# Crime Detection using MRCC based Hybrid (GA-FCM) clustering approach

## G. Maheswari

*Assistant Professor, Mangayarkarasi College of Arts and Science for Women, Madurai, Tamil Nadu, India*
*(E-Mail: gm291276@gmail.com)*

## K. Chitra

*Assistant Professor, Government Arts College, Melur, Madurai, Tamil Nadu, India*

## *ABSTRACT*

*In any country, state, or district, crime is a huge issue. Relevant, timely information must be collected so that the problem can be monitored. Analysis of crime is defined by the exploration and detection of crime and their relation between different criminals and the heads of crime. Various analytical or predictive techniques in data mining have been developed and used in various fields. Many researchers use different forms of data mining techniques to detect and track crime. We suggest a method in this paper for the detection of crime using data mining techniques. Next, it is necessary to collect data from online sites, which are unstructured criminality. In the preprocessing phase, the data can be normalized. Crime detection is analyzed using MRCC based hybrid (GA-FCM)  clustering approach, which iteratively generates crime clusters that are based on similar crime attributes. Criminal identification and prediction are analyzed using Ensemble-based classification. Crime verification of our results is done using .NET Framework. The strategy aims to strengthen safety by aiding offender details classification and reporting departments and thus easily identifies criminal. For comparison purpose the same procedure can be applied on PIMA diabetic dataset for diabetic prediction.*

*Keywords: Crime Data mining, MRCC based hybrid (GA-FCM) clustering, Ensemble Classification, and .NET Framework*

## 1. INTRODUCTION

Crime is an offense that is sometimes punished and enforced by statute against society. The reality that offenders conduct offenses anywhere and in some manner, has been noticed. In the world today, offenders are becoming more technically competent. As a consequence, law enforcement officials must keep pace with them. Even the perpetrators now utilize the system widely from massive organized militant networks like Al-Qaeda, and huge cocaine gangs such as the Medellin Cartel to keep up to date with the rules. In an era in which the usage of technology for a great many reasons is so common, this initiative illustrates how the police will overcome

certain technical barriers by analyzing the large volume of data relating to crime created to identify crime that may arise in a particular place or overtime in the future. Due to the spiraling rates of crime across the board, in recent years, there has to be a mechanism for understanding future patterns of criminality, so that we could at least be prepared to deal with them if we can not avoid certain crimes. The issue, though, is ultimately to recognize and forecast crime trends accurately with sufficient accuracy such that possible future criminal acts are identified and finally thwarted. Standard methodologies use machine learning and deep learning methodologies for the prevention and nipping of illegal activity in the Bud. But this has not deterred or prevent any substantial increase in crime in a region. This article was carried out as a measurement of the crime patterns that will occur in a given area or region or a certain interval of time in the future by a novel classification approach. These novel methodologies can predict the crime using crime data to prevent the seed of crime from rooting itself. This work contributes primarily to the future determination of the incidence and pattern of crime, which would include its predictions. This is done by first collecting crime data from crime records then preprocessing this data by normalization, and then clustering can be done by using the MRCC based hybrid (GA-FCM) clustering approach. We intend to gauge crime patterns by applying Ensemble classification algorithms. The rest of the paper can be structured as follows; section 2 depicts the related works, The problem on crime detection in existing methodologies was pointed on section 3, the suggested methodology was implemented in section 4, section 5 describes the evaluation of the suggested methodology performance finally section 6 concludes the paper.

## 2. RELATED WORKS

Information of prior research relevant to crime prevention, which has certain incorporated shortcomings in the literary study. Many of the scholars mentioned key distinction ((Chen, et al. 2004);(Kulis and Jordan 2011);(Malathi and Baboo 2011); (Lee and Estivill-Castro 2011)); (Anitha 2020) and classification techniques for crime detection, criminal identification theoretically; however, none of them provides a sound implementation for the same. Although (Li and Juhola 2014) It is an field where better learning can gain as the crime research says, DMT. CDCI provides a streamlined and visualized method of tracking, recognizing, and forecasting crime and confirmed crime to shield India from abominable crimes. (Bandekar and Vijayalakshmi 2020) The scholars explored the study of crime in India. Crime is analyzed by classification based on machine learning. (Yerpude and Gudur 2017) To forecast factors causing the high crime rate, the writers' clarified algorithms such as the decision tree, Naïve-Bayes was used on the data collection. Output is achieved through supervised and unattended learning. (Uzlov, et al. 2018) Authors dedicated themselves to the review and perspective of the use of methods of data mining in the work of the national police criminal analysts through a process of proactive crime prevention and investigation polices and the implementation thereof. It also defines data processing techniques to improve the efficiency of law enforcement agencies' intelligence management by creating advanced intelligent technical capabilities. (Sevri, et al. 2017) The goal is to disclose the relationship between independent record attributes. (Pramanik,

