

## A New Method for Data Clustering based on The Combination of Genetic Optimization and Firefly Algorithms

Mahsa Afsardeir<sup>1</sup>, Mansureh Afsardeir<sup>2\*</sup>

1. MSc, Department of Computer Engineering, Science and Research Branch, Islamic Azad University, Tehran, Iran.
2. MSc, Department of Biomedical Engineering, Dezful Branch, Islamic Azad University, Dezful, Iran.  
(Corresponding Author) [afsardeir2007@gmail.com](mailto:afsardeir2007@gmail.com)

### Abstract

**Introduction:** With the progress of technology and increasing the volume of data in databases, the demand for fast and accurate discovery and extraction of databases has increased. Clustering is one of the data mining approaches that is proposed to analyze and interpret data by exploring the structures using similarities or differences. One of the most widely used clustering methods is the k-means. In this algorithm, cluster centers are randomly selected and each object is assigned to a cluster that has maximum similarity to the center of that cluster. Therefore, this algorithm is not suitable for outlier data since this data easily changes centers and may produce undesirable results. Therefore, by using optimization methods to find the best cluster centers, the performance of this algorithm can be significantly improved. The idea of combining firefly and genetics algorithms to optimize clustering accuracy is an innovation that has not been used before.

**Method:** In order to optimize k-means clustering, in this paper, the combined method of genetic algorithm and firefly worm is introduced as the firefly genetic algorithm.

**Findings:** The proposed algorithm is evaluated using three well-known datasets, namely, Breast Cancer, Iris, and Glass. It is clear from the results that the proposed algorithm provides better results in all three datasets. The results confirm that the distance between clusters is much less than the compared approaches.

**Discussion and Conclusion:** The most important issue in clustering is to correctly determine the cluster centers. There are a variety of methods and algorithms that performs clustering with different performance. In this paper, based on firefly metaheuristic algorithms and genetic algorithms a new method has been proposed for data clustering. Our main focus in this study was on two determining factors, namely the distance within the data cluster (distance of each data to the center of the cluster) and the distance that the headers have from each other (maximum distance between the centers of the clusters). In the k-means algorithm, clustering is not accurate since the cluster centers are selected randomly. Employing firefly algorithms and genetics, we try to obtain more accurate centers of the clusters and, as a result, correct clustering.

**Keywords:** Data Mining, Genetic Algorithms, Firefly Algorithm, Firefly -Genetic Algorithm, K-Means Algorithm.

## روشی نوین جهت خوشه‌بندی داده مبتنی بر ترکیب الگوریتم‌های بهینه‌سازی ژنتیک و کرم شب‌تاب

سال دوم، زمستان ۱۴۰۰  
شماره چهارم، صص: ۳۵ - ۴۳

تاریخ دریافت: ۱۴۰۰/۰۵/۲۶  
تاریخ پذیرش: ۱۴۰۰/۰۷/۱۲

مهسا افسردیر<sup>۱</sup> منصوره افسردیر<sup>۲\*</sup>

۱. کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده فنی مهندسی، واحد علوم تحقیقات، دانشگاه آزاد اسلامی، تهران، ایران.

[M.afsardeir2012@yahoo.com](mailto:M.afsardeir2012@yahoo.com)

۲. کارشناسی ارشد، گروه مهندسی پزشکی، دانشکده فنی مهندسی، واحد دزفول، دانشگاه آزاد اسلامی، دزفول، ایران. (نویسنده مسئول)

[afsardeir2007@gmail.com](mailto:afsardeir2007@gmail.com)

**چکیده:** خوشه‌بندی یکی از مسائل مهم در داده‌کاوی است که بدون هدف از پیش تعیین شده، داده‌ها را براساس شباهت درون خوشه‌ها تقسیم‌بندی می‌کند. از روش‌های متداول خوشه‌بندی الگوریتم  $k$ -means است که با دریافت ورودی، داده‌ها را به  $k$  خوشه تقسیم‌بندی می‌کند. یکی از معایب این روش حساسیت به شرایط اولیه است که منجر به کاهش دقت در خوشه‌بندی می‌شود. از روش‌های بهبود عملکرد  $k$ -means می‌توان استفاده از الگوریتم‌های فراابتکاری را نام برد. این پژوهش به دو روش بهینه‌سازی ژنتیک و کرم شب‌تاب پرداخته و الگوریتم جدیدی با عنوان الگوریتم ژنتیکی کرم شب‌تاب جهت بهینه‌سازی خوشه‌بندی  $k$ -means ارائه کرده‌است. الگوریتم کرم شب‌تاب از الگوریتم‌های هوش جمعی است که از ویژگی نور چشمک‌زن کرم شب‌تاب الهام گرفته‌است؛ الگوریتم ژنتیک نیز نوعی الگوریتم فراابتکاری است که از تکنیک‌های زیست‌شناسی مانند وراثت و جهش استفاده می‌کند. در الگوریتم  $k$ -means چون مراکز خوشه‌ها به صورت تصادفی انتخاب می‌شوند، خوشه‌بندی دقت لازم را ندارد. لذا ما با استفاده از الگوریتم‌های فراابتکاری سعی در به دست آوردن مراکز دقیق خوشه‌ها و در نتیجه، دستیابی به خوشه‌بندی صحیح داریم. در روش پیشنهادی، ابتدا الگوریتم  $k$ -means را روی داده‌های ورودی اجرا کرده و خوشه‌بندی انجام می‌شود. سپس مضرری از مراکز خوشه که در این الگوریتم به دست آمده را به عنوان حد پایین و حد بالای الگوریتم پیشنهادی استفاده می‌کنیم. جمعیت اولیه به صورت تصادفی بین حد پایین و حد بالا تولید می‌شود. در حلقه اصلی الگوریتم جمعیت را به دو دسته جمعیت مساوی تقسیم می‌کنیم، بر روی دسته اول الگوریتم ژنتیک را اجرا می‌کنیم، بر روی دسته دوم براساس الگوریتم کرم شب‌تاب، موقعیت‌های جدید را به دست می‌آوریم. حال جمعیت قبلی و جمعیت جدید به دست آمده از الگوریتم ژنتیک و جمعیت جدید به دست آمده از الگوریتم کرم شب‌تاب را تلفیق کرده و آن‌ها را از خوب به بد مرتب می‌کنیم و تعداد مورد نیاز از آن‌ها را انتخاب و به ابتدای حلقه می‌رویم. این فرایند را تا برقراری شرط توقف ادامه می‌دهیم. در پایان الگوریتم  $k$ -means، الگوریتم کرم شب‌تاب، الگوریتم ژنتیک و الگوریتم پیشنهادی بر روی سه مجموعه داده اعمال شده و نتایج مقایسه شد. نتایج شبیه‌سازی نشان می‌دهد که الگوریتم ژنتیکی کرم شب‌تاب عملکرد بهتری در مقایسه با سایر روش‌ها داشته‌است.

