

Improving the efficiency of data clustering with chaotic evolutionary algorithms

Sajad Manteghi *

Head of Technology and Information
Department, General Directorate of Education,
Kohgiluyeh and Boyer-Ahmad Province, Iran

Sara Khosravani Pour

Instructor, Computer Group, General Directorate
of Education and Training, Kohgiluyeh and
Boyer-Ahmad Province, Iran.

Abstract

Nowadays, clustering plays an important role in most research fields such as engineering, medicine, biology, data mining, etc. In fact, clustering means unsupervised division. By using it, the data are divided into categories that are more similar to each other in terms of the parameters of interest. One of the famous methods in this field is k-means. In this method, despite the dependence on initial conditions and convergence to local optimal points, N numbers of data are grouped into k clusters with high speed. In this article, to solve the existing problems, the combined method is used based on evolutionary algorithms, chaos theory and k-means; that is in addition to solving the mentioned problems, it will also be independent of the number of variables. In this article, for the purpose of validation, the proposed methods are implemented on 13 different famous collections, and the results are compared with genetic algorithm, particle community, bee colony, simulated refrigeration, differential evolution, harmony search, and k-means methods. The high ability and robustness of these methods will be evident based on the results.

Keywords: clustering, k-means algorithm, evolutionary algorithms, chaos, chaotic evolutionary algorithm

Received: 27/January/2023

Accepted: 15/August/2023

ISSN: 2980-8936

* Corresponding Author: S.manteghi@eng.ui.ac.ir

بهبود کارآیی خوشه‌بندی داده‌ها با الگوریتم‌های تکاملی آشوب‌گونه

ریاست فناوری و اطلاعات، اداره کل آموزش و پرورش، کهگیلویه و بویراحمد،
ایران.

سجاد منطقی*

مریی گروه کامپیوتر، اداره کل آموزش و پرورش، کهگیلویه و بویراحمد، ایران.

سارا خسروانی پور

چکیده

امروزه، خوشه‌بندی نقش مهمی را در اغلب زمینه‌های تحقیقاتی مانند مهندسی، پزشکی، زیست‌شناسی، داده‌کاوی و ... ایفا می‌نماید. در واقع خوشه‌بندی به معنای تقسیم‌بندی بدون نظارت می‌باشد. داده‌ها با استفاده از آن به دسته‌هایی که از نظر پارامترهای موردعلاقه، شباهت بیشتری به یکدیگر دارند، تقسیم می‌گردند. یکی از روش‌های معروف در این زمینه k -means می‌باشد. در این روش علی‌رغم وابستگی به شرایط اولیه و همگرایی به نقاط بهینه محلی، تعداد N داده به k خوشه با سرعت بالا، دسته‌بندی می‌شوند. در این مقاله جهت رفع مشکلات موجود از روش ترکیبی مبتنی بر الگوریتم‌های تکاملی و تئوری آشوب و k -means بهره گرفته خواهد شد؛ که علاوه بر رفع مشکلات ذکرشده، مستقل از تعداد متغیرها نیز خواهد بود. در این مقاله به منظور اعتبارسنجی، روش‌های پیشنهادی بر روی ۱۳ مجموعه متفاوت مشهور پیاده‌سازی می‌گردد و نتایج با روش‌های الگوریتم ژنتیک، اجتماع ذرات، کلونی زنبور عسل، تبرید شبیه‌سازی‌شده، تکاملی تفاضلی، جستجوی هارمونی و k -means مقایسه خواهند گردید. توانایی بالا و مقاوم بودن این روش‌ها بر اساس نتایج مشهود خواهد بود.

کلیدواژه‌ها: خوشه‌بندی، الگوریتم k -means، الگوریتم‌های تکاملی، آشوب، الگوریتم تکاملی آشوب‌گونه

مقدمه

پردازش داده، یکی از شاخص‌های بسیار مهم در دنیای اطلاعات است. خوشه‌بندی یکی از بهترین روش‌هایی است که برای کار با داده‌ها ارائه شده است. خوشه‌بندی قابلیت ورود به فضای داده و تشخیص ساختار داده‌ها را امکان‌پذیر می‌نماید. لذا به عنوان یکی از ایده‌آل‌ترین روش‌ها، برای کار با دنیای عظیم داده‌ها محسوب می‌شود. خوشه‌بندی، یافتن ساختاری در مجموعه‌ای از داده‌ها است که طبقه‌بندی نشده‌اند. به بیان دیگر می‌توان گفت که خوشه‌بندی قراردادن داده‌ها در گروه‌هایی است که اعضای هر گروه از جهات خاصی شبیه یکدیگرند. معیار شباهت در این جا، فاصله می‌باشد؛ یعنی نمونه‌هایی که به یکدیگر نزدیکند در یک خوشه قرار می‌گیرند. با محاسبه فاصله بین دو داده می‌توان فهمید که چقدر این دو داده به هم نزدیک هستند و بر این اساس در یک خوشه قرار داده می‌شوند. توابع ریاضی مختلفی برای محاسبه فاصله وجود دارند؛ فاصله اقلیدسی، فاصله همینگ و

خوشه‌بندی یک فرآیند بدون ناظر است که یک مجموعه از اشیاء را به گروه‌های متجانس تقسیم می‌کند. خوشه بندی یکی از مهم‌ترین روش‌های داده کاوی می‌باشد ولی در زمینه‌های دیگری مثل شناخت الگو، فن آوری اطلاعات، پردازش تصاویر، زیست‌شناسی، روانشناسی و بازاریابی نیز به کار می‌رود. هدف اصلی خوشه‌بندی پیدا کردن یک ساختار معنادار برای داده‌ها می‌باشد. الگوریتم k-means یکی از معروف‌ترین الگوریتم‌های خوشه‌بندی است که در سال ۱۹۷۹ انتشار یافت و هنوز هم به طور گسترده‌ای از آن استفاده می‌شود. متأسفانه این الگوریتم وابسته به مقادیر اولیه مراکز خوشه‌هاست و به همین دلیل همیشه خوشه‌بندی به طور کاملاً صحیح انجام نمی‌شود. با توجه به قرارگیری خوشه‌بندی در دسته مسائل غیرقطعی - سخت، لزوم استفاده از الگوریتم‌های فرامکاشفه‌ای و تکاملی جهت دستیابی به پاسخ مناسب در زمان منطقی مشخص می‌باشد.

بهینه‌سازی شاخه‌ای از علوم ریاضیات است که در آن سعی می‌شود نقطه بهینه توابع را با توجه به تعدادی محدودیت به دست آورد. امروزه بسیاری از مسائل بهینه‌سازی اغلب از نوع مسائل غیرچندجمله‌ای سخت هستند. از جمله راه - حل‌های موجود در برخورد با این گونه مسائل، استفاده از الگوریتم‌های تقریبی یا مکاشفه‌ای است. این گونه الگوریتم‌ها تضمینی نمی‌دهند که جواب به دست آمده بهینه باشد. امروزه الگوریتم‌های تکاملی نظیر الگوریتم ژنتیک، بهینه‌سازی انبوه ذرات، کلونی مورچگان و غیره در حل بسیاری از مسائل پیچیده که روش‌های عددی قادر به حل آن‌ها نیستند و یا در حل آن دچار مشکلاتی می‌شوند، کاربرد دارند. مسئله خوشه‌بندی نیز یکی از این چنین مسائل است و با توجه به کاربرد بالای این موضوع، محققین همواره در پی یافتن روشی کامل و دقیق برای حل این مسئله بوده‌اند؛ اما با توجه به ماهیت این موضوع، روش‌های عددی موجود قابلیت حل این مسئله را ندارند و معضل پاسخ‌های محلی در این روش‌ها مشکل ساز بوده است.

الگوریتم‌های تکاملی الهام گرفته از پدیده‌های طبیعی هستند. این الگوریتم‌ها از مولدهای تصادفی برای تولید جمعیت اولیه استفاده می‌کنند؛ ولی در طبیعت پدیده‌ها تصادفی نیستند و از نظم خاصی پیروی می‌کنند. در این مقاله در الگوریتم‌های تکاملی از پدیده آشوب به منظور تولید اعداد شبه تصادفی برای خوشه‌بندی استفاده شده است. آشوب نوعی بی‌نظمی منظم است. بی‌نظمی از آن جهت که نتایج آن غیرقابل پیش‌بینی است و منظم بدان جهت که از نوعی قطعیت برخوردار است؛ بنابراین می‌توان راه کاری جامع تر و کامل تر ارائه نمود. لذا هدف از این مقاله ارائه الگوریتم - های خوشه‌بندی بر مبنای الگوریتم‌های تکاملی آشوب گونه می‌باشد که بتواند نسبت به الگوریتم‌های مورد مقایسه بهبودی را حاصل نماید. در این مقاله چندین الگوریتم تکاملی آشوب گونه همانند الگوریتم ژنتیک آشوب گونه، الگوریتم بهینه‌سازی انبوه ذرات آشوب گونه و الگوریتم کلونی زنبور آشوب گونه برای مسئله خوشه‌بندی پیشنهاد شده است.

