
**DEVELOPING UNCERTAINTY METAHEURISTIC ALGORITHM FOR CLUSTERING
SIMILAR PATTERNS IN AUTISM DATASET****C.Bhuvaneshwari**

Ph.D. Research scholar, Department of Computer Science, L.R.G. Government Arts College for Women, Tiruppur, Tamilnadu, India. bhuvana.chinna@gmail.com

Dr.A.Anbarasi

Assistant Professor, Department of Computer Science, L.R.G. Government Arts College for Women, Tiruppur, Tamilnadu, India. anbarasi2@gmail.com

***Corresponding Author:** C.Bhuvaneshwari

*Ph.D. Research scholar, Department of Computer Science, L.R.G. Government Arts College for Women, Tiruppur, Tamilnadu, India. bhuvana.chinna@gmail.com

Abstract

An individual with autism has substantially impaired intellectual abilities, linguistic abilities, recognition of objects, interaction, and interpersonal abilities. Although early discovery of autism can help with diagnosis and the implementation of appropriate measures for effect mitigation, this condition cannot be cured. Though researchers developed many machine learning models in prediction of autism, still the problem exist when the hospitals, therapy facilities, and smartphone apps generate a lot of data about autism, but they lack in producing rich data with prior representation of categories or labels as well. Most of the existing unsupervised learning models are not focusing on the uncertainty conditions like vagueness, inconsistency, incompleteness, hesitancy and ambiguity issues commonly occurred in autism dataset. Thus, they failed to achieve highest rate of accuracy in presence of these existing issues. Hence, in this paper an empowered uncertainty based unsupervised learning model coined as neutrosophic clustering (ENU) is constructed to handle the issue of noisy and outlier instances presented in autism dataset, that affects the clustering accuracy and making difficult to understand the pattern of autism exhibiting unknown patterns. The process of clustering and the centroid selection is done by adopting the knowledge of sparrow search optimization algorithm (SSOA) instead of performing random selection. With the optimized cluster centroid selection using highest fitness value and the kernel function enabled neutrosophic clustering discriminates the autism existing instances from normal instances more precisely. The simulation results is conducted of autism dataset collected from Kaggle repository and the performance outcome prove that the proposed ENC-SSOA produced highest rate of autism detection compared with other state of art.

Keywords: Autism, uncertainty, unsupervised learning, neutrosophic clustering, sparrow search optimization algorithm, kernel function

Introduction

The main focus of health care providers in recent times has been digitized medical care assessment. In general, the medical sector is one of the most significant and critical systems in a person's life [1]. The unified presentation of data is assisted by automated analysis, which also offers scientists a simple and adaptable methodology for data collecting. The

troubling aspect is that an autistic person needs an early diagnosis in order to reach their life developments. One of the biggest challenges is making an early diagnosis of the autism condition because it might take consultants up to six months to identify the key symptoms[2]. Data mining techniques were utilized to speed up this process to avoid these drawn-out processes and traditional techniques used to determine the characteristics of persons with autism, which have in effect led to a significant advancement in this sector to better understand the condition of individuals with autism.

The machine learning plays a vital role in detection of autism at its early stage using supervised learning models, unsupervised learning models and artificial intelligence. [3] When the autism dataset is available with the class labels the usage of supervised or classification algorithms are more opted. When the quality of raw data is not satisfactory, then it is primary to perform preprocessing to overcome the issue of missing values and dimensionality curse [4]. To discover the hidden pattern in autism dataset and to understand the latent structure without the class labels or in presence of unknown instances, unsupervised learning accomplishes it. Unsupervised learning has the advantage of enabling researchers to employ unlabeled data that might not contain existing categories.

While using autism dataset, the presence of inconsistency, vagueness, indeterminacy in determining the similarity among the unlabeled instances leads to produce least accuracy in prediction process. This issue is considered as the main challenge in this paper, to develop a empowered uncertainty model for understanding the depth pattern of unknown instances in autism dataset. While using conventional clustering algorithms, the selection of centroids among the population is done in an arbitrary manner, during such process there is high chance of selecting insignificant instances as centroids. Thus, the metaheuristic model is adopted in this work to perform clustering and cluster centroid selection to improve the accuracy in computing similar instances and categorizing them according to identify the variation of pattern among autism and healthy children.

