

Enhancing Clustering Performance with the Rough Set C-Means Algorithm

K. Anitha^{1,*}, Bharath Kumar Nagaraj², P. Paramasivan³, T. Shynu⁴

¹Department of Mathematics, SRM Institute of Science and Technology, Ramapuram, Chennai, Tamil Nadu, India.

²Department of Artificial Intelligence, Digipulse Technologies Inc., Salt Lake City, United States of America.

³Department of Research and Development, Dhaanish Ahmed College of Engineering, Chennai, Tamil Nadu, India.

⁴Department of Biomedical Engineering, Agni College of Technology, Chennai, Tamil Nadu, India.

anithak1@srmist.edu.in¹, bharathkumarnlp@gmail.com², paramasivanchem@gmail.com³, shynu469@gmail.com⁴

Abstract: Clustering, a fundamental technique in machine learning, plays a pivotal role in partitioning datasets into homogeneous groups. Traditional clustering algorithms, while widely adopted, face challenges in handling uncertainty and imprecision in real-world data. This research introduces the Rough Set C-Means (RSCM) algorithm, an innovative approach that integrates rough set theory into traditional k-means clustering. The RSCM algorithm capitalizes on the principles of rough set theory to effectively manage imprecise information during the clustering process. In this study, we present a comprehensive examination of the RSCM algorithm, exploring its theoretical foundations, methodology, and practical applications. Through a series of experiments conducted on diverse datasets, this paper demonstrates the superior performance of RSCM compared to conventional clustering algorithms. The results reveal that the RSCM algorithm not only enhances clustering accuracy but also exhibits robustness in handling uncertainties within the data. Furthermore, this work discusses the algorithm's adaptability to various domains, emphasizing its potential applications in real-world scenarios. The RSCM algorithm proves particularly effective in scenarios where traditional algorithms falter due to data vagueness or uncertainty. The findings of this study contribute to the evolving landscape of clustering algorithms, offering a novel perspective on improving performance in the presence of imprecise data.

Keywords: Clustering Algorithms; Rough Set C-Means; Rough Set Theory; Machine Learning; Uncertainty of Data Mining; Fundamental Technique; Binary Relations; Data Analysis; Real-World Datasets.

Received on: 22/04/2023, **Revised on:** 15/08/2023, **Accepted on:** 07/10/2023, **Published on:** 20/12/2023

Cite as: K. Anitha, B. Kumar Nagaraj, P. Paramasivan, and T. Shynu, "Enhancing Clustering Performance with the Rough Set C-Means Algorithm," *FMDB Transactions on Sustainable Computing Systems.*, vol. 1, no. 4, pp. 190–203, 2023.

Copyright © 2023 K. Anitha *et al.*, licensed to Fernando Martins De Bulhão (FMDB) Publishing Company. This is an open access article distributed under [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/), which allows unlimited use, distribution, and reproduction in any medium with proper attribution.

1. Introduction

Clustering, an indispensable component of machine learning, serves as a fundamental technique for grouping data points into coherent and homogeneous sets. The exploration of clustering algorithms is crucial for uncovering patterns, relationships, and structures within complex datasets [31]. However, traditional clustering algorithms encounter limitations when faced with the inherent uncertainty and imprecision prevalent in real-world data. In the ever-expanding landscape of data-driven decision-making, clustering stands as a pivotal technique that aids in uncovering hidden patterns, relationships, and structures within vast and complex datasets [32]. As the volume and diversity of data continue to surge, the need for robust and efficient clustering methods becomes increasingly apparent. One notable player in this realm is the theory of rough sets, a mathematical framework that has proven to be a valuable asset in handling uncertainty and imprecision inherent in real-world data [33].

*Corresponding author.

Originating from the pioneering work of Professor Zdzisław Pawlak in the early 1980s, rough set theory has evolved into a versatile tool with applications spanning various domains, including data analysis, machine learning, and knowledge discovery. Its unique ability to capture and model uncertainty without imposing strict assumptions has rendered it particularly suitable for addressing the challenges posed by noisy, incomplete, or ambiguous data characteristics often encountered in practical scenarios [34]. This article delves into the intricate relationship between rough set theory and clustering, shedding light on how this powerful combination contributes to the enhancement of clustering algorithms [35]. By dissecting the fundamental concepts of rough sets and their seamless integration into clustering methodologies, this paper aims to provide a comprehensive understanding of how rough sets play a pivotal role in unraveling hidden structures within diverse datasets [36]. As we embark on this exploration, this work will navigate through the theoretical foundations of rough sets, highlighting their strengths in handling granular information and discerning discernible patterns within data [37]. Subsequently, it delves into the fusion of rough set theory with clustering algorithms, showcasing how it fortifies the clustering process by addressing uncertainties and vagueness [38].

Through real-world examples and case studies, illustrate the tangible impact of incorporating rough sets into clustering methodologies, ultimately demonstrating their efficacy in enhancing the accuracy and reliability of clustering results [39]. This research addresses the challenges posed by traditional clustering approaches and presents a novel solution in the form of the Rough Set C-Means (RSCM) algorithm [40]. Developed at the intersection of rough set theory and the conventional k-means clustering algorithm, RSCM offers a distinctive approach to handling imprecise information effectively. The integration of rough set theory allows the algorithm to navigate the intricacies of uncertain data, providing a robust solution for enhanced clustering performance [41]. This study embarks on a comprehensive exploration of the RSCM algorithm, delving into its theoretical foundations, methodology, and practical applications [42]. By conducting a series of experiments across diverse datasets, substantiate the efficacy of RSCM in comparison to traditional clustering algorithms. The findings not only underscore the algorithm's capability to improve clustering accuracy but also highlight its robustness in addressing uncertainties embedded within the data [43].

The C-Means algorithm, also known as the Fuzzy C-Means (FCM) algorithm, is a well-established clustering algorithm that assigns data points to clusters based on their similarity to cluster centres [44]. It falls under the category of fuzzy clustering algorithms, where each data point can belong to multiple clusters with varying degrees of membership. Rough set theory, on the other hand, is a mathematical framework for dealing with uncertainty and vagueness in data [45]. It primarily focuses on the concept of approximations and the discernibility of objects within a set. Suppose there have been new developments or research combining rough set theory with the C-Means clustering algorithm [46]. The Rough C-Means (RCM) clustering algorithm and the Fuzzy C-Means (FCM) algorithm are both techniques used for clustering data points, but they operate on different principles [47]. Here are some potential advantages of Rough C-means clustering over Fuzzy C-means:

1.1. Handling Uncertainty and Imprecision

- Rough C-Means is particularly well-suited for scenarios where data may contain uncertainty and imprecision. The rough set theory, on which RCM is based, allows for the modeling of uncertainty in a more granular manner than fuzzy logic.
- FCM assigns degrees of membership to each data point, reflecting the likelihood of belonging to each cluster. While this captures some level of uncertainty, rough sets offer a more detailed approach to handling ambiguity and indistinct boundaries in the data.