الگوریتم‌های تکاملی در مقایسه با سایر الگوریتم‌های بهینه‌سازی برتری‌هایی دارند که موجب شده است به طور گسترده مورد استفاده قرار بگیرند. هم‌چنین محدودیتی در مورد تابع شایستگی ندارند و لزومی ندارد که این تابع مثلاً مشتق‌پذیر باشد. علاوه بر این موارد، چون الگوریتم‌های تکاملی دارای جمعیتی از موجودات هستند و روی بخش‌های مختلفی از جمعیت به طور موازی کار می‌کنند، احتمال کمتری برای قرار گرفتن در بهینه‌های محلی دارند. این قابلیت الگوریتم‌های تکاملی اجازه می‌دهد که کار بهینه‌سازی را به طور موازی روی چندین بخش جمعیت انجام داد. از بین الگوریتم‌های تکاملی می‌توان به موارد زیر اشاره کرد:

- الگوریتم ژنتیک^۱ (Holland, 1975).
- الگوریتم تبرید شبیه‌سازی‌شده^۲ (Alfonzetti et al., 2006).
- الگوریتم بهینه‌سازی ازدحام ذرات^۳ (Kennedy & Eberhart, 1995).
- الگوریتم رقابت استعماری^۴ (Atashpaz-Gargari & Lucas, 2007).
- الگوریتم کلونی مورچگان^۵ (Dorigo & Gambardella, 1997).
- الگوریتم کلونی زنبور عسل^۶ (Karaboga & Basturk, 2007).

طبق تحقیقاتی که انجام شده الگوریتم‌های تکاملی برای خوشه‌بندی مناسب هستند (Kao et al., 2008; Niknam et al., 2011)؛ ولی تا به حال الگوریتم‌های تکاملی آشوب‌گونه برای خوشه‌بندی استفاده نشده‌اند که بدانیم تأثیرات آن‌ها در مسئله خوشه‌بندی چیست. آشوب در لغت به معنی درهم‌ریختگی، آشفتگی و بی‌نظمی است و مترادف آن در مکانیک، تلاطم است. این واژه به معنی فقدان هر گونه ساختار بانظم است و معمولاً در محاورات، آشوب و آشفتگی نشانه بی‌نظمی و سازمان‌نیافتگی به نظر آمده و جنبه منفی در بر دارد؛ اما در واقع با پیدایش نگرش جدید و روشن شدن ابعاد علمی و نظری آن، امروزه دیگر بی‌نظمی و آشوب به مفهوم سازمان‌نیافتگی، ناکارآیی و درهم‌ریختگی تلقی نمی‌گردد؛ بلکه وجود جنبه‌های غیرقابل پیش‌بینی و اتفاقی در پدیده‌های پویا است که ویژگی‌های خاص خود را دارند (Schuster & Just, 2006). تئوری آشوب مطالعه رفتار سیستم‌های پویا است که به شرایط اولیه حساسیت بالایی دارند. تغییر اندکی در شرایط اولیه چنین سیستم‌هایی باعث اثرات بسیار در آینده خواهد شد (Kellert, 1993).

هدف در این پژوهش این است که با بررسی الگوریتم‌های موجود در زمینه خوشه‌بندی، راهکاری برای ایجاد تغییرات در الگوریتم‌های موجود برای بالابردن سرعت هم‌گرایی و بهبود کیفیت خوشه‌بندی ارائه گردد به گونه‌ای که محدودیت‌های موجود را پوشش دهد. از نمونه محدودیت‌های موجود می‌توان موارد زیر را برشمرد:

- کارآیی برای پایگاه داده‌های با حجم بالا.
- کشف خوشه‌ها با اشکال مختلف.
- عدم حساسیت به ترتیب داده‌های ورودی.
- قابلیت تفسیر و استفاده.

کارهای گذشته

در سال ۱۹۹۹ یک الگوریتم ژنتیک ترکیبی جدید به منظور تقسیم‌بندی بهینه سراسری یک داده دریافتی به تعداد مشخصی خوشه توسط کریشنا^۱ و همکارانش معرفی شده است (Krishna & Murty, 1999). در این روش یک الگوریتم ژنتیک ترکیب شده با گرادیان نزولی کلاسیک در خوشه‌بندی یعنی k-means استفاده شده است. این الگوریتم k-means با ژنتیک (GKA) نامیده می‌شود. در این روش، از الگوریتم k-means در GKA به عنوان یک عملگر جستجو به جای عملگر آمیزش استفاده می‌شود. همچنین یک عملیات جهش مغرضانه خاص برای خوشه‌بندی تعیین می‌گردد که فاصله بر مبنای جهش نامیده می‌شود. از تئوری مارکوف محدود شده نیز در روش مذکور استفاده شده و نشان داده شده است که GKA به بهینه سراسری همگرا می‌شود. جهش به کار برده شده در روش GKA به این صورت است که تغییرات جهش یک مقدار به فاصله نقطه داده از مراکز خوشه بستگی دارد. در بسیاری از فراخوانی‌ها هر آلل به نقطه داده بستگی دارد و این مقدار نشان‌دهنده خوشه‌ای است که به آن تعلق دارد. یک عملیات مثل احتمال تغییرات یک آلل برای تعداد اعضای یک خوشه در صورتی که مراکز خوشه به داده نزدیک‌تر باشد؛ تعریف شده است. برای اعمال عملگر جهش برای آلل $s_w(i)$ با توجه به الگوی x_i از $d_j = d(x_i, c_j)$ که از فاصله اقلیدسی بین x_i و c_j است استفاده می‌شود.

موالیک^۲ و همکارانش در سال ۲۰۰۰ یک الگوریتم به نام خوشه‌بندی ژنتیک برای مسئله خوشه‌بندی ارائه کرده‌اند. مزیت این روش استفاده از الگوریتم k-means در ژنتیک است. روش مذکور بر روی ۴ مجموعه داده آزمایش شده است (Maulik & Bandyopadhyay, 2000). در این روش هر رشته دنباله‌ای از اعداد حقیقی است که نشان‌دهنده k مرکز خوشه است. برای یک فضای n بعدی طول یک کروموزوم $N \times K$ است که N موقعیت اول نشان‌دهنده N بعد مرکز خوشه اول است و N موقعیت بعدی نیز مرکز خوشه بعدی و بقیه نیز به همین صورت هستند. جمعیت اولیه به صورت تصادفی تولید می‌شود. برازندگی بر حسب معیار فاصله محاسبه می‌شود به گونه‌ای که ابتدا خوشه‌ها با توجه به مراکز کدگذاری شده در کروموزوم مورد نظر تشکیل می‌شوند. پس از تشکیل خوشه‌ها مراکز خوشه‌ها در کروموزوم با میانگین نقاط درون خوشه جایگزین می‌شود.

الگوریتمی مبتنی بر کلونی مورچگان برای مسئله خوشه‌بندی توسط شلوکار^۳ و همکارانش در سال ۲۰۰۴ ارائه شد. در این الگوریتم عوامل توزیعی به کار گرفته می‌شود که در آن از روش واقعی مورچه‌ها برای پیدا کردن کوتاه‌ترین مسیر از لانه خود به سمت منبع غذا و بازگشت استفاده شده است (Shelokar et al., 2004). هدف این الگوریتم اختصاص بهینه از N شیء عضو \mathcal{R}^n در K خوشه است به گونه‌ای که مجموع مربعات فاصله اقلیدسی بین شیء و مراکز خوشه‌ای که به آن تعلق دارد کمینه باشد. الگوریتم برای ایجاد راه‌حل‌ها R عامل در نظر می‌گیرد. یک عامل با یک رشته راه‌حل خالی S با طول N شروع می‌شود به گونه‌ای که هر عنصر از رشته با یکی از نمونه‌های آزمایشی مطابق است. مقدار اختصاص یافته به یک نمونه از رشته راه‌حل S بیانگر شماره خوشه‌ای است که نمونه آزمایشی به S اختصاص داده است. برای ساخت یک راه‌حل، عامل از اطلاعات دنباله‌ای فرومون برای اختصاص هر یک از المان‌های رشته S به یک خوشه مناسب استفاده می‌کند.

در ابتدای الگوریتم، ماتریس فرومون T، با مقدار کوچک τ_0 مقداردهی اولیه می‌گردد. مقدار دنباله $\tau_{i,j}$ در مکان (i,j) نشان‌دهنده فرومون جمع شده از نمونه iام از خوشه jام است. برای مسئله تقسیم N نمونه درون K خوشه، اندازه ماتریس فرومون $N \times K$ است. در هر مرحله از تکرار، هر یک از عامل‌ها یا مورچه‌های نرم‌افزار همانند دنباله‌ای از

1. Krishna
2. Mualik
3. Shelokar

را حل‌ها را که از فرآیند ارتباطات فرومون برای به دست آوردن یک بخش نزدیک به بهینه از N نمونه دریافتی به K گروه به گونه‌ای که هدف مورد نظر حاصل گردد استفاده می‌کنند. پس از تولید یک نسل از R راه‌حل‌های آزمایشی، یک جستجوی محلی برای بهبود برازندگی این راه‌حل‌ها اجرا می‌شود. ماتریس فرومون با توجه به کیفیت راه‌حل‌های تولیدشده توسط عامل‌ها به‌روز می‌شود. هدایت توسط ماتریس فرومون بهبودیافته صورت می‌گیرد. عامل‌ها راه‌حل‌های بهبودیافته را می‌سازند و مراحل بالا را تا رسیدن به تعداد از قبل معین‌شده تکرار می‌کنند.