et al. 2016) authors Illustrate a modern paradigm for large data research that uses several big data tools to predict crime trends. (Zulfadhilah, et al. 2016) The author using the clustering using the K-Means algorithm has been analyzed and categorized into 3 clusters: large, medium, and small.

# 3. PROBLEM STATEMENT

Recent approaches use machine learning and deep learning techniques for detecting, forecasting, and addressing crimes even more efficiently, but evenly the crime rate is still increasing. This occurs because of the poor performance of the classifier does not appear to be highly effective. Therefore, an efficient classification system is required to solve the current classification problems.

# 4. PROPOSED METHODOLOGY

The crime detection approach is employed by using MRCC (Measuring Modularity and extensibility) based hybrid (GA-FCM) clustering mechanism was represented in figure 1. Initially, the crime dataset is taken, and the data is preprocessed with the use of the normalization approach. Then the preprocessed data is clustered with the use of the MRCC based Hybrid (GA-FCM) clustering technique in which coupling and cohesion process is carried out for the clustering mechanism. The cohesion and the coupling mechanism are generally used in the field of data retrieval ,data analytics etc. Here it is used for the purpose of clustering. The features are extracted, and the attributes are selected and matched based on training, and testing samples (i.e., Attribute selection based feature extraction) technique is used. In feature extraction approach, the reusability and extensibility mechanism is employed. Then, the extracted attributes are classified and matched using the Ensemble-based classifier, which in turn generates the questionnaire, and the input data will be matched for similarity, and the prediction is made. If the similarity instances are greater than the threshold limit, then the system retrieves the crime details/crime data. If the similarity is less than the threshold limit, then it will be stored in the database.

### 4.1 CDCI dataset

First of all, we generate crime detection and criminal identification (CDCI) crime dataset through two chronological steps

(1) Data extraction extracts the unstructured crime data from various crime Web sources, namely—NCRB, CPJ, etc. Web sources during the period of 2014–2020.

(2) In the stage of data processing, the data can be cleansed, and data can be reduced, then it can be organized and sort out into 5,000 crime incidents. 35 crime features reflect the organized CDCI crime dataset. The collected data can be applied in the. NET Framework.
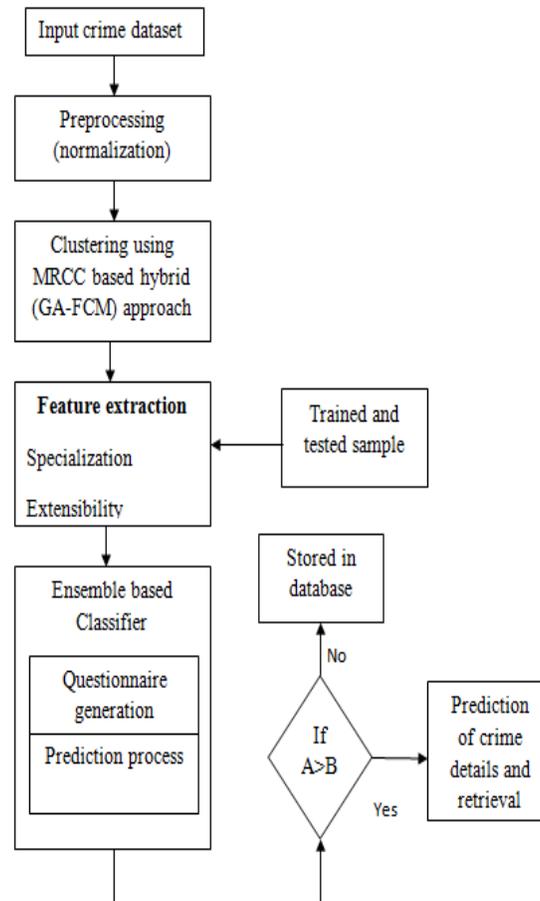
**Figure 1 Schematic representation of the suggested methodology**