**واژه‌های کلیدی:** داده‌کاوی، الگوریتم ژنتیک، الگوریتم کرم شب‌تاب، الگوریتم ژنتیکی کرم شب‌تاب، خوشه‌بندی  $k$ -means.

## ۱. مقدمه

با پیشرفت فناوری و افزایش حجم داده در پایگاه داده‌ها، ضرورت کشف و استخراج سریع و دقیق از پایگاه داده‌ها بیش از پیش نمایان شده‌است. بنابراین نیاز به طراحی سیستم‌هایی که قادر به کشف دانش و اطلاعات مورد نیاز از سان باشند، به‌خوبی اح‌ساس می‌شود. داده‌کاوی مهم‌ترین فناوری برای بهره‌وری مؤثر، صحیح و سریع از داده‌های حجیم است و اهمیت آن رو به فزونی است [۱]. داده‌کاوی پل ارتباطی میان علم آمار، علم کامپیوتر، هوش مصنوعی، الگوشناسی، فراگیری ماشین داده می‌باشد. داده‌ها اغلب حجیم‌اند و به‌تنهایی قابل‌استفاده نیستند، اما اطلاعات نهفته در داده‌ها قابل‌استفاده است. بنابراین بهره‌گیری از فرآیند داده‌کاوی جهت شناسایی الگوها و مدل‌ها و نیز ارتباط عناصر مختلف در پایگاه داده جهت کشف دانش نهفته در داده‌ها و نهایتاً تبدیل داده به اطلاعات، روزبه‌روز ضروری‌تر می‌شود. خوشه‌بندی یکی از رویکردهای داده‌کاوی است که در جهت تحلیل و تفسیر داده‌ها مطرح می‌شود. در واقع خوشه‌بندی، به‌دنبال کشف ساختار در داده‌های جمع‌آوری شده می‌باشد [۲]. این روش با استفاده از شباهت‌ها یا تفاوت‌ها (مانند فاصله) موجود میان نقاط داده در یک مجموعه به کشف ساختار می‌پردازد. از آنجا که خوشه‌بندی یک روش بدون نظارت محسوب می‌شود، در زمینه‌های گوناگون از جمله پردازش تصویر، مطالعات زمین‌لرزه، بازاریابی، زیست‌شناسی، داده‌کاوی از اهمیت ویژه‌ای برخوردار است. یکی از روش‌های پرکاربرد خوشه‌بندی روش k-means است. هدف اصلی k-means افزایش شباهت داخلی خوشه‌ها و کاهش شباهت اشیا بیرون خوشه‌ها است. در این الگوریتم مراکز خوشه‌ها به صورت تصادفی انتخاب می‌شوند و هر شی را با توجه به بیشترین شباهت آن به مراکز خوشه‌ها، به خوشه‌ها تخصیص می‌دهند. از این‌رو، این الگوریتم برای داده‌های پرت مناسب نیست زیرا این داده‌ها به راحتی مراکز را تغییر می‌دهند و ممکن است نتایج نامطلوبی حاصل شود. بنابراین با استفاده از روش‌های بهینه‌سازی جهت یافتن بهترین مراکز خوشه‌ها می‌توان تا حد بسیار زیادی عملکرد این الگوریتم را بهبود بخشید. در سال‌های اخیر الگوریتم‌های هوشمند بسیاری در زمینه بهینه‌سازی خوشه‌بندی صورت گرفته‌است. هدف اصلی در این گونه تحقیقات، انتخاب صحیح مرکز خوشه‌ها برای جدا کردن نمونه‌ها از یکدیگر و قراردادن نمونه‌های مشابه در یک خوشه و در نتیجه گروه‌بندی صحیح می‌باشد. بنابراین با توجه به اهمیت روش‌های خوشه‌بندی، با ارائه الگوریتم‌های فراابتکاری و ترکیب این الگوریتم‌ها با یکدیگر می‌توان مسائل خوشه‌بندی را نسبت به حل آن‌ها به‌تنهایی با جواب بهینه‌تری تولید کرد [۳] [۴] [۵] [۱۱]. محاسبه راه‌حل‌های بهینه برای اکثر مسائل بهینه‌سازی که در بسیاری از زمینه‌های کاربردی و عملی مشاهده می‌شوند، کار دشواری است. روش‌های حل مسائل بهینه‌سازی مشتمل بر دو دسته روش‌های دقیق و روش‌های ابتکاری می‌باشد. روش‌های دقیق راه‌حل‌های بهینه را به‌دست آورده و شرایط بهینگی را تضمین می‌کنند. روش‌های فراابتکاری راه‌حل‌های با کیفیت

بالا را در زمان معقولی، تولید می‌کنند، اما تضمینی برای یافتن راه‌حل بهینه سراسری ندارند.

از جمله الگوریتم‌های خوشه‌بندی می‌توان به الگوریتم ژنتیک و الگوریتم کرم شب‌تاب که امروزه به‌طور گسترده در مسائل مختلف استفاده می‌شوند، اشاره کرد. با ارائه الگوریتم‌های فراابتکاری و ترکیب این الگوریتم‌ها با یکدیگر می‌توان مسائل خوشه‌بندی را نسبت به حل آن‌ها به‌تنهایی با جواب بهینه‌تری تولید کرد [۱۱].

در کوششی الگوریتم کرم شب‌تاب بهبود یافته را با الگوریتم K-means ترکیب کرده‌اند تا بتواند با این الگوریتم جدید، مراکز صحیح و با دقت بالا را برای خوشه‌ها به‌صورت بهینه پیدا کند. [۱۶] همچنین الگوریتم کرم شب‌تاب را برای خوشه‌بندی به‌منظور کاهش مشکلاتی که در خوشه‌بندی وجود دارد به‌کار برده است و در نهایت با الگوریتم زنبور عسل و الگوریتم بهینه‌سازی ازدحام ذرات مقایسه کرده‌است. [۱۷] از جمله الگوریتم‌های دیگری که برای حل مشکلات خوشه‌بندی ارائه شده، الگوریتم K-FA است که از الگوریتم کرم شب‌تاب برای یافتن مراکز خوشه بهینه استفاده می‌کند سپس الگوریتم k-means از این مراکز برای اصلاح و مقداردهی مراکز خوشه استفاده می‌کند و موجب بهره‌وری بیشتر می‌شود. [۱۸]