Related Work

Ming Jiang et al [5] developed an unsupervised learning model by tracking the eye movement of child to detect autism spectrum disease. The features are reduced using cluster fix n they are classified into two categories. To train the model, they used support vector machine to predict the autism. Using SVM works for the dataset with known patters and with class variables.

Stevens et al [6] introduced a hierarchical and gaussian mixture method to discover the autism based on phenotypes behavior. Mapping ASD subtypes and their biological links with unsupervised machine learning. Each subgroup is examined using regression method. The uncertainty due to noise and outliers are not focused in this work.

Suman Rajet al [7] devised a deep learning model to predict the presence of autism by collecting three various autism dataset which contains non-clinical information. The dataset comprised of three different age groups like children, adult and adolescent subjects. The results show CNN produced best results compared with other classification models due to the problem of overfitting.

Rajabet al [8] performed a detailed analysis on feature subset selection using different statistical methods. The impacts of significant feature selection to improve the ASD prediction is discussed in their work. To perform validation different empirical methods are used for classification and training. The class imbalance affects the accuracy of the models.

Rasool et al [9] designed a novel early prediction model using machine learning and data mining algorithms. The autism dataset underwent feature subset selection using correlation among the attributes, the prediction was done using conventional supervised learning models.

Chelseat al [10] conducted detailed survey on ASD prediction using unsupervised models and to discover the importance of clustering models. Different clustering algorithms are used for detecting autism and analyzed how these algorithms work in determining the similarity and grouping them.

All the above existing models either work on the supervised learning or unsupervised learning, the problem of uncertainty is not well handled. As the hesitation affects the performance of the prediction accuracy this proposed work aims to overcome the above-mentioned limitations by developing uncertainty-based clustering model to understand the pattern of autism children and predict them at its early stage.

Methodology: Empowered Neutrosophic Clustering with sparrow search optimization Algorithm for Autism pattern discovery

In this work the proposed neutrosophic clustering with sparrow search optimization algorithm for autism pattern discovery is constructed. Initially, the autism dataset [11] collected from the Kaggle repository comprised of raw dataset with 1054 instances with 17 attributes and 1 class variable. The dataset is preprocessed in the previous work [12] using the boosted regression tree and fuzzy backward feature elimination for data imputation to convert the missing values to complete dataset and the relevant features which are highly correlated with the dependent class variable is used for further processing. The reduced dataset comprised of 11 attributes instead of using 17 attributes to overcome the time and computation complexity. The pattern of autism and healthy toddlers are discriminated in this work by devised an uncertainty unsupervised neutrosophic clustering. The process of clustering is optimized by selecting the potential features using a novel metaheuristic algorithm known as sparrow search optimization algorithm. The detailed description of the proposed model for autism pattern discovery is depicted in the figure 1.