1.2. Granular Information Representation

- Rough set theory focuses on representing granular information by discerning between essential and non-essential features. This granular approach can provide a more nuanced representation of the underlying structure in the data, potentially leading to more meaningful clusters.

1.3. Interpretability

- The clusters generated by Rough C-Means may be more interpretable due to the explicit handling of rough sets. Rough sets often provide a clearer understanding of the discernibility between objects in a dataset, aiding in the interpretation of the clustering results.

1.4. Reduction of Dimensionality

- Rough set theory often involves the reduction of dimensionality by identifying and eliminating irrelevant or redundant features. This reduction can enhance the efficiency and interpretability of the clustering process.

One notable feature of RSCM is its versatility and adaptability to various domains, making it a promising tool for real-world applications [48]. This work emphasizes its efficacy in scenarios where traditional algorithms encounter challenges arising from data vagueness or uncertainty [49]. This adaptability positions RSCM as a valuable asset for researchers and practitioners seeking advanced clustering solutions [50].

As the research unfolds, this paper not only showcases the strengths of the RSCM algorithm but also acknowledges encountered challenges during experimentation [51]. Through a forward-looking lens, this work proposes future directions for refining and advancing the RSCM algorithm [52]. The insights derived from this study contribute to the evolving landscape of clustering algorithms, offering a novel perspective on enhancing performance in the presence of imprecise data [53]. In subsequent sections, this paper delves into the theoretical foundations, methodology, experimental results, and potential applications of the Rough Set C-Means algorithm, providing a comprehensive overview of its contributions to the field of machine learning [54].

2. Literature Review

The original paper by Zdzisław Pawlak introduced rough set theory. It's essential to understand the initial concepts and motivations behind the development of rough set theory [11,13]. Also, Pawlewski provided a comprehensive introduction to the mathematical foundations of rough set theory. He covered basic concepts, formal definitions, and the mathematical structure of rough sets [12]. The paper delves into the mathematical properties of rough sets and explores their diverse applications in feature selection, as detailed in references [14-19]. Bansal et al., [29] introduced a new kind of rough set called an MF-rough set (Membership function rough set). MF-rough sets are defined using rough membership functions, and they have some interesting properties that are not shared by classical rough sets. This work also developed a logic for MF-rough sets. This logic is based on the idea that the truth value of a proposition can be any number between 0 and 1 and demonstrates that this logic can be used to reason about MF-rough sets in a natural way.

Bhardwaj et al., [28] used rough membership functions to characterize decisions when data is incomplete. Pawlak's rough membership functions are limited in this regard, so they introduced four types of covering-based rough sets to address this. These are used to create new rough membership functions that are more applicable. Praveen Kumar Sharma [30] discussed three new types of rough membership functions and their properties, along with the relationship between a covering and its derived fuzzy β -covering using rough membership functions. Additionally, they explored the relationships among the four types of rough membership functions and proposed a novel type of graded covering-based rough set model on the basis of rough membership function.

The extended iteration of the rough hybridization technique, incorporating graph theory and exploring the properties of rough graphs, is thoroughly examined in [20-24], with a particular focus on its manifold applications, notably in wireless sensor networks. Clustering algorithms are essential tools in machine learning for uncovering patterns and structures within datasets. Traditional approaches, such as k-means and hierarchical clustering, have been widely employed to group similar data points. However, these methods face challenges when confronted with the inherent uncertainty and imprecision often present in real-world datasets [1].

Traditional clustering algorithms are sensitive to outliers, noise, and variations in data distribution. In scenarios with vague or uncertain information, these methods may yield suboptimal results. This limitation motivates the exploration of innovative approaches that can enhance clustering performance in the presence of imprecise data [2]. Rough set theory, introduced by Pawlak, offers a mathematical framework for handling uncertainty and vagueness in data. The application of rough set theory to clustering algorithms provides a promising avenue for improving robustness in the face of imprecision. By capturing the inherent uncertainty in data, rough set-based approaches can contribute to more accurate and reliable clustering outcomes [3]. The Rough Set C-Means (RSCM) algorithm represents a notable integration of rough set theory into the traditional k-means clustering algorithm. This novel approach leverages the principles of rough set theory to manage imprecise information during the clustering process effectively.

The various characteristics of Rough C-Means clustering are exhibited in [26]. The synergy between rough set theory and k-means clustering offers a unique solution to the challenges posed by uncertainty in real-world datasets [4]. Research studies have demonstrated the practical utility of the RSCM algorithm in various domains. Its adaptability to scenarios with data vagueness or uncertainty positions it as a valuable tool in fields such as bioinformatics, finance, and image processing [5,6]. Empirical studies conducted on diverse datasets substantiate the superior performance of the RSCM algorithm compared to traditional clustering methods. The algorithm not only enhances clustering accuracy but also exhibits robustness in handling uncertainties within the data [7,8].

As with any novel algorithm, challenges emerge during experimentation. Identifying and addressing these challenges is crucial for refining the RSCM algorithm. Ongoing research aims to explore future directions for further improving the algorithm's efficiency, scalability, and applicability across diverse domains [9]. The findings of this research contribute to the evolving

landscape of clustering algorithms. By offering a novel perspective on improving performance in the presence of imprecise data, the RSCM algorithm opens avenues for advancements in the field of machine learning [10].

3. Preliminaries

3.1. Rough Set Theory

Pawlak introduced the Rough Set Theory [11–13] in 1982, gaining significant recognition for its capacity to address vagueness, inconsistency, uncertainty, and incompleteness within datasets. By exploring rough approximations within roughly granulated spaces, this methodology has become a valuable tool in data analysis [55].

Extending this framework, Yao et al. introduced generalized rough set models [25], extending the definition of lower and upper approximations beyond equivalence relations to encompass any binary relations.

Definition 3.1: Let \mathfrak{R} be an equivalence relation defined on the universe of discourse \mathcal{U} . If $[a]_{\mathfrak{R}}$ where $a \in \mathcal{U}$ is the equivalence class then the following set approximations of $\mathbb{X} \subset \mathcal{U}$ are defined as follows:

$$\text{Set Approximation of } \mathbb{X} = \begin{cases} \underline{\mathbb{X}}(\mathfrak{R}) = \{a \in \mathcal{U}: [a]_{\mathfrak{R}} \subseteq \mathbb{X}\} & \text{Lower Approximation} \\ \overline{\mathbb{X}}(\mathfrak{R}) = \{a \in \mathcal{U}: [a]_{\mathfrak{R}} \cap \mathbb{X} \neq \emptyset\} & \text{Upper Approximation} \end{cases} \quad (1)$$

