge Discovery*”, 2008, pp. 429–433 .

11. **Y. Zhou, H. Yu, X. Cai**, A novel k-means algorithm for clustering and outlier detection, in: “*Second International Conference on Future Information Technology and Management Engineering*”, 2009, pp. 476–480 .
12. **G.S.D.S. Jayakumar, B.J. Thomas**, A new procedure of clustering based on multivariate outlier detection, *J. Data Science*, 11 (2013) 69–84 .
13. **M. Ahmed, A. Naser**, A novel approach for outlier detection and clustering improvement, in: “*Proceedings of the 8th IEEE Conference on Industrial Electronics and Applications (ICIEA)*”, 2013, pp. 577–582 .
14. **S. Kumar, Tommy W. S. Chow and Michael Pecht**, Approach to Fault Identification for Electronic Products Using Mahalanobis Distance, “*IEEE transactions on instrumentation and measurement*”. 59(2009) ,2055-2064.
15. **Guojun Gan , Michael Kwok-Po Ng**, k -means clustering with outlier removal ,in: *Pattern Recognition Letters* 90 (2017) 8–14 .