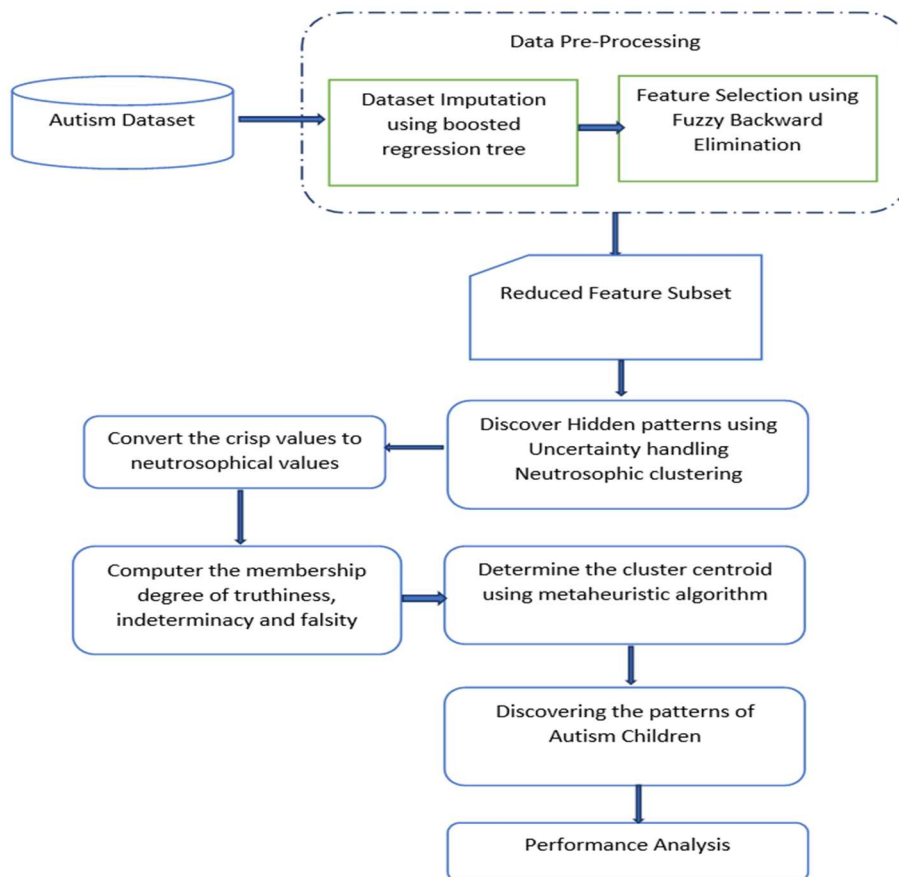


Figure 1: Overall Framework of the proposed uncertainty-based clustering in autism prediction

Determining Similarity using Neutrosophic Clustering for Autism Prediction

Unlike other uncertainty theories, that have only the ability of handing few factors of uncertainty by defining the elements either in terms of degree of membership or non-membership only with complete details. This proposed work utilizes the generalization of uncertainty theories even in the presence of incompleteness, inconsistency, vagueness and indeterminacy is termed as Neutrosophic theory introduced by Smarandache [13, 14]. Most of the real time datasets, especially while dealing with medical or disease diagnosis often comprised of vague, inconsistent and incomplete datasets which affects the prediction model's performance. Hence, in this proposed work neutrosophical representation and determining the unknown pattern of autism dataset by neutrosophical clustering are constructed to improve the pattern discovery in autism presence or absence children. The neutrosophic logic is denoted with the tristate elements which all belong to the degree of membership, independent to each other and their interval lies between $]0,1[$. The each element (or) record in the autism dataset is represented neutrosophical towards the degree of belongingness as autism presence in the triplet form truthiness, falsity and indeterminacy.

$$NL(T) = \{ (\tilde{\mu}, \tilde{\nu}, \tilde{\delta} \in [0,1]^3) \}$$

Clustering is a key component of data mining and machine learning. Based on the proximity metric used, the input data are divided into categories. Every instance in the dataset are handled equally in K-Means, DBSCAN and fuzzy C-means, without factoring into consideration outlier and noisy scenarios. However, genuine datasets like the one for the

autism disease may contain noise and outliers, which call for firm decision-making. Neutrosophic clustering, which manages both noise and outliers, can solve this problem. With a new objective function, it considers both the degree of identity to determined and indeterminate clusters. The multifaceted attribute field is handled using kernel as its objective function, that is specified in the equation below.

$$J_{MNLC}(\ddot{\mu}, \vartheta, \ddot{\pi}, C) = \sum_{i=1}^M \sum_{j=1}^C (\varphi_1 \ddot{\mu}_{ij})^p \|\alpha(z_i) - \alpha(C_j)\|^2 + \sum_{i=1}^N (\varphi_2 \mathfrak{S}_{ij})^p \|\alpha(z_i) - \alpha(C_{imx})\|^2 + \tau^2 \sum_{i=1}^N (\varphi_3 \vartheta_i)^p$$

and $\|\alpha(z_i) - \alpha(C_j)\|^2 = \mathfrak{U}(z_i, z_j) - 2\mathfrak{U}(z_i, C_j) + \mathfrak{U}(C_i, C_j)$