### 4.2 Preprocessing

This move includes the further phase of translating data into a numerical type that is simple for the computer to comprehend. A major step in pre-processing is the omission of invaluable words. This can be achieved with standardization. The first step in the normalization process is to obtain the z-score.

$$Z=[(n-\mu)/\sigma] \qquad (1)$$

Where $\mu$ is the mean of the data amount, and $\sigma$ is the standard deviation of the data. While the data mean and standard deviation are not known, then the standard Z- score will be calculated using the sample mean and standard deviation.

$$Z = \frac{n-\bar{n}}{M} \qquad (2)$$

n̄ is the mean of the sample, and M is the standard deviation of the sample.

Then the Hat matrix can be calculated.

$$H = n * (n^T n)^{-1 n^T} \tag{3}$$

The variance for the Hat matrix is,

$$Var(\widehat{H}_i) = \sigma^2(1 - h_{ii}) \tag{4}$$

$$Var(\widehat{H}_i) = \sigma^2(1 - \frac{1}{i} - [(n_i - \bar{n}^2)/\sum_{j=1}^{i}(n_j - \bar{n}^2)]) \tag{5}$$

Then the residual which can be calculated by

$$r_{i=} \frac{\widehat{h}_i}{\bar{\sigma}}\sqrt{1 - h_{ii}} \tag{6}$$

Where $\bar{\sigma}$ is an estimate if the σ

$$\widehat{\sigma}^2 = \frac{1}{i-j}\sum_{j=1}^{i}\overline{mh^2}_j \tag{7}$$

Where m is the number of parameters.

Then the process of the feature scaling can be done to bring all the values in between 0 to 1.This method is called as the concord based normalization.

$$n' = \frac{(Z - n_{min\ Var(\widehat{H}_i)})r_i}{(n_{max} - n_{min})\widehat{\sigma}^2} \tag{8}$$

### 4.3 Clustering using MRCC based hybrid (GA-FCM) approach

The key concept of MRCC hybrid (GA-FCM) clustering in which related elements may be clustered together into a series of clusters, such that similarities (cohesion) intracluster are strong and that differences between the inter-clusters (couples) are small. The suggested approach can be used for improving the process of cohesion and coupling. In software design, its purpose is similar to high coherence and low coupling. A feature attribute data matrix is an input data package. Components are the entities on the basis of their similarities we want to group. The qualities of the elements are features. Although FCM needs less technical assessment, FCM is typically locally oriented. Although FCM needs less technical assessment, FCM is typically locally oriented. Each portion is incorporated in the FCM algorithm into a genetic algorithm, which conserves the benefits of the FCM and GA algorithms. This hybrid approach incorporates

GA-FCM for training and uses global GA FCM exploration to find an acceptable initial clustering method for FCM and local exploration in order to prevent slipping into the local maximum. In a clear and efficient inertia weight adjustment Technique with a new release of GA, the global search and local search power are adjusted and balanced. A feature attribute data matrix is an input data package. Components are the entities on the basis of their similarities we want to group. The qualities of the elements are features. This model represents the natural selection mechanism in which the most suitable features related to the crime are chosen to replicate and create offspring in the next generation. As follows, the new strategy function is defined:

$$\beta = \beta_{max}\text{-}(N_{iter}/N_{iter(total)}{}^{*}\beta_{max} - \beta_{min} \tag{9}$$

where $\beta$ is the current iteration and $N$ is the regulatory factor for fine-tuning ability of GA
As expressed in equation (9), data size will decrease with an increase in the iteration number. According to the MRCC hybrid (GA-FCM) algorithm mutation probability is associated with the fitness,

$$F_f(\text{i})=\begin{cases} F_{f\,min} + \left(F_{f_{max}}\text{-}F_{f\,min}\right) * K, K < 1 \\ F_{f_{max}}/N_{iter}, \end{cases}_{\text{i=1...n}} \tag{10}$$