یک روش مبتنی بر الگوریتم بهینه‌سازی زنبور عسل برای هدف خوشه‌بندی در سال ۲۰۰۸ توسط فتحیان و همکارانش پیشنهاد شده است (Fathian et al., 2007). این الگوریتم در چند گام در ادامه آورده شده است:

گام ۱: نمایش رشته: برای راه‌حل انتخابی مسئله از یک کروموزوم استفاده می‌شود که هر ژن درون کروموزوم بیانگر یک پارامتر از راه‌حل انتخابی است. در این روش، یک کروموزوم به عنوان مجموعه‌ای از k مرکز خوشه ابتدایی و هر ژن به عنوان ابعاد مرکز خوشه در نظر گرفته شده است.

گام ۲: تعیین پارامترهای ورودی مدل: الگوریتم با سه پارامتر تعیین‌شده توسط کاربر و یک الگوریتم از پیش تعیین‌شده شروع می‌شود.

گام ۳: تولید تصادفی مجموعه‌ای از راه‌حل‌های اولیه: یک مجموعه از مراکز خوشه‌های اولیه به صورت تصادفی از نقاط مجموعه داده تولید می‌شوند. هر راه‌حل نشان‌دهنده‌ی مراکز خوشه k است.

گام ۴: انتخاب ملکه: بعد از تولید راه‌حل‌های اولیه مذکور، راه‌حل‌های اولیه بر اساس تابع کارآیی تخمین زده می‌شود و بهترین راه‌حل به عنوان ملکه قرار می‌گیرد.

گام ۵: جفت‌گیری: از تبرید شبیه‌سازی‌شده برای انتخاب مجموعه‌ای از راه‌حل‌های موجود در فضای جستجو برای انجام جفت‌گیری بین بهترین راه‌حل و راه‌حل‌های انتخابی استفاده می‌شود.

گام ۶: فرآیند تولید مثل.

گام ۷: نوزادان انتخابی و ملکه توسط کارگران با ژل‌های سلطنتی تغذیه می‌شوند.

گام ۸: اگر بهترین راه‌حل جدید از ملکه بهتر باشد، ملکه با راه‌حل تولیدی جایگزین می‌شود.

گام ۹: معیارهای پایانی بررسی می‌گردد. اگر معیارهای پایانی ارضاء شدند الگوریتم پایان می‌یابد در غیر این صورت راه‌حل‌های قبلی دور ریخته می‌شوند و راه‌حل‌های جدید تولید شده و به مرحله ۵ برمی‌گردند.

یک الگوریتم تکاملی ترکیبی برای ترکیب الگوریتم تبرید شبیه‌سازی‌شده با بهینه‌سازی انبوه ذرات اصلی توسط نیک‌نام و همکارانش در سال ۲۰۰۸ پیشنهاد شده است (Niknam & Amir, 2010). در الگوریتم پیشنهادی از تبرید شبیه‌سازی‌شده به عنوان یک جستجوگر محلی اطراف Gbest استفاده می‌شود؛ به عبارت دیگر Gbest در هر تکرار تغییر می‌کند. این الگوریتم PSO-SA نامیده می‌شود. متغیرهای کنترلی مراکز خوشه‌ها هستند.

یک تکنیک ترکیبی مبتنی بر ترکیب الگوریتم‌های k -means، جستجوی ساده Nelder-Mead و بهینه‌سازی ازدحام ذرات توسط کائو و همکارانش در سال ۲۰۰۸ پیشنهاد شده است. این الگوریتم پیشنهادی K-NM-PSO نامیده می‌شود. جستجوی K-NM-PSO برای مراکز خوشه از یک مجموعه داده دلخواه برای الگوریتم k -means است؛ اما با استفاده از این الگوریتم می‌توان به صورت مؤثر بهینه‌ی سراسری را یافت (Kao et al., 2008).

یک الگوریتم جدید بر اساس ترکیب دو الگوریتم خوشه‌بندی K -means و رقابت استعماری بهبودیافته توسط نیک‌نام و همکارانش در سال ۲۰۱۲ پیشنهاد می‌گردد که این ترکیب K-MICA نامیده می‌شود و از روش حداکثرسازی انتظار اصلاح‌شده (EM) برای تعیین تعداد خوشه‌ها استفاده شده است (Niknam et al., 2011).

روش پیشنهادی

الگوریتم‌های بهینه‌سازی به روش هوشمند، پیشرفت قابل‌ملاحظه‌ای در مقابل روش‌های ریاضی از خود نشان داده‌اند. در این تحقیق از روش جدیدی مبتنی بر ترکیب الگوریتم‌های تکاملی (هم‌چون ژنتیک، بهینه‌سازی انبوه ذرات و کلونی زنبور)، پدیده آشوب و الگوریتم k-means، برای خوشه‌بندی اطلاعات استفاده خواهد شد. در واقع هدف این است که داده‌های موجود را بتوان به‌نحوی با حداقل خطا خوشه‌بندی نمود. زمانی خوشه‌بندی ایده‌آل میسر می‌گردد و هدف نهایی برآورده شود و حداقل شباهت بین داده‌های موجود در دسته‌های مجزا، هم‌چنین حداکثر شباهت بین داده‌های موجود در یک دسته برقرار باشد.

در این راستا از مزایای الگوریتم k-means، آشوب و الگوریتم‌های تکاملی بهره گرفته خواهد شد. روش k-means علی‌رغم معایبش که در ادامه توضیح داده می‌شود، به دلیل سادگی و سهولت شبیه‌سازی آن بسیار کاربرد دارد.

- وابستگی به شرایط اولیه: مقادیر اولیه در این روش پاسخ نهایی را دست‌خوش تغییرات قرار می‌دهد و خوشه‌های به‌دست آمده بنا به مقادیر اولیه بسیار متفاوت خواهند بود.
- همگرایی به نقاط بهینه محلی: در صورتی که تابع هدف دارای نقاط بهینه زیادی باشد استفاده از این روش جواب اصلی را نتیجه نخواهد داد.

جهت مرتفع‌نمودن این معایب، در تحقیقات اخیر ایده‌های متفاوتی بیان گردیده است که تا حد قابل‌قبولی کارساز بوده است. الگوریتم‌های تکاملی ارائه گردیده در این زمینه را می‌توان به شرح زیر برشمرد؛ ژنتیک، اجتماع ذرات و زنبور مصنوعی.

با ترکیب نمودن الگوریتم k-means و الگوریتم‌های تکاملی آشوب‌گونه، می‌توان به هدف مورد نظر که خوشه‌بندی اطلاعات می‌باشد رسید. در این مسیر دقت جواب نهایی به طور چشمگیری افزایش یافته است. همچنین نگرانی در خصوص همگرایی پاسخ به بهینه محلی تا حدود بسیار زیادی برطرف گردیده است و بدون در نظر گرفتن تعداد متغیرها و مقادیر اولیه جواب مطلوب به دست خواهد آمد. در ادامه الگوریتم‌های تکاملی آشوب‌گونه شرح داده خواهد شد.

الگوریتم تکاملی مبتنی بر آشوب

در این بخش الگوریتمی ترکیبی از الگوریتم‌های تکاملی و پدیده آشوب ارائه شده است. در روش‌های پیشنهادی از الگوریتم‌های تکاملی ژنتیک، بهینه‌سازی گروه ذرات، تکاملی تفاضلی، کلونی زنبور و رقابت کشورهای استعماری استفاده شده است. آشوب شاخه‌ای از حوزه‌های تحقیقات غیرخطی پویا است که به طور گسترده مورد مطالعه قرار گرفته است. موارد کاربردی زیادی در سیستم‌های واقعی، چه ساخت انسان و چه طبیعی، با استفاده از روش‌های نوین در حوزه غیرخطی است که به‌ظاهر اتفاقی به نظر می‌رسد. هرچند این ظاهر تصادفی، منشأ اتفاقی ندارد و به کلی نتیجه یک پردازش قطعی و تعریف شده است. بیش از صد سال است که زیست‌شناسان نمونه‌هایی از جمعیت گونه‌های مختلف را مورد بررسی قرار داده‌اند. به‌تازگی دانشمندان قادر به ایفا کردن مدل‌های آشوب‌گونه در جمعیت‌های خاصی بوده‌اند (Liz & Ruiz-Herrera, 2012). به عنوان مثال، مطالعه روی سیاه‌گوش کانادایی نشان داد رفتار آشوب‌گونه در رشد جمعیت وجود دارد (Lai, 1996). به همین منظور در این پژوهش تغییری در الگوریتم‌های تکاملی مبتنی بر جمعیت انجام گرفته و پدیده آشوب با الگوریتم‌های تکاملی ترکیب شده تا از مولدهای آشوبی در تولید جمعیت الگوریتم‌ها استفاده شود. در این مقاله از الگوریتم‌های تکاملی آشوب‌گونه برای خوشه‌بندی، با هدف

دست‌یابی به خوشه‌بندی با پایداری بالا استفاده شده است. فرمول تولید جمعیت با پدیده آشوب در ادامه ذکر شده است (Neal, 1993).

$$f(x) = 1 - 2 \times \left| u - \frac{1}{2} \right| \quad (1)$$

در فرمول بالا u یک عدد تصادفی در محدوده [۰ و ۱] است.

