

CrossRef DOI of original article:

Towards Optimized K means Clustering using Nature-Inspired Algorithms for Software Bug Prediction

Tameswar Kajal¹, Geerish Suddul² and Kumar Dookhitram³

¹ University of technology

Received: 1 January 1970 Accepted: 1 January 1970 Published: 1 January 1970

Abstract

In today's software development environment, the necessity for providing quality software products has undoubtedly remained the largest difficulty. As a result, early software bug prediction in the development phase is critical for lowering maintenance costs and improving overall software performance. Clustering is a well-known unsupervised method for data classification and finding related patterns hidden in datasets.

Index terms— data clustering, K-means algorithm, Nature-inspired algorithms, software bug detection, coral reefs.

1 Introduction

In an era of technological disruption, the demand for software adoption has accelerated. They are a part of our society and play an important role in shaping it. Our modern society is becoming increasingly reliant on complex software systems. Thus, it is critical to build reliable and trustworthy systems in a cost-effective and timely manner. The presence of defective modules in a software drives up development and maintenance expenses, leading to customer dissatisfaction. The need for quality assurance has inevitably remained the biggest challenge in today's software development environment. Hence, software bug prediction is an important task to help developers locate bugs more efficiently.

Software bug prediction is an imperative task in Software Development Life cycle (SDLC) as it pertains to the overall success of software. One method in this direction is to use machine learning (ML) methods to predict defects in software. In addition, implementing this method earlier in the SDLC process enhances quality of the product and lowers the cost of software maintenance. Many researchers have applied different theories and methodologies in the field of software bug prediction. Two things are clear from the literature when it comes to defect prediction. Initially, no single prediction approach dominates (Lessmann et al., 2008), and next, the employment of various set of data, data pre-processing, validation systems, and performance statistics makes it challenging to make sense of the multiple prediction outcomes (Myrtveit et al., 2005). There are two common ML model used for prediction based on dataset availability. The first, known as supervised approach, in which a software defect prediction model is built from training set of data and then tested on a testing dataset. Secondly, unsupervised approach, in which the defect prediction model for software is built from scratch using the present testing dataset without training the dataset.

Clustering algorithms have been commonly used to evade the lack of training datasets available being a constraint. Cluster analysis groups things into clusters based on their similarity to create a visual representation of data (Jain and Dubes, 1998). As pointed out by Kaur, 2010, one of the better instances of unsupervised learning is K-means clustering. Clustering is beneficial because it makes it easier to obtain or locate relevant information at a faster rate. Among the different clustering approaches that already exist, the Kmeans methodology is obviously fairly popular. (Gayathri et al., 2015). The preliminary values of the initial centroids, which are generated randomly each time the algorithm is run, have a significant impact on the performance of k-means. K-means frequently fall into local optima that produce poor clustering results. Obtaining a globally optimal clustering result involves a time-consuming, exhaustive approach that tests all partitioning choices. A heuristic

3 RELATED WORKS

45 approach to the problem is to use an optimization algorithm to search for global optima in each computer
46 iteration.

47 Our unsupervised approach uses the k-means approach to divide the unlabeled dataset into defective and
48 non-defective non-overlapped clusters for bug prediction. The goal of this research is to verify the hybrids'
49 efficacy as well as to quantify the quality of results produced by each clustering hybrid model. In this study, we
50 have applied the k-means clustering algorithm, an unsupervised algorithm with different NIAs including Genetic
51 algorithm (GA), Bat algorithm (BA), Particle Swarm Optimization (PSO), Coral Reefs Optimization (CRO),
52 Cuckoo Search optimization (CSO) algorithm, Ant colony optimization (ACO), Firefly algorithm (FA) and Grey
53 Wolf Optimizer (GWO) for software bug prediction. The rest of this paper is organized as follows. Section 2
54 presents a discussion of the related work in software bug prediction. An overview of the methodology, consisting
55 of the algorithms used are presented in Section 3. Section 4 describes the proposed method. Section 5 describes
56 the Dataset and Data Processing method. The evaluation methodology is discussed in section 6. The results and
57 discussion part is discussed in Section 7. Section 8 discusses the practical implications followed by conclusions
58 and future works in section 9.

2 II.

3 Related Works

61 K-means clustering is a well-known partitioned clustering algorithm that has been used in a variety of applications.
62 In the literature, several variations of Kmeans have been proposed to improve its performance for the broad
63 clustering problem. ??ong et al. (2012) studied the integration of bio-inspired optimization methods into K-
64 means clustering for software bug prediction in order to assess clustering performance. The main optimization
65 algorithms tested include the Firefly algorithm, Cuckoo search algorithm, Bat algorithm, Wolf and Ant Colony
66 Optimization (ACO) algorithms. Results show that the combination of these algorithms acquired improved
67 performance accuracy compared with ordinary k-means, at the same time accelerating the search process and
68 avoid local optima. Zhong et al.,2004 compared the k-means algorithm to natural-gas algorithms. The natural
69 gas algorithm outperformed the k-means algorithm in terms of mean square error values. However, this method
70 necessitates the use of a software expert to determine whether the software is appropriate.

71 Annisa et al., 2020, came up with an improved version of k-means algorithm for software bug prediction, that
72 locate the initial centroid of the k-means algorithm and determine the number of clusters present. Because it
73 produces better accuracy than the simple K-Means method, this proposed method could be useful for clustering
74 other data types. Seliya and Khoshgoftaar, 2007 proposed K-means for software failure prediction. Their method
75 iteratively labels clusters as fault-prone or not using expert domain knowledge as a restriction.