Where f= fitness(i) $- f_{min}$, which differs from current ith solution quality; fitness(i) and fmin reflect established ith solution health and current global optimum population fitness; $F_{f_{max}}$ and $F_{f_{min}}$ represent the maximum and minimum of mutation probability based on fuzzy respectively. The fitness of the solution and the likelihood of mutation, relative to the K, can be seen in the equation ( 10). The chance of mutation will typically vary as far as replication is concerned. The optimal location is evaluated while the function is initialized.

$$\sigma_j=(\mathbf{\gamma}(1+\beta_p))^{*}\cos(\mathbf{\pi}^{*}\beta_p/2)/(\mathbf{\gamma}\,(1+\beta_p\,2)^{*}\,\beta_p{}^{\beta_p-1/2})^{*}\,(\beta_p - 1/2)^{1/(\beta_p-1/2)} \tag{11}$$

Where $\sigma_j$ represents the random size of the nest.

The equation (12) can be rewritten as,

$$n_p=\text{rand}(\tau_{san},1)^{*}(\,u_b - l_b)+l_b \tag{12}$$

Where $n_p$ represents the random position of the data.
After that, the objective function needs to be computed to find out the fitness solution for coupling and cohesion.

$$\sigma^2= \gamma^{2\,*}c_{ij^2} \tag{13}$$

Where $\sigma^2$ is the error variance, and $c_{ij^1}$ *is* the original distance of the unknown node, the standard deviation from equation (13) is found to be equal to the error variance. The target feature, which is a mean mistake for the unity and coupling technique, will then be determined.

$$f(c_i y_i) = \frac{1}{n}\sum_{j=1}^{n}(e_{ij} - e_{ij'})^2 \tag{14}$$

When utilizing the GA-FCM-Algorithm, MRCC hybrid can be used to approximate the undefined optimal fitness approach for unified and combined applications. The optimal work out approach has been obtained where the identical features are find out (cohesion) they can be grouped out (coupling) .

$$obj_{fn} = -20*\exp(-2*\sqrt{\textstyle\sum \sigma_v})/2 - \exp(\textstyle\sum \cos(2\pi * \sigma_v)/d_b) + 20\exp \tag{15}$$

### 4.4 Feature extraction

The attributes that can then be derived for classification purposes using the Extensibility Specialization-based Gray level co-occurrence matrix (ES-GLCM). Extensibility can calculate the potential to expand a framework and the amount of effort needed to introduce an expansion, and the specialization may even process the data for function extraction procedures. This includes a single data size value, and another data value in the Ø direction is known as l and the neighborhood division of m. Typically m gets a single word meaning, and Ø will gain directionally. And the received directional value will override the data attributes used in the classification process.

The ES- GLCM process can be determined as follows will

$$R(n,p) = G(n,p,m, Ø)/\sum_{k=1}^{H}\sum_{l=1}^{H} G(n,p,m, Ø) \tag{16}$$

Where G(n,p,m, Ø) is the frequency of the particular component having the word score values of l and m, r (n,p) was the component of the n and l. By implementing the ES-GLCM, the different attributes can be obtained.

### Extensible features

It includes the general details on the data required for the criminal data detection process.

$$EF = -\sum_{n=1}^{H-1}\sum_{p=1}^{H-1} R(p,m)*\log(R(p,m)) \tag{17}$$

## Specialized features

Both values collected with the ES-GLCM will be determined by summing up whether it can determine whether the data is large or small. If the homogeneity is small, the angular moment is weak. The datas will usually be measured for uniformity.Then the specialized features can be sort out.

$$SF=\sum_{n=1}^{H-1}\sum_{p=1}^{H-1} R( p,m)\verb|^|2 \tag{18}$$

### 4.5 Adaptive Ensemble classification

We should create a questionnaire for the clearance of the test. Information will be identified, and information from the perpetrator seen based on the answers given. It consists of an automated grouping of outputs and the predictive precision of many algorithms in machine learning. This blends a random forest grouping with a variety of decision tree approaches. Introducing the ensemble definition is given as follows.:

Phase 1: Use K of random training set data points.

Phase 2: Decision tree can be built upon the K points

Phase 3: Choose the exact replicate from phases 1 and 2  from that can be used to build a decision tree.

Phase 4: Create a new sample for each and every decision tree groups and ultimately predict the majority-voting group.