where $\mathfrak{U}(z_i - z_j) = \alpha^{\ddot{\mu}}(z_i)\ddot{\mu}(z_i)$ signifies the function of inner dot product, and if it used gaussian objective function then $\mathfrak{U}(z_i, z_j)$ and $\mathfrak{U}(C_i, C_j)$, then the modified neutrosophic objective function would be formulated a

$$J_{MNLC}(\ddot{\mu}, \vartheta, \ddot{\pi}, C) = \sum_{i=1}^M \sum_{j=1}^C ((\varphi_1 \ddot{\mu}_{ij})^p (1 - \mathfrak{U}(z_i, C_j))) + \sum_{i=1}^M ((\varphi_2 \mathfrak{S}_{ij})^p (1 - \mathfrak{U}(z_i, C_{imx}))) + \tau^2 \sum_{i=1}^M (\varphi_3 \vartheta_i)^p$$

$$C_j = \frac{\sum_{i=1}^M (\varphi_1 \ddot{\mu}_{ij})^p \mathfrak{U}(z_i, z_j)}{\sum_{i=1}^N (\varphi_1 \ddot{\mu}_{ij})^p}$$

Membership degree of truthiness is mathematically formulated as

$$\ddot{\mu}_{ij} = \frac{\varphi_2, \varphi_3 \mathfrak{U}(z_i, C_j)^{-\left(\frac{2}{p-1}\right)}}{\sum_{j=1}^C \mathfrak{U}(z_i, C_j)^{-\left(\frac{2}{p-1}\right)} + \mathfrak{U}(z_i, C_{imx})^{-\left(\frac{2}{p-1}\right)} + \tau^{-\left(\frac{2}{p-1}\right)}}$$

Membership degree of indeterministic is determined by

$$\mathfrak{S}_{ij} = \frac{\varphi_1, \varphi_3 \mathfrak{U}(z_i, C_j)^{-\left(\frac{2}{p-1}\right)}}{\sum_{j=1}^C \mathfrak{U}(z_i, C_j)^{-\left(\frac{2}{p-1}\right)} + \mathfrak{U}(z_i, C_{imx})^{-\left(\frac{2}{p-1}\right)} + \tau^{-\left(\frac{2}{p-1}\right)}}$$

Membership Degree of falsity is formulated as

$$\vartheta_{ij} = \frac{\varphi_1, \varphi_2 (\tau)^{-\left(\frac{2}{p-1}\right)}}{\sum_{j=1}^C \mathfrak{U}(z_i, C_j)^{-\left(\frac{2}{p-1}\right)} + \mathfrak{U}(z_i, C_{imx})^{-\left(\frac{2}{p-1}\right)} + \tau^{-\left(\frac{2}{p-1}\right)}}$$

Sparrow search Algorithm for Centroid Selection and Optimized Neutrosophic Clustering

To discover the best instances in autism dataset as centroids, virtual sparrow as search agent which involved in discovering the centroid instances and each sparrow is represented as the matrix [15]

$$P = \begin{bmatrix} P_{1,1} & P_{1,2} & \dots & \dots & P_{1,f} \\ P_{2,1} & P_{2,2} & \dots & \dots & P_{2,f} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ P_{n,1} & P_{n,2} & \dots & \dots & P_{n,f} \end{bmatrix}$$

Where 'm' denotes number of features in the autism dataset, total number of sparrows involved as search agent as 'f' for optimizing the neutrosophic clustering. The fitness value of the sparrows are computed and displayed in the vectors as

$$Fz_x = \begin{bmatrix} fz([P_{1,1} & P_{1,2} & \dots & \dots & P_{1,f}]) \\ fz([P_{2,1} & P_{2,2} & \dots & \dots & P_{2,f}]) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ fz([P_{v,1} & P_{v,2} & \dots & \dots & P_{v,f}]) \end{bmatrix}$$

Fz_x is the fitness score of every row of an individual, where 'v' stands for the number of sparrows. When searching for food, the producers with the highest fitness values are given top consideration [16]. The producers have a wider search region than the remainder of the scavengers in the group because they are in charge of finding the best supply of food and directing the entire population's movement. Each cycle updates the workers' position according to the estimation.

