

## خوشه بندی داده‌ها به روش الگوریتم فاخته

<sup>1</sup> مهری ملالو<sup>5</sup> و محمد ابراهیم شیری<sup>1</sup>

m.mollalo@aut.ac.ir دانشجوی کارشناسی ارشد علوم کامپیوتر، دانشگاه صنعتی امیرکبیر،

<sup>1</sup> عضو هیئت علمی گروه علوم کامپیوتر، دانشگاه صنعتی امیرکبیر ، shiri@aut.ac.ir

چکیده - خوشه بندی از جمله روش‌های پرکاربرد در تجزیه و تحلیل داده‌ها است که بدون هیچ دانش قبلی، داده‌ها را به گروه‌های معنی‌داری تقسیم می‌کند که یکی از موضوعات اساسی در داده کاوی است. در این مقاله یک الگوریتم خوشه بندی جدید بر اساس سبک خاص زندگی جمعی یک دسته پرمنده بنام فاخته ارائه شده است. الگوریتم فاخته که یکی از روش‌های جدید بهینه‌سازی تکاملی می‌باشد که با الهام گرفتن از روش تخم‌گذاری فاخته‌ها، پرورش تخم‌ها، تلاش برای زنده ماندن در میان دیگر فاخته‌ها و مهاجرت به سمت محیط بهتر برای مساله خوشه بندی داده‌ها استفاده شده است. الگوریتم معروف خوشه بندی k-means به مقداردهی اولیه بسیار حساس است و به راحتی در بهینه‌های محلی گیر می‌افتد در صورتی که الگوریتم پیشنهاد شده به خوشه‌های با اندازه و ابعاد متفاوت حساس نیست، برای مجموعه داده‌های چند بعدی مناسب است و همچنین سرعت بسیار خوبی در همگرایی به نقطه بهینه سراسری و دقت بالا در حل مسائل از خود نشان داده است. کارایی الگوریتم پیشنهادی بر روی مجموعه‌های داده‌ای پایگاه UCI مورد ارزیابی قرار گرفته شده است. نتایج بدست آمده از آزمایشات نشان دهنده بیبود عملکرد این الگوریتم در مقایسه با الگوریتم‌های معروف خوشه بندی همچون k-means و PSO است.

کلید واژه - الگوریتم بهینه‌سازی فاخته، الگوریتم تغییریافته فاخته، خوشه بندی، داده کاوی

### مدل [6] و الگوریتم‌های مبتنی بر مش [7] ارائه شده است.

در بسیاری از روش‌های خوشه بندی، فاصله (میزان شباهت) داده‌های متعلق به یک خوشه نسبت به فاصله آنها با داده‌های دیگر خوشه‌ها، معیار تعیین خوشه در نظر گرفته شده است. در ساده‌ترین روش، هر داده می‌تواند با در نظر گرفتن میزان شباهت یا کمترین فاصله اش تا خوشه‌ها، دقیقاً در یکی از خوشه‌ها عضویت داشته باشد. [9-8].

به غیر از روش‌های فوق، الگوریتم خوشه بندی k-means برای شناسایی خوشه‌هایی که در فضای ورودی غیر خطی از هم جدا هستند، استفاده می‌شود [51-51].

الگوریتم k-means [9] به دلیل سادگی پیاده سازی و هزینه کم محاسباتی و کارایی بالا، محبوب‌ترین روش خوشه بندی داده‌ها است. با این حال یکی از مشکلات این الگوریتم این است که به مقداردهی اولیه بسیار حساس است و به راحتی در بهینه‌های محلی گیر می‌افتد. بطوریکه با افزایش ابعاد مجموعه داده‌ها، پیدا کردن جواب بهینه به یک مساله NP-hard تبدیل می‌شود [59]. لذا به دلیل اهمیت استراتژی خوشه بندی، جدیداً از الگوریتم‌های بهینه سازی سراسری تصادفی همچون الگوریتم ژنتیک (GA)،

### ۱- مقدمه

خوشه بندی، الگوها و نمونه‌ها را بدون دخالت هیچ ناظری، در گروه‌های مختلفی دسته‌بندی می‌کند [5]. در پنجاه سال گذشته از نقطه نظرهای زیادی به مسئله خوشه بندی توجه شده است. این مسئله در زمینه‌های مختلف از قبیل تشخیص الگو، تجزیه و تحلیل داده‌ها، پردازش تصویر، علوم اقتصادی، زیست‌شناسی، بیوانفورماتیک و فشرده‌سازی داده‌ها مورد استفاده قرار گرفته شده است. بنابراین مطالعه در مورد الگوریتم‌های جدید خوشه بندی در زمینه‌های تحقیقاتی از جمله داده کاوی، یادگیری ماشین، آمار، و زیست‌شناسی موضوع مهمی است.