76 The k-means algorithm based on quad tree was proposed by Bishnu and Bhattacharjee, 2012 and it was
77 compared to some clustering algorithms. Their proposed algorithm has error rates that are comparable to k-
78 means, Linear Discriminant Analysis and Naive Bayes. Catal et al. 2009 used the x-means clustering algorithm
79 to create faulty and non-faulty clusters based on software metrics. Lines of code, cyclomatic complexity, operand
80 and operator are the metrics. If the metric values are complex than the threshold, the software entity is predicted
81 to be defective, and vice versa. Almayyan, 2021 used dataset from the NASA repository and used three clustering
82 algorithms, Farthest First, X-means and Self-organizing map. This article presents a comparison of software defect
83 prediction algorithms based on Bat, Cuckoo, Grey Wolf Optimizer (GWO), and Particle Swarm Optimization
84 (PSO) in order to evaluate different feature selection algorithms. The Farthest First clustering algorithm was
85 found to be effective in predicting software faultiness, and Bat and Cuckoo were found to be useful in comparison
86 to all other metaheuristic algorithms.

87 Though several academics have sought to merge K-means clustering with nature-inspired algorithms (NIAs),
88 their efforts have been restricted to almost identical group movements, such as the Firefly, Artificial Bee Colony
89 (ACO), and Particle Swarm Optimization (PSO) algorithms ??Jensi and Jiji, 2015). In addition, only a few
90 bio-inspired optimization methods that are integrated with K-means are provided in the previous studies. Only
91 7 of the 28 NIAs hybridized with K-means (Genetic Algorithm, Particle Swarm Optimization, Bat Algorithm,
92 Artificial Bee Colony, Differential Evolution, Harmony Search, and Symbiotic Organism Search) dedicated their
93 hybridization to solving automatic clustering problems, accounting for 20.6 percent of the total (Ikotun et al.,
94 2021). In general, it can be seen that the rate of publishing on K-means hybridization with specific NIAAs is
95 minimal. More research is needed in this area to see if there are any other ways to improve the performance
96 of the existing hybridization algorithm. This suggests that combining Kmeans with these other NIAs to solve
97 automatic clustering problems should be investigated.

98 The purpose of this research is to look into the mechanics of incorporating certain NIAs into the Kmeans
99 clustering algorithm. The optimization function adds to the existing best solution by progressively improving
100 it with a new solution from an unknown fragment of the search space. When a new solution is identified to be
101 better than the present one, the searching agents replace the solutions and continue searching until some stopping
102 criteria are fulfilled.

4 III.

5 Methodology a) K means Clustering Algorithm

The K-means clustering algorithm is a partitioned clustering technique that divides a dataset into k number of clusters using a certain fitness measure. Due to the large amount of data objects in real-world datasets, distributing data items into appropriate clusters to obtain an ideal cluster outcome is computationally expensive and time-consuming (Ikotun et al.2021).

Given a dataset $X = \{x_i\}$, where $i = 1, 2, \dots, n$ of d -dimension data points of size n , X is partitioned into ' k ' clusters such that $J(c, k) = \sum_{k=1}^k \sum_{i \in C_k} \|x_i - \mu_k\|^2$ (1)

With the objective function: minimize the sum of the square error over all the k clusters. That is, minimize $J(C) = \sum_{k=1}^k \sum_{i \in C_k} \|x_i - \mu_k\|^2$ (2)

When assigning N objects to k clusters, the purpose of the clustering algorithm is to limit the number of potential possibilities. This can be expressed numerically as: $S(N, K) = \frac{1}{K!} \sum_{i=0}^{K-1} \binom{K-1}{i} i^i (K-i)^{N-i}$ (3)

b) Nature-inspired algorithms (NIAs) Nature-inspired computation has gained popularity in the previous two decades and has been used in practically every field of research and engineering (Yang et al.2013). NIAs are global optimization strategies for solving difficult real-world issues (Okwu et al. 2020). NIAs have successfully provided suboptimal solutions to automatic clustering problems in a reasonable amount of time (Hruschka et al. 2009). The population is used for the exploration of search space in the nature-inspired metaheuristic, ensuring a higher possibility of finding optimal cluster partitions (Nanda and Panda, 2014). It has been discovered that combining K-means with NIAs for automatic clustering improves the performance of algorithms when dealing with cluster analysis. In most circumstances, the automatic cluster number determination aids in the selection of near-optimal starting cluster centroids for the clustering process rather than the normal random selection (Zhou et al. 2017).

6 c) Combination of k-means with Nature-Inspired Algorithms (NIAs)

Clustering using NIAs is now as simple as assigning combinations of centroids to the searching agents, allowing them to heuristically find the best answer. Though the specifics of conducting a heuristic search vary depending on which nature-inspired optimization algorithm technique is used, the initialization stage and the finishing step, where the quality of the discovered solution is evaluated as a stopping condition, are both comparable.

S is defined as the solution space that contains a finite number of x_i , where i is the solution's index, in the initialization construct. The search agents represent the solutions x , each of which holds a set of centroids, regardless of the types of bio-inspired optimization methods used. Typically, a large population of searching agents, N , is utilized to collaboratively search for the best feasible cluster configurations (as expressed by the locations of the optimal centroids). K is the number of clusters that must be formed, which is generally a userdefined figure. D is the dimension of the search space, which is the number of attributes a data point possesses.

To find the optimal configuration of centroids we let $cen_{j,v}$ be the centroids at the j th cluster and the v th attribute. To obtain the centroid location, the following formula is used: $cen_{j,v} = \frac{\sum_{i=1}^N w_{i,j} x_{i,v}}{\sum_{i=1}^N w_{i,j}}$ (4)

In our concept, the matrix $cen_{j,v}$ contains all of the cluster centers and is a two-dimensional matrix with $K \times D$ characteristics. $F(cen) = \sum_{k=1}^k \sum_{i=1}^N W_{i,j} \sum_{v=1}^D (X_{i,v} - cen_{j,v})^2$ (5)

The calculation method loops $K \times D$ times to analyze the values of all the attributes of x in each cluster v to calculate the distance between each x and the centroid.