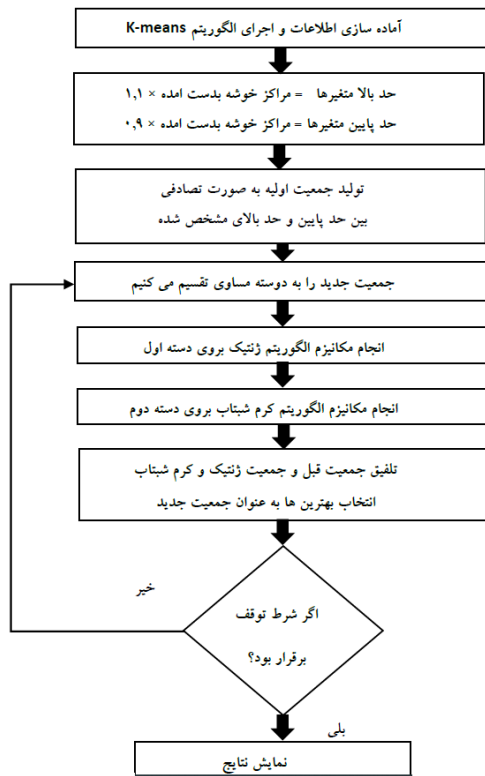
در [۱۹] از ترکیب الگوریتم k-means و الگوریتم ژنتیک برای بهبود کیفیت خوشه‌بندی، ثابت در خوشه‌بندی و امکان جستجوی سراسری بهتر استفاده شده‌است.

بنابراین ایده ترکیب الگوریتم کرم شب‌تاب و ژنتیک به‌منظور بهینه‌سازی دقت خوشه‌بندی یک ابتکار و نوآوری است که تاکنون استفاده نشده‌است. بنابراین هدف از این مقاله ترکیب دو الگوریتم کرم شب‌تاب و الگوریتم ژنتیک جهت افزایش کارایی خوشه‌بندی k-means است. در ادامه، به مبانی k-means و روش‌های بهینه‌سازی ژنتیک، کرم شب‌تاب و شرح روش پیشنهادی این مقاله و نحوه پیاده‌سازی آن پرداخته شده‌است. سپس، نتایج خوشه‌بندی با الگوریتم پیشنهادی را بر روی داده‌های استاندارد مشابه مقایسه کرده و در نهایت به نتیجه‌گیری پرداخته‌ایم.

## ۲. روش پیشنهادی

هدف از الگوریتم k-means بخش‌بندی داده درون k خوشه به روشی است که شباهت درون خوشه‌ای کم و فاصله بین خوشه‌ای بالا باشد. اما به دلیل حساس بودن به شرایط اولیه دقت خوشه‌بندی را کاهش می‌دهد و نمی‌تواند جواب بهینه‌ای ایجاد کند. در این الگوریتم مراکز خوشه‌ها به صورت تصادفی انتخاب می‌شوند و هر شی را با توجه به بیشترین شباهت آن به مراکز خوشه‌ها، به خوشه‌ها تخصیص می‌دهند. از این‌رو، این الگوریتم برای داده‌های پرت مناسب نیست زیرا این داده‌ها به راحتی مراکز را تغییر می‌دهند و ممکن است نتایج نامطلوبی حاصل شود.

در سال‌های اخیر الگوریتم‌های هوشمند بسیاری در زمینه بهینه‌سازی خوشه‌بندی صورت گرفته‌است. هدف اصلی در این گونه تحقیقات، انتخاب صحیح مرکز خوشه‌ها به‌منظور جدا کردن نمونه‌ها از



شکل ۱: فرآیند روش پیشنهادی

## ۲.۲. الگوریتم بهینه‌سازی کرم شبتاب

الگوریتم کرم شبتاب الگوریتمی برگرفته از طبیعت است که رفتار اجتماعی کرم‌های شبتاب را شبیه‌سازی می‌کند. این الگوریتم توسط یانگ در سال ۲۰۰۹ معرفی گردید [۸]. کرم‌های شبتاب نورهایی تولید می‌کنند که الگوی نوری هر کدام با دیگری متفاوت است. آن‌ها به‌منظور جذب جفت و شکار از این نور استفاده می‌کنند، میزان این نور رابطه مستقیم با جاذبیت کرم شبتاب دارد با در نظر گرفتن میزان نور هر کرم به‌عنوان مقدار تابع هدف، می‌توان رفتار کرم‌های شبتاب را به‌صورت یک الگوریتم بهینه‌ساز مدل نمود [۸] [۱۳]. الگوریتم کرم شبتاب بر پایه جذب و روشنایی است. این موضوع موجب تقسیم‌بندی خودکار کل جمعیت در زیر گروه‌ها با یک متوسط فاصله معین می‌شود و هر گروه می‌تواند در اطراف یک بهینه محلی ازدحام کند. در میان همه این‌ها، بهترین بهینه سراسری می‌تواند یافت شود. دوم، این تقسیم‌بندی اجازه یافتن همه بهینه‌های زمان را می‌دهد. این قابلیت الگوریتم کرم شبتاب را مخصوصاً برای مسائل بهینه‌سازی چند کیفیتی غیرخطی مناسب می‌سازد. به‌علاوه پارامترها در الگوریتم کرم شبتاب می‌تواند برای کنترل تصادفی به نسبت تکرار تنظیم شود، به‌گونه‌ای که همگرایی نیز می‌تواند با میزان‌سازی این پارامترها سرعت یابد. این مزایا الگوریتم کرم شبتاب را برای قابلیت سروکار داشتن با مسائل پیوستگی، خوشه‌بندی، کلاس‌بندی و به‌علاوه بهینه‌ساز ترکیبی مناسب می‌سازد [۱۱].

یکدیگر و قراردادن نمونه‌های مشابه در یک خوشه و در نتیجه گروه‌بندی صحیح می‌باشد.

از آن‌جا که روش‌های موجود نمی‌توانند مرکز خوشه را با دقت بالایی محاسبه‌کنند بنابراین به خوشه‌بندی دقیق‌تر با استفاده از ترکیب روش‌های بهینه‌سازی نیاز داریم. از این‌رو، می‌توان از یک خوشه‌بندی دقیق‌تر برای بهبود شناسایی اشیاء در تصویر و لبه‌یابی دقیق‌تر استفاده کرد.