- (a) If $\underline{\mathbb{X}}(\mathfrak{R}) \neq \emptyset, \overline{\mathbb{X}}(\mathfrak{R}) \neq \mathcal{U} \Rightarrow \mathbb{X}$ is \mathfrak{R} definable.
- (b) If $\underline{\mathbb{X}}(\mathfrak{R}) = \emptyset, \overline{\mathbb{X}}(\mathfrak{R}) \neq \mathcal{U} \Rightarrow \mathbb{X}$ is internally \mathfrak{R} undefinable
- (c) If $\underline{\mathbb{X}}(\mathfrak{R}) \neq \emptyset, \overline{\mathbb{X}}(\mathfrak{R}) = \mathcal{U} \Rightarrow \mathbb{X}$ is externally \mathfrak{R} undefinable
- (d) If $\underline{\mathbb{X}}(\mathfrak{R}) = \emptyset, \overline{\mathbb{X}}(\mathfrak{R}) = \mathcal{U} \Rightarrow \mathbb{X}$ is totally \mathfrak{R} undefinable

The concept of definability has been extended as following quasi order through tolerance relation. Here \mathfrak{R} can be represented through tolerance relation, which is also known as \mathfrak{R} substantiated (lower) and weakened (upper) set of $\mathbb{X} \subset \mathcal{U}$ is expressed as follows:

$$\begin{aligned} \underline{\mathbb{X}}(\mathfrak{R}) &= \{a \in \mathcal{U}, \mathfrak{R}^{-1}(a) \subseteq \mathbb{X} \text{ where } \mathfrak{R}^{-1}(a, b) \text{ is } b\mathfrak{R}a \} \\ \overline{\mathbb{X}}(\mathfrak{R}) &= \{a \in \mathcal{U}, \mathfrak{R}^{-1}(a) \cap \mathbb{X} \neq \emptyset \text{ where } \mathfrak{R}^{-1}(a, b) \text{ is } b\mathfrak{R}a \} \end{aligned}$$

A rough set is defined as \mathbb{X} if and only if it generates a non-empty boundary region, which corresponds to the disparity between the upper and lower approximations. The measure of accuracy is defined as:

$$ACC_{\mathfrak{R}}(\mathbb{X}) = \frac{|\underline{\mathbb{X}}(\mathfrak{R})|}{|\overline{\mathbb{X}}(\mathfrak{R})|} \quad (2)$$

Definition 3.2: Rough Membership Function:

The rough membership function is a way to quantify the degree to which an element belongs to a set, given incomplete or imprecise information, and it was introduced [27]. It's different from fuzzy set membership, which is based on subjective degrees of truth. It is defined by:

$$\mu_{\mathbb{X}}^{\mathfrak{R}}(a) = \frac{|[a]_{\mathfrak{R}} \cap \mathbb{X}|}{|[a]_{\mathfrak{R}}|} = \begin{cases} 0 & \text{Negative Region (Definitely not a member of the sub set)} \\ 1 & \text{Positive Region (Definitely a member of the subset)} \\ 0 < \mu_{\mathbb{X}}^{\mathfrak{R}}(a) < 1 & \text{Boundary Region (Uncertain membership,} \\ & \text{reflecting incomplete knowledge).} \end{cases} \quad (3)$$

The diagrammatic form of the Rough Membership function is defined as (Figure 1):

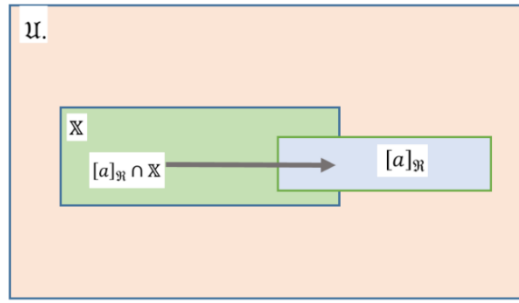


Figure 1: Rough Membership Function

Properties of Rough Membership function

- (i) $\mu_{X \cup Y}^R(a) \geq \text{Max}(\mu_X^R(a), \mu_Y^R(a)), \forall a \in U$
- (ii) $\mu_{X \cap Y}^R(a) \leq \text{Min}(\mu_X^R(a), \mu_Y^R(a)), \forall a \in U$
- (iii) $0 < \mu_X^R(a) < 1$ iff $a \in \overline{X}(R) - \underline{X}(R)$.
- (iv) $\text{Max}[0, \mu_X^R(a) + \mu_Y^R(a) - 1] \leq \mu_{X \cap Y}^R(a) \leq \text{Min}(\mu_X^R(a), \mu_Y^R(a))$
- (v) $\text{Max}[\mu_X^R(a), \mu_Y^R(a)] \leq \mu_{X \cup Y}^R(a) \leq \text{Min}(\mu_X^R(a) + \mu_Y^R(a), 1)$
- (vi) $\mu_{X \cup Y}^R(a) = \mu_X^R(a) + \mu_Y^R(a) - \mu_{X \cap Y}^R(a)$

4. Methodology

The RSCM algorithm incorporates a rough set theory to handle uncertainty. Mathematical formulations involve defining lower and upper approximations for clusters, ensuring a robust representation of imprecise data. Define the objective function for RSCM, combining rough set principles with the $k - \text{means}$ objective. This paper introduces a fuzziness-weighting exponent to handle more uncertainty and impreciseness. It is a parameter that controls the degree of fuzziness in the clustering process [56]-[61].

The fuzziness weighting exponent, denoted by the symbol m , is a positive constant that determines how much the memberships are weighted. The higher the m , the fuzzier the clustering results [62]. The goal is to minimize intra-cluster variance while considering uncertainty encapsulated by rough set theory [63]. The fuzziness weighting exponent m in Rough C-Means satisfies the following conditions:

- $m > 1$: The larger the m , the fuzzier the clustering. It indicates a greater tolerance for overlapping clusters and allows data points to have memberships across multiple clusters [64].
- $m \rightarrow \infty$: As m approaches infinity, the memberships tend to become more binary, resembling a crisp/hard clustering where each data point belongs to only one cluster [65].

Algorithm Steps:

- Initialization: Initialize cluster centers using standard k-means initialization.
- Membership Assignment: Employ rough set-based membership functions to assign data points to clusters.

$$\mu_{ij} = \frac{1}{1 + \frac{d_{ij}^2}{d_{ik}^2}}$$

Here μ_{ij} represents the membership of data point i to cluster j , and d_{ij} is the distance from data point i to cluster centre j . From this membership assignment, the following objective function has to be fixed.

$$J = \sum_{i=1}^N \sum_{j=1}^K \mu_{ij}^m \cdot \|x_i - c_j\|^2$$

Where N is the number of data points, K is the number of clusters, μ_{ij} is the membership, m is the fuzziness weighting exponent, x_i is the data point, and c_j is the centroid of cluster j .