Cluster centers can be designated by data points. For example, in a two-cluster clustering task, the objective function requires three variables. As a result, there are three dimensions.

Three variables, and hence three-dimensional spaces, are required, and the i th data point may be written as $x_i = (i, [x_{i,1}, x_{i,2}, x_{i,3}, x_{i,4}, x_{i,5}, x_{i,6}])$. The clustering strategy can be formulated as follows: $clmat_{i,j} = \min_k \{ \|x_i - cen_k\| \}$ (6)

Sets of functional parameters must be defined in order to execute the bio-inspired optimization algorithms. Despite the fact that some of their parameters are shared, each set of parameters for the hybrid bio-inspired clustering algorithms is designed independently. The six models investigated are K means with Genetic Algorithm, K means with Bat algorithm, K means with Ant colony algorithm, K means with Cuckoo Search Algorithm, K means with Firefly Algorithm and K means with Coral reefs algorithm. The most significant variations are in how the global optimal exploration is carried out for all these algorithms. The evaluation stage comes right after the exploration construct, and it compares if the new solution is better than the current best one.

7 d) Genetic Algorithm

Genetic Algorithm (Ga) are randomized heuristic search algorithms that are based on natural selection and genetic principles (Goldberg, 1989). The genetic operators used in the combination of K-means and GA are selection, distance-based mutation, and the K-means operator. The parameters have been set according to the

161 study of Bouhmal et al. 2015. $P(0)$ is chosen at random as the starting population. Each allele in the population
 162 can be given a cluster number from the uniform distribution over the set $\{1, \dots, K\}$ at random.

163 According to the distribution given by, the selection operator selects a chromosome from the preceding
 164 population at random as follows: $P(s_i) = F(s_i) / \sum_{j=1}^K F(s_j)$ (7)

165 The possibility of solutions surviving in the future population is ranked in the current population. Each
 166 solution in the population must be assigned a figure of merit or a fitness value. $F(s_W) = \begin{cases} g(s_W) & \text{if } g(s_W) < 0 \\ 0 & \text{otherwise} \end{cases}$ (8)

168 e) Bat Algorithm (BA) Bat echolocation is used in the bat algorithm (BA), which is a heuristic optimization
 169 tool (Yang, 2010). The four basic parameters of a BA are pulse frequency, pulse rate, velocity, and a constant.
 170 The parameters have been set according to the study ??Huang and Ma, 2020).

171 The frequency, velocity, and position for each bat are initialized. The virtual bats' movement is described
 172 by updating their velocity and position using the equations below for each time step t , where T is the iteration
 173 limit. $f_i = f_{min} + (f_{max} - f_{min}) \cdot r$ (9) $V_{i,t+1} = v_{i,t} + [X_{i,t} + X^*] \cdot f_i$ (10) $X_{i,t+1} = X_{i,t} + v_{i,t}$ (11)

174 A random number is generated when the bat positions are updated; if the random number is greater than the
 175 pulse emission rate, a new location is formed around the current best solutions, as shown in the equation below.

176 $X_{new} = x_{old} + EA \cdot t$ (12)

177 8 f) Ant Colony Optimization (ACO)

178 The ACO heuristic was inspired by investigations of ant foraging behavior in real colonies, which indicated that
 179 ants can often figure out the shortest path between food source and nest (Zheng et al. 2003). The parameters
 180 have been set according to the study ??Tang et al. 2012).

181 When the ant moves from i to j , the path node at the start can set as A , $A = \{0, 1, \dots, n-1\}$. This reflects the
 182 role of pheromones accumulated by ants during exercise during ant migration and reveals the relative relevance
 183 of the trajectory. The larger τ is, it indicates the high probability for subsequent ants to choose this path.

184 The probability of the ant moving from I to j is computed using the following formula: $P_{ij}^k(t) = \tau_{ij}^k(t) / \sum_{l \in A} \tau_{il}^k(t)$ (13)

186 9 g) Firefly Algorithm (FA)

187 Firefly algorithm is a very strong technique for solving restricted optimization and NP-hard problems (Apostolopoulos and Vlachos, 2011). The parameters have been set according to the study ??Tang et al. 2012).

189 The attractiveness of a firefly I on a firefly j is determined by the degree of the firefly i 's brightness and the
 190 distance r_{ij} between the firefly I and the firefly j , as shown below: $I(r) = I_s / r^2$ (14)

191 Consider the case when there are n fireflies and the solution for firefly I is x_i . The brightness of the firefly I is
 192 linked to the objective function $f(x_i)$. $I = f(x_i)$ (15)

193 Each firefly has an attraction value, and the less dazzling (attractive) one is drawn to the brighter one and
 194 transferred there. The attractiveness value α is relative based on the distance between fireflies. Where pheromone
 195 is τ , which is a constant that represents weight. The time of iteration is N_c and the initial setting is τ . The
 196 predicted heuristic factor is τ , which demonstrates the relevance of visibility relative to other factors. It also
 197 represents the significance of the heuristic component in the entire path of the ant's movement. $\tau(r) = \tau_0 e^{-\alpha r^2}$ (16)

199 Where $V_{i,t}$ and $X_{i,t}$ are the velocity and position at time t , $V_{i,t+1}$ and $X_{i,t+1}$ are the velocity and position
 200 at time $t+1$, and r is a random number between 0 and 1.

201 Where τ_0 is the firefly attraction value at $r = 0$ and α is the media light absorption coefficient.

202 Where E is a random number A_t represents the average loudness of all bats at time t .

203 An initial population of n nests is randomly generated at the positions, $X = \{x_{0,1}, x_{0,2}, \dots, x_{0,n}\}$, to evaluate
 204 the objective values to find the current global best $g_{t,0}$.