The adaptive ensemble  Classifier generates a collection of randomly selected learning dataset decision-making zones. This combines votes of several Decision Trees into a single word and then settles on the final test entity category.  For calculating the trusted classified value, the correlation of crime scenes has to be calculated.

$$\text{Correlation (i,j)}= \sqrt{(C_{n_{fea}}(i,1)\text{-}V_{n_{fea}}(j-1))+(C_{n_{fea}}(i,1)-(C_{n_{fea}}(j,1)^2} \tag{19}$$

As we conditioned, the prototypes on normal and crime scenes depend upon its features. Here the crime scene features can be calculated by using the equation (20)

$$C_{n\_fea}=[C_{n\_fea}\ \text{Corr}] \tag{20}$$

Hence the classifier can evaluate the features of the crime scene and can discriminate the crime data and normal data. Then the trust value of the crime data was calculated.

$$C_{tv}= (C_{n\_fea}\ dist\ corr) \tag{21}$$

Where   $C_{tv}$ represents the trusted value of the crime data, depends upon the trusted value the normal and the abnormal crime scene can be identified.

All of the above procedure can be applied to the Pima Indian Diabetes Dataset for comparison purpose . The dataset will provides details from a community near Phoenix, Arizona, the USA on 768 patients (268 tested positive cases and 500 tested negative cases). The Tested positive and Tested negative show whether or not the individual has a diabetic substance. Every instance consists of 8 integer attributes. These data contain personal health information and medical test

results. After undergoing preprocessing and classification steps, the dataset revealed some significant results, which are demonstrated in Result and discussion section.

## 5. RESULT AND DISCUSSION

Analysis of Crime Dataset has been carried out using MRCC based Hybrid (GA-FCM) clustering approach. In essence, high cohesion (0.91) and low coupling (0.45) can be obtained, which means keeping parts of a code base that are related to each other in a single place. Low coupling, at the same time, is about separating unrelated parts of the code base as much as possible. The obtained data analysis results can be reported below,

$$\text{Cohesion} = \frac{a}{kl} \qquad (22)$$

Here a be the summing of the distance attributes, K be the suggested method and l be the number of the attributes,

Ratio of the summation of the similarities between all pairs to the total number of pairs called as coupling.

$$\text{Coupling} = \langle I_i \cap I_j | I_i \cup I_j \rangle \qquad (23)$$



**Figure 2 Witness Information**

Throughout investigating cases and delivering deterrence, witnesses play a crucial function. Here the evidence produced details on the accused as seen in Figure 2.

(a)



(b)

**Figure 3 Question generation**

For crime mapping in a particular zone, the questionnaire was administered for crime management, as depicted in figure 3.

**Figure 4 Criminal records classification**

The classifier can analyze the questionnaire data depend upon them. The offender details were classified from the data as depicted in figure 4.The performance of the proposed classification method was analyzed by comparing it with some recent existing methods on data classification. (Chiang, et al. 2020) (Gavrilescu and Vizireanu 2019). True positives (TP) are the criterion for the precision of the number of instances. False Positive (FP) is the state under which the majority of instances was marked as right when in reality inaccurate. False Negative (FN) is the state under which, although valid, the number of instances are marked as incorrect. True Negative (TN) is the state under which the amount of documents listed as accurate though in reality correct.

*A. Accuracy*

This defines the ordered errors, calculate the composition's arithmetical field. Low accuracy creates a disparity between a measurement and a "real" value. This ensures that the exceptional data samples are checked using the same algorithm many times and the computer or device produces correct tests. The exactness of the final details is the percentage of the real results.

$$Accuracy\ (A)\ =\ (TP + TN)/(TP\ +\ TN\ +\ FP\ +\ FN) \tag{24}$$

*B. Sensitivity*

Sensitivity is often referred to as the actual positive degree identification . The percentage of real positives accurately detected is calculated by certain sectors.

$$Sensitivity = TP/(TP + FN) \tag{25}$$

## C. Specificity

Specificity, also known as the real negative score, calculates the amount of individual negatives accurately defined.