در این مقاله با ترکیب دو الگوریتم ژنتیک و کرم شبتاب سعی در بهبود عملکرد الگوریتم k-means داشته‌ایم. با توجه به این که هر یک از الگوریتم‌های فراابتکاری دارای نقاط قوتی هستند با ترکیب این الگوریتم‌ها تحت عنوان روشی به نام الگوریتم ژنتیکی کرم شبتاب می‌توان کارایی را افزایش داد. به همین منظور، در این تحقیق از ترکیب مزایای الگوریتم ژنتیک و کرم شبتاب برای بهبود عملکرد الگوریتم k-means استفاده شده است. شکل ۱ فرآیند روش پیشنهادی را نمایش می‌دهد. در ادامه الگوریتم پیشنهادی به تفصیل شرح داده خواهد شد.

## ۱.۲. خوشه‌بندی k-means

الگوریتم خوشه‌بندی k-means یکی از تقسیم‌بندی‌های اساسی در تجزیه و تحلیل خوشه‌ها است و به دلیل سادگی، پیاده‌سازی آسان و همگرایی سریع، کاربرد مؤثر و پیچیدگی زمانی خطی به‌طور گسترده‌ای استفاده می‌شود [۶]. این الگوریتم ابتدا k نقاط داده را به‌صورت تصادفی به‌عنوان مراکز ثقل اولیه در نظر می‌گیرد و سپس هر داده را به نزدیکترین خوشه اختصاص داده تا زمانی که معیار همگرایی صورت‌پذیرد [۷]. مراحل این الگوریتم به شرح زیر است:

- ۱) به صورت تصادفی k شیء دلخواه را از مجموعه داده D به عنوان مرکز خوشه اولیه انتخاب می‌کند.
- ۲) هر شیء را با توجه به بیشترین شباهت آن به مراکز خوشه‌ها، به خوشه‌ها اختصاص می‌دهد.
- ۳) به‌روزرسانی میانگین خوشه‌های جدی، مثلاً محاسبه مقدار میانگین اشیاء برای هر خوشه.
- ۴) با توجه به مراکز جدید خوشه‌ها تا زمانی که هیچ تغییری رخ ندهد به گام دوم برگرد.

میزان جذب کرم‌های شب‌تاب ( $\beta$ ) نسبی بوده و به فاصله بین دو کرم شب‌تاب ( $r$ ) و ضریب جذب نور ( $\gamma$ ) بستگی دارد که از رابطه (۱) قابل-محاسبه است.

$$\beta(r) = \beta_0 e - e^{-\gamma r^2} \quad (1)$$

در رابطه (۱)،  $\beta_0$  نشان‌دهنده میزان جذب در منبع نور می‌باشد. فاصله بین هر دو کرم شب‌تاب  $i$  و  $j$  در نقطه  $x_i$  و  $x_j$  با رابطه (۲) محاسبه می-شود.

$$r_{ij} = \|x_i - x_j\| = \sqrt{\sum_{k=1}^d (x_{i,k} - x_{j,k})^2} \quad (2)$$

در رابطه (۲)،  $x_{i,k}$  نشان‌دهنده  $k$  امین بخش از مختصات فضایی کرم شب‌تاب  $i$  می‌باشد. همچنین حرکت کرم شب‌تاب و جذب  $i$  به کرم شب‌تاب  $j$  که درخشان‌تر است از رابطه (۳) محاسبه می‌شود.

$$x_n^i = x_n^{i-1} + \beta_0 e^{-\gamma r_{ij}^2} (x_{nj}^{i-1} - x_{ni}^{i-1}) + \alpha e_i^{n-1} \quad (3)$$

در رابطه (۳)،  $x_i$  موقعیت کرم شب‌تاب کم‌نورتر،  $x_j$  موقعیت کرم درخشان‌تر،  $n$  شماره تکرار،  $\alpha$  عددی تصادفی و  $e_i^{n-1}$ ، یک بردار از اعداد تصادفی است که می‌تواند دارای توزیع یکنواخت یا گوسی باشد.

### ۳.۲. الگوریتم ژنتیک

از الگوریتم ژنتیک در مسائل جستجو و بهینه‌سازی استفاده می‌شود. ابتدا یک نسل اولیه (به‌صورت تصادفی) ایجاد می‌گردد، که در واقع کروموزوم‌های اولیه هستند [۹]. هر یک از این کروموزوم‌ها جوابی برای مسئله هستند اما، نه جواب اصلی که ما به دنبال آن هستیم. سپس پدیده جهش با احتمال خیلی کم ممکن است رخ دهد. در نهایت کروموزوم‌ها از نظر امتیاز رتبه‌بندی می‌شوند، این امتیازدهی معمولاً بر اساس مقدار تابع هدف است. برخی کروموزوم‌ها با هم ترکیب شده و نسل بعد را به‌وجود می‌آورند، احتمال انتخاب کروموزوم‌های با امتیاز بالاتر، بیشتر است اما در عین حال، احتمال انتخاب شدن برای تمام کروموزوم‌ها حتی کروموزوم‌های با کمترین امتیاز وجود دارد. با نسل جدید به‌وجود آمده این مراحل را تکرار می‌شود تا به جواب مطلق برسیم [۱۰][۱۲][۱۳].

از نظر الگوریتمی، مراحل اصلی GA به شرح زیر است:

- (۱) به‌طور تصادفی یک گروه اولیه از راه‌حل‌ها ایجاد می‌شود.
- (۲) ارزش تناسب هر راه‌حل، در گروه، ارزیابی می‌شود.
- (۳) با تکرار و اعمال مراحل زیرگروه جدیدی از راه‌حل‌ها را ایجاد می‌شود:

- دو والد از گروه راه‌حل‌ها بر اساس مکانیسم انتخاب، انتخاب می‌شوند.
- عمل تقاطع و جهش انجام می‌شود.
- کروموزوم‌های شایسته‌تر جایگزین می‌شوند.

(۴) در صورتی که شرط پایان برقرار شده باشد، پایان می‌یابد. در غیر این صورت به مرحله قبل برمی‌گردد. [۱۳][۱۴][۱۵]

در الگوریتم‌های ژنتیکی، در طی مرحله تولیدمثل از عملگرهای ژنتیکی استفاده می‌شود. با تأثیر این عملگرها بر روی یک جمعیت، نسل بعدی آن جمعیت تولید می‌شود. این عملگرها شامل:

**عملگر تقاطع:** بر روی یک زوج کروموزوم از نسل مولد عمل کرده و یک زوج کروموزوم جدید تولید می‌کند. به عبارت دیگر، عملگر تقاطع در یک لحظه بر روی دو کروموزوم اعمال شده و دو نوزاد به‌وسیله ترکیب ساختار دو کروموزوم ایجاد می‌کند. روش‌های تقاطع مختلفی وجود دارد: تقاطع تک‌نقطه‌ای (نقطه‌ای به صورت تصادفی به‌عنوان نقطه برش در دو کروموزوم والد انتخاب شده و از آن نقطه به بعد با هم جابجا می‌شوند)، تقاطع دو نقطه‌ای (به‌طور تصادفی دو نقطه از کروموزوم را انتخاب کرده و تمام ژن‌های بین این دو نقطه را در دو کروموزوم تعویض می‌کنیم) و ...