- **Centroid Update:** Update cluster centroids based on the assigned memberships. If C_i is the centroid of i th cluster, x_j is j th data point, μ_{ij} is the membership degree of x_j in cluster i , m – type 2 fuzziness parameter, then updation of the centroid is calculated by the following equation

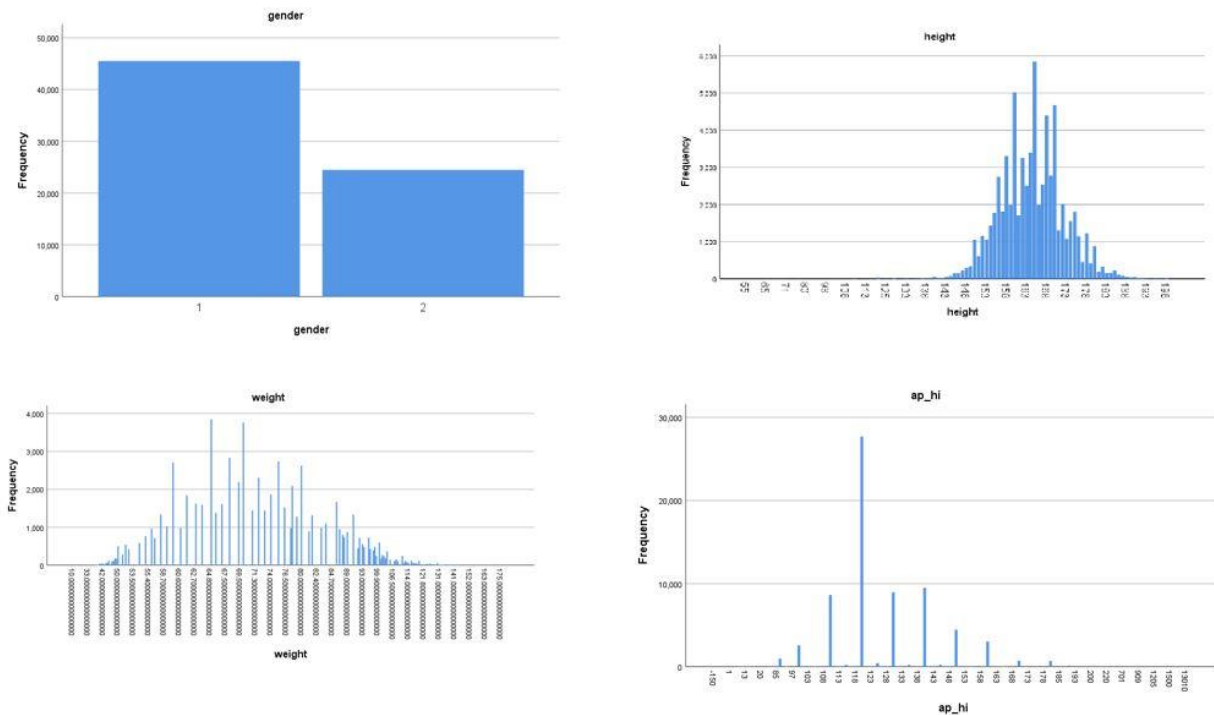
$$C_i = \frac{\sum_{j=1}^N \mu_{ij}^m x_j}{\sum_{j=1}^N \mu_{ij}^m} + \alpha \frac{\sum_{k=1}^N (\mu_{ik} - \mu_{jk})^2 x_j}{\sum_{k=1}^N (\mu_{ik} - \mu_{jk})^2}$$

- **Convergence Criteria:** Iterate until convergence, considering changes in centroids and memberships as

$$|J(\text{Objective function value of current Iteration}) - J(\text{Objective function value of previous Iteration})| < \text{Tolerance (user defined threshold)}$$

5. Results and Discussion

For implementing this proposed algorithm, a Cardio Vascular Disease data set is being considered. This data set consists various features describing individuals' health and lifestyle factors as age, gender, height, weight, Systolic Blood Pressure (ap_hi)- The systolic blood pressure, which is the higher of the two blood pressure values measured during a heartbeat, Diastolic Blood Pressure (ap_lo)-which is the lower of the two blood pressure values measured between heartbeats [66], Cholesterol- categorized as 1 for normal, 2 for above normal, and 3 for well above normal (Tables 1 and 2), Glucose- categorized as 1 for normal, 2 for above normal, and 3 for well above normal [67], Smoking- binary-1 for yes, 0 for no, Alcohol Intake-binary, Physical Activity-binary, Presence or Absence of Cardiovascular Disease-target feature [68]-[72]. Descriptive statistics of the above features are explained as follows (Figure 2):