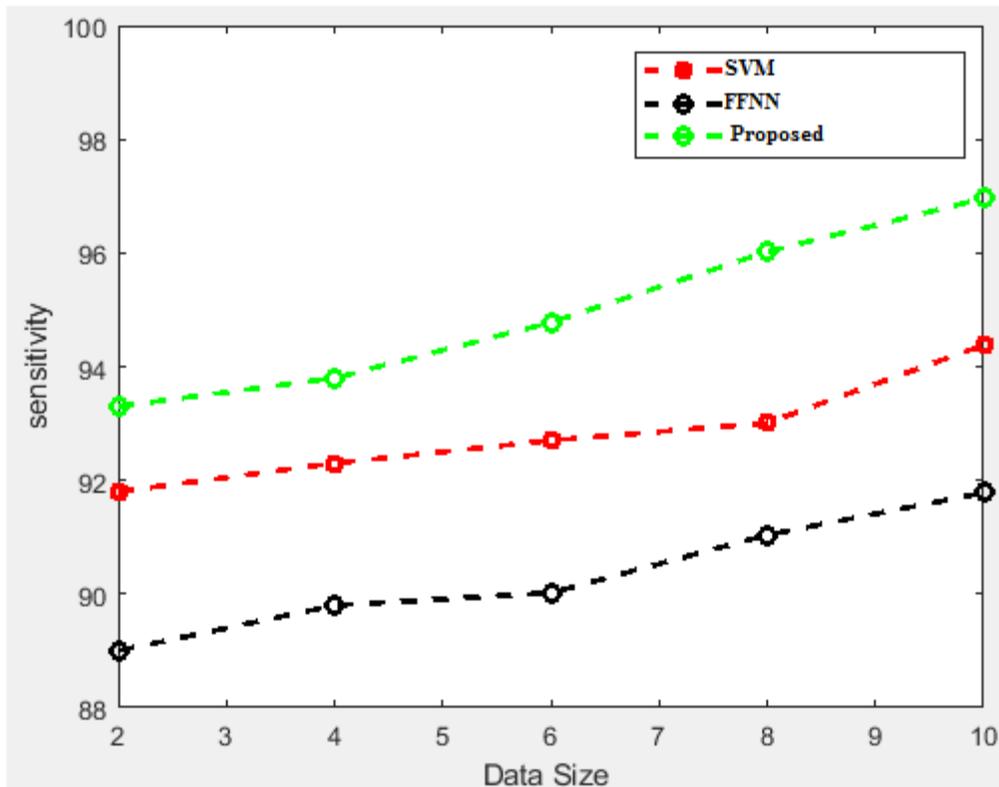$$Specificity \; = \; TP/(TP + FP) \qquad\qquad (26)$$

The public CDCI and PIMA Diabetes data collection has been analyzed. It has achieved the highest predictive accuracy of 95.9 and 95.7% as displayed in Table 1 .

**Table 1 Comparative analysis**

| Performance metrics | CDCI dataset | PIMA diabetic dataset |
|---|---|---|
| Accuracy | 0.959 | 0.957 |
| Sensitivity | 0.975 | 0.972 |
| Specificity | 0.87 | 0.87 |