**عملگر جهش:** این عملگر در کروموزوم‌های متفاوت تغییرات تصادفی برنامه‌ریزی نشده‌ای ایجاد کرده و ژن‌هایی را که در جمعیت اولیه وجود نداشته‌اند را وارد جمعیت می‌کند. نقش جهش در الگوریتم ژنتیک بازگرداندن مواد ژنتیکی گم شده و یا پیدایش شده داخل جمعیت است. تا از همگرایی زودرس الگوریتم به جواب‌های بهینه محلی جلوگیری شود. عملگرهای جهش مختلفی وجود دارد، از میان آن‌ها جهش یکنواخت را به‌صورت مختصر توضیح خواهیم داد. در این عملگر، ژنی از کروموزوم به‌طور تصادفی انتخاب شده است و مقدار آن به مقدار تصادفی دیگری تبدیل می‌شود. ابتدا یک عدد تصادفی در بازه  $[0, 1]$  که طول کروموزوم مورد نظر است، انتخاب شده و ژن موجود در آن مکان از کروموزوم تغییر می‌کند.

**عملگر انتخاب:** برای انتخاب بهترین جواب‌ها در تولید مجدد نسل (تولید جمعیت جدید) باید از روشی استفاده کرد که تا حد ممکن بهترین جواب را انتخاب کند. عملگر انتخاب چرخ‌رولت یک روش انتخاب است که در آن عنصری که عدد برازش (تناسب) بیشتری داشته باشد، انتخاب می‌شود. در واقع به نسبت عدد برازش برای هر عنصر یک احتمال تجمعی نسبت می‌دهیم با این احتمال که شانس انتخاب هر عنصر تعیین شود.

### ۴.۲. الگوریتم ژنتیکی کرم شب‌تاب

هدف از الگوریتم k-means بخش‌بندی داده درون  $k$  خوشه به‌طوری که شباهت درون خوشه‌ای بالا، فاصله بین خوشه‌ای بالا باشد اما به دلیل حساس بودن به شرایط اولیه دقت خوشه‌بندی را کاهش داده و جواب بهینه‌ای به ما نمی‌دهد. جهت بهینه‌سازی خوشه‌بندی k-means در

این مقاله روش ترکیبی الگوریتم ژنتیک و کرم شب تاب با عنوان الگوریتم ژنتیکی کرم شب تاب معرفی شده است. با توجه به این که هریک از الگوریتم های فراابتکاری دارای نواقص و نقاط قوتی هستند با ترکیب این الگوریتم ها می توان کارایی را افزایش داد. در الگوریتم پیشنهادی ابتدا اطلاعات مورد نظر با استفاده از روش k-means خوشه بندی می شوند. سپس مضربی از مراکز خوشه که در این روش به دست آمده است را به عنوان حد پایین و حد بالای الگوریتم پیشنهادی استفاده می کنیم. که به ترتیب با رابطه (۴) و رابطه (۵) محاسبه می شود.

$$(4) \text{ مراکز خوشه به دست آمده } = 1/1 \times \text{حد بالای خوشه ها}$$

$$(5) \text{ مراکز خوشه به دست آمده } = 0/9 \times \text{حد پایین خوشه ها}$$

پس از این مرحله جمعیت اولیه به صورت تصادفی بین حد پایین و حد بالا تولید می شود. در حلقه اصلی الگوریتم جمعیت را به دو دسته جمعیت مساوی تقسیم می کنیم. بر روی دسته اول الگوریتم ژنتیک را اجرایی کنیم:

- با احتمال  $p_c$  تقاطع تک نقطه ای را اجرایی کنیم.

عملگر تقاطع در یک لحظه بر روی دو کروموزوم اعمال شده و دو نوزاد به وسیله ترکیب ساختار دو کروموزوم ایجاد می کند. مفهوم مهمی که در ارتباط با این عملگر مطرح است نرخ تقاطع  $p_c$  است. اگر تعداد کروموزوم های تولید شده را با  $x$  و اندازه جمعیت اولیه را با  $pop - size$  نشان دهیم، خواهیم داشت:

$$(6) p_c = \frac{x}{pop - size}$$

با احتمال  $p_m$  جهش یکنواخت را اجرایی کنیم.

این عملگر در کروموزوم های متفاوت تغییرات تصادفی بر نامهریزی نشده ای ایجاد کرده و ژن هایی را که در جمعیت اولیه وجود نداشته اند، وارد جمعیت می کند. درباره این عملگر مفهوم مهمی مطرح شده که به آن نرخ جهش  $p_m$  گویند.

با استفاده از عملگر انتخاب چرخ رولت برای ورود به جمعیت جدید انتخاب می شود.

برای انتخاب بهترین جواب ها و تولید مجدد نسل (تولید جمعیت جدید) باید از روشی استفاده کرد که تا حد ممکن بهترین جواب را انتخاب کند. در میان روش های مختلف انتخاب ما چرخ رولت پیشنهادی هالند است.

بر روی دسته دوم بر اساس الگوریتم کرم شب تاب موقعیت های جدید را به دست می آوریم.

میزان جاذبیت هر کرم شب تاب را با توجه به ضریب جذب معادله (۷) محاسبه می کنیم.

$$\beta(r) = \beta_0 e - e^{-\gamma r^2}$$

که  $\beta_0$  ،  $\gamma$  و  $r$  به ترتیب میزان جذب از پیش تعریف شده، ضریب جذب نور و فاصله بین دو کرم شب تاب  $i$  و  $j$  می باشند.

سپس حرکت هر کرم شب تاب به سمت کرم شب تاب جاذب تر با استفاده از فرمول (۸) محاسبه خواهد شد.

$$(8) x_n^i = x_n^{i-1} + \beta_0 e^{-\gamma r_{ij}^2} (x_{nj}^{i-1} - x_{ni}^{i-1}) + \alpha e_i^{n-1}$$

که  $x_i$  موقعیت کرم شب تاب کم نورتر،  $x_j$  موقعیت کرم شب تاب درخشان تر،  $n$  شماره تکرار  $\alpha$  عددی تصادفی و  $e_i^{n-1}$  یک بردار از اعداد تصادفی است که می تواند دارای توزیع یکنواخت یا گوسی باشد.