205 The new position is updated accordingly by performing a Levy flight: $x_i(t+1) = x_i(t) + \tau \cdot L \cdot \text{vy}(\tau)$, (17)

206 10 i) Coral Reefs Optimization Algorithm (CRO)

207 CRO is another nature-inspired algorithm, based on an artificial simulation of the process of coral reef formation
 208 and reproduction (Sanz et al. 2014). The CRO algorithm has never been utilized in the realm of software bug
 209 detection to our knowledge. Corals reproduce at each iteration step in the CRO algorithm, producing new
 210 individuals. The parameters have been set according to the study (Medeiros et al., 2015).

211 By allocating a coral to each square (i, j) , the CRO algorithm generates a $N \times M$ square grid in which each
 212 square (i, j) may represent an alternate solution to a problem (or colony of corals). The formation of coral is the
 213 second phase. After three phases, the entire collection of existing corals in the reef is graded according to their
 214 level of healthiness (broadcast spawning, brooding, and larvae setting).

215 11 j) Particle Swarm Optimization (PSO)

216 The behavior of particles in a swarm is the central concept of the PSO. Each particle has its own location in a
 217 multidimensional space and communicates with the others. To move about in space, the particles employ social

218 and cognitive information. When the algorithm comes to a halt, the best solution has been discovered (Koochi
219 and Groza, 2014). The parameters have been set according to the study (Rana et al., 2010).

220 The inertia weight balances the algorithm’s local and global search abilities. The proportional contribution of
221 the prior velocity to the current velocity is defined by the inertia weight. $V_{i,k+1} = wv_{i,k} + c1 \text{rand}(p_{\text{best}_i} - x_{i,k}) + c2 \text{rand}(g_{\text{best}} - X_{i,k})$ (18) $X_{i,k+1} = X_{i,k} + v_{i,k+1}$ (19)
222

223 12 k) Grey Wolf Optimizer (GWO)

224 The Grey Wolf Optimizer (GWO) is a simple, population-based, flexible, and derivative-free metaheuristic
225 optimization method that intelligently avoids stagnation in local optima spots of the search space. It simulates
226 the social behaviors of grey wolves in the aspects of their hierarchical leadership and hunting movement ??Mirjalili
227 et al., 2013). Grey wolves’ leadership and haunting mechanism help to design a new metaheuristic algorithm
228 with three steps: searching prey, encircling prey, and attacking prey.

229 During the GWO operation, the position of the wolves is continuously updated, with appropriate mathematical
230 formulas ??Hou et al., 2022). The parameters have been set according to the study (Wang et al., 2019).

231 IV.

232 13 Proposed Method

233 The purpose of clustering is to discover a proper set of centroids using the metaheuristic of the nature-inspired
234 method as a guide. The metaheuristic will always insist on centroids being moved in a progressive manner in
235 each phase, with the goal of finding the best grouping. The ideal group’s ultimate result should be that the
236 data points inside each cluster are closest to their centroid. During the search, the centroids move around in
237 the search space, following the swarming pattern of the nature-inspired optimization method, until no further
238 progress is seen. It comes to a halt when there is no other possible relocation that will yield a better result. Along
239 with the success of employing nature-inspired metaheuristic algorithms to solve automatic clustering problems,
240 it has been discovered that combining two or more metaheuristics for the same objective improves clustering
241 performance. The performance of hybrid algorithms, according to Nanda and Panda 2014, is superior to that of
242 separate algorithms in terms of robustness, effectiveness, and accuracy.

243 V.

244 14 Dataset and Data Processing

245 The dataset was collected from the online PROMISE repository. AR1, AR3, AR4, AR5, AR6, KC1, KC2,
246 JM1, CM1, PC1 and PC5 were used respectively. With reference to the paper, by Shepperd et al. 2013, data
247 cleaning is mandatory before using any datasets available. Indeed, we noted a huge class imbalance issue with the
248 available datasets (faulty, non-faulty), and all data inconsistencies, missing and null values were removed. Each
249 dataset selected represents a NASA software system that includes various metrics. Each dataset is made up of a
250 number of software modules and attributes. Modules with defects are classified as prone to faults, whereas those
251 without defects are To address the curbs of the K-means clustering approach in generating globally optimum
252 clusters, the suggested method uses the k-means algorithm together with a range of NIAs for software bug
253 prediction. By adding an exploration function to the k-means algorithm, the combination of these strategies may
254 improve the model. The exploration function improves the existing solution by examining regions outside of its
255 immediate vicinity, and if a new, better solution than the current best one is discovered, the search agents will
256 move toward it. The exploring procedure will continue until certain stopping criteria are met. Nature-inspired
257 algorithms are metaheuristic algorithms, which means they have the ability to explore the combinatorial search
258 space heuristically rather than exhaustively. The integration methods are based on representing the search agents
259 as a combination of centroid locations, then the search agents explore the search space for the best solution.

260 Where $\delta > 0$ denotes the step size, which should be connected to the problem’s scales. In most circumstances,
261 we can use $\delta = 1$. classified as non-fault prone. For the training purpose, the entire dataset is used except for
262 the last column (output column), only columns consisting of numerical values were considered. VI.

263 15 Evaluation a) Experimental Setup

264 The main goal of this research is to demonstrate the utility of the k-means algorithm with different NIAs, which
265 we accomplished using Tensorflow to train the model. TensorFlow is an open-source machine learning platform
266 to build and deploy prediction models. Google Colab was also used to run the results, which allowed the code
267 to run with no configuration and free GPU access. Each dataset is performed 10 times in the trials to find the
268 average CPU time and objective function values/best fitness value.