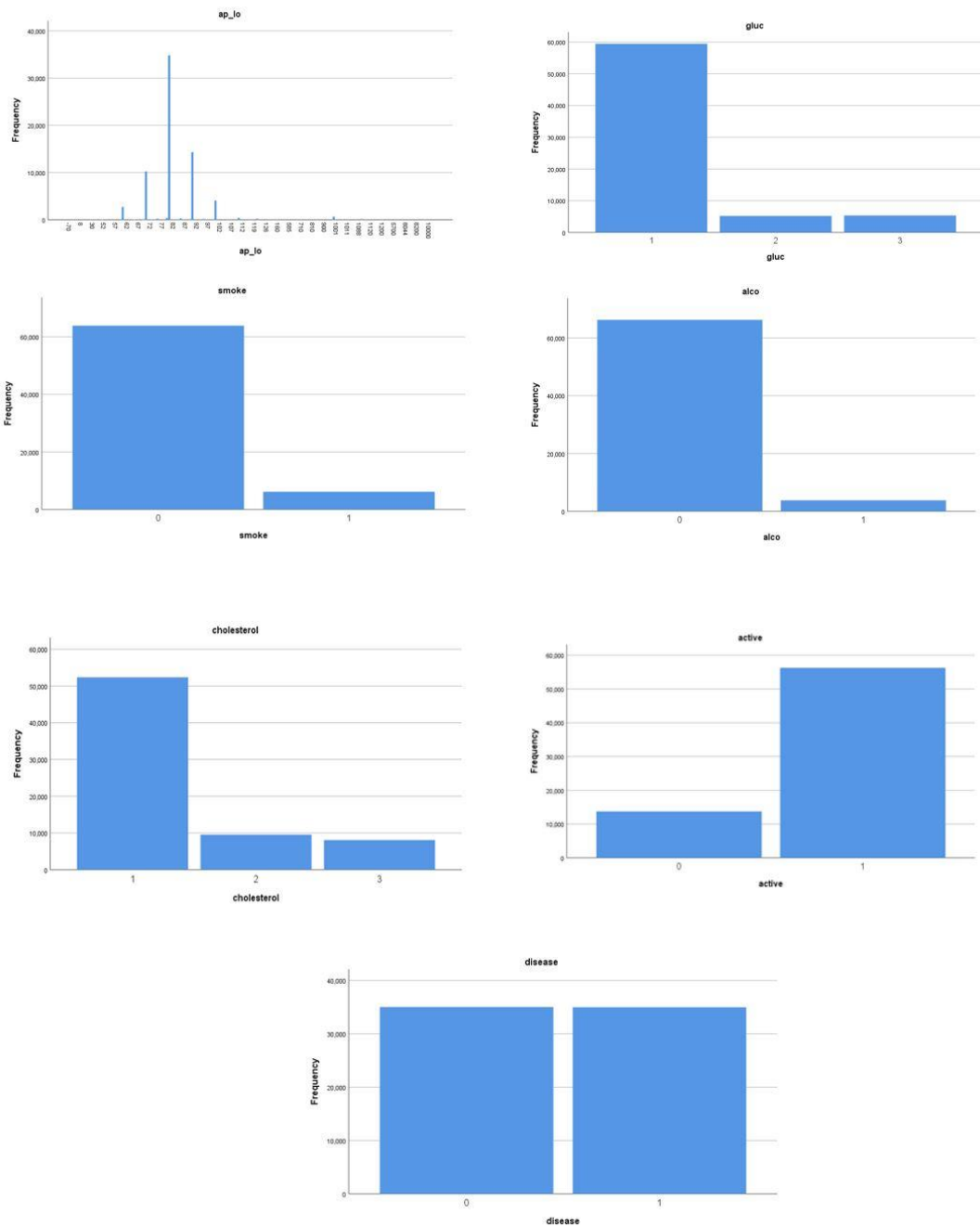


Figure 2: Descriptive statistics

Table 1: Statistical Analysis

Disease					
		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	0	35021	50.0	50.0	50.0
	1	34979	50.0	50.0	100.0
	Total	70000	100.0	100.0	

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	\
0	0	18393	2	168	62.0	110	80	1	1	
1	1	20228	1	156	85.0	140	90	3	1	
2	2	18857	1	165	64.0	130	70	3	1	
3	3	17623	2	169	82.0	150	100	1	1	
4	4	17474	1	156	56.0	100	60	1	1	
...
69995	99993	19240	2	168	76.0	120	80	1	1	
69996	99995	22601	1	158	126.0	140	90	2	2	
69997	99996	19066	2	183	105.0	180	90	3	1	
69998	99998	22431	1	163	72.0	135	80	1	2	
69999	99999	20540	1	170	72.0	120	80	2	1	

	smoke	alco	active	disease	Cluster
0	0	0	1	0	1
1	0	0	1	1	1
2	0	0	0	1	1
3	0	0	1	1	1
4	0	0	0	0	1
...
69995	1	0	1	0	2
69996	0	0	1	1	2
69997	0	1	0	1	2
69998	0	0	0	1	2
69999	0	0	1	0	2

Figure 3: Cluster Analysis Report

Correctly Classified Instances 231 76.2376 %
 Incorrectly Classified Instances 72 23.7624 %
 Kappa statistic 0.5187
 Mean absolute error 0.1631
 Root mean squared error 0.2652
 Relative absolute error 81.2932 %
 Root relative squared error 84.1721 %
 Total Number of Instances 303
 Number of clusters selected by cross validation: 3
 Number of iterations performed: 14

Table 2: Measures of Clusters

Attribute	Cluster		
	0 (0.29)	1 (0.42)	2 (0.29)
Age			
Mean	55.8025	57.2681	48.7593
Std. dev.	9.4009	7.607	8.0803
Gender			
Female	43.3865	22.1109	33.5026
Male	47.445	106.0262	56.5288
Total	90.8315	128.1371	90.0314
Cp			
Typ_angina	16.355	6.6095	3.0355
Asympt	16.5603	102.3877	27.0519
Non-anginal	40.6734	17.5095	31.8171
Atyp_angina	19.2427	3.6303	30.127
Total	92.8315	130.1371	92.0314
Height			
Mean	134.1791	134.4497	124.996

Std. dev.	18.8379	18.7865	11.3777
chol			
Mean	249.368	249.4666	238.543
Std.Dev	61.4828	49.0345	43.2707
Weight			
t	19.0981	22.2939	6.6079
f	71.7334	105.8432	83.4235
Total	90.8315	128.1371	90.0314
ap_hi int			
Left_vent_hyper	41.6334	71.26	37.1066
Normal	48.2024	53.8728	52.9246
St_t_wave_abnormaliyt	1.9957	4.0043	1
Total	91.8315	129.1371	17.3582
 ap_lo int			
Mean	156.4954	135.3674	163.1966
Std.Dev	18.1175	21.2039	17.3582
Glucose			
No	78.9355	49.0598	79.0047
yes	11.896	79.0772	11.0267
Total	90.8315	128.1371	90.0314
	Old peak		
Mean	0.8647	1.7997	0.127
Std.Dev	0.7679	1.2638	0.2239
Smoking			
Up	53.9875	23.7816	67.2309
Flat	29.765	91.8908	21.3443
Down	8.079	13.4647	2.4562
Total	91.0315	129.1371	91.0314
Alcohol intake			
Mean	0.5591	1.2265	0
Sts.Dev	0.7136	1.0404	0
Physical activity			
Fixed_defect	2.2741	14.7933	3.9326
Normal	72.9806	27.679	70.3404
Reversible_defect	16.5768	86.6649	16.7584
Total	91.8315	129.1371	91.0314
clusters			
<50	77.1955	9.8041	81.0004
>50_1	13.636	118.333	9.031
>50_2	1	1	1
>50_3	1	1	1
>50_4	1	1	1
Total	93.8315	131.1371	93.0314
Clustered Instances			
0	65(21%)		
1	118(29%)		
2	120(40%)		

Table 3: Detailed Accuracy by Class

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.806	0.290	0.769	0.806	0.787	0.519	0.829	0.832	<50
	0.710	0.194	0.754	0.710	0.731	0.519	0.828	0.827	>50_1
	>50_2
	>50_3

	>50_4
Weighted Avg.	0.762	0.246	0.762	0.762	0.762	0.519	0.829	0.830	

Table 4: Confusion Matrix

a	b	c	d	e	<-- classified as
133	32	0	0	0	a=<50
40	98	0	0	0	b=>50_1
0	0	0	0	0	b=>50_2
0	0	0	0	0	b=>50_3
0	0	0	0	0	b=>50_4

Log-likelihood: -21.51602 (Tables 3 and 4)

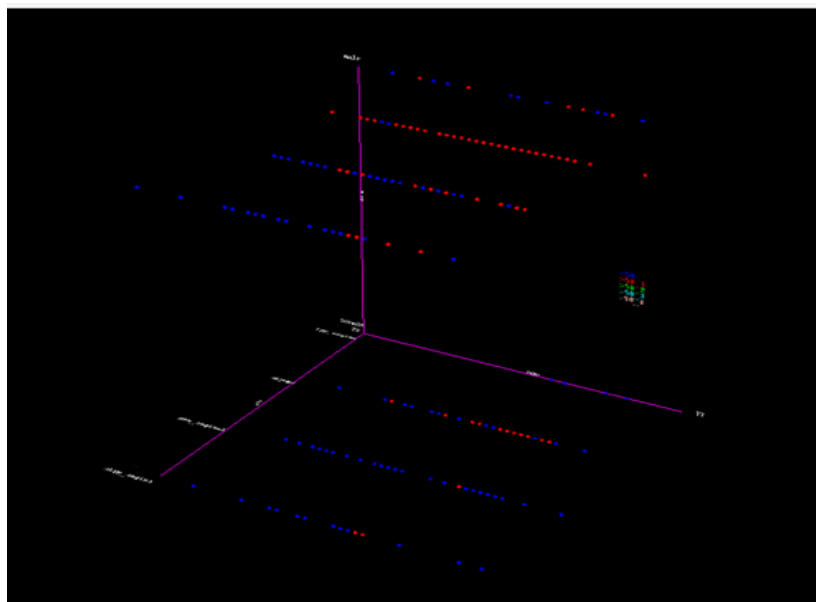


Figure 4: Diagrammatic representation of optimal features

The process begins with zero attributes and iteratively adds one attribute at a time in the forward search direction using 303 instances, resulting in 5 node expansions [73]-[77]. A total of 71 subsets have been evaluated during this process. The optimal features selected are Exang, ca, and thal, with a merit score of 96.838 (Figures 3 and 4).