الگوریتم‌های پیشنهادی

در روش پیشنهادی از الگوریتم‌های تکاملی آشوب گونه به منظور خوشه‌بندی داده‌ها استفاده شده است. در روش‌های پیشنهادی پس از تولید جمعیت به روش آشوب گونه، الگوریتم k -means روی آخرین عضو جمعیت، جهت بهبود عملکرد اجرا می‌گردد. در شکل‌های ۱ تا ۴ شبه‌کد الگوریتم‌های پیشنهادی نشان داده شده است.

— ذره را با استفاده از فرمول زیر تولید کن:

$$Paticle(i) = 1 - 2 \times \left| u - \frac{1}{2} \right|$$

— $Paticle(i)$ ذره نام جمعیت است و u عددی بین [۰ و ۱] است.

— اجرای الگوریتم k -means بر روی آخرین فرد جمعیت.

— برای تمام ذرات، سرعت و موقعیتی تصادفی ایجاد کن.

— تا زمانی که شرایط خاتمه محقق نشده است:

○ یک واحد به t اضافه کن.

○ مقدار تابع هدف را به ازای هر ذره محاسبه کن.

○ به ازای i از یک تا n (تعداد جمعیت):

▪ بهینه محلی را بیاب و در $x^{i,best}[t]$ ذخیره کن.

○ مقدار بعدی i .

○ بهینه سراسری را محاسبه کن و در $x^{gbest}[t]$ قرار بده.

○ به ازای i از یک تا n :

▪ به ازای j از یک تا d :

$$v_j^i[t+1] = wv_j^i[t] + c_1r_1(x_j^{i,best}[t] - x_j^i[t]) + c_2r_2(x_j^{gbest}[t] - x_j^i[t])$$

$$x_j^i[t+1] = x_j^i[t] + v_j^i[t+1]$$

▪ مقدار بعدی j .

مقدار بعدی i .

شکل (۱): الگوریتم خوشه‌بندی بهینه‌سازی ازدحام ذرات آشوب گونه.

- تولید جمعیت اولیه^۱ شامل n کروموزوم به صورت آشوب گونه.
- اجرای الگوریتم k -means روی آخرین فرد جمعیت.
- بررسی تابع ارزیابی $f(x)$ برای هر کروموزوم x در جمعیت.
- ایجاد یک جمعیت جدید بر اساس تکرار گام‌های زیر:
 - انتخاب دو کروموزوم والد از یک جمعیت بر اساس میزان برازندگی آن‌ها.
 - در نظر گرفتن مقدار مشخصی برای احتمال اعمال عملگر ترکیب^۲ (تقاطع) و سپس انجام عملیات ترکیب بر روی والدین به منظور ایجاد فرزندان (اگر هیچ ترکیب جدیدی صورت نگیرد، فرزندان همان والدین خواهند بود).
 - در نظر گرفتن احتمال جهش و سپس تغییر فرزندان در هر مکان.
 - جایگزینی فرزندان جدید در جمعیت.
- استفاده از جمعیتی جدید برای اجراهای بعدی الگوریتم.
- توقف اجرای الگوریتم در صورت مشاهده شرایط توقف و برگرداندن بهترین جواب در جمعیت فعلی.
- رفتن به مرحله ۲.

شکل (۲): الگوریتم خوشه‌بندی ژنتیک آشوب گونه.

- تولید جمعیت اولیه به صورت آشوب گونه.
- اعمال الگوریتم k -means بر روی آخرین فرد جمعیت.
- **تا زمانی که** یک امپراتوری باقی مانده است مراحل زیر را تکرار کن:
 - **تا زمانی که** هیچ مستعمره‌ای در امپراتوری وجود ندارد:
 - انتقال مستعمرات به سمت کشورهای امپریالیستی آن‌ها.
 - اگر هزینه مستعمره‌ای از امپریالیست خود بیشتر بود:
 - موقعیت مستعمره با امپریالیست را تعویض کن.
 - تابع هزینه را برای تمام امپراتوری‌ها را محاسبه کن.
 - ضعیف‌ترین مستعمره از ضعیف‌ترین امپراتوری را به امپراتوری قوی بده.
 - ضعیف‌ترین امپراتوری را حذف کن.

شکل (۳): الگوریتم خوشه‌بندی رقابت استعماری آشوب گونه.

- تولید جمعیت اولیه به صورت آشوب‌گونه.
- اجرای الگوریتم k-means بر روی آخرین فرد جمعیت.
- محاسبگی شایستگی هر فرد جمعیت و انتخاب بهترین فرد جمعیت.
- تا زمانی که شرایط توقف برقرار نشده مراحل زیر را تکرار کن.
- **جهش:** برای هر فرد جمعیت مثل x یک بردار مانند u از جهش بردار x تولید کن.
- **ترکیب:** برای هر فرد جمعیت مثل x آن را با بردار u ترکیب کن و یک فرزند با نام z تولید کن.
- **انتخاب:** اگر شایستگی فرزند یعنی z بهتر از شایستگی والد یعنی x باشد آن گاه z را جایگزین x کن.

شکل (۴): الگوریتم خوشه‌بندی تکامل تفاضلی آشوب‌گونه

- تولید جمعیت اولیه به صورت آشوب‌گونه.
- اجرای الگوریتم k-means روی آخرین فرد جمعیت.
- محاسبگی شایستگی هر یک از زنبورهای جمعیت اولیه.
- تا زمانی که شرایط توقف برقرار نشده گام‌های زیر تکرار می‌گردد:
- انتخاب مکان‌های مناسب جهت جستجو در همسایه‌ها
- استفاده از زنبورها جهت یافتن مکان مناسب و تخمین برآزندگی
- انتخاب بهترین زنبورهای در هر مسیر
- زنبورهای باقی‌مانده به جستجوی تصادفی می‌پردازند و برآزندگی هر یک از آن‌ها محاسبه می‌گردد.

شکل (۵): الگوریتم خوشه‌بندی کلونی زنبور آشوب‌گونه

روش‌های محاسبه کیفیت خوشه‌بندی

روش‌های مختلفی برای محاسبه کیفیت خوشه‌بندی وجود دارد، در این پژوهش از معیارهای فیشر (FM)، اطلاعات متقابل نرمال شده (NMI) و AR به منظور ارزیابی کیفیت خوشه‌بندی روش‌های مختلف استفاده شده است. در ادامه معیارهای استفاده شده در پژوهش مورد بررسی قرار می‌گیرند:

معیار F-Measure از ایده‌های دقت و فراخوانی به منظور بازیابی اطلاعات استفاده می‌کند. هر کلاس i به صورت مجموعه‌ای از n_i بخش مطلوب برای یک تحقیق و پرس‌وجو در نظر گرفته می‌شود. هر خوشه j (که توسط الگوریتم تولید شده) به صورت مجموعه‌ای از n_j بخش بازیابی شده توسط یک پرس‌وجو در نظر گرفته می‌شود (Coley, 1999). n_{ij} تعداد عناصر کلاس i در خوشه j را می‌دهد. برای هر کلاس i و خوشه j معیارهای دقت و فراخوانی به صورت زیر تعریف می‌شوند:

$$R(i, j) = \frac{n_{ij}}{n_j} \quad (2)$$

$$P(i, j) = \frac{n_{ij}}{n_i} \quad (3)$$

مقدار F-Measure متناظر نیز به صورت زیر محاسبه می‌گردد (Perruchet & Peereman, 2004):

$$FM = 2 \times \frac{P \cdot R}{P + R} \quad (۴)$$

واضح است که هر چقدر مقدار F-Measure بیشتر باشد، کیفیت خوشه‌بندی نیز بیشتر می‌شود. برای استفاده از معیار NMI، ابتدا تمام نمونه‌های دیگر متعلق به مجموعه داده D که در خوشه C_i قرار ندارند؛ به صورت یک خوشه مستقل D/C_i نمایش داده می‌شود. حال یک خوشه‌بندی شامل دو خوشه C_i و D/C_i ایجاد شده است که آن را P_1 می‌نامیم که $P_1 = \{C_i, D/C_i\}$ می‌باشد. اکنون خوشه‌بندی $P(D)$ که روی داده‌های نمونه‌برداری شده اعمال شده است، نیز باید به صورت دو خوشه‌ای ارائه شود تا در نهایت نتایج حاصل از این دو خوشه‌بندی طی فرآیندی با هم منطبق شوند. برای این منظور همه خوشه‌ها در $P(D)$ به دو خوشه C^* و D/C^* تقسیم می‌شوند. خوشه C^* از اجتماع همه خوشه‌هایی که بیش از ۵۰٪ از نمونه‌هایشان در خوشه C_i وجود دارند، تشکیل می‌شود و مابقی خوشه‌ها نیز در خوشه D/C^* قرار می‌گیرند. این خوشه‌بندی را P_2 می‌نامیم که $P_2 = \{C^*, D/C^*\}$ است. حال از اطلاعات متقابل^۱ نرمال شده (NMI) [۲۱] که معیار متداول برای ارزیابی شباهت بین دو افراز (نتیجه خوشه‌بندی) است، برای اندازه‌گیری شباهت بین دو خوشه‌بندی P_1 و P_2 استفاده می‌شود. از آنجایی که معیار اطلاعات متقابل نرمال نشده (MI)، وابسته به اندازه خوشه‌هاست، معمولاً از معیار NMI استفاده می‌شود. رابطه NMI بین دو خوشه‌بندی P_1 و P_2 به صورت زیر محاسبه می‌شود:

$$NMI(P_1, P_2) = \frac{MI(P_1, P_2)}{\frac{-1}{2m} (\sum_{i=0}^l P_i \log \frac{P_i}{m} + \sum_{j=0}^l P_j \log \frac{P_j}{m})} \quad (۵)$$