هدف اصلی روش‌های ارائه شده برای خوشه بندی، پیدا کردن گروه‌هایی از داده‌ها است که دارای بیشترین شباهت بین داده‌های درون هر خوشه و کمترین شباهت بین داده‌های بین هر دو خوشه است. در سال‌های اخیر الگوریتم‌های خوشه بندی زیادی از جمله، سلسه مراتبی [1]، الگوریتم‌های بخش‌بندی [9-4]، الگوریتم‌های مبتنی بر چگالی [1]، الگوریتم‌های مبتنی بر



## 2- خوشبندی

هدف از انجام خوشبندی، طبقه‌بندی اشیا مطابق با میزان تشابه بین آنها و سازماندهی داده‌ها در دو یا چند گروه است. تکنیک‌های خوشبندی از جمله روش‌های بدون نظارت است و هیچ شناخت اولیه‌ای از کلاس‌ها و یا در واقع خوشبندی‌ها در آن وجود ندارد.

ایده‌ی خوشبندی برای اولین بار در دهه‌ی 5991 ارائه شد و امروزه مورد توجه بسیاری از محققان قرار گرفته است. خوشبندی یک گروه از اشیا اطلاق می‌شود که شباهت بیشتری با هم نسبت به اشیا در خوشبندی‌های دیگر دارند. عبارت "شباهت" باید به معنای شباهت ریاضی و قابل محاسبه تفسیر گردد.

در فضای داده‌ای، شباهت دو شیء معمولاً به معنای فاصله آن دو شیء در فضای برداری بیان می‌شود. فاصله می‌تواند در بین بردار داده‌ها با هم و یا فاصله‌ی یک داده تا یک نمونه‌ی معرف محاسبه شود. نمونه‌ی معرف معمولاً از پیش شناخته شده نیست و در طی انجام خوشبندی و جداسازی داده‌ها تعیین می‌شود.

مسئله‌ی خوشبندی را بصورت کلی می‌توان به شرح زیر بیان کرد: در مجموعه داده‌ی  $\{x_1, x_2, \dots, x_N\}$ ،  $X$ ،  $N$  شی داده شده به  $K$  خوشبندی ( $C_1, C_2, \dots, C_K$ ) تقسیم می‌شود. بطوریکه، شباهت بین داده‌های درون هر خوشبندی حداکثر و شباهت بین داده‌های بین خوشبندی‌های متفاوت حداقل شود. این مسئله در معادله (5) نمایش داده شده است.

$$\left\{ \begin{array}{l} \bigcup_{i=1}^K C_i = X \\ C_i \cap C_j = \emptyset \quad i, j = 1, 2, \dots, K \quad i \neq j \\ C_i \neq \phi \quad i = 1, 2, \dots, K \end{array} \right. \quad (5)$$

که  $K$  تعداد خوشبندی‌ها و  $C_i$  خوشبندی  $i$  ام می‌باشد.

از نقطه نظر ریاضی خوشبندی  $C_i$  می‌تواند به صورت معادله (1) مشخص شود:

(1)

$$\left\{ \begin{array}{l} C_i = \left\{ x_j \mid \|x_j - z_i\| \leq \|x_j - z_p\|, x_j \in X \right\}, p \neq i, p = 1, \dots, K \\ z_i = \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j \quad i = 1, 2, \dots, K \end{array} \right.$$

||

برای بیانگر فاصله اقلیدسی یا فاصله کسینوسی بین هر دو نقطه در مجموعه نمونه است.  $z_i$  مرکز خوشبندی  $C_i$  است که در واقع میانگین تمامی نقطه‌ها در خوشبندی  $C_i$  نشان می‌دهد.

الگوریتم کلونی مورچه‌ها (ACO) و الگوریتم اجتماع ذرات (PSO) برای خوشبندی داده‌ها استفاده شده است [54-57]. در همه این روش‌ها، الگوریتم با یک جمعیت اولیه شروع بکار کرده و سپس با تعدادی تکرار، فضای راه حل را برای رسیدن به جواب تقریباً بهینه جستجو می‌کند.