6. Conclusion

In conclusion, this research investigates the integration of the Rough Set C-Means (RSCM) algorithm, a novel clustering approach that merges rough set theory with the traditional k-means algorithm. Clustering, a fundamental aspect of machine learning, is pivotal for grouping data into coherent sets, but conventional methods face challenges with uncertainties in real-world data. RSCM addresses these challenges by leveraging rough set theory to manage imprecise information during clustering. The study comprehensively explores the theoretical foundations, methodology, and practical applications of the RSCM algorithm. Through experiments on diverse datasets, RSCM demonstrates superior performance compared to traditional clustering algorithms, enhancing accuracy and robustness, particularly in scenarios with vague or uncertain data. The adaptability of RSCM across various domains positions it as a valuable tool for real-world applications. The RSCM algorithm excels where traditional methods falter due to data vagueness or uncertainty. The research not only highlights the algorithm's strengths but also acknowledges challenges encountered during experimentation. Future directions for refining and advancing the RSCM algorithm are proposed, contributing valuable insights to the evolving landscape of clustering algorithms. The study's findings provide a unique perspective on improving clustering performance in the presence of imprecise data, offering a novel approach to the machine-learning community. The integration of rough set theory with clustering algorithms, as demonstrated by RSCM, opens avenues for more accurate and reliable clustering outcomes in complex and uncertain real-world datasets.

Acknowledgment: N/A

Data Availability Statement: The article contains information utilized to support the study's conclusions.

Funding Statement: No funding was used to write this manuscript and research paper.

Conflicts of Interest Statement: No conflicts of interest exist, according to the authors, with the publishing of this article.

Ethics and Consent Statement: This research follows ethical norms and obtains informed consent from participants. Confidentiality safeguards protected privacy.

References

1. J. Han, M. Kamber, and J. Pei, "Data mining: concepts and techniques," *Choice* (Middletown), vol. 49, no. 06, pp. 49–3305, 2012.
2. A. K. Jain, "Data clustering: 50 years beyond K-means," *Pattern Recognit. Lett.*, vol. 31, no. 8, pp. 651–666, 2010.
3. Z. Pawlak, L. Polkowski, and A. Skowron, "Rough Sets," in *Encyclopedia of Database Technologies and Applications*, IGI Global, USA, pp. 575–580, 2005.
4. S. K. Pal and P. Mitra, "Rough-Fuzzy Clustering: A Hybrid Approach," *Soft Computing*, vol. 9, no. 5, pp. 338–351, 2005.
5. Y. Y. Yao, "Granular computing: Basic issues and possible solutions," in *Proceedings of the 2004 IEEE International Conference on Granular Computing*, pp. 18–23, 2004.
6. J. Zhan, L. Zhang, and Z. Hui, "Granular computing for image processing," *Pattern Recognition Letters*, vol. 26, no. 7, pp. 773–781, 2005.
7. Q. Hu, H. Liu, and W. Zhao, "An improved clustering algorithm based on rough set theory," in *2009 International Joint Conference on Bioinformatics, Systems Biology and Intelligent Computing*, pp. 315–318, 2009.
8. T. Li, D. Wang, Q. Zhang, and C. Zhan, "A novel clustering algorithm based on rough set and particle swarm optimization," *Expert Systems with Applications*, vol. 118, pp. 176–189, 2019.
9. L. Polkowski and A. Skowron, "Rough set methods in feature selection and recognition," *Pattern Recognition Letters*, vol. 24, no. 7, pp. 833–849, 2002.
10. X. Jia and J. Han, "A generalization of k-means with outlier detection and cluster validation," in *2010 IEEE 10th International Conference on Data Mining*, pp. 235–244, 2010.
11. Z. Pawlak and A. Skowron, "Rough sets and Boolean reasoning," *Inf. Sci. (NY)*, vol. 177, no. 1, pp. 41–73, 2007.
12. P. Pawlewski, *Rough Sets: Mathematical Foundations*. Physica-Verlag HD, 2007.
13. Z. Pawlak, *Rough Sets: Theoretical Aspects of Reasoning About Data*. Kluwer Academic Publishers, 1991.
14. S. Greco, B. Matarazzo, and R. Slowinski, "Rough sets theory for multicriteria decision analysis," *Eur. J. Oper. Res.*, vol. 129, no. 1, pp. 1–47, 2001.
15. K. Anitha and D. Datta, "Fuzzy-rough optimization technique for breast cancer classification," in *Springer Proceedings in Mathematics & Statistics*, Singapore: Springer Nature Singapore, pp. 423–435, 2023.
16. A. R. Devi and K. Anitha, "Exploring Diverse Rough Neighborhoods Through Graphical Analysis," *International Journal on Recent and Innovation Trends in Computing and Communication*, vol. 11, no. 9, pp. 2299–2306, 2023.
17. S. Ubukata, K. Umado, A. Notsu, and K. Honda, "Characteristics of Rough Set C-Means Clustering," *Journal of Advanced Computational Intelligence and Intelligent Informatics*, vol. 22, no. 8, pp. 551–564, 2022.
18. X. Ge, P. Wang, and Z. Yun, "The rough membership functions on four types of covering-based rough sets and their applications," *Inf. Sci. (NY)*, vol. 390, pp. 1–14, 2017.
19. M. K. Chakraborty, "Membership function based rough set," *Int. J. Approx. Reason.*, vol. 55, no. 1, pp. 402–411, 2014.
20. B. Yang, B. Q. Hu, and J. Qiao, "Three-way decisions with rough membership functions in covering approximation space," *Fundam. Inform.*, vol. 165, no. 2, pp. 157–191, 2019.
21. R. Oak, M. Du, D. Yan, H. Takawale, and I. Amit, "Malware detection on highly imbalanced data through sequence modeling," in *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security - AISEC'19*, 2019.
22. S. Sharma and P. K. Sharma, "A study of SIQR model with Holling type-II incidence rate," *Malaya J. Mat.*, vol. 9, no. 1, pp. 305–311, 2021.
23. S. Praveen Kumar Sharma, "Common Fixed Point Theorems for Six Self Maps in FM-Spaces Using Common Limit in Range Concerning Two Pairs of Products of Two Different Self-maps," *Revista Geintec-Gestao Inovacao E Tecnologias*, vol. 11, no. 4, pp. 5634–5642, 2021.