که اطلاعات متقابل، $MI(P_1, P_2)$ از رابطه زیر به دست می‌آید.

$$MI(P_1, P_2) = \sum_{i=0}^l \sum_{j=0}^l \frac{r_{ij}}{m^2} \cdot \log \frac{mp_{ij}}{r_{ij}}, \quad (۶)$$

که در این رابطه، p_{11} نشان‌دهنده تعداد نمونه‌های مشترک موجود در خوشه‌های C_i و C^* است. p_{10} نشان‌دهنده تعداد نمونه‌های مشترک موجود در D/C^* و C_i است. p_{01} نشان‌دهنده تعداد نمونه‌های مشترک موجود در D/C^* و D/C_i است. p_{00} نشان‌دهنده تعداد نمونه‌های مشترک موجود در D/C_i و D/C^* است. هم‌چنین m تعداد کل نمونه‌هاست. در واقع p_i ، $p_{.j}$ به ترتیب بیانگر کل نمونه‌های موجود در C_i و C^* هستند. معیار NMI نشان‌دهنده میزان شباهت خوشه‌بندی P_1 و P_2 می‌باشد.

نتایج تجربی

در این بخش نتایج به دست آمده از الگوریتم پیشنهادی بر روی داده‌های معروف، مورد بررسی قرار خواهد گرفت. آنگاه، پاسخ نهایی الگوریتم‌های پیشین با الگوریتم مورد بحث در این تحقیق، با هم مقایسه شده و با سه معیار NMI، FM، AR مورد ارزیابی قرار می‌گیرد.

به طور کلی مجموعه داده‌های مختلفی برای انجام فرآیند خوشه‌بندی وجود دارد؛ که هر کدام دارای ویژگی‌های خاص خود هستند. لذا برای دستیابی به اهداف موجود؛ مجموعه داده‌های Breast Cancer، Glass، Wine، Iris، Galaxy، Inosphere و Half-ring برگزیده می‌شوند. این مجموعه داده از سایت معروف یادگیری ماشین (UCI) انتخاب شده‌اند.

مجموعه داده Wine از آزمایشگاه MCI گرفته شده است که بر اساس فعل و انفعالات شیمیایی شراب به دست آمده است. این مجموعه داده شامل ۱۷۸ نمونه است که هر نمونه ۱۳ ویژگی پیوسته دارد. این مجموعه داده در ۳ کلاس متفاوت طبقه‌بندی می‌شود.

با توجه به جدول ۱ برای مجموعه داده Wine، جواب بهینه برای الگوریتم‌های Chaotic GA مقدار ۱۶۲۹۲.۳۲، Chaotic ICA مقدار ۱۶۲۹۲.۳۵ و Chaotic ABC مقدار ۱۶۲۹۲.۳۸ می‌شود؛ در حالی که الگوریتم‌های دیگر تحت هیچ شرایطی توانایی رسیدن به این مقدار را ندارند. همان‌طور که می‌دانیم انحراف معیار یکی از شاخص‌های پراکندگی است که نشان می‌دهد به طور متوسط داده‌ها چه مقدار از مقدار متوسط فاصله دارند. اگر مقدار انحراف معیار کوچک باشد، نشانه این است که داده‌ها به میانگین نزدیک هستند و پراکندگی اندکی دارند و انحراف معیار بزرگ بیانگر پراکندگی قابل توجه در داده‌ها است. انحراف معیار الگوریتم Chaotic ICA و Chaotic ABC از الگوریتم‌های دیگر کمتر است. همان‌طور که در جدول قابل مشاهده است الگوریتم‌های Chaotic PSO، Chaotic ICA، Chaotic DE، Chaotic GA و Chaotic ABC در مقایسه با الگوریتم‌های PSO، ICA، DE، GA و ABC زمان اجرای کمتری دارند.

جدول (۱): نتایج میانگین ۱۰۰ بار اجرای الگوریتم‌های تکاملی موجود بر روی مجموعه داده Wine.

زمان اجرا	انحراف معیار	تابع هزینه			الگوریتم‌های تکاملی
		حداکثر هزینه	متوسط هزینه	حداقل هزینه	
1.8298s	$\pm(15.98344)$	16381.54	16352.85	16314.17	Chaotic PSO
1.6710s	$\pm(1.10363)$	16296.62	16294.32	16292.35	Chaotic ICA
1.7786s	$\pm(11.34668)$	16346.55	16324.65	16303.97	Chaotic DE
1.8626s	$\pm(3.97853)$	16308.73	16298.65	16292.32	Chaotic GA
3.1025s	$\pm(2.28436)$	16301.51	16295.39	16292.38	Chaotic ABC
2.0696s	$\pm(23.68992)$	16443.59	16366.69	16322.79	PSO
1.7847s	$\pm(2.56136)$	16307.78	16298.28	16294.0	ICA
0.6742s	$\pm(12856.19664)$	73115.70	34923.27	18777.46	SA
1.8319s	$\pm(25.18457)$	16440.56	16359.06	16309.41	DE
3.8392s	$\pm(3.84541)$	16313.53	16300.18	16294.97	GA
6.2975s	$\pm(12.39669)$	16379.03	16339.87	16311.50	ABC
5.9578s	$\pm(94.67190)$	16828.79	16626.69	16400.99	HS
0.0130s	$\pm(685.714)$	18294.85	16886.12	16555.68	K-means

در جدول ۲ نتایج ارزیابی ۱۰۰ بار اجرای مستقل، هر یک از الگوریتم‌های مورد بررسی با معیارهای FM، NMI، AR قرار داده شده است. با توجه به مطالب ذکر شده بالاتر بودن مقدار معیارهای FM، NMI و AR بیانگر کیفیت بالاتر خوشه‌بندی است. با توجه به جدول، الگوریتم Chaotic ICA نتایج بهتری را نسبت به دیگر الگوریتم‌های آشوب‌گونه دارد و یکی از الگوریتم‌های با کیفیت بالا با توجه به معیارهای ذکر شده محسوب می‌شود.

جدول (۲): نتایج ارزیابی داده‌های Wine با معیارهای NMI، FM و AR.

میانگین	معیار AR	معیار FM	معیار NMI	الگوریتم‌های تکاملی
0.506694	0.377163	0.718891	0.424027	Chaotic PSO
0.506703	0.377488	0.718469	0.424152	Chaotic ICA
0.506643	0.377294	0.718569	0.424066	Chaotic DE
0.502052	0.37068	0.716046	0.419431	Chaotic GA
0.503871	0.376033	0.717771	0.417808	Chaotic ABC
0.507545	0.378295	0.719348	0.424991	PSO
0.505384	0.375627	0.717690	0.422833	ICA
0.468752	0.303719	0.725092	0.377444	SA
0.507414	0.377949	0.719281	0.425012	DE
0.503498	0.373345	0.716142	0.421006	GA
0.505108	0.3749483	0.717911	0.422463	ABC
0.506748	0.3764747	0.717424	0.426345	HS
0.498304	0.364192	0.705746	0.424974	K-means

مجموعه داده Breast Cancer مربوط به اطلاعات غده سرطانی بیماران مبتلا به سرطان سینه است که در دانشگاه Wisconsin گردآوری شده است. در این مجموعه داده ۶۸۳ نمونه در دو کلاس خوش خیم و بدخیم کلاس‌بندی شده است. هر یک از نمونه‌های این مجموعه داده شامل ۹ ویژگی با مقادیر پیوسته است. مقدار پاسخ بهینه به دست آمده از حل مسئله بر روی مجموعه داده Breast Cancer برای الگوریتم Chaotic ICA، ۱۱۱۹.۱۲۴ و برای الگوریتم Chaotic GA، ۱۱۱۹.۸۰۱ خواهد بود. الگوریتم‌های دیگر نسبت به آن دارای مقدار هزینه بیشتری هستند؛ به نحوی که انحراف معیار کمتری از الگوریتم‌های دیگر دارند. با مقایسه نتایج به دست آمده از جدول ۳، این مسئله واضح می‌گردد که عملکرد الگوریتم‌های پیشنهادی مناسب بوده است و سرعت اجرایی الگوریتم‌های آشوب گونه نیز بالاتر است.