بیشتر الگوریتم‌های ارائه شده برای حل مسائل بهینه‌سازی، از رفتار حشرات و پرندگان در طبیعت الهام گرفته شده‌اند. الگوریتم بهینه‌سازی فاخته (COA<sup>1</sup>) یک الگوریتم بهینه‌سازی تکاملی می‌باشد که از سبک خاص زندگی یک نوع پرندگان بنام فاخته‌آ الهام گرفته شده است. تخم‌گذاری و تولیدمثل خاص فاخته‌ها اساس این الگوریتم بهینه‌سازی می‌باشد. در این الگوریتم فاخته‌ها در دو مدل استفاده شده‌اند: فاخته‌های بالغ و فاخته‌هایی که در تخم هستند. فاخته‌های بالغ در لانه‌ی بعضی پرندگان میزبان شناسایی نشوند و از بین نرونده، رشد کرده و به فاخته‌های بالغ تبدیل می‌شوند. ویژگی‌های محیطی و مهاجرت گروه‌های فاخته‌ها، باعث همگرا شدن الگوریتم و پیدا کردن جواب بهینه سراسری می‌شود. در این مقاله الگوریتم بهینه‌سازی فاخته برای خوشبندی داده‌ها گسترش داده شده است و کارایی آن با الگوریتم‌هایی همچون k-means و PSO بر روی مجموعه‌های داده‌ای مختلفی از پایگاه داده UCI مورد مقایسه قرار گرفته شده است. نتایج بدست آمده، نشان دهنده کارایی عملکرد این روش نسبت به روش‌های دیگر خوشبندی می‌باشد.

در ادامه در بخش 1 مسئله خوشبندی داده‌ها را مطرح می‌شود. سپس در بخش 9 ویژگی زندگی پرندگی فاخته و جزئیات مربوط به تخم‌گذاری و مهاجرت آن و همچنین الگوریتم پیشنهاد شده با جزئیات بیشتری مورد بررسی قرار خواهد گرفت. در ادامه در بخش 4 نتایج بدست آمده از شبیه‌سازی الگوریتم پیشنهاد شده بر روی مجموعه‌های داده‌ای نشان داده شده است. در نهایت در بخش 1 نتیجه‌گیری کارها ارائه شده است.

<sup>1</sup> Cuckoo Optimization Algorithm  
<sup>2</sup> Cuckoo



اندازند و در نهایت در هر لانه یک تخم شانس رشد را پیدا می-کنند.

### 2-3- الگوریتم پیشنهاد شده جهت خوشبندی داده‌ها

روش پیشنهاد شده جهت خوشبندی داده‌ها به روش الگوریتم فاخته (CCA<sup>4</sup>) برای مسال با ابعاد زیاد کارایی خوبی از خود نشان داده است. همانند سایر الگوریتم‌های تکاملی، CCA

هم با یک جمعیت اولیه کار خود را شروع می‌کند.

این جمعیت از فاخته‌ها هر کدام بدنیال پیدا کردن K مرکز خوشبندی می‌باشد که برای رسیدن به چنین هدفی باید تولید مثل و تخم‌گذاری کنند تا محیط توسط فاخته‌های پیشتری جستجو شود. فاخته‌ها تخم‌های خود را در لانه تعدادی پرنده‌ی میزبان قرار می‌دهند. تعدادی از تخم‌ها توسط پرنده میزبان شناسایی شده و از بین می‌روند و برخی دیگر شانس زنده ماندن خواهند داشت. موقعیتی که در آن بیشترین تعداد تخم‌ها نجات یابند پارامتری خواهد بود که COA قصد بهینه‌سازی آن را دارد.

مقدار تابع برازش برای مسائل خوشبندی داده‌ها برابر میانگین فاصله داده‌ها از مراکز خوشبندی می‌باشد و به صورت معادله (4) تعریف می‌شود:

$$\text{fitness}(z_i) = \left( \frac{1}{n} \sum_{m=1}^K \sum_{x_j \in c_m} \|x_j - z_{im}\| \right)^{-1} \quad (4)$$

که  $\{z_{i1}, z_{i2}, \dots, z_{iK}\}$  نشان دهنده فاخته i است که شامل K مرکز خوشبندی است و  $z_{im}$  نمایانگر مرکز خوشبندی m از فاخته i است. همچنین  $\{x_1, x_2, \dots, x_N\}$  تعداد داده‌های مجموعه داده‌ای می‌باشد. مقدار تابع برازش فوق نشان می‌دهد که مراکز خوشبندی با میانگین کمترین فاصله از مقدار برازش بیشتری برخوردارند.

در طبیعت هر فاخته بین 1 تا 11 تخم می‌گذارد. عادت دیگر هر فاخته این است که در یک دامنه مشخص تخم‌های خود را می‌گذارد که حداقل دامنه تخم‌گذاری (ELR<sup>5</sup>) می‌باشد. در

از معادله (1) در می‌یابیم که  $C_i$  توسط داده‌هایی که نزدیکترین عناصر به  $z_i$  هستند ساخته می‌شود. بنابراین وظیفه خوشبندی انجام فرایندی است که K مرکز خوشبندی ( $C_1, C_2, \dots, C_K$ ) را تعیین کند. به طوریکه N شی داده شده به K خوشبندی می‌شود به طوریکه مجموع مجذور فاصله اقلیدسی بین هر شی و مرکز خوشبندی می‌شود که به آن شی اختصاص داده شده است حداقل باشد. این مساله در معادله (9) نمایش داده شده است که جمع مربع خطاهای (SSE<sup>3</sup>) می‌باشد [58].

$$SSE = \sum_{i=1}^K \sum_{x_j \in c_i} \|x_j - z_i\|^2 \quad (9)$$

که  $\{z_1, z_2, \dots, z_K\}$  نشان دهنده K مرکز خوشبندی است و  $x_j \in \{x_1, x_2, \dots, x_N\}$  تعداد داده‌های مجموعه داده‌ای می‌باشد.