$$P_{ij}^{t+1} = \begin{cases} P_{ij}^t \cdot e\left(\frac{-i}{\beta \cdot itr_{max}}\right) & \text{if } AV_2 \geq T \\ P_{ij}^t + R \cdot K & \text{if } AV_2 < T \end{cases}$$

Where $j = 1, 2, \dots, f$ is used to indicate the size and t stands for the current iteration. P_{ij}^t stands for the i th sparrow's j th size at iteration 't' is the present iteration. A random number is represented by $\beta \in [0, 1]$. The alarm value is $AV_2 \in [0, 1]$ and the hazard threshold value is $T \in [0.5, 1.0]$. The standard distribution is signified as Q and K is the value of the matrix. When $AV_2 < T$, the producers enter wide search mode because there are no competitors in the area. When a few sparrows spot a predator and $AV_2 > T$, all of the sparrows must take off rapidly for a safe area. A few of the sparrows attempt to outbid the energy producer by being depicted as

$$P_{ij}^{t+1} = \begin{cases} R \cdot e\left(\frac{P_{wrst}^t - P_{ij}^t}{i^2}\right) & \text{if } i > v/2 \\ P_p^{t+1} + |P_{ij}^{ct} - P_{pd}^{t+1}| \cdot B^+ \cdot K & \text{otherwise} \end{cases}$$

P_{pd} refers to producer's potential location and the P_{wrst} refers to the worst global present location. B refers to the matrix with f variables.

. When $i > v/2$, it is advised that the scrounger with the lowest fitness value prefer to be the hungriest. The main location of the sparrow population is referred to as

$$P_{ij}^{t+1} = \begin{cases} P_{bst}^t + D \cdot |P_{ij}^t - P_{bst}^t| & \text{if } fz_i > fz_g \\ P_{ij}^t + H \cdot \left(\frac{|P_{ij}^t - P_{wrst}^t|}{(fz_i - fz_w) + \epsilon} \right) & \text{if } fz_i = fz_g \end{cases}$$

Where the current global optimal position is denoted by P_{bst} , H is the arbitrary values [0.1], by applying normal distribution 'D' it controls the movement step size. global best fitness value and worst fitness value are fz_g and fz_w respectively. The sparrow is at the edge of the group when $fz_i > fz_g$. When the sparrow is at middle it is denoted as $fz_i = fz_g$.

The autism dataset is given as the input dataset and it is converted into the neutrosophic values with the computation of membership degree of truthiness, falsity and indeterminacy. Then each attribute is preprocessed by using the imputation model and feature selection model proposed in our previous work []. In this proposed work, instances in the autism dataset is clustered to determine the cluster centroid and to determine the cluster centroids metaheuristic model is applied. Depending on the objective function expressed, the cluster centroids which has the highest neighborhood instances and the best fitness value is determined by applying the empowered neutrosophic clustering with sparrow search optimization. By representing each instance with neutrosophic values, and the kernel based neutrosophic handles the outliers and noisy instances more effectively.

With the clustered instances the pattern of children with autism and normal children are discovered more effectively using the proposed empowered neutrosophic clustering with sparrow search optimization algorithm (ENC-SSOA).

Experimental Results and Discussions

In this section, the evaluation of the proposed empowered neutrosophic clustering with sparrow search optimization algorithm (ENC-SSOA) for discovering the pattern of autism children at its early stage is discussed. The python software is used to construct ENC-SSOA for detecting autism. The existing clustering algorithms used for comparison are K-Means, DBSCAN and Fuzzy C means.