The crime dataset result are depicted graphically below,

**Figure 5 Sensitivity comparisons with proposed work**

Figure 5 shows the comparative sensitivity study of current and proposed approaches. Data scale is seen on the x-axis and sensitivity value on the y-axis. As a consequence, the method suggested was evidently more reactive than many conventional strategies.
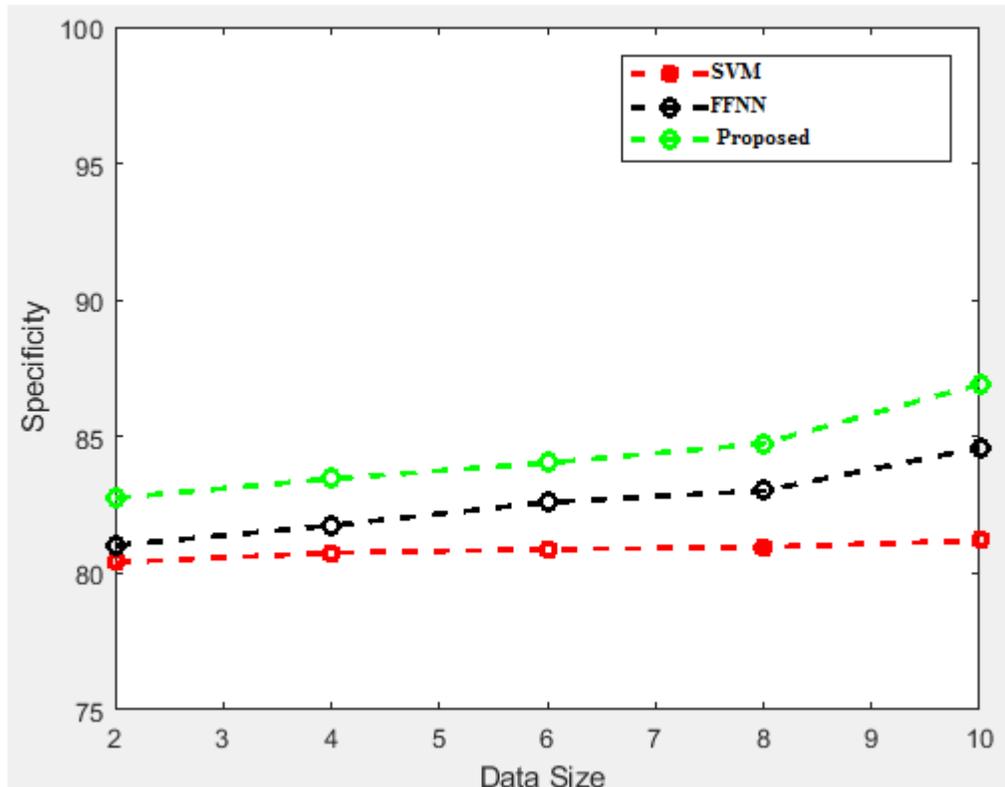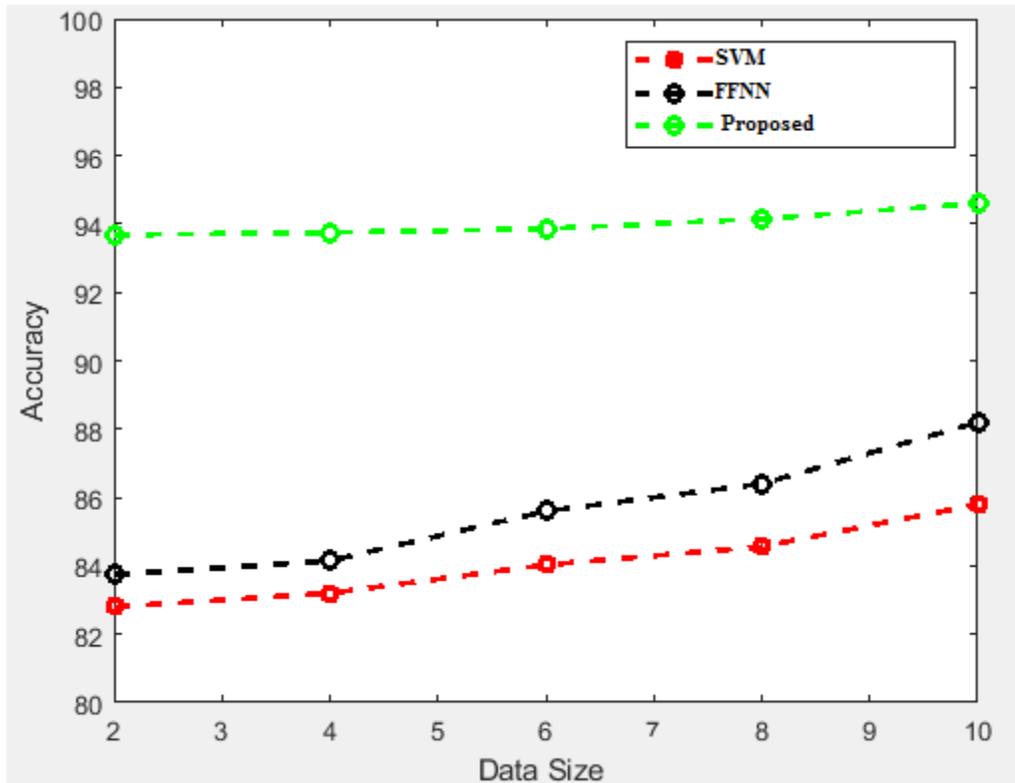


**Figure 6 Specificity comparisons with proposed work**

Figure 6 offers a qualitative overview with respect to specificity of the suggested and current processes. The x-axis shows the data scale and the y-axis the specificity value. The result shows that the suggested method is stronger than other conventional approaches to produce a higher accuracy score.