برای هر یک از داده‌ها در مجموعه ارائه شده، میزان فاصله تا نزدیکترین خوشبندی خطای می‌باشد.

### 3- روش پیشنهاد شده: خوشبندی داده‌ها به روش الگوریتم فاخته

#### 3-1- شیوه‌ی خاص زندگی و تخم‌گذاری فاخته‌ها

فاخته‌ها پرنده‌گانی هستند که خود را از دردسر هرگونه لانه سازی و وظایف والدین رهانیده‌اند و به نوعی زیرکی جهت پرورش جوجه‌های خود متولّ شده‌اند. این پرنده‌گان تخم‌های خود را در لانه سایر انواع پرنده‌گان قرار می‌دهند و صبر می‌کنند تا آنها در کنار تخم‌های خود به تخم‌های این پرنده‌گان نیز رسیدگی کنند. در این بین هستند پرنده‌گانی که تخم‌های فاخته‌فاخته را از لانه بیرون پرت می‌کنند [59].

جوچه‌های فاخته زودتر از تخم‌های پرنده میزبان از تخم بیرون می‌آیند و زودتر هم رشد می‌کنند. در اکثر موارد جوچه فاخته تخم‌ها و یا جوچه‌های پرنده میزبان را از لانه بیرون می-

<sup>4</sup> Cuckoo Clustering Algorithm

<sup>5</sup> Egg Laying Radius

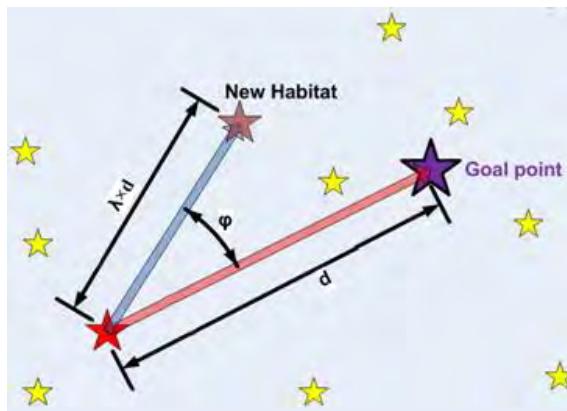
<sup>3</sup> Sum of squared error

به سمت محل هدف طی نمی‌کنند. آنها فقط قسمتی از مسیر را طی کرده و در آن مسیر نیز مقداری انحراف دارند. این نحوه حرکت در شکل 1 بوضوح قابل مشاهده است.

همانطور که در شکل 2 دیده می‌شود هر فاخته فقط درصدی از کل مسیر را به سمت هدف ایده‌آل فعلی طی می‌کند ( $\lambda\%$ ) و یک انحراف  $\varphi$  نیز دارد. این دو پارامتر به فاخته‌ها کمک می‌کنند تا محیط بیشتری را جستجو کنند و از گیرافتادن در بهینه‌های محلی اجتناب کنند.  $\lambda$  عددی تصادفی بین 0 و 1 است و  $\varphi$  عددی بین  $\pi/6$  و  $\pi/2$  می‌باشد.

وقتی تمام فاخته‌ها به سمت نقطه هدف مهاجرت کردند و نقاط سکونت جدید هر کدام مشخص شد، دوباره مراحل تخمگذاری و محاسبه ELR برای هر یک از فاخته‌ها انجام می‌شود.

با توجه به این واقعیت که همیشه تعادلی بین جمعیت پرنده‌گان در طبیعت وجود دارد عددی مثل  $N_{max}$  حداکثر تعداد فاخته‌هایی که می‌توانند در یک محیط زندگی کنند را کنترل می‌کند. این تعادل بدلیل محدودیت‌های غذایی، شکار شدن توسط شکارچیان و نیز عدم امکان پیدا کردن لانه‌های مناسب برای تخم‌ها وجود دارد.