269 The clustering results of the new hybrid clustering algorithms are compared to the K-means, which serve as a
270 benchmarking reference. The full dataset is used for training, and cluster formation is referred to until perfection
271 is attained using the entire set of data. The ultimate clustering result’s quality is determined by each cluster’s
272 integrity, which is represented by the objective function’s final fitness value.

273 The hardware configuration used for all experiments in this study is as follows: Corei7-6500U CPU @2.50 GHz
274 2.60 GHz, Windows 10, 64-bit operating system, x64 based processor, RAM: 8 GB DDR4, and Hard Disk: SSD.

16 b) Performance Evaluation Measures

In order to assess the effectiveness of combining the k-means algorithm and optimization algorithms in the prediction of software bugs, the evaluation metrics, accuracy and F-measure have been calculated accordingly as shown in the Equation (1): $Accuracy = (TP + TN) / (TP + TN + FN + FP)$, (20)

Where TP = true positive, TN = true negative, FN = false negative and FP = false positive.

On the other hand, the external metric used to determine the accuracy of the clustering findings, known as the F-measure, is also computed.

The F-measure, which is the average of precision and sensitivity performance, is calculated as follows: $F = 2 * P * Sensitivity / P + Sensitivity$, (21) Where P refers to precision and sensitivity is calculated by finding the non-defective modules that were accurately categorized. From the table above, K-means clustering is optimized using the various NIAs. We can see that all of the proposed algorithms perform better than the traditional standalone k-means algorithm. K-means appears to take the shortest computation time in any of the tests, maybe because it stops early in local optima (Table ??). This is evident from the accuracy obtained from the table above. NIAs speed up the process of clustering centroids and illustrate that all partitioning clustering methods can be linked with the natural search process to prevent local optima. Secondly, simple Kmeans were applied to the robust nature of GA, which shows adequate prediction accuracy for all datasets. Even though GA may converge to the global optimum due to mutation, GA faces the issue in terms of computational challenges. The application of k means with the Bat algorithm apparently yields the same accuracy. This hybrid algorithm improves the convergence speed of BA and helps the k means algorithm independent of the initial centers. Next, K means is combined with PSO. The PSO method is used to start the process because of its fast convergence, and then the K-Means algorithm is used to refine the PSO algorithm's outcome to near-optimal solutions. The hybridization of these two methods yields effective results in terms of efficiency and precision. The PSO algorithm can be used to generate good initial cluster centroids for the K-Means.

17 VII.

18 Results and Discussions

19 Practical Implications

Metaheuristics algorithms have shown to be effective optimizers. This research found that each of the hybrid K means based-nature-inspired optimization algorithm models outperformed the standalone K means algorithm in terms of accuracy and F1 score. Following the intrinsic limitations of K-means design and the virtues of Nature-inspired optimization techniques, it seems feasible to integrate them, allowing them to complement and function together. The algorithms' successful integration gives reason to believe that more advanced optimization mining techniques can be developed. This study can be used as a roadmap for researchers who want to incorporate other new emerging NIAs into improved clustering methods in the field of software bug detection.

20 IX.

21 Conclusion and Future Works

Prediction of defect-prone software modules is an important goal in software engineering. The traditional clustering algorithm usually gets trapped in the problem of local optima. As a result, the natureinspired method provides an alternative technique for solving clustering problems using its searching capabilities. This study's main contribution is combining the clustering algorithm with the different NIAs for software bug detection. To the authors' knowledge, only PSO, Cuckoo, Bat, and GWO (Grey Wolf Optimizer) algorithms were applied with clustering algorithms for software bug detection (Almayyan, 2021). The results are improved significantly when clustering algorithms are combined with bio-inspired optimization methods, apparently for the hybrid model of k means clustering withCoral reefs algorithm, achieving an accuracy of 96%.For future work, this work can be replicated with other related datasets for the analysis of bug prediction in software. ¹

¹ Towards Optimized K means Clustering using Nature-Inspired Algorithms for Software Bug Prediction © 2023 Global Journals

1

| Dataset | Modules | Defective modules | Software metrics (Attributes) |
|---------|---------|-------------------|-------------------------------|
| AR1 | 121 | 9 | 29 |
| AR3 | 63 | 8 | 29 |
| AR4 | 107 | 20 | 29 |
| AR5 | 36 | 8 | 29 |
| AR6 | 101 | 15 | 29 |
| KC1 | 2109 | 1783 | 22 |
| KC2 | 522 | 107 | 21 |
| JM1 | 7782 | 1672 | 21 |
| CM1 | 327 | 42 | 37 |
| PC1 | 705 | 61 | 37 |
| PC5 | 1711 | 471 | 38 |

Figure 1: Table 1 :

2

| Datasets | AR1 | AR3 | AR4 | AR5 | AR6 | KC1 | KC2 | JM1 | CM1 | PC1 | PC5 |
|----------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| k-Means | 88.90 | 88.00 | 89.01 | 88.85 | 88.43 | 89.10 | 89.00 | 88.80 | 89.00 | 89.19 | 89.99 |
| K-Means +GA | 90.50 | 90.58 | 91.28 | 91.55 | 90.11 | 90.00 | 90.54 | 90.53 | 91.25 | 90.00 | 90.05 |
| K-Means +BAT | 90.00 | 91.59 | 91.00 | 92.34 | 92.00 | 92.98 | 91.34 | 90.00 | 91.25 | 92.56 | 92.00 |
| K-Means +PSO | 92.50 | 92.65 | 92.87 | 93.01 | 93.00 | 92.99 | 94.10 | 92.67 | 92.89 | 93.10 | 93.58 |
| K-Means +Coral Reefs | 94.00 | 94.54 | 94.56 | 94.87 | 94.00 | 95.96 | 95.66 | 96.88 | 95.01 | 95.04 | 95.54 |
| K-Means +Cuckoo | 94.50 | 94.58 | 94.58 | 94.00 | 94.56 | 95.45 | 95.88 | 95.67 | 95.44 | 94.56 | 94.78 |
| K-Means +ACO | 94.00 | 93.56 | 93.50 | 94.10 | 93.78 | 93.03 | 93.56 | 93.44 | 93.89 | 94.01 | 94.52 |
| K-Means +Firefly | 92.56 | 92.67 | 93.00 | 93.44 | 93.02 | 93.56 | 94.78 | 93.67 | 94.88 | 94.34 | 94.54 |
| K-Means +GWO | 90.09 | 92.47 | 94.65 | 93.22 | 92.00 | 92.60 | 93.00 | 92.50 | 94.50 | 94.12 | 94.13 |