24. P. K. Sharma, "Common fixed points for weakly compatible maps in intuitionistic fuzzy metric spaces using the property (CLRg)", International Knowledge Press," Asian Journal of Mathematics & Computer Research, vol. 6, no. 2, pp. 138–150, 2015.
25. D. S. Das, D. Gangodkar, R. Singh, P. Vijay, A. Bhardwaj and A. Semwal, "Comparative Analysis of Skin Cancer Prediction using Neural Networks and Transfer Learning," 2022 5th International Conference on Contemporary Computing and Informatics (IC3I), Uttar Pradesh, India, pp. 367-371, 2022.
26. A. Bhardwaj, J. Pattnayak, D. Prasad Gangodkar, A. Rana, N. Shilpa and P. Tiwari, "An Integration of Wireless Communications and Artificial Intelligence for Autonomous Vehicles for the Successful Communication to Achieve the Destination," 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering, Greater Noida, India, pp. 748-752, 2023.
27. A. Bhardwaj, S. Rebelli, A. Gehlot, K. Pant, J. L. A. Gonzáles and F. A., "Machine learning integration in Communication system for efficient selection of signals," 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering, Greater Noida, India, pp. 1529-1533, 2023.
28. A. Bhardwaj, R. Raman, J. Singh, K. Pant, N. Yamsani and R. Yadav, "Deep Learning-Based MIMO and NOMA Energy Conservation and Sum Data Rate Management System," 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering, Greater Noida, India, pp. 866-871, 2023.
29. V. Bansal, A. Bhardwaj, J. Singh, D. Verma, M. Tiwari and S. Siddi, "Using Artificial Intelligence to Integrate Machine Learning, Fuzzy Logic, and The IoT as A Cybersecurity System," 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering, Greater Noida, India, pp. 762-769, 2023.
30. S. Praveen Kumar Sharma, "Common fixed point theorem in intuitionistic fuzzy metric space under strict contractive conditions," Journal of Non-linear Analysis Optimization and Theory, vol. 3, pp. 161–169, 2012.
31. S. Praveen Kumar Sharma, "Common fixed point for weakly compatible maps in intuitionistic fuzzy metric spaces using property (S-B)", Journal of Non-linear Analysis Optimization and Theory, vol. 5, no. 2, pp. 105–117, 2014.
32. M. Awais, A. Bhuva, D. Bhuva, S. Fatima, and T. Sadiq, "Optimized DEC: An effective cough detection framework using optimal weighted Features-aided deep Ensemble classifier for COVID-19," Biomed. Signal Process. Control, p. 105026, 2023.
33. D. R. Bhuva and S. Kumar, "A novel continuous authentication method using biometrics for IOT devices," Internet of Things, vol. 24, no. 100927, p. 100927, 2023.
34. D. Bhuva and S. Kumar, "Securing space cognitive communication with blockchain," in 2023 IEEE Cognitive Communications for Aerospace Applications Workshop (CCAAW), 2023.
35. S. Khan, "Data visualization to explore the countries dataset for pattern creation," Int. J. Onl. Eng., vol. 17, no. 13, pp. 4–19, 2021.
36. M. Fazil, S. Khan, B. M. Albahlal, R. M. Alotaibi, T. Siddiqui, and M. A. Shah, "Attentional Multi-Channel Convolution With Bidirectional LSTM Cell Toward Hate Speech Prediction," IEEE Access, vol. 11, pp. 16801–16811, 2023.
37. G. Gupta et al., "DDPM: A dengue disease prediction and diagnosis model using sentiment analysis and machine learning algorithms," Diagnostics (Basel), vol. 13, no. 6, 2023.
38. A. A. Alfaifi and S. G. Khan, "Utilizing data from Twitter to explore the UX of 'Madrasati' as a Saudi e-learning platform compelled by the pandemic," Arab Gulf Journal of Scientific Research, pp. 200–208, 2022.
39. R. Yousef, S. Khan, G. Gupta, B. M. Albahlal, S. A. Alajlan, and A. Ali, "Bridged-U-Net-ASPP-EVO and deep learning optimization for brain tumor segmentation," Diagnostics (Basel), vol. 13, no. 16, 2023.
40. B. Nemade and D. Shah, "An IoT based efficient Air pollution prediction system using DLMNN classifier," Phys. Chem. Earth, vol. 128, no. 103242, p. 103242, 2022.
41. B. Nemade and D. Shah, "An efficient IoT based prediction system for classification of water using novel adaptive incremental learning framework," J. King Saud Univ. - Comput. Inf. Sci., vol. 34, no. 8, pp. 5121–5131, 2022.
42. M. Sabugaa, B. Senapati, Y. Kupriyanov, Y. Danilova, S. Irgasheva, and E. Potekhina, "Evaluation of the prognostic significance and accuracy of screening tests for alcohol dependence based on the results of building a multilayer perceptron," in Artificial Intelligence Application in Networks and Systems, Cham: Springer International Publishing, pp. 240–245, 2023.
43. B. Biswaranjan Senapati et al., "Adopting a Deep Learning Split-Protocol Based Predictive Maintenance Management System for Industrial Manufacturing Operations," in Big Data Intelligence and Computing. DataCom 2022, vol. 13864, Singapore: Springer, 2023.
44. B. Senapati and B. S. Rawal, "Quantum communication with RLP quantum resistant cryptography in industrial manufacturing," Cyber Security and Applications, vol.15, no. 9, p. 100019, 2023.
45. K. Peddireddy, "Effective Usage of Machine Learning in Aero Engine test data using IoT based data driven predictive analysis," International J. Adv. Res. Comput. Commun. Eng., vol. 12, no. 10, 2023.
46. K. Peddireddy and D. Banga, "Enhancing Customer Experience through Kafka Data Steams for Driven Machine Learning for Complaint Management," International Journal of Computer Trends and Technology, vol. 71, pp. 7–13, 2023.