شکل 2. مهاجرت یک نمونه فاخته به سمت محل اقامت هدف

شبیه کد الگوریتم خوشبندی داده‌ها به روش زندگی فاخته (CCA) به صورت زیر می‌باشد (الگوریتم 1):

**الگوریتم 1:** شبیه کد الگوریتم خوشبندی داده‌ها به روش الگوریتم (CCA)  
**Require:** Data set,  $X = \{x_1, x_2, \dots, x_N\}$ ;  
Cluster number, K;  
**Ensure:** Clusters:  $\{C_1, C_2, \dots, C_K\}$ ;

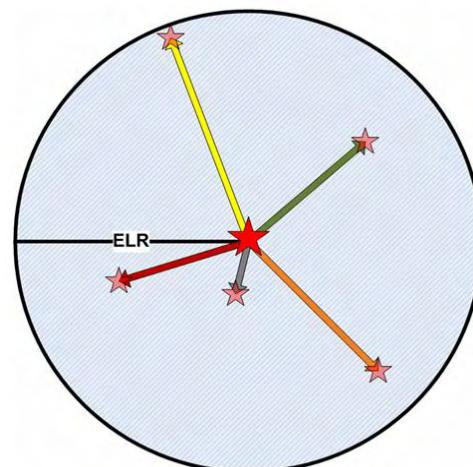
یک مساله بهینه‌سازی به حد بالای متغیرها  $var_{hi}$  و به حد پایین متغیرها  $var_{low}$  می‌گویند. هر فاخته دارای ELR مخصوص به خود خواهد بود که با تعداد کل تخم‌های موجود در محیط، تعداد تخم‌های فعلی فاخته و همچنین حد بالا و پایین متغیرهای مساله متناسب است.

ELR به صورت معادله (1) محاسبه می‌شود:

$$ELR = \alpha \times \frac{\text{Number of current eggs}}{\text{Total number of eggs}} \times (var_{hi} - var_{low}) \quad (1)$$

$\alpha$  متغیری است که حداقل شعاع تخم‌گذاری با آن تنظیم می‌شود.

هر فاخته بصورت تصادفی تخم‌هایی را در لانه پرنده‌گان می‌زیبان که در ELR خود قرار دارد، می‌گذارد (شکل 1). نکته جالب در مورد جوجه فاخته‌ها این است که فقط یک تخم در هر لانه امکان رشد دارد. چرا که وقتی جوجه‌های فاخته از تخم بیرون می‌آیند تخم‌های پرنده می‌زیبان از لانه بیرون می‌کنند و اگر جوجه‌های پرنده می‌زیبان زودتر از تخم خارج شده باشند جوجه فاخته بیشترین مقدار غذا را که پرنده می‌زیبان می‌آورد، خواهد خورد و پس از چند روز جوجه‌های خود پرنده می‌زیبان از گرسنگی می‌میرند و فقط جوجه فاخته زنده می‌ماند.



شکل 1. تخمگذاری تصادفی در محدوده ELR. ستاره قرمز رنگی که در وسط قرار دارد محل سکونت اولیه فاخته‌ای با 5 تخم مرغ می‌باشد؛ ستاره‌های صورتی لانه‌ی جدید تخم‌مرغ‌ها هستند.

وقتی جوجه فاخته‌ها رشد کردند و تبدیل به فاخته بالغ شدند، پرنده‌های فاخته دیگر با میزانی فاصله و انحراف به سمت بهترین منطقه فعلی در بین تمام فاخته‌ها که در آنجا شansas زنده ماندن تخم‌ها بیشتر است مهاجرت می‌کنند.

هنگام مهاجرت به سمت نقطه هدف، فاخته‌ها تمام مسیر را



الگوریتم پیشنهادی بر روی سه مجموعه دادهای ارزیابی شده است. مجموعه دادهای Iris شامل داده‌هایی در مورد سه نوع گل زنبق است که دارای 151 نمونه داده است که به سه کلاس تقسیم شده‌اند که در هر کلاس 11 نمونه قرار دارد و هر نمونه 578 چهار متغیر ورودی دارد. مجموعه دادهای Wine دارای 154 نمونه داده است که به سه دسته یا کلاس با 9 متغیر و بعد تقسیم شده است. مجموعه دادهای Glass یکی دیگر از مجموعه-هاست که در رابطه با اجزاء شیمیایی تشکیا دهنده 6 نوع شیشه است و دارای 6 کلاس است و از 154 نمونه تشکیل شده است و همچنین دارای 9 فیلد و ورودی است.

جدول ۵: مشخصات مجموعه‌های دادهای

خوشه	مجموعه داده	نمونه‌ها	ویژگی / بعد داده‌ها
9	Iris	511	4
9	Wine	578	59
6	Glass	154	9

## 2-4- تنظیم پارامترها

پارامترهای الگوریتم‌های خوشبندی مختلف به این صورت تنظیم شده‌اند که برای PSO 11 ذره اولیه و  $w=0.72$  و  $C_1 = C_2 = 1.49$  را تنظیم شده است.