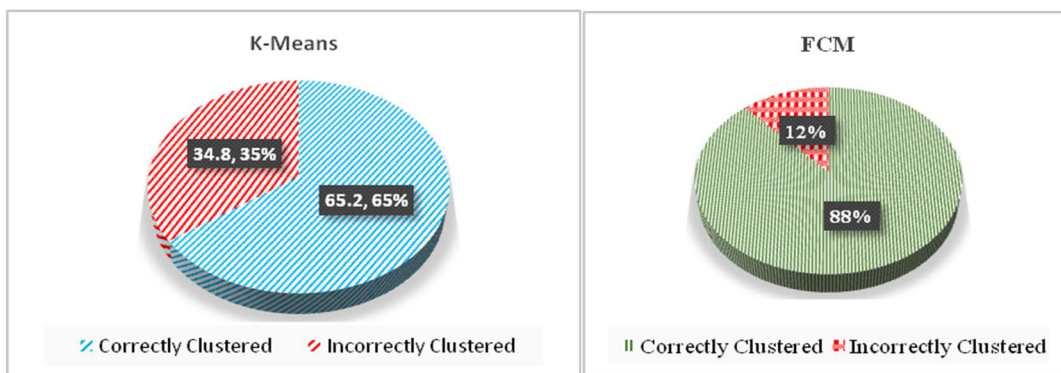




Figure 2a)-2d) Correctly and Incorrectly Clustered patterns of autism dataset

Figure 2a)-2d) shows the four different clustering models performance based on correctly and incorrectly clustered instances which determines the autism patterns. The uncertainty clustering model defines each element in terms of truth, falsity and indeterminacy. From the obtained results, it is proved that uncertainty based neutrosophic clustering model intelligently identifies the border lying instances and noisy instances to overcome the problem of inconsistency prevailed in autism detection. The other conventional clustering models, with its ability of local searching due to the presence of border and noisy instances the correctly clustering rate is low compared with the proposed ENC-SSOA algorithm.

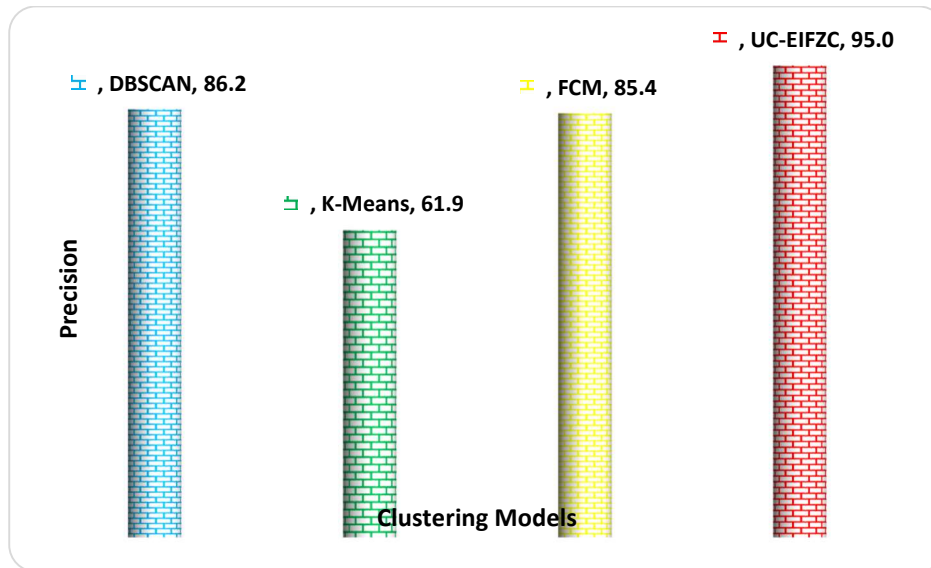


Figure 3 Comparative result based on precision

The incompleteness and vagueness in understanding the unknown instances of autism dataset is a very challenging task while using the conventional clustering models. Figure 3 displays the precision value of four different clustering models involved in grouping the similar patterns of instances in autism dataset, to understand the discrimination among healthy and autism children. The uncertainty based ENC-SSOA with its knowledge of indeterminacy index defines each unknown pattern of autism present children more

precisely with highest percentage rate of 95.0% while DBSCAN, K-Means and FCM produced 8.2%, 61.9% and 85.4% respectively.