Figure 2: Table 2 :

4

| Datasets | AR1 | AR3 | AR4 | AR5 | AR6 | KC1 | KC2 | JM1 | CM1 | PC1 | PC5 |
|------------------|------|------|------|------|------|------|------|------|------|------|------|
| k-Means | 0.66 | 0.79 | 0.82 | 0.80 | 0.75 | 0.81 | 0.80 | 0.81 | 0.82 | 0.82 | 0.80 |
| K-Means+GA | 0.84 | 0.83 | 0.83 | 0.80 | 0.83 | 0.84 | 0.84 | 0.85 | 0.82 | 0.81 | 0.85 |
| K-Means +BAT | 0.83 | 0.81 | 0.83 | 0.86 | 0.86 | 0.86 | 0.85 | 0.85 | 0.85 | 0.87 | 0.85 |
| K-Means +PSO | 0.85 | 0.85 | 0.87 | 0.87 | 0.86 | 0.85 | 0.87 | 0.85 | 0.87 | 0.87 | 0.87 |
| K-Means +Coral | 0.86 | 0.86 | 0.86 | 0.85 | 0.86 | 0.86 | 0.87 | 0.88 | 0.86 | 0.87 | 0.88 |
| Reefs | | | | | | | | | | | |
| K-Means | 0.89 | 0.85 | 0.88 | 0.89 | 0.86 | 0.89 | 0.86 | 0.89 | 0.89 | 0.87 | 0.88 |
| +Cuckoo | | | | | | | | | | | |
| K-Means+ ACO | 0.84 | 0.83 | 0.86 | 0.85 | 0.84 | 0.86 | 0.85 | 0.85 | 0.86 | 0.85 | 0.86 |
| K-Means +Firefly | 0.86 | 0.85 | 0.83 | 0.87 | 0.87 | 0.85 | 0.83 | 0.85 | 0.86 | 0.88 | 0.85 |
| K-Means+ GWO | 0.82 | 0.82 | 0.81 | 0.86 | 0.84 | 0.83 | 0.79 | 0.85 | 0.84 | 0.84 | 0.85 |

VIII.

Figure 3: Table 4 :

.1 Acknowledgments

This study received no formal support from public, private, or not-for-profit funding organizations.

.2 ()

Year 2023 C Furthermore, K means and Coral reefs algorithm are combined. The results for this combined method are quite promising since they show that using the CRO method for a clustering application can produce better results to using hybrid genetic algorithms, which is the most often used clustering optimization technique. To best of our knowledge, CRO has not been used with clustering for software bug detection. The hybrid model of k means with Cuckoo Search algorithm shows significant accuracy, likewise CRO algorithm. Cuckoo search is used to provide a robust initialization, whereas K-means is utilized to construct solutions faster. K means is also combined with Ant Colony Optimization algorithm. The suggested method's learning mechanism is based on the use of a defined parameter termed pheromone, which eliminates undesirable K-means algorithm solutions. The suggested method improves the K-means algorithm by making it less reliant on starting parameters such as randomly picked beginning cluster centers, resulting in a more stable algorithm. K means with firefly also produce near accuracy with CRO and Cuckoo search algorithm. This is because fireflies with high similarity are dispersed, resulting in a more diverse distribution of the entire swarm in search space. K means with GWO has also shown rapid convergence. This improvement is caused by the fact that K-means significantly affects the GWO population and separates it into two clusters. Because GWO often operates as three clusters and has three wolves in the search space, K-means is advantageous for GWO. As a result, it can be concluded that K-means combined with GWO increased GWO's effectiveness to some extent.

High clustering accuracy and efficiency were obtained from the hybrid clustering of Coral reefs and Cuckoo Search Algorithm. Hybrid clustering of Coral reefs algorithm has never been applied in the field of software bug detection and has indeed shown promising results. Hybrid clustering of Coral reefs algorithm locate cluster centroids without causing premature convergence. The findings of the evaluation results add evidence that NIAs can indeed speed up the process and avoid local optima. Because fewer iterations are required to achieve the best cluster outcome, selecting the number of clusters enhances the hybridized clustering method's convergence speed. The computational time for each algorithm is computed as shown in Table ???. Less computational time was noted when K means was integrated with Coral reefs and Cuckoo Search algorithm respectively. For statistical performance, the F1 score has been calculated for all the algorithms as shown in Table ???. Again, the F1 Score shows that K-means with Coral reefs resulted in dependable and significant performance that can be used to predict software defects. When a good validity measure is applied, most metaheuristic algorithms can automatically divide datasets into an appropriate number of clusters, according to Gbaje et al.2019.

[Engineering ()] , Engineering . 2008. 34 p. .

[Yang ()] , X S Yang . *Swarm Intelligence and Bio-Inspired Computation: Theory and Applications* 2013. V. Amsterdam, The Netherlands: Elsevier Science Publishers B.