برای CCA جمعیت اولیه 11 فاخته و حداقل فاخته‌هایی که می‌توانند در محیط زنده باشند نیز 11 در نظر گرفته شده است. حداقل و حداقل مقادیر تخمی که هر فاخته می‌تواند 9 و 8 است.

## 3-4- نتایج شبیه‌سازی

در این آزمایشات، تعداد تکرارهای الگوریتم برابر 91 تکرار در نظر گرفته شده است و از معیار فاصله اقلیدسی برای به دست آوردن فاصله بین نمونه‌ها نسبت به مرکز خوشبندی استفاده شده است. شکل 9 مجموعه دادهای Iris را قبل از خوشبندی نشان می‌دهد و در شکل 4 مجموعه دادهای Iris بعد از خوشبندی نمایش داده است. این مجموعه دادهای دارای چهار بعد است. همانطور که به روشنی در تصویر قابل مشاهده است، این الگوریتم به درستی، مجموعه نمونه‌های ورودی را به سه خوش و دسته تقسیم‌بندی کرده است.

```

1: Initialize:
numCuckooS: number of initial population;
minNumberOfEggs: minimum number of eggs for each
cuckoo;
maxNumberOfEggs: maximum number of eggs for each
cuckoo;
maxIter: maximum iteration of the Cuckoo Clustering
Algorithm;
maxNumOfCuckoos: maximum number of cuckoos that
can live at the same time;
varlow;
varhi;
Generate the Position of Initialize Population randomly
which each single  $z = \{z_1, z_2, \dots, z_{numCuckooS}\}$ ;
cuckoo
contains  $\mathbf{K}$  randomly generated centroid vectors  $z_i$ 
 $z_i = \{z_{i1}, z_{i2}, \dots, z_{iK}\}$ 
for t=1:maxIter do
    for i=1:numCuckooS do
        Initialize number of eggs for each  $z_i$ 
    end for
    for i=1:numCuckooS do
        Calculate ELR(egg laying radius) according to
        Eq.(5) for current  $z_i$ 
    end for
    Lay eggs in different nests
    Remove the eggs on the same positions for all
    Cuckoos
    for i=1:numCuckooS do
        Evaluate fitness function for each  $z_i$  and his
        eggs and put them under each other
    end for
    if Population > maxNumOfCuckoos then
        Kill cuckoos in worst position
    end if
    Update the global best position  $z_{best_i}$ : Select the
    best  $z_{best_i}$  from  $\{z_1, z_2, \dots, z_{numCuckooS}\}$ 
    Move all cuckoos toward  $z_{best_i}$  (Immigration of
    cuckoo)
end for

```

## 4- نتایج آزمایشات

### 1-4- مجموعه داده

در این بخش نتایج شبیه‌سازی روش پیشنهادی خوشبندی به روش الگوریتم فاخته (CCA) نسبت به روش‌های k-means و PSO ارائه داده می‌شود.

در این مقاله، برای بررسی و ارزیابی صحت عملکرد الگوریتم پیشنهادی از سه مجموعه Iris، Wine و Glass از مجموعه داده‌های مربوط به پایگاه داده UCI استفاده شده است. جزئیات مربوط به این مجموعه داده‌ها در جدول 5 ارائه شده است.

فاصله نمونه‌ها از مراکز، با هم مقایسه شده‌اند.

معیار دقت با استفاده از برچسب کلاس‌ها، میزان دقت روش خوشبندی در انتساب داده‌ها به خوشبندی را نشان می‌دهد. لازم به ذکر است که برچسب کلاس‌ها تنها برای ارزیابی نتایج استفاده می‌شوند و در فرآیند خوشبندی دخالتی ندارند. رابطه (6) این معیار را نشان می‌دهد.

$$\text{Accuracy}(\pi) = \frac{\sum_{m=1}^p \text{majority}(C_m | L_m)}{N} \quad (6)$$

$\text{majority}(C_m | L_m)$  در رابطه (6) تعداد داده‌ای هایی از

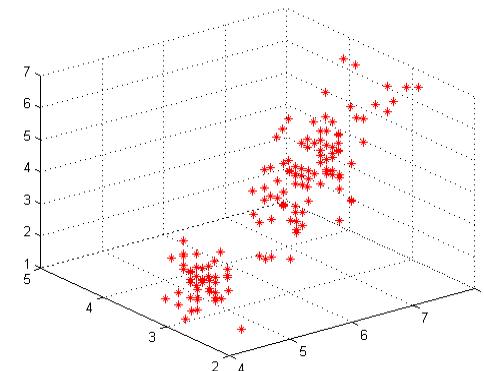
یک کلاس است که دارای اکثریت در خوشبندی  $C_m$  می‌باشد.