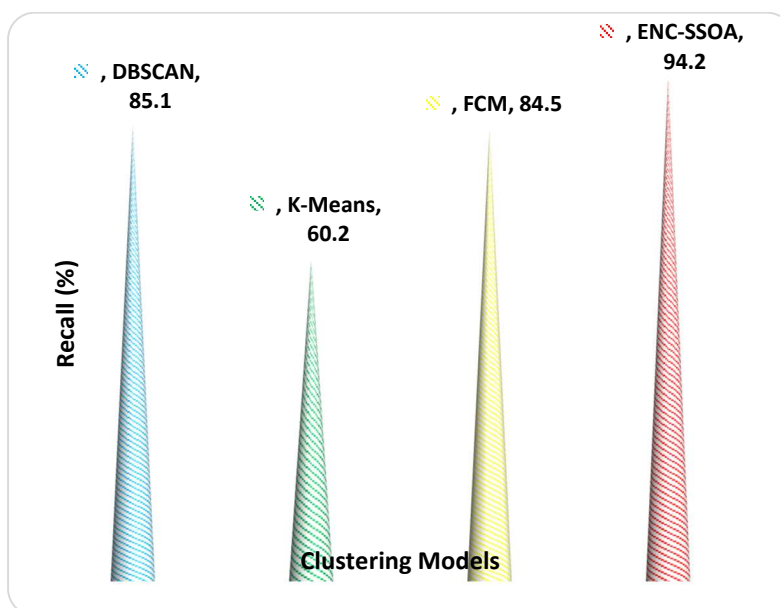


Figure 4 Comparative result based on recall

Figure 4 illustrates the comparative analysis based on recall of four clustering models involved in autism prediction when the patterns of autism dataset is unknown. The dataset used in this work is unlabeled, by computing the similarity among the instances, the process of clustering is done using the uncertainty-based clustering. Each attribute is represented in terms of neutrosophic values with the factors of membership towards belongingness, non-belongingness and indeterminacy to assess whether it belongs to autism or healthy child. The discrimination of autism child from the healthy children group assessment helps to investigate their problem at its early stage to improve their living style. Hence the proposed ENC-SSOA achieves highest recall rate of 94.2% while other state of arts produced less values due to lack of knowledge about uncertainty.

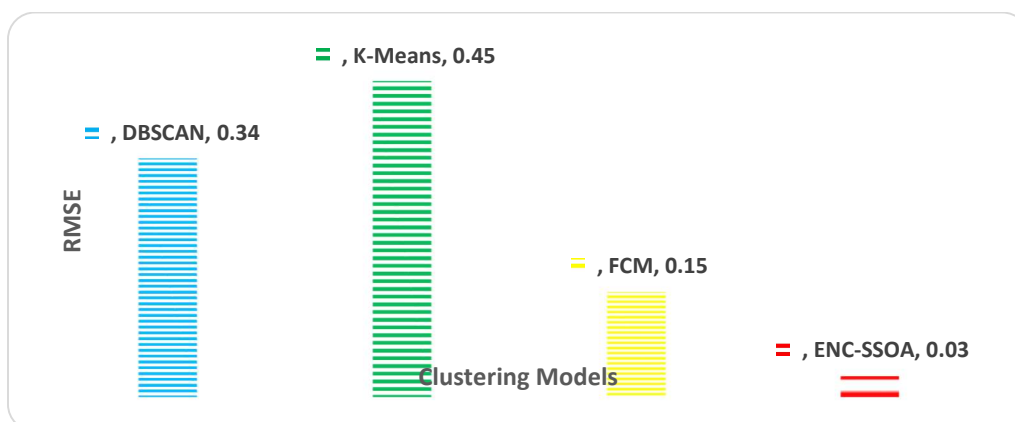


Figure 5 Comparative result based on RMSE

Figure 5 explores error rate of the DBSCAN, K-Means, Fuzzy C Means and proposed Uncertainty handling enriched Neutrosophic Clustering. The error rate of ENC-SSOA is considerably reduced while compared with other three clustering models. The reason is, fuzzy c means defines each element in the form of membership grade only, the issue of vagueness in determining the border lying instances and outliers are not focused hence its error rate in detection of autism is increased. The K-Means and DBSCAN models suffers from local optima and thus the autism with uncertainty conditions are not greatly dealt by them.