حال جمعیت قبلی و جمعیت جدید حاصل از الگوریتم ژنتیک و جمعیت جدید حاصل از الگوریتم کرم شب تاب را تلفیق کرده و آن ها را از خوب به بد مرتب می کنیم و به تعداد مورد نیاز از آن ها انتخاب و به ابتدای حلقه می رویم. این فرآیند را تا برقراری شرط توقف ادامه می دهیم که در این تحقیق حداکثر تعداد تکرار در نظر گرفته شده است.

## ۵.۲. ارزیابی الگوریتم پیشنهادی

در این پژوهش از فاصله درون خوشه ای و فاصله بین مراکز خوشه به عنوان تابع هزینه برای الگوریتم ها استفاده می شود و تابع هزینه نهایی با استفاده از جمع دو این الگوریتم با ضرایب اهمیت به ترتیب  $w_1$  و  $w_2$  محاسبه خواهد شد. فاصله درون خوشه ای با توجه به فرمول فاصله اقلیدسی به صورت رابطه (۹) محاسبه خواهد شد.

$$(9) \text{ فاصله درون خوشه ای} = \sum data_i - cluster_i$$

که در رابطه (۹)  $data_i$  داده نام،  $cluster_i$  مرکز خوشه نام و  $n$  تعداد جمعیت است. شباهت بین مراکز خوشه به صورت رابطه (۱۰) محاسبه می شود.

$$(10) \text{ فاصله برون خوشه ای} = \sum_{i=1}^c \sum_{j=1}^c (cluster_i - cluster_j)$$

که در این رابطه  $cluster_i$  مرکز خوشه نام،  $cluster_j$  مرکز خوشه نام و  $c$  تعداد مراکز خوشه است. ارزیابی نهایی با تابع هدف نهایی نیز با استفاده از رابطه (11) و ضریب  $w_1$  و  $w_2$  به ترتیب برای فاصله درون خوشه ای و فاصله بین مراکز خوشه محاسبه خواهد شد.

$$(11) \text{ Cost function} = w_1 \times (\text{فاصله درون خوشه ای}) + w_2 \times (\text{فاصله برون خوشه ای})$$

## ۶.۲. پارامترهای روش پیشنهادی

برای یک مقایسه مناسب، باید شرایط حل برای همه الگوریتم ها یکسان باشد از این رو، پارامترهای مربوط به الگوریتم ها در جدول (۲) نشان داده شده است. لازم به ذکر است که الگوریتم پیشنهادی تمام پارامترهای مربوط به الگوریتم های ژنتیک و کرم شب تاب را دارا خواهد بود و به دلیل

مجموعه داده Glass	۲۱۹	۷	۴
-------------------	-----	---	---

احتمالی بودن فضای جستجوی الگوریتم‌ها، هریک از الگوریتم‌ها را ۳۱ بار اجرا کرده‌ایم و نتایج در ادامه آمده‌است.

**پایگاه داده iris:** این مجموعه داده مربوط به اطلاعات سه نمونه گل زنبق است، که در سال ۱۹۸۷ توسط فیشر تهیه شده‌است. این مجموعه داده دارای ۱۵۱ نمونه و ۳ کلاس است که هرکدام از این نمونه‌ها دارای ۹ ویژگی هستند. ویژگی‌های مربوط به این مجموعه داده عبارتند از: طول کاسبرگ، پهنای کاسبرگ، طول گلبرگ، پهنای گلبرگ.

**پایگاه داده Glass:** این مجموعه داده مربوط به اطلاعات شناسایی شیشه است، که در سال ۱۹۸۷ تهیه شده‌است. این مجموعه داده دارای ۲۱۹ نمونه و ۷ کلاس است که هرکدام از این نمونه‌ها دارای ۴ ویژگی هستند. ویژگی‌های مربوط به این مجموعه داده عبارتند از: سدیم، منیزیم، آلومینیوم، سیلیکون، پتاسیم، کلسیم، باریوم، آهن و شاخص انکساری.

### ۸.۲. نتایج و تجزیه و تحلیل داده‌ها

در این قسمت نتایج شبیه سازی مربوط به الگوریتم‌های ارائه شده در قسمت قبل را بررسی خواهیم کرد. الگوریتم‌های k-mean، ژنتیک، کرم شب‌تاب و الگوریتم ژنتیکی کرم شب‌تاب بر روی ۳ مجموعه داده اجرا و نتایج مقایسه شده است. همان‌طور که در قسمت ۲-۵ اشاره شد تابع شایستگی مجموع وزن دار فاصله درون خوشه‌ای و فاصله مراکز خوشه‌ها با ضریب  $w_1$  و  $w_2$  است. حال برای تحلیل تابع شایستگی؛ هر سه حالت را به صورت کلی بررسی می‌کنیم. جدول (۳) تا (۵) به ترتیب نتایج را بر اساس فاصله بین خوشه‌ای و درون خوشه‌ای، فاصله درون خوشه‌ای و شباهت بین مراکز خوشه‌ای نمایش می‌دهد.

در ادامه الگوریتم‌های ژنتیک، کرم شب‌تاب و ژنتیکی کرم شب‌تاب برای هریک از مجموعه داده‌های موردنظر، اجرا شده و نتایج به صورت نمودار هم‌گرایی برای ۱۰۰ بار تکرار در شکل (۲) تا (۵) نشان داده شده‌است.

### جدول ۳: نتایج بر اساس فاصله بین خوشه‌ای و درون خوشه‌ای

پایگاه داده / الگوریتم	k-means	الگوریتم ژنتیک	الگوریتم کرم شب‌تاب	ژنتیکی کرم شب‌تاب
مجموعه داده Breast Cancer	۳۲۳/۸۷	۳۲۲/۷۴	۳۲۲/۵۴	۳۲۱/۹۸
مجموعه داده Iris	۳۱/۲۸	۳۰/۹۹	۳۰/۸۵	۳۰/۴۱
مجموعه داده Glass	۶۸/۵۴	۶۸/۲۳	۶۸/۳۲	۶۷/۸۷

### جدول ۴: نتایج بر اساس شباهت بین مراکز خوشه‌ای

پایگاه داده / الگوریتم	k-means	الگوریتم ژنتیک	الگوریتم کرم شب‌تاب	ژنتیکی کرم شب‌تاب
مجموعه داده Breast Cancer	۳۲۸/۶۵	۳۲۷/۶۵	۳۲۷/۶۳	۳۲۶/۹۹