نتایج حاصل از شبیه‌سازی الگوریتم‌ها بر روی مجموعه‌های داده‌ای Iris، Wine و Glass در جدول 1 مشاهده می‌شود. با مقایسه نتایج گزارش شده در جدول 1 مشاهده می‌شود که الگوریتم پیشنهادی، نسبت به الگوریتم‌های خوشبندی k-means و PSO از دقت و کارایی بالاتری برخوردار است.

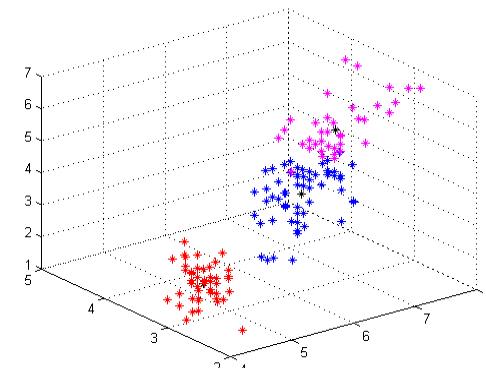
جدول 1: مقایسه کارایی الگوریتم‌های k-means، PSO و CCA بر روی مجموعه‌های داده‌ای موجود در جدول 5

	میانگین فاصله تا مرکز	دقت	روش	مجموعه داده
97/9461	1/881	k-means		
97/9119	1/999	PSO	Iris	
<b>97/0510</b>	<b>0/998</b>	CCA		
5661/6	1/711	k-means		
5614/8	1/699	PSO	Wine	
<b>1627/4</b>	<b>0/776</b>	CCA		
141/1749	1/477	k-means		
198/9111	1/481	PSO	Glass	
<b>213/5036</b>	<b>0/508</b>	CCA		

همانطور که در مجموعه داده‌ای Iris مشاهده می‌شود، الگوریتم CCA میانگین فاصله تا مرکز را 97/1151 تا 97/1119 تشخیص داده است و این در صورتی است که الگوریتم‌های k-means و PSO به ترتیب مقادیر 97/9461 و 97/9119 را برای این معیار تشخیص داده‌اند که با مقایسه این مقادیر، برتری روش CCA نسبت به الگوریتم‌های k-means و PSO قابل مشاهده است. زیرا هرچه مقدار بدست آمده از معیار میانگین فاصله تا مرکز کمتر و مقدار معیار دقت بزرگتر باشد نشان دهنده عملکرد بهتر الگوریتم



شکل 9: مجموعه داده‌ای Iris قبل از خوشبندی



شکل 4: مجموعه داده‌ای Iris بعد از خوشبندی که 9 خوشبندی موجود با رنگ‌های متفاوت نشان داده شده‌اند.

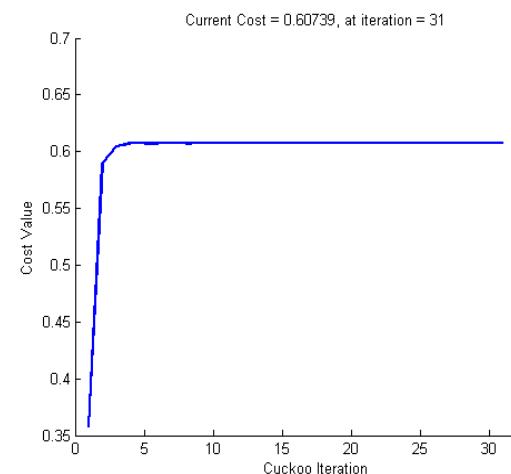
نتایج شبیه‌سازی‌های انجام شده بر روی مجموعه‌های داده‌ای Iris و Wine به صورت میانگین بعد از 91 مرتبه تکرار الگوریتم، در جدول (1) ارائه شده است. یکی از ویژگی‌های الگوریتم پیشنهادی، سرعت قابل توجه آن در همگرا شدن به جواب نزدیک به جواب بهینه است. این الگوریتم در تکرارهای پایین (تا تکرار 51) به جواب نزدیک به جواب بهینه همگرا می‌شود. همانطور که در تصویر 1 قابل مشاهده است این الگوریتم در هر 91 تکراری که انجام شده است به صورت میانگین تا تکرار 51 به جواب تقریباً بهینه همگرا شده است که این نشان از سرعت بالای الگوریتم در رسیدن به جواب در تکرارهای پایین می‌باشد. الگوریتم پیشنهاد شده (CCA) با الگوریتم‌های خوشبندی PSO و k-means بر اساس دو معیار دقت<sup>6</sup> و میانگین<sup>7</sup>