[Rana et al. ()] 'A hybrid sequential approach for data clustering using K-Means and particle swarm optimization algorithm'. Sandeep Rana , Sanjay Jasola , Rajesh Kumar . *International Journal of Engineering, Science and Technology* 2010. 2 (6) p. .

[Kao and Kao ()] 'A hybridized approach to data clustering'. Yi-Tung Kao , Erwiezahara , I-Wei Kao . *Expert Systems with Applications* 2008. 34 (3) p. .

[Yunlongzhu et al. ()] 'A new approach for data clustering using hybrid artificial bee colony algorithm'. Xiaohui Yunlongzhu , Wenping Yan , Liang Zou , Wang . *Neuro computing* 2012. 97 p. .

[Goldberg and Yang ()] 'A new metaheuristic bat-inspired algorithm'. D E Goldberg , ; X-S Yang . *Genetic Algorithms in Search, Optimization, and Machine Learning*, (New York) 1989. 2010. Addison-Wesley. 284 p.

[Gayathri et al. (2015)] 'A Novel Approach for Clustering Based On Bayesian Network'. R Gayathri , A Cauveri , R Kanagapriya , V Nivetha , P Tamizhselvi , K P Kumar . *Proceedings of the 2015 International Conference on Advanced Research in Computer Science Engineering & Technology*, (the 2015 International Conference on Advanced Research in Computer Science Engineering & Technology) 2015. March. 2015. ACM. p. 60.

[Tóth et al. ()] 'A Public Bug Database of GitHub Projects and Its Application in Bug Prediction'. Z Tóth , P Gyimesi , R Ferenc . *Computational Science and Its Applications -ICCSA 2016*, (Cham) 2016. Springer International Publishing. p. .

[Hruschka et al. ()] 'A Survey of Evolutionary Algorithms for Clustering'. E Hruschka , R J G B Campello , A A Freitas , A De Carvalho . *IEEE Trans. Syst. Man Cybern. Part C Appl. Rev* 2009. 39 p. .

[Nanda and Panda ()] 'A survey on nature inspired metaheuristic algorithms for partitionial clustering'. S J Nanda , G Panda . *Swarm Evol. Comput* 2014. 16 p. .

[Jain and Dubes ()] *Algorithms for clustering data*, A K Jain , R C Dubes . 1988. Prentice-Hall, Inc.

[Zhou et al. ()] 'An Automatic K-Means Clustering Algorithm of GPS Data Combining a Novel Niche Genetic Algorithm with Noise and Density'. X Zhou , J Gu , S Shen , H Ma , F Miao , H Zhang , H Gong . *ISPRS Int. J. Geo-Inf* 2017. 6 p. 392.