Conclusion

In this paper, the uncertainty prevailed in autism prediction is handled by representing each instance in the dataset using the triplet factors known as truthiness, indeterminacy and falsity towards the belongingness of autism. The neutrosophic clustering is empowered by the sparrow search optimization algorithm (SSOA). The clustering and the cluster centroids are done by applying SSOA to determine the best centroids by computing fitness value of each instances based on their influence in determining the patterns of autism and non-autism categories. The kernel function overcomes the issue of handling vague information in defining the membership degree of neutrosophic factors. The experimental results show that the performance of empowered neutrosophic clustering integrated with sparrow search optimization algorithm (ENC-SSOA) produced highest clustering accuracy compared with other conventional clustering models in autism prediction.

References

1. Fadi Thabtah, Firuz Kamalov, Khairan Rajab, A new computational intelligence approach to detect autistic features for autism screening, *International Journal of Medical Informatics*, Volume 117, September 2018, Pages 112-124
2. Usta, M.B, Karabekiroglu K, Sahin, B.Aydin, M. Bozkurt, A. Karaosman, T. Aral, A. Cobanoglu, C. Kurt, A.D Kesim, N, Use of machine learning methods in prediction of short-term outcome in autism spectrum disorders. *Psychiatry Clin. Psychopharmacol.*29, 320–325, 2019.
3. AnirudhR., ThiagarajanJ. J, Bootstrapping graph convolutional neural networks for autism spectrum disorder classification, *International Conference on Acoustics, Speech and Signal Processing*, 3197–3201, 2019.
4. TovarA.E.,Rodriguez-GranadosA.,Arias-TrejoN, Atypical shape bias and categorization in autism: Evidence from children and computational simulations. *Developmental Science*,23(2), 2020.
5. Ming Jiang, Qi Zhao, Learning Visual Attention to Identify People with Autism Spectrum Disorder *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3267-3276
6. Stevens, E., Dixon, D. R., Novack, M. N., Granpeesheh, D., Smith, T, Linstead, E, Identification and analysis of behavioral phenotypes in autism spectrum disorder via unsupervised machine learning. *International journal of medical informatics*, 129, 29–36, 2019
7. Suman Raj,Sarfaraz Masood,Analysis and Detection of Autism Spectrum Disorder Using Machine Learning Techniques, *International Conference on Computational Intelligence and Data Science*, *Procedia Computer Science* 167 (2020) 994–1004.

8. Rajab K.D., Padmavathy A., Thabtah F, Machine Learning Application for Predicting Autistic Traits in Toddlers. Arab J Sci Eng 46, 3793–3805 (2021).
9. Rasool Azeem Musa, Mehdi Ebady Manaa, Ghassan Abdul-Majeed, Predicting Autism Spectrum Disorder (ASD) for Toddlers and Children Using Data Mining Techniques Journal of Physics: Conference Series 1804 (2021) 012089.
10. Chelsea M. Parlett-Pelleriti¹ · Elizabeth Stevens, Dennis Dixon, Erik J. Linstead, Review Journal of Autism and Developmental Disorders, Applications of Unsupervised Machine Learning in Autism Spectrum Disorder Research: a Review, Review Journal of Autism and Developmental Disorders, Vol.:(0123456789), 2022
11. <https://www.kaggle.com/fabdelja/autism-screening-for-toddlers>
12. C.Bhuvanewari, Dr.A.Anbarasi (2021), An Improvised Multilayer Perceptron Network Using Boosted Regression Tree Based Missing Value Imputation And Fuzzy Backward Elimination Feature Selection For Autism Disease Prediction , Volume: 5 Issue 2 , 2119 – 2131,
13. Smarandache, Florentin; Mohamed Abdel-Basset; and Said Broumi, Neutrosophic Sets and Systems, vol. 47, 2021, *Neutrosophic Sets and Systems* 47, 1 (2021).
14. Smarandache, Florentin, Broumi, Said (eds.) (2020). *Neutrosophic Theories in Communication, Management and Information Technology*. New York: Nova Science Publishers.
15. Gharehchopogh, F.S., Namazi, M., Ebrahimi L, Advances in Sparrow Search Algorithm: A Comprehensive Survey, Arch Computat Methods Eng 30, 427–455 (2023)
16. Awadallah, M.A., Al-Betar, M.A., Doush, I.A, Recent Versions and Applications of Sparrow Search Algorithm. Arch Computat Methods Eng 30, 2831–2858 (2023)