<sup>6</sup> Iterations  
<sup>7</sup> Accuracy

## مراجع

- [1] A.K. Jain, M.N. Murty, P.J. Flynn, Data clustering: a review, ACM Computing Surveys
- [2] S. Guha, et al., "CURE: an efficient clustering algorithm for large databases," SIGMOD Rec., vol. 27, pp.73\_84, 1998.
- [3] J. MacQueen, "Some Methods for Classification and Analysis of Multivariate Observations," presented at the Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability- Vol. 1, 1967.
- [4] R. T. Ng and J. Han, "Efficient and Effective Clustering Methods for Spatial Data Mining," presented at the Proceedings of the 20th International Conference on Very Large Data Bases, 1994.
- [5] A. Hinneburg, et al., "An Efficient Approach to Clustering in Large Multimedia Databases with Noise," 1998, pp.58\_65.
- [6] A. P. Dempster, et al., "Maximum Likelihood from Incomplete Data via the EM Algorithm," Journal of the Royal Statistical Society. Series B (Methodological), vol. 39, pp. 1\_38, 1977.
- [7] G. Sheikholeslami, et al., "WaveCluster: a wavelet-based clustering approach for spatial data in very large databases," The VLDB Journal, vol. 8, pp. 289\_304, 2000.
- [8] Duda, R., Hart, P., (1973), Pattern Classification and Scene Analysis. J. Wiley & Sons, Inc. New York, NY, USA.
- [9] Krishnapuram, R. and Keller, J., (1996), The possibilistic c-means algorithm: insights and recommendations. IEEE Trans. Fuzzy Systems 4,385-393.
- [10] I.S. Dhillon, Y. Guan, B. Kulis, Weighted graph cuts without eigen vectors: a multilevel approach, IEEE Transactions on Pattern Analysis and Machine Intelligence 29 (11) (2007) 1944\_1957.
- [11] I.S. Dhillon, Y. Guan, B. Kulis, A unified view of kernel k-means, spectral clustering and graph partitioning, Technical Report TR\_04\_25, UTCS, 2005.
- [12] I.S. Dhillon, Y. Guan, B. Kulis, Kernel k-means: spectral clustering and normalized cuts, in: Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM Press, NY, USA, 2004, pp.551\_556.
- [13] X.Rui and D. Wunsch, II, "Survey of clustering algorithms, "Neural Networks, IEEE Transactions on, vol , 16 .pp. 645\_678, 2005.
- [14] J. Handl and B. Meyer, "Ant-based and swarm-based clustering," Swarm Intelligence, vol. 1, pp. 95-113, 2007/12/01 2007.
- [15] E. R. Hruschka, et al., "Evolving clusters in gene-expression data, " Information Sciences, vol. 176 ,pp. 1898-1927, 2006.
- [16] P. S. Shelokar, et al., "An ant colony approach for clustering," Analytica Chimica Acta, vol. 509, pp. 187-195, 2004.
- [17] D. W. van der Merwe and A. P. Engelbrecht, "Data clustering using particle swarm optimization," in Evolutionary Computation, 2003. CEC '03. The 2003 Congress on, 2003, pp. 215-220 Vol.1.
- [18] P-N. Tan, M. Steinbach, V. Kumar, Introduction to Data Mining, Pearson Addison-Wesley, 2006.
- [19] TheLife of Birds,Parenthood. <http://www.pbs.org/lifeofbirds/home/index.html> (retrieved 05.09.09)

در خوشبندی داده‌ها می‌باشد.



تصویر ۵. همگرا شدن الگوریتم به جواب تقریباً بهینه در تکرارهای پایین

## 5- نتیجه‌گیری

در این مقاله یک الگوریتم جدید، مبتنی بر سبک خاص زندگی پرندگان فاخته جهت خوشبندی داده‌ها ارائه شد. در این الگوریتم از روش تخم‌گذاری و مهاجرت فاخته‌ها برای خوشبندی داده‌ها الهام گرفته شد. خوشبندی داده‌ها در واقع پیدا کردن بهینه‌ترین مرکز خوشبندی می‌باشد که این کار بوسیله‌ی بهینه کردن تابع برازش انجام می‌شود. یکی از مزایای این الگوریتم این است که به راحتی از گیرافتان در بهینه‌های محلی اجتناب می‌کند و همچنین سرعت اجرای الگوریتم نسبت به الگوریتم‌های خوشبندی PSO و k-means مقایسه شد. نتایج شبیه‌سازی‌های انجام شده بر روی مجموعه‌های داده‌ای مختلف نشان از بهبود کارایی و عملکرد CCA در مقایسه با الگوریتم‌های خوشبندی k-means و PSO است.

## سپاسگزاری

در این قسمت جا دارد که از استاد فرزانه و داشمند، جناب آقای دکتر پرهام مرادی دولت آبادی، به پاس حمایت‌ها و راهنمایی‌های ارزنده‌شان تشکر و قدردانی کنم و برای ایشان از خداوند منان عزت روز افزون را خواستارم.