21 CONCLUSION AND FUTURE WORKS

- 378 [Wang and Li ()] ‘An Improved Grey Wolf Optimizer Based on Differential Evolution and Elimination
379 Mechanism’. Jie-Sheng Wang , Shu-Xia Li . 10.1038/s41598-019-43546-3. [https://doi.org/10.1038/
380 s41598-019-43546-3](https://doi.org/10.1038/s41598-019-43546-3) *SciRep* 2019. 9 p. 7181.
- 381 [Apostolopoulos and Vlachos ()] ‘Application of the Firefly Algorithm for Solving the Economic Emissions
382 Load Dispatch Problem’. T Apostolopoulos , A Vlachos . 10.1155/2011/523806. *International journal of
383 Combinatorics* 2011.
- 384 [Inacio et al. (2015)] ‘Applying the Coral Reefs Optimization Algorithm to Clustering Problems’. G Inacio , Joao
385 C Medeiros , Anne M P Xavier-Junior , Canuto . 10.1109/IJCNN.2015.7280845. Conference Paper July 2015.
- 386 [Agbaje et al. ()] ‘Automatic Data Clustering Using Hybrid Firefly Particle Swarm Optimization Algorithm’. M
387 B Agbaje , A E Ezugwu , R Els . *IEEE Access* 2019. 7 p. .
- 388 [Huang and Ma] *Bat Algorithm Based on an Integration Strategy and Gaussian Distribution*, Jianqiang Huang ,
389 Yan Ma . 10.1155/2020/9495281. <https://doi.org/10.1155/2020/9495281> 2020.
- 390 [Lessmann et al.] ‘Benchmarking classification models for software defect prediction: A proposed framework and
391 novel findings’. S Lessmann , B Baesens , C Mues , S Pietsch . *IEEE Transactions on Software*
- 392 [Yang and Deb (2009)] ‘Cuckoo search via Lévy flights’. X.-S Yang , S Deb . *Proceedings of the World Congress
393 on Nature & Biologically Inspired Computing (NABIC '09)*, (the World Congress on Nature & Biologically
394 Inspired Computing (NABIC '09)Coimbatore, india) December 2009. p. .
- 395 [Bouhmala et al. (2015)] ‘Enhanced Genetic Algorithm with K-Means for the Clustering Problem’. N Bouhmala
396 , A Viken , J B Lønnum . *International Journal of Modeling and Optimization* April 2015. 5 (2) .
- 397 [Kaur and Kaur ()] ‘Fault Prediction using K-Canberra Means Clustering’. Deepinder Kaur , Arashdeep Kaur .
398 *CNC* 2010. (in Press)
- 399 [Mirjalili et al. ()] ‘Grey wolf optimizer’. S Mirjalili , S M Mirjalili , A Lewis . 10.1016/j.advengsoft.2013.12.007.
400 *Adv. Eng. Softw* 2014. 69 p. .
- 401 [Jensi and Wiselinjiji (2015)] ‘HYBRID DATA CLUSTERING APPROACH USING K-MEANS AND
402 FLOWER POLLINATION ALGORITHM, Advanced Computational Intelligence’. R Jensi , G Wiselinjiji
403 . *An International Journal (ACII)* April 2015. 2 (2) .
- 404 [Hou et al.] ‘Improved Grey Wolf Optimization Algorithm and Application’. Y Hou , H Gao , Z Wang , C Du .
405 10.3390/s22103810. <https://doi.org/10.3390/s22103810> *Sensors* 2022 p. 3810.
- 406 [Riskiannisa and Riana (2020)] ‘Improved point center algorithm for k-means clustering to increase software
407 defect prediction’. Didirosiyadi Riskiannisa , Dwiza Riana . *International Journal of Advances in Intelligent
408 Informatics* November 2020. 6 (3) p. .
- 409 [Tang et al. ()] ‘Integrating nature-inspired optimization algorithms to K-means clustering’. R Tang , S Fong ,
410 X Yang , S Deb . 10.1109/ICDIM.2012.6360145. *Seventh International Conference on Digital Information
411 Management (ICDIM 2012)*, 2012. p. .
- 412 [Abiodun et al.] ‘K-Means-Based Nature-Inspired Metaheuristic Algorithms for Automatic Data Clustering
413 Problems: Recent Advances and Future Directions’. M Abiodun , Mubarak S Ikotun , Absalom E Almutari ,
414 Ezugwu . 10.3390/app112311246. <https://doi.org/10.3390/app112311246> *Appl. Sci* 2021 p. 11246.
- 415 [Mousa et al. ()] ‘Local Search Based Hybrid Particle Swarm Optimization for Multiobjective Optimization’. A
416 A Mousa , M A El-Shorbagy , Abd El-Wahed , WF . *Swarm and Evolutionary Computation* 2012. 3 p. .
- 417 [Okwu and Tartibu ()] ‘Metaheuristic Optimization: Nature-Inspired Algorithms Swarm and Computational
418 Intelligence’. M O Okwu , L K Tartibu . *Theory and Applications* 2020. Germany: Springer Nature:
419 Berlin/Heidelberg. 927.
- 420 [Koochi and Groza ()] ‘Optimizing Particle Swarm Optimization algorithm’. I Koochi , V Z Groza .
421 10.1109/CCECE.2014.6901057. *2014 IEEE 27th Canadian Conference on Electrical and Computer Engi-
422 neering (CCECE)*, 2014. p. .
- 423 [Myrtveit et al. ()] ‘Reliability and validity in comparative studies of software prediction models’. I Myrtveit , E
424 Stensrud , M Shepperd . *IEEE Transactions on Software Engineering* 2005. 31 (5) p. .
- 425 [Shepperd et al. ()] M Shepperd , Q Song , Z Sun , C Mair . *Data quality: Some Eng*, 2013. 39 p. .
- 426 [Catal et al. ()] *Software fault prediction of unlabeled program modules*, C Catal , U Sevim , B Diri . 2009. 2009.
427 p. .
- 428 [Bishnu and Bhattacharjee ()] ‘Software fault prediction using quad tree-based kk-means clustering algorithm’.
429 P S Bishnu , V Bhattacharjee . 10.1109/TKDE.2011.163. <https://doi.org/10.1109/TKDE.2011.163>
430 *IEEE Trans Knowl Data Eng* 2012. 24 (6) p. .
- 431 [Macqueen ()] ‘Some methods for classification and Analysis of Multivariate Observations’. J B Macqueen .
432 *Proceedings of 5 th Berkeley Symposium on Mathematical Statistics and Probability*, (5 th Berkeley Symposium
433 on Mathematical Statistics and Probability) 1967. University of California Press. p. .

- 434 [Zheng et al. ()] ‘The application of ant colony system to image texture classification’. H Zheng , Z Zheng ,
435 Y Xiang . *Proceedings of the 2nd International Conference on Machine Learning and Cybernetics*, (the
436 2nd International Conference on Machine Learning and Cybernetics Xi’an, China) 2003. 3 p. . (textute read
437 texture)
- 438 [Salcedo-Sanz et al. ()] ‘The Coral Reefs Optimization Algorithm: A Novel Metaheuristic for Efficiently Solving
439 Optimization Problems’. S Salcedo-Sanz , J Del Ser , S Gil-Lpez , I Landa-Torres , J A Portilla-Figueras .
440 *The Scientific World Journal* 2014. Hindawi Publishing Corporation. 2014.
- 441 [Fong et al. ()] ‘Towards Enhancement of Performance of K-Means Clustering Using Nature-Inspired Opti-
442 mization Algorithms’. Simon Fong , Suash Deb , Xin-She Yang , Yan Zhuang . 10.1155/2014/564829.
443 <https://doi.org/10.1155/2014/564829> *Computational Intelligence and Metaheuristic Algorithms with*
444 *Applications* 2014.
- 445 [Almayyan (2021)] ‘Towards Predicting software defects with clustering techniques’. Waheeda Almayyan .
446 *International Journal of Artificial Intelligence and Application (IJAIA)* January 2021. 12 (1) .
- 447 [Zhong et al. ()] ‘Unsupervised learning for expert-based software quality estimation’. S Zhong , T M Khosh-
448 goftaar , N Seliya . 10.1109/HASE.2004.1281739. <https://doi.org/10.1109/HASE.2004.1281739>
449 *Proceedings of the eighth IEEE international conference on high assurance systems engineering HASE*, (the
450 eighth IEEE international conference on high assurance systems engineering HASE) 2004. 2004. p. .
- 451 [Tang and Fong (2012)] *Xin-she Yang and Suash Deb, Integrating Nature inspired Optimization algorithms to*
452 *k-means clustering*, Rui Tang , Simon Fong . Aug 2012. University of Macau