جدول (۳): نتایج میانگین ۱۰۰ بار اجرای الگوریتم‌های تکاملی موجود بر روی مجموعه داده Breast Cancer

زمان اجرا	انحراف معیار	تابع هزینه			الگوریتم‌های تکاملی
		حداکثر هزینه	متوسط هزینه	حداقل هزینه	
1.9617s	$\pm (5.13322)$	1157.590	1143.954	1129.214	Chaotic PSO
1.7867s	$\pm (0.12401)$	1119.704	1119.124	1119.017	Chaotic ICA
1.7970s	$\pm (21.2836)$	1199.244	1152.211	1121.562	Chaotic DE
1.9979s	$\pm (0.19226)$	1120.435	1119.801	1119.413	Chaotic GA
3.4208s	$\pm (1.44462)$	1129.270	1127.810	1123.524	Chaotic ABC
2.1781s	$\pm (4.51902)$	1154.574	1144.140	1130.852	PSO
2.0071s	$\pm (14.96884)$	1174.003	1128.943	1119.177	ICA
1.2121s	$\pm (273.36210)$	2695.564	1948.092	1438.788	SA
1.9954s	$\pm (22.48397)$	1226.464	1182.066	1132.111	DE
4.1158s	$\pm (2.58264)$	1134.412	1125.219	1121.378	GA
7.7236s	$\pm (4.14139)$	1143.261	1133.078	1124.143	ABC
4.4225s	$\pm (80.60777)$	1813.978	1621.6620	1333.659	HS
0.0143s	$\pm (0.46024)$	1129.270	1128.709	1128.335	K-means

با توجه به نتایج کارآیی روش‌های موردبررسی با استفاده از معیارهای NMI، FM و AR که در جدول ۴ نشان داده شده است، میانگین مقدار معیارهای ذکر شده برای الگوریتم‌های Chaotic ICA و Chaotic GA، به ترتیب مقادیر ۰.۸۷۸۶۵۵ و ۰.۸۷۵۲۱۶ خواهد بود که بیانگر کارآیی بالاتر الگوریتم‌های آشوب گونه در مقایسه با الگوریتم‌های دیگر است. همچنین با توجه به نتایج رتبه‌بندی بر اساس ارتباطات معنادار الگوریتم‌ها با یکدیگر الگوریتم‌های Chaotic ICA با تمام الگوریتم‌های موجود ارتباط معنادار دارد و پس از آن الگوریتم Chaotic GA با تعداد زیادی از الگوریتم‌های موجود ارتباط معنادار دارد. نتایج ذکر شده نشان‌دهنده عدم تصادفی بودن روش‌های پیشنهادی است.

جدول (۴): نتایج ارزیابی داده‌های Breast Cancer با معیارهای NMI، FM و AR.

الگوریتم‌های تکاملی	معیار NMI	معیار FM	معیار AR	میانگین
Chaotic PSO	0.765217	0.960358	0.859962	0.861846
Chaotic ICA	0.790138	0.966176	0.879651	0.878655
Chaotic DE	0.750243	0.956735	0.847722	0.851567
Chaotic GA	0.785005	0.964976	0.875667	0.875216
Chaotic ABC	0.731376	0.954300	0.839575	0.841750
PSO	0.756517	0.958201	0.852685	0.855801
ICA	0.777578	0.963084	0.869173	0.869945
SA	0.412089	0.851566	0.429269	0.564308
DE	0.754986	0.957753	0.851088	0.854609
GA	0.776086	0.962916	0.868648	0.869217
ABC	0.766346	0.960636	0.860927	0.862636
HS	0.578951	0.899503	0.655523	0.711326
K-means	0.717068	0.948598	0.820352	0.828673

مجموعه داده Iris در واقع مجموعه‌ای از داده‌ها است که شامل سه نمونه گل زنبق است. این مجموعه شامل ۵۰ نمونه است که هر نمونه دارای ۴ ویژگی است و در ۳ کلاس متمایز طبقه‌بندی می‌گردد. جدول ۵ نتایج میانگین ۱۰۰ بار اجرای الگوریتم‌های تکاملی موجود بر روی مجموعه داده Iris را نشان می‌دهد.

جدول (۵): نتایج میانگین ۱۰۰ بار اجرای الگوریتم‌های تکاملی موجود بر روی مجموعه داده Iris

الگوریتم‌های تکاملی	تابع هزینه			انحراف معیار	زمان اجرا
	حداقل هزینه	متوسط هزینه	حداکثر هزینه		
Chaotic PSO	97.3298	98.7982	99.5574	$\pm (0.49263)$	1.7454s
Chaotic ICA	96.6555	96.6629	96.9222	$\pm (0.02699)$	1.5330s
Chaotic DE	96.8732	97.6374	99.0106	$\pm (0.58680)$	1.5875s
Chaotic GA	96.6664	96.6801	96.6970	$\pm (0.00628)$	1.7363s
Chaotic ABC	96.7286	97.3214	101.5813	$\pm (0.845477)$	2.9370s
PSO	97.8824	98.9718	100.3142	$\pm (0.52493)$	1.8018s
ICA	96.6554	96.7205	97.4632	$\pm (0.18222)$	1.7502s
SA	130.7615	181.1466	280.4573	$\pm (25.10238)$	1.0737s
DE	149.3569	106.4994	309.2958	$\pm (32.22541)$	1.7631s
GA	96.6661	97.3955	98.3262	$\pm (0.38596)$	3.9793s
ABC	97.8159	100.0714	103.8752	$\pm (1.46867)$	5.9294s
HS	107.2247	136.8219	148.4858	$\pm (7.54184)$	3.4983s
K-means	97.3259	101.5442	123.9695	$\pm (9.71673)$	0.7682s

جدول (۶): نتایج ارزیابی داده‌های Iris با معیارهای FM، NMI و AR.

میانگین	معیار AR	معیار FM	معیار NMI	الگوریتم‌های تکاملی
0.809868	0.753492	0.908629	0.767482	Chaotic PSO
0.801736	0.739832	0.903772	0.761604	Chaotic ICA
0.801490	0.740055	0.903494	0.760922	Chaotic DE
0.795313	0.732701	0.899840	0.753398	Chaotic GA
0.796512	0.736773	0.902291	0.750472	Chaotic ABC
0.806792	0.748582	0.906767	0.765028	PSO
0.805876	0.745936	0.906324	0.765369	ICA
0.710969	0.564961	0.888349	0.679596	SA
0.808365	0.748006	0.907483	0.769606	DE
0.825738	0.777250	0.918301	0.781664	GA
0.810464	0.752736	0.908755	0.769900	ABC
0.745283	0.643119	0.879649	0.713079	HS
0.762444	0.681335	0.875371	0.730627	K-means

مجموعه داده Ionosphere مربوط به اتمسفر است، این داده‌های راداری توسط یک سیستم در لابراتوار گوس بی^۱ جمع‌آوری شده است. این مجموعه داده دارای ۳۵۱ نمونه است که در ۲ کلاس مختلف طبقه‌بندی شده است. هر یک از نمونه‌های این مجموعه داده شامل ۳۴ ویژگی پیوسته است.

جدول (۷): نتایج میانگین ۱۰۰ بار اجرای الگوریتم‌های تکاملی موجود بر روی مجموعه داده Ionosphere

زمان اجرا	انحراف معیار	تابع هزینه			الگوریتم‌های تکاملی
		حداکثر هزینه	متوسط هزینه	حداقل هزینه	
2.0872s	$\pm (7.25921)$	346.983	326.7905	319.202	Chaotic PSO
1.9458s	$\pm (0.55444)$	1533.331	1531.951	1530.967	Chaotic ICA
2.0105s	$\pm (14.64452)$	1656.076	1611.289	1581.566	Chaotic DE
2.1175s	$\pm (0.42419)$	1533.745	1532.607	1531.832	Chaotic GA
3.5230s	$\pm (9.14E-13)$	1536.158	1536.158	1536.158	Chaotic ABC
2.3357s	$\pm (71.27108)$	1842.476	1657.443	1593.872	PSO
2.0606s	$\pm (32.92895)$	1919.147	1834.948	1731.838	ICA
1.2607s	$\pm (90.13086)$	2419.121	2164.404	1992.604	SA
2.1532s	$\pm (41.20099)$	2063.174	1958.693	1856.828	DE
4.3329s	$\pm (29.52986)$	1762.222	1655.09	1609.814	GA
7.3836s	$\pm (23.8655)$	2006.783	1951.96	1899.199	ABC
7.4590s	$\pm (49.96369)$	2507.927	2422.111	2267.1851	HS
0.0142s	$\pm (18.34671)$	2311.822	2274.583	2217.282	K-means

بعد از بررسی پاسخ الگوریتم‌های آشوب گونه بر روی مجموعه داده Ionosphere بر اساس جدول ۷ این امر مسلم می‌گردد که روش ارائه گردیده قابل انطباق با سیستم‌های واقعی نیز خواهد بود. لذا می‌توان از آن در اکثر مسائل بهینه‌سازی استفاده نمود. به عبارتی دیگر این اطمینان وجود دارد که الگوریتم پیشنهادی، محدودیت‌های الگوریتم‌های گذشته را تا حد قابل قبولی پوشش خواهد داد.