**Figure 7 Accuracy comparisons with proposed work**

The comparative accuracy study for both proposed and existing approaches is shown in Figure 7. The X-axis indicates the scale of details and the Y-axis displays the exactness. The outcome indicates that the system introduced is able to have greater classification precision than conventional approaches.

## 6. CONCLUSION

An innovative classification method method for managing the crime data was implemented in the proposed study. This research utilizes the MRCC-based hybrid (GA-FCM) algorithm, which is implemented into an adaptive ensemble classification to efficiently achieve the forecast performance with a qualified persistent crime data set classifier. This allows the architecture for the application classification process to deliver a stronger outcome. In the field of sentimental research the success of classification methods is assessed by their precision, their specificity, and their accuracy, which is equally crucial. The results of this analysis reached a data classification accuracy of 95.9 percent in crime dataset. Also the same algorithm can be applied for predicting the diabetes. The performance clearly shows that the implemented algorithm can predicts diabetes effectively with high range of accuracy. The algorithm 's performance quality is substantially increased. In contrast to other data mining methods, the performance may be more reliable.

## *REFERENCES*

1.    *Anitha A 2020, Prediction of Crime Rate Using Data Clustering Technique, in Soft Computing for Problem Solving, Springer, pp. 443-454.*

2. *Bandekar SR & Vijayalakshmi C 2020, 'Design and Analysis of Machine Learning Algorithms for the reduction of crime rates in India,' Procedia Computer Science, vol. 172, pp. 122-127.*

3. *Chen H, Chung W, Xu JJ, Wang G, Qin Y & Chau M 2004, 'Crime data mining: a general framework and some examples,' computer, vol. 37, no. 4, pp. 50-56.*

4. *Chiang H-S, Sangaiah AK, Chen M-Y & Liu J-Y 2020, 'A novel artificial bee colony optimization algorithm with SVM for bio-inspired software-defined networking,' International Journal of Parallel Programming, vol. 48, no. 2, pp. 310-328.*

5. *Gavrilescu M & Vizireanu N 2019, 'Feedforward Neural Network-Based Architecture for Predicting Emotions from Speech,' Data, vol. 4, no. 3, p. 101.*

6. *Kulis B & Jordan MI 2011, 'Revisiting k-means: New algorithms via Bayesian nonparametrics,' arXiv preprint arXiv:1111.0352*

7. *Lee I & Estivill-Castro V 2011, 'Exploration of massive crime data sets through data mining techniques,' Applied Artificial Intelligence, vol. 25, no. 5, pp. 362-379.*

8. *Li X & Juhola M 2014, 'Country crime analysis using the self-organizing map, with special regard to demographic factors,' AI & society, vol. 29, no. 1, pp. 53-68.*

9. *Malathi A & Baboo SS 2011, 'Evolving data mining algorithms on the prevailing crime trend– an intelligent crime prediction model,' Int J Sci Eng Res, vol. 2, no. 6*

10. *Pramanik MI, Lau RY & Chowdhury MKH 2016, 'Automatic Crime Detector: A Framework for Criminal Pattern Detection in Big Data Era,' PACIS, p. 311.*

11. *Sevri M, Karacan H & Akcayol MA 2017, 'Crime analysis based on association rules using apriori algorithm,' International Journal of Information and Electronics Engineering, vol. 7, no. 3, pp. 99-102.*

12. *Uzlov D, Vlasov O & Strukov V 2018, 'Using Data Mining for Intelligence-Led Policing and Crime Analysis,' 2018 International Scientific-Practical Conference Problems of Infocommunications. Science and Technology (PIC S&T), pp. 499-502.*

13. *Yerpude P & Gudur V 2017, 'Predictive modelling of crime dataset using data mining,' International Journal of Data Mining & Knowledge Management Process (IJDKP) Vol, vol. 7*

14. *Zulfadhilah M, Prayudi Y & Riadi I 2016, 'Cyber profiling using log analysis and k-means clustering,' International Journal of Advanced Computer Science and Applications, vol. 7, no. 7, pp. 430-435.*