### جدول ۱: مقادیر اولیه پارامترهای الگوریتم‌های فراابتکاری

نام الگوریتم	نام پارامتر	اندازه پارامتر
k-means	تعداد تکرار	۱۰۰
	تعداد مراکز خوشه	۳
	اندازه جمعیت	۱۰
	تعداد مراکز خوشه	۳
	درصد جهش	۰/۱
	درصد تقاطع	۰/۹
الگوریتم ژنتیک	اندازه جمعیت	۱۰
	تعداد مراکز خوشه	۳
	ضریب جذب نور	۰/۹
	میزان جذب	۲
	شدت نور	۱
	تعداد تکرار	۱۰۰
الگوریتم کرم شب‌تاب	تعداد تکرار	۱۰۰
	تعداد مراکز خوشه	۳
	اندازه جمعیت	۱۰
	درصد جهش	۰/۱
	درصد تقاطع	۰/۹
	ضریب جذب نور	۰/۹
الگوریتم ژنتیکی کرم شب‌تاب	ضریب جذب نور	۰/۹
	میزان جذب	۲
	شدت نور	۱

### ۷.۲. پایگاه داده

برای ارزیابی مناسب الگوریتم‌های فراابتکاری از ۳ مجموعه داده شناخته شده برای مقایسه استفاده شده‌است؛ نام و مشخصات این پایگاه داده‌ها در جدول (۲) آمده‌است.

**پایگاه داده Breast Cancer:** این مجموعه داده مربوط به سرطان سینه است که در سال ۱۹۹۱ از دانشگاه ویسکانسین توسط دکتر ویلیام اچ. ولبرگ در بیمارستان مدیسون جمع‌آوری شده‌است. این مجموعه داده دارای ۶۴۴ نمونه و ۲ کلاس است که هرکدام از این نمونه‌ها دارای ۴ ویژگی هستند. ویژگی‌های مربوط به این مجموعه داده عبارتند از: ضخامت توده، یکنواختی اندازه سلول، یکنواختی شکل سلول، اندازه سلول اپیتelial، کروماتین مطلوب، هسته آشکار، هسته عادی، چسبندگی حاشیه‌ای و تقسیم غیرمستقیم هسته سلول.

### جدول ۲: مقادیر اولیه پایگاه داده‌ها

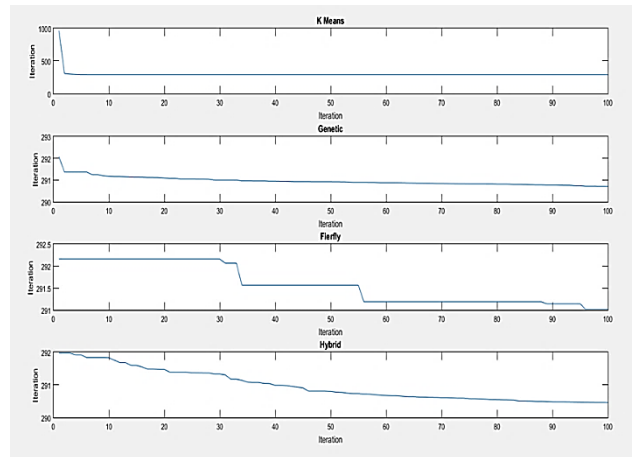
نام پایگاه داده	تعداد نمونه‌ها	تعداد کلاس‌ها	تعداد ویژگی‌ها
مجموعه داده Breast Cancer	۶۴۴	۲	۹
مجموعه داده Iris	۱۵۱	۳	۹

۲۸/۶۵	۲۹/۵۱	۲۹/۶۵	۲۹/۷۴	مجموعه داده Iris مجموعه داده Glass
۶۵/۹۱	۶۶/۷۸	۶۶/۲۳	۶۶/۵۴	

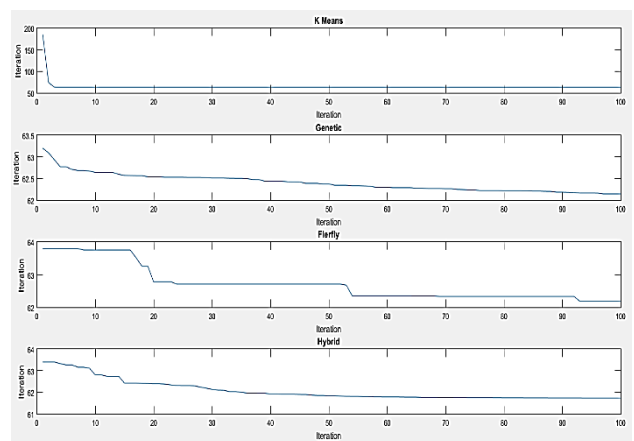
جدول ۵: نتایج بر اساس فاصله بین خوشه‌های و درون خوشه‌های

الگوریتم	الگوریتم کرم	الگوریتم ژنتیک	k-means	پایگاه داده / الگوریتم
ژنتیکی کرم	شب‌تاب	ژنتیک		
شب‌تاب				
۳/۵۴	۳/۸۷	۴/۲۱	۴/۳۹	مجموعه داده Breast Cancer
۲/۰۹	۲/۳۵	۲/۳۱	۲/۴۳	مجموعه داده Iris
۲/۰۱	۲/۳۲	۲/۴۵	۲/۵۶	مجموعه داده Glass

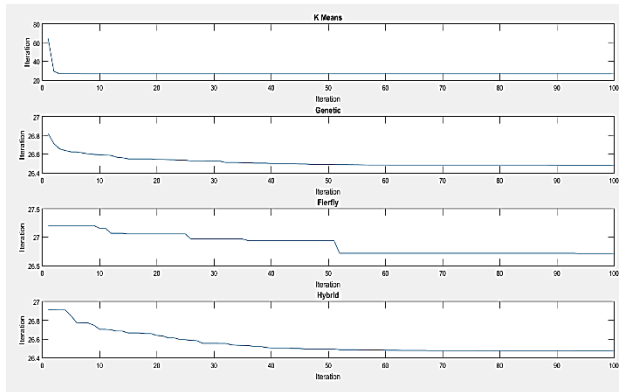
همان‌طور که از نتایج جدول (۳) تا (۵) مشخص است الگوریتم پیشنهادی نتایج بهتری در هر سه نوع تابع هدف به‌دست می‌دهد. همچنین از نتایج کاملاً پیداست که فاصله بین خوشه‌های بسیار کمتر از دو نوع دیگر می‌باشد.