جدول (۸): نتایج ارزیابی داده‌های Ionosphere با معیارهای NMI، FM و AR.

میانگین	معیار AR	معیار FM	معیار NMI	الگوریتم‌های تکاملی
0.346415	0.150641	0.576561	0.312042	Chaotic PSO
0.337670	0.173175	0.707732	0.132102	Chaotic ICA
0.321428	0.154845	0.694127	0.115313	Chaotic DE
0.337383	0.173329	0.707315	0.131504	Chaotic GA
0.332469	0.167902	0.702429	0.127076	Chaotic ABC
0.313036	0.136161	0.701124	0.101822	PSO
0.295523	0.014900	0.839445	0.032224	ICA
0.296213	0.012832	0.849208	0.026599	SA
0.276291	0.023602	0.770898	0.034372	DE
0.303919	0.142661	0.675297	0.093799	GA
0.299716	0.017010	0.853086	0.029051	ABC
0.299514	0.058938	0.775883	0.063718	HS
0.176829	0.000324	0.527865	0.002298	K-means

مجموعه داده Half-ring شامل ۴۰۰ نمونه است که در ۲ کلاس طبقه‌بندی می‌گردد. هر یک از نمونه‌های این مجموعه داده دارای ۲ ویژگی است. در جدول ۹ کارآیی الگوریتم‌ها بر اساس تابع هزینه آورده شده‌اند که با توجه به نتایج این جدول، الگوریتم Chaotic ABC هزینه کمتری در مقایسه با الگوریتم‌های دیگر دارد.

جدول (۹): نتایج میانگین ۱۰۰ بار اجرای الگوریتم‌های تکاملی موجود بر روی مجموعه داده Half-ring.

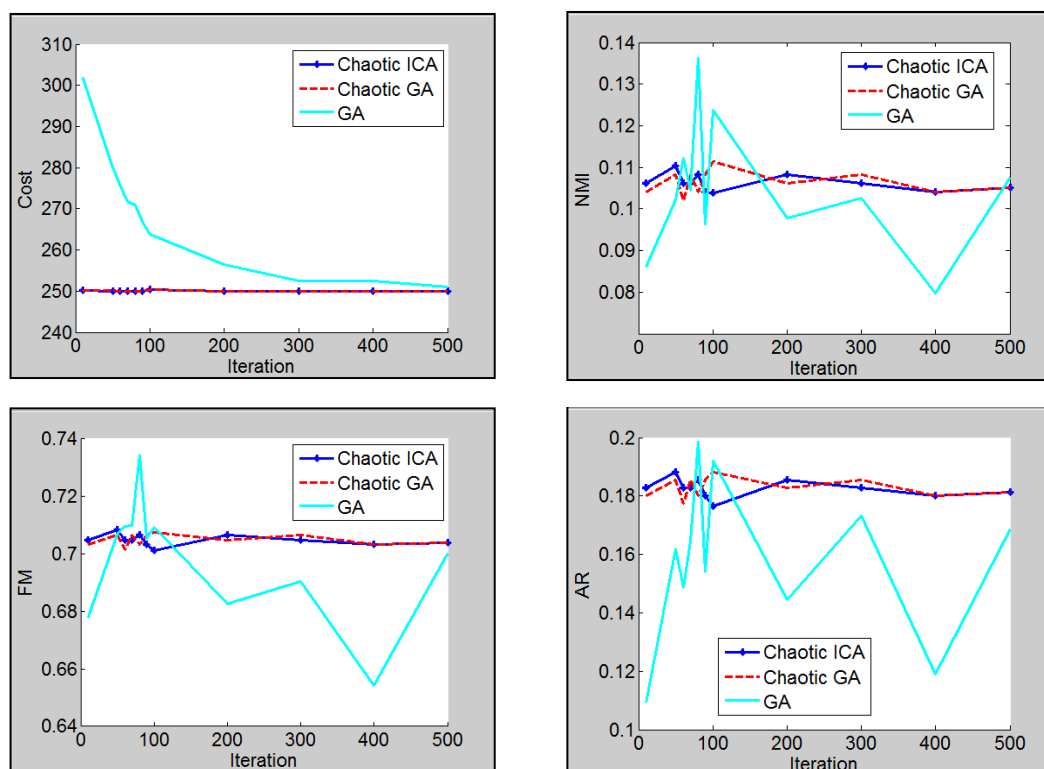
زمان اجرا	انحراف معیار	تابع هزینه			الگوریتم‌های تکاملی
		حداکثر هزینه	متوسط هزینه	حداقل هزینه	
1.8360s	$\pm (0.04052)$	193.2442	193.0917	193.0345	Chaotic PSO
1.6198s	$\pm (0.00327)$	193.0368	193.0299	193.0277	Chaotic ICA
1.6667s	$\pm (0.02313)$	193.1191	193.0454	193.0277	Chaotic DE
1.81264s	$\pm (0.00412)$	193.0515	193.036	193.0277	Chaotic GA
3.19743s	$\pm (0.00125)$	193.0368	193.0288	193.0277	Chaotic ABC
1.9353s	$\pm (0.03844)$	193.2037	193.0928	193.0402	PSO
1.7545s	$\pm (0.00250)$	193.0373	193.029	193.0277	ICA
1.0461s	$\pm (54.01495)$	463.9958	286.5792	203.4527	SA
1.8381s	$\pm (0.02580)$	193.1486	193.0504	193.0277	DE
3.6748s	$\pm (0.00935)$	193.071	193.031	193.0277	GA
6.8335s	$\pm (1.86E-05)$	193.0278	193.0277	193.0277	ABC
3.2754s	$\pm (0.53657)$	195.8395	193.8955	193.0975	HS
0.0112s	$\pm (7.49415)$	287.9756	268.1763	244.9351	K-means

الگوریتم‌های آشوب‌گونه با توجه به جدول ۱۰ کارآیی بالایی بر اساس معیارهای NMI، FM و AR بر روی مجموعه داده Half-ring دارند.

جدول (۱۰): نتایج ارزیابی داده‌های Half-ring با معیارهای NMI، FM و AR.

میانگین	معیار AR	معیار FM	معیار NMI	الگوریتم‌های تکاملی
0.419188	0.193460	0.770180	0.293924	Chaotic PSO
0.422985	0.199610	0.772064	0.297281	Chaotic ICA
0.421410	0.197055	0.771288	0.295886	Chaotic DE
0.427841	0.207491	0.774446	0.301586	Chaotic GA
0.420970	0.196340	0.771075	0.295495	Chaotic ABC
0.419321	0.193683	0.770234	0.294047	PSO
0.421963	0.197952	0.771563	0.296375	ICA
0.466485	0.291851	0.806328	0.301275	SA
0.421418	0.197070	0.77129	0.295894	DE
0.419999	0.194767	0.770594	0.294636	GA
0.419228	0.196110	0.771006	0.290568	ABC
0.413602	0.183890	0.767231	0.289683	HS
0.175657	0.000499	0.524445	0.002028	K-means

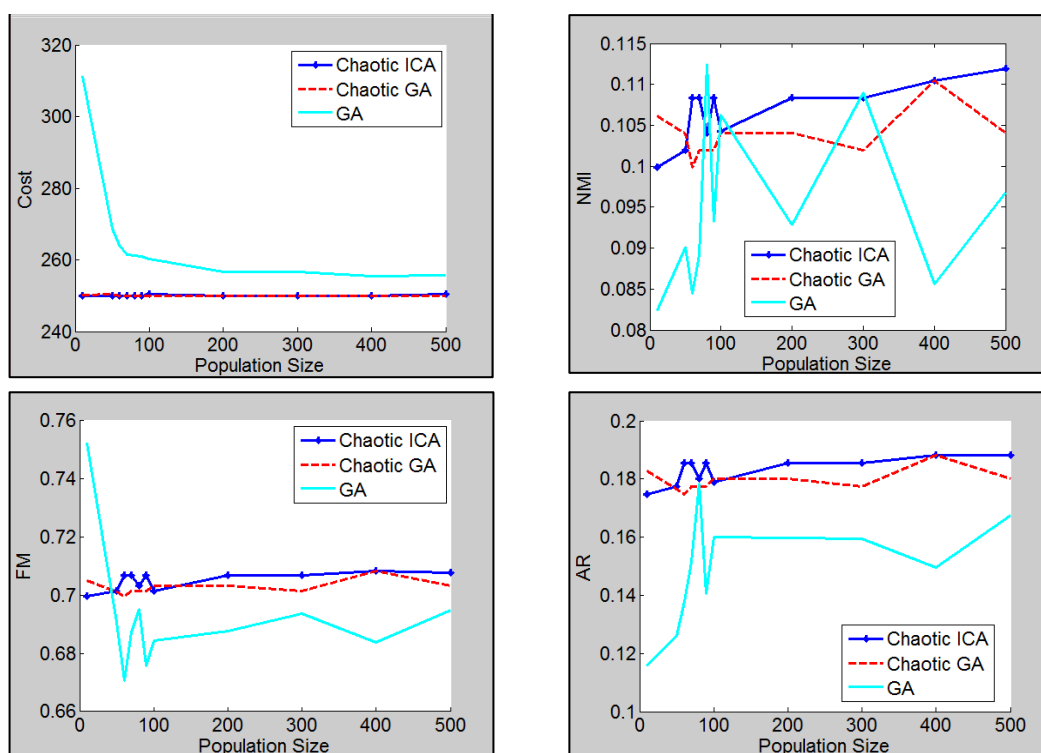
شکل ۶ دقت و کارایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA که بهترین پاسخ را در میان الگوریتم‌های مورد بررسی داشته‌اند را نشان می‌دهند.