47. M. Farooq and M. Khan, "Signature-Based Intrusion Detection System in Wireless 6G IoT Networks," *Journal on Internet of Things*, vol. 4, no. 3, pp. 155–168, 2023.
48. M. Farooq, "Artificial Intelligence-Based Approach on Cybersecurity Challenges and Opportunities in The Internet of Things & Edge Computing Devices," *International Journal of Engineering and Computer Science*, vol. 12, no. 07, pp. 25763–25768, 2023.
49. A. Peddireddy and K. Peddireddy, "Next-Gen CRM Sales and Lead Generation with AI," *International Journal of Computer Trends and Technology*, vol. 71, no. 3, pp. 21–26, 2023.
50. K. Peddireddy, "Kafka-based Architecture in Building Data Lakes for Real-time Data Streams," *International Journal of Computer Applications*, vol. 185, no. 9, pp. 1–3, 2023.
51. H. M. Albert et al., "Crystal formation, structural, optical, and dielectric measurements of l-histidine hydrochloride hydrate (LHHCLH) crystals for optoelectronic applications," *J. Mater. Sci.: Mater. Electron.*, vol. 34, no. 30, 2023.
52. N. Sirisha, M. Gopikrishna, P. Ramadevi, R. Bokka, K. V. B. Ganesh, and M. K. Chakravarthi, "IoT-based data quality and data preprocessing of multinational corporations," *The Journal of High Technology Management Research*, vol. 34, no. 2, 2023.
53. M. A. Tripathi, K. Madhavi, V. S. P. Kandi, V. K. Nassa, B. Mallik, and M. K. Chakravarthi, "Machine learning models for evaluating the benefits of business intelligence systems," *The Journal of High Technology Management Research*, vol. 34, no. 2, 2023.
54. D. S. Kumar, A. S. Rao, N. M. Kumar, N. Jeebaratnam, M. K. Chakravarthi, and S. B. Latha, "A stochastic process of software fault detection and correction for business operations," *The Journal of High Technology Management Research*, vol. 34, no. 2, 2023.
55. B. Prasanth et al., "Maximizing Regenerative Braking Energy Harnessing in Electric Vehicles Using Machine Learning Techniques," *Electronics*, vol. 12, no. 5, 2023.
56. R. Jain et al., "Internet of Things-based smart vehicles design of bio-inspired algorithms using artificial intelligence charging system," *Nonlinear Eng.*, vol. 11, no. 1, pp. 582–589, 2022.
57. D. Ganesh, S. M. S. Naveed, and M. K. Chakravarthi, "Design and implementation of robust controllers for an intelligent incubation Pisciculture system," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 1, no. 1, pp. 101–108, 2016.
58. S. Buragadda, K. S. Rani, S. V. Vasantha, and K. Chakravarthi, "HCUGAN: Hybrid cyclic UNET GAN for generating augmented synthetic images of chest X-ray images for multi classification of lung diseases," *Int. J. Eng. Trends Technol.*, vol. 70, no. 2, pp. 229–238, 2022.
59. N. Venkatesan and M. K. Chakravarthi, "Adaptive type-2 fuzzy controller for nonlinear delay dominant MIMO systems: an experimental paradigm in LabVIEW," *Int. J. Adv. Intell. Paradig.*, vol. 10, no. 4, p. 354, 2018.
60. G. S. Sudheer, C. R. Prasad, M. K. Chakravarthi, and B. Bharath, "Vehicle Number Identification and Logging System Using Optical Character Recognition," *International Journal of Control Theory and Applications*, vol. 9, no. 14, pp. 267–272, 2015.
61. M. Chakravarthi and N. Venkatesan, "Design and implementation of adaptive model based gain scheduled controller for a real time non linear system in LabVIEW," *Research Journal of Applied Sciences, Engineering, and Technology*, vol. 10, no. 2, pp. 188–196, 2015.
62. M. Chakravarthi and N. Venkatesan, "Design and Implementation of Lab View Based Optimally Tuned PI Controller for A Real Time Non Linear Process," *Asian Journal of Scientific Research*, vol. 8, no. 1, 2015.
63. K. Peddireddy, "Streamlining enterprise data processing, reporting and realtime alerting using Apache Kafka," in *2023 11th International Symposium on Digital Forensics and Security (ISDFS)*, 2023.
64. R. Boina, A. Achanta, and S. Mandvikar, "Integrating data engineering with intelligent process automation for business efficiency," *International Journal of Science and Research*, vol. 12, no. 11, pp. 1736–1740, 2023.
65. S. Mandvikar and A. Achanta, "Process automation 2.0 with generative AI framework," *Int. J. Sci. Res. (Raipur)*, vol. 12, no. 10, pp. 1614–1619, 2023.
66. S. Venkatasubramanian and R. Mohankumar, "DDoS Attack Detection in WSN Using Modified BGRU With MFO Model," in *Advanced Applications of Generative AI and Natural Language Processing Models*, A. J. Obaid, Ed. IGI Global, USA, pp. 286–305, 2024.
67. S. Venkatasubramanian and S. Hariprasath, "Aquila Optimization-Based Cluster Head Selection and Honey Badger-Based Energy Efficient Routing Protocol in WSN," in *Proceedings of the International Conference on Intelligent Computing, Communication, and Information Security. ICICIS 2022. Algorithms for Intelligent Systems*, Singapore: Springer, pp. 273–290, 2022.
68. R. S. Gaayathri, S. S. Rajest, V. K. Nomula, R. Regin, "Bud-D: Enabling Bidirectional Communication with ChatGPT by adding Listening and Speaking Capabilities," *FMDB Transactions on Sustainable Computer Letters.*, vol. 1, no. 1, pp. 49–63, 2023.
69. V. K. Nomula, R. Steffi, and T. Shynu, "Examining the Far-Reaching Consequences of Advancing Trends in Electrical, Electronics, and Communications Technologies in Diverse Sectors," *FMDB Transactions on Sustainable Energy Sequence*, vol. 1, no. 1, pp. 27–37, 2023.

70. P. S. Venkateswaran, F. T. M. Ayasrah, V. K. Nomula, P. Paramasivan, P. Anand, and K. Bogeshwaran, "Applications of artificial intelligence tools in higher education," in *Advances in Business Information Systems and Analytics*, IGI Global, USA, pp. 124–136, 2023.
71. D. Johnson, A. Menezes, and S. Vanstone, "The elliptic curve digital signature algorithm (ECDSA)," *Int. J. Inf. Secur.*, vol. 1, no. 1, pp. 36–63, 2001.
72. S. Zhang and M. A. Karim, "Color image encryption using double random phase encoding," *Microw. Opt. Technol. Lett.*, vol. 21, no. 5, pp. 318–323, 1999.
73. C. Butpheng, K.-H. Yeh, and H. Xiong, "Security and privacy in IoT-cloud-based e-health systems-A comprehensive review," *Symmetry (Basel)*, vol. 12, no. 7, p. 1191, 2020.
74. E. Mosqueira-Rey, D. Alonso-Ríos, V. Moret-Bonillo, I. Fernández-Varela, and D. Álvarez-Estévez, "A systematic approach to API usability: Taxonomy-derived criteria and a case study," *Inf. Softw. Technol.*, vol. 97, pp. 46–63, 2018.
75. R. Sikder, M. S. Khan, M. S. Hossain, and W. Z. Khan, "A survey on android security: development and deployment hindrance and best practices," *TELKOMNIKA*, vol. 18, no. 1, p. 485, 2020.
76. G. S. Chhabra, V. P. Singh, and M. Singh, "Cyber forensics framework for big data analytics in IoT environment using machine learning," *Multimed. Tools Appl.*, vol. 79, no. 23–24, pp. 15881–15900, 2020.
77. N. Sultana, N. Chilamkurti, W. Peng, and R. Alhadad, "Survey on SDN based network intrusion detection system using machine learning approaches," *Peer Peer Netw. Appl.*, vol. 12, no. 2, pp. 493–501, 2019.