شکل ۲: نمودار همگرایی اجرای الگوریتم‌ها برای مجموعه داده BC



شکل ۳: نمودار همگرایی اجرای الگوریتم‌ها برای مجموعه داده Glass



شکل ۴: نمودار همگرایی اجرای الگوریتم‌ها برای مجموعه داده Iris

### ۳. نتیجه‌گیری

در خوشه‌بندی مهمترین موضوع، تعیین صحیح مراکز خوشه، برای یک خوشه‌بندی صحیح است. روش‌ها و الگوریتم‌های بسیاری در این زمینه ارائه شده که هر کدام با عملکرد متفاوت، خوشه‌بندی را انجام می‌دهند. در این مقاله روش پیشنهادی جدیدی بر پایه دو الگوریتم فراابتکاری کرم شب‌تاب و الگوریتم ژنتیک برای خوشه‌بندی داده ارائه شده است. بیشترین تمرکز ما در این پژوهش بر دو عامل تعیین‌کننده فاصله درون خوشه‌های (فاصله هر داده تا مرکز خوشه) و میزان فاصله‌ای که سرخوشه‌ها از یکدیگر دارند (بیشترین فاصله مراکز خوشه‌ها) بوده است. در الگوریتم k-means چون مراکز خوشه به‌صورت تصادفی انتخاب می‌شود، خوشه‌بندی دقت لازم را ندارد. با استفاده از الگوریتم‌های فراابتکاری کرم شب‌تاب و ژنتیک سعی در به‌دست آوردن مراکز دقیق خوشه‌ها و در نتیجه، خوشه‌بندی صحیح بودیم. همان‌گونه که از نتایج شبیه‌سازی‌ها بر روی داده‌های استاندارد UCI پیداست، الگوریتم ژنتیک و ترکیب آن با الگوریتم کرم شب‌تاب، نسبت به الگوریتم k-means، الگوریتم کرم شب‌تاب و الگوریتم ژنتیک، خوشه‌بندی را بهبود بخشیده است و دقت پاسخ ارائه شده در مسئله بالاتر از سایر روش‌ها است. همچنین از نتایج مشخص است که روش پیشنهادی فاصله درون خوشه‌های را به صورت بهینه‌تری نسبت به دیگر روش‌ها انجام داده است. بررسی نتایج حاصل از تکرارهای بلندمدت و کوتاهمدت نشان می‌دهد که الگوریتم پیشنهادی به تعداد تکرارهای بهینه‌سازی بلندمدت وابستگی چندانی ندارد. در مجموع ترکیب حاصل از الگوریتم کرم شب‌تاب و الگوریتم ژنتیک، نسبت به الگوریتم‌های کرم شب‌تاب، ژنتیک و k-means نتایج بهتر و قابل قبول‌تری در خوشه‌بندی ارائه داده است.



- [18] Hassanzadeh, Tahereh, Meybodi, Mohammad Reza, A New Hybrid Approach For Data Clustering Using Firefly Algorithm And K-Means, The 16th CSI International Symposium On Artificial Intelligence And Signal Processing (AISP), 2012, 007-011.
- [19] Binlu, Fangyuan Ju, An Optimized Genetic K-Means Clustering Algorithm, Computer Science and Information Processing (CSIP), 2012, 1296-1299.
- [1] Yaghini, M., and Ghazanfari, N. "Tabu-KM: a hybrid clustering algorithm based on tabu search approach." *International Journal of Industrial Engineering & Production Research* 21, no. 2 (2010)
- [2] Jiawei, H., Kamber, M., Han, J., Kamber, M. and Pei, J. "Data Mining: Concepts and Techniques Elsevier." (2006).
- [3] Taber, R. "Clustering (Xu, R. and Wunsch II, DC; 2009) [Book review]." *IEEE Computational Intelligence Magazine* 4, no. 3 (2009): 92-95.
- [4] Berzal, F. and Matín, N. "Data mining: concepts and techniques by Jiawei Han and Micheline Kamber." *ACM Sigmod Record* 31, no. 2 (2002): 66-68.
- [5] Jain, A. K., and Dubes, R. C. Jain, Anil K., and Richard C. Dubes. "Algorithms for clustering data." Prentice-Hall, Inc., 1988.
- [6] Jain, A.K., "Data clustering: 50 years beyond K-means." *Pattern recognition letters* 31, no. 8 (2010): 651-666.
- [7] Hassanzadeh, T. and Meybodi, M.R. "A new hybrid approach for data clustering using firefly algorithm and K-means." In *The 16th CSI international symposium on artificial intelligence and signal processing (AISP 2012)*, pp. 007-011. IEEE, 2012.
- [8] Yang, X. S. "Firefly algorithms for multimodal optimization." In *International symposium on stochastic algorithms*, pp. 169-178. Springer, Berlin, Heidelberg, 2009.
- [9] Moscato, P. "On evolution, search, optimization, genetic algorithms and martial arts: Towards memetic algorithms." *Caltech concurrent computation program, C3P Report 826* (1989): 1989.
- [10] Dianati, M., Song, I. and Treiber, M. *An introduction to genetic algorithms and evolution strategies*. Technical report, University of Waterloo, Ontario, N2L 3G1, Canada, 2002.
- [11] Wahid, F., Ghazali, R. & Ismail, L.H. Improved Firefly Algorithm Based on Genetic Algorithm Operators for Energy Efficiency in Smart Buildings. *Arab J Sci Eng* 44, 4027-4047 (2019).
- [12] Mahshwar, Keshva Kaushik, Vikram Arora. A Hybrid Data Clustering Using Firefly Algorithm Based Improved Genetic Algorithms. *Sciedirect. Procedia Computer Science* 58 (2015) 249-256.
- [13] M A El-Shorbagy, Adel M El-Refaey, A hybrid genetic–firefly algorithm for engineering design problems, *Journal of Computational Design and Engineering*, Volume 9, Issue 2, April 2022, Pages 706-730,
- [14] Mustafa Servet Kiran, Ahmet Babalik. Improved Artificial Bee Colony Algorithm for Continuous Optimization Problems. *Journal of Computer and Communications*, 2014, 2, 108-116.
- [15] Abdullah, A., Deris, S., Mohamad, M.S., Hashim, S.Z.M. (2012). A New Hybrid Firefly Algorithm for Complex and Nonlinear Problem. In: Omatu, S., De Paz Santana, J., González, S., Molina, J., Bernardos, A., Rodríguez, J. (eds) *Distributed Computing and Artificial Intelligence*. *Advances in Intelligent and Soft Computing*, vol 151. Springer, Berlin, Heidelberg.
- [16] Hassanzadeh, t, meybodi.m, "A new hybrid Approach for Data clustering using firefly Algorithm and k-means"
- [17] J.senthilnath, S.N.omkar, "clustering using firefly algorithm:performance study", *swarm and Evolutionary Computation*, volume1,ISSue3, pp 164-171,September 2011.