شکل ۶: مشخصه همگرایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA بر اساس تعداد تکرارها

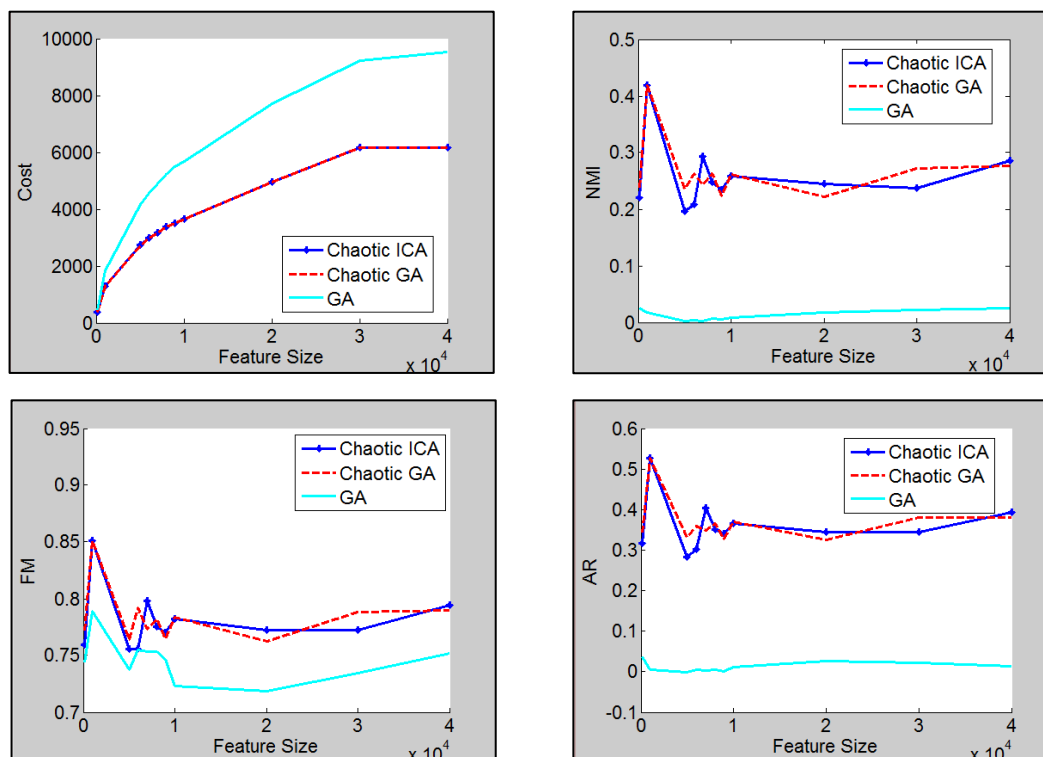
شکل ۶ همگرایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA بر اساس تغییرات در تعداد تکرارها است. نمودارهای همگرایی الگوریتم‌ها با توجه به معیارهای تابع هزینه، NMI، FM و AR است. همان‌طور که مشاهده می‌گردد تغییرات در تکرارها تأثیری بر روی الگوریتم‌های آشوب گونه ندارد.

در شکل ۷ همگرایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA نشان داده شده است. الگوریتم‌های آشوب‌گونه به تغییرات در اندازه جمعیت نیز غیرحساس هستند.



شکل ۷: مشخصه همگرایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA بر اساس اندازه جمعیت

در شکل ۸ مشخصه‌های همگرایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA بر اساس تغییر در تعداد ویژگی‌های مجموعه داده است. همان‌طور که مشاهده می‌گردد الگوریتم‌های آشوب‌گونه به تغییرات در اندازه ویژگی‌های مجموعه داده نیز حساس نیستند.



شکل ۸: مشخصه همگرایی الگوریتم‌های Chaotic ICA، Chaotic GA و GA بر اساس تعداد ویژگی‌ها

نتیجه‌گیری

در این مقاله چندین روش برای بهبود کارایی خوشه‌بندی داده‌ها ارائه گردید. این روش‌ها مبتنی بر الگوریتم‌های تکاملی آشوب‌گونه می‌باشند. نتایج تجربی بر روی مجموعه داده‌های مختلف نشان می‌دهد که روش‌های پیشنهادی در مقایسه با حالات استاندارد خوشه‌بندی دارای کارایی مناسب‌تری می‌باشند. شبیه‌سازی‌ها بر روی مجموعه داده‌های متنوع اجرا گردیده است؛ لذا پاسخ‌های نهایی به‌دست آمده گویای این مطلب می‌باشد که وابستگی به تعداد نمونه، تعداد خوشه‌ها، تعداد پارامترها و تعداد تکرار وجود ندارد و الگوریتم‌های آشوب‌گونه عملکرد مناسبی دارند.

منابع

- Alfonzetti, S., Dilettoso, E., & Salerno, N. (2006). Simulated annealing with restarts for the optimization of electromagnetic devices. *IEEE transactions on magnetics*, 42(4), 1115-1118.
- Atashpaz-Gargari, E., & Lucas, C. (2007). Imperialist competitive algorithm: an algorithm for optimization inspired by imperialistic competition. 2007 IEEE congress on evolutionary computation,
- Coley, D. A. (1999). *An introduction to genetic algorithms for scientists and engineers*. World Scientific Publishing Company.
- Dorigo, M., & Gambardella, L. M. (1997). Ant colonies for the travelling salesman problem. *biosystems*, 43(2), 73-81.
- Fathian, M., Amiri, B., & Maroosi, A. (2007). Application of honey-bee mating optimization algorithm on clustering. *Applied Mathematics and Computation*, 190(2), 1502-1513.
- Holland, J. H. (1975). *Adaptation in natural and artificial systems*. University of Michigan Press, Ann Arbor, MI.
- Kao, Y.-T., Zahara, E., & Kao, I.-W. (2008). A hybridized approach to data clustering. *Expert Systems with Applications*, 34(3), 1754-1762.

- Karaboga, D., & Basturk, B. (2007). A powerful and efficient algorithm for numerical function optimization: artificial bee colony (ABC) algorithm. *Journal of global optimization*, 39, 459-471.
- Kellert, S. H. (1993). *In the wake of chaos: Unpredictable order in dynamical systems*. University of Chicago press.
- Kennedy, J., & Eberhart, R. (1995). Particle swarm optimization. Proceedings of ICNN'95-international conference on neural networks.
- Krishna, K., & Murty, M. N. (1999). Genetic K-means algorithm. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 29(3), 433-439.
- Lai, D. (1996). Comparison study of AR models of the Canadian lynx data: A close look at BDS statistic. *Computational statistics & data analysis*, 22(4), 409-423.
- Liz, E., & Ruiz-Herrera, A. (2012). Chaos in discrete structured population models. *SIAM Journal on Applied Dynamical Systems*, 11(4), 1200-1214.
- Maulik, U., & Bandyopadhyay, S. (2000). Genetic algorithm-based clustering technique. *Pattern recognition*, 33(9), 1455-1465.
- Neal, R. W. (1993). *The logistic Lattice in Random Number Generation*, Division of Mathematics University of Texas at San Antonio].
- Niknam, T., & Amiri, B. (2010). An efficient hybrid approach based on PSO, ACO and k-means for cluster analysis. *Applied soft computing*, 10(1), 183-197.
- Niknam, T., Fard, E. T., Pourjafarian, N., & Roust, A. (2011). An efficient hybrid algorithm based on modified imperialist competitive algorithm and K-means for data clustering. *Engineering Applications of Artificial Intelligence*, 24(2), 306-317.
- Perruchet, P., & Peereman, R. (2004). The exploitation of distributional information in syllable processing. *Journal of Neurolinguistics*, 17(2-3), 97-119.
- Schuster, H. G., & Just, W. (2006). *Deterministic chaos: an introduction*. John Wiley & Sons.
- Shelokar, P., Jayaraman, V. K., & Kulkarni, B. D. (2004). An ant colony approach for clustering. *Analytica chimica acta*, 509(2), 187-195.