

# یک روش مبتنی بر خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده جهت شاخص‌گذاری اطلاعات تصویری

محمد رضا کیوان پور<sup>۱</sup>، سعیده رنجبران<sup>۲</sup> و نجوا ایزدپناه<sup>۳</sup>  
<sup>۱</sup>دانشکده فنی و مهندسی، دانشگاه الزهراء(س)، تهران، ایران  
<sup>۲</sup>و<sup>۳</sup>دانشکده برق، رایانه و فناوری اطلاعات، دانشگاه آزاد اسلامی قزوین، قزوین، ایران

## چکیده

در سامانه‌های رایج بازیابی تصویر مبتنی بر محتوا از ساختارهای شاخص‌گذاری چندبعدی برای سرعت‌بخشیدن به عملیات جستجو استفاده می‌شود. در اکثر حوزه‌های کاربردی، ابعاد بالای از بردارهای ویژگی چندبعدی برای توصیف تصاویر مورد نیاز است؛ اما ساختارهای شاخص‌گذاری چندبعدی رایج کارایی خود را با افزایش ابعاد فضای ویژگی از دست می‌دهند. افزایش ابعاد فضای داده موجب افزایش نمای اندازه فضای جستجو و تعداد گره‌ها در ساختارهای شاخص‌گذاری چندبعدی و همچنین موجب افزایش هم‌پوشانی بین گره‌های ساختارهای شاخص‌گذاری چندبعدی می‌شود. این مسائل منجر به افزایش هزینه جستجو از طریق ساختارهای شاخص‌گذاری چندبعدی رایج می‌شود. هدف این پژوهش ارائه یک ساختار شاخص‌گذاری تصویر مبتنی بر خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده، جهت مدیریت بردارهای ویژگی متناظر با تصاویر، در فضای ویژگی با ابعاد بالاست که در ضمن از هم‌پوشانی نواحی گره‌ها نیز جلوگیری به عمل آورد. آزمون‌های مختلف نشان داده که ساختار شاخص‌گذاری پیشنهادی در فضاهای با ابعاد بالا کارایی مناسبی داشته و نسبت به رویکردهای پیشین دارای برتری است.

واژگان کلیدی: بازیابی تصویر مبتنی بر محتوا، ساختارهای شاخص‌گذاری چندبعدی، خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده، روش‌های پیگرد افکنش.

## ۱- مقدمه

در سال‌های اخیر با رشد فناوری و افزایش توانایی رایانه‌ها در مدیریت داده‌های تصویری و حجم وسیع تصاویر تولیدشده در حوزه‌های مختلف، طراحی و پیاده‌سازی پایگاه‌داده‌های تصویری و بازیابی از آنها به‌عنوان یک الزام مطرح است (Datta et al., 2008, Srinivasa Rao et al., 2010, Bohm et al., 2001). به‌طور متداول، اغلب سامانه‌های بازیابی تصویر به‌عنوان یک راهکار فعل و انفعالی برای بازیابی محتوا و تفسیر نیمه خودکار تصویر به‌کار می‌روند (Asbaghi et al., 2008). انواع مختلفی از کاربردهای بازیابی تصویر مطرح است؛ برای مثال، بخشی از تلاش‌های صورت‌گرفته توسط محققان به بازیابی تصاویر متنی<sup>۱</sup> اختصاص یافته است (Keyvanpour and Tavoli, 2013, Keyvanpour and

1 Document Image Retrieval

(Tavoli, 2012). روش‌های بازیابی تصویر مبتنی بر محتوا و مبتنی بر متن از روش‌های بنیادین برای بازیابی اطلاعات تصویری محسوب می‌شوند (Keyvanpour and Asbaghi, 2008). در بازیابی تصویر مبتنی بر محتوا، ویژگی‌های بصری<sup>۲</sup> شامل ویژگی رنگ<sup>۳</sup>، ویژگی بافت<sup>۴</sup>، ویژگی شکل<sup>۵</sup> و ویژگی‌های محلی<sup>۶</sup> از هر تصویر استخراج می‌شوند و در قالب بردارهای ویژگی چندبعدی<sup>۷</sup> سازماندهی می‌شوند. در مرحله جستجو، وقتی کاربر پرس‌وجوی خود را با انتخاب تصویر مورد نظرش مطرح می‌کند، یک بردار ویژگی برای تصویر مورد جستجو<sup>۸</sup> محاسبه می‌شود؛ سپس، با استفاده از

<sup>2</sup>Visual features

<sup>3</sup>Color feature

<sup>4</sup>Texture feature

<sup>5</sup>Shape feature

<sup>6</sup>Local features

<sup>7</sup>Multi-dimensional Feature vectors

<sup>8</sup>Image query

مبتنی بر تراکم<sup>۱</sup>، طبقه‌بندی کرده‌اند (Xu et al., 2007). الگوریتم‌های سلسله‌مراتبی به دو نوع سلسله‌مراتبی تقسیم‌کننده<sup>۲</sup> شامل الگوریتم PDDP (Boley, 1998) و سلسله‌مراتبی تجمیعی<sup>۳</sup> شامل الگوریتم‌های BIRCH (Zhang et al., 1996) و CURE (Guha et al., 1998)، تقسیم شده‌اند که با هدف تجزیه فضای داده در یک ساختار سلسله‌مراتبی توسعه یافته‌اند.

هدف از این پژوهش، ارائه یک ساختار شاخص‌گذاری چندبعدی مبتنی بر خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده<sup>۴</sup>، جهت مدیریت بردارهای ویژگی متناظر با تصاویر، در فضای ویژگی با ابعاد بالاست که در ضمن از هم‌پوشانی نواحی گره‌ها نیز جلوگیری به عمل می‌آورد. بدین منظور از الگوریتم FastICA (Hyvarinen and Oja, 2000)، با هدف پیدا کردن یک مؤلفه غیرگوسی معنادار از داده‌ها، که افکنش داده‌ها<sup>۵</sup> روی آن تقریب آشفستگی منفی<sup>۶</sup> را بیشینه کرده و در نتیجه افکنش داده‌ها روی آن برای تجزیه داده‌ها به دو زیرخوشه می‌تواند حاوی اطلاعات مفیدی باشد، استفاده شده است. در ادامه ابتدا در بخش ۲، نحوه عملکرد سامانه بازیابی تصویر براساس ساختار شاخص‌گذاری پیشنهادی، و سپس در بخش ۳، ساختار شاخص‌گذاری چندبعدی پیشنهادی (NO-NGP-tree<sup>۷</sup>) مطرح شده است. در بخش ۴، پیاده‌سازی و آزمون مطرح شده است. و در نهایت بخش ۵ به نتایج و فعالیت‌های آینده پرداخته است.

## ۲- نحوه عملکرد کلی سامانه بازیابی تصویر براساس ساختار شاخص‌گذاری پیشنهادی

سامانه بازیابی تصویر براساس ساختار شاخص‌گذاری پیشنهادی، مانند معماری کلی سامانه‌های بازیابی تصویر مبتنی بر محتوای ارائه‌شده در (Rui and Huang, 1999)، از چهار قسمت اصلی استخراج ویژگی، مجموعه‌های داده، موتور بازیابی و طرح شاخص‌گذاری چندبعدی تشکیل شده

معیارهای شباهت<sup>۱</sup>، بردار ویژگی تصویر مورد جستجو با بردارهای ویژگی متناظر با تصاویر پایگاه داده مقایسه شده و تصاویری که بیشترین شباهت را با بردار تصویر مورد جستجو دارند به‌عنوان نتایج برگردانده می‌شوند (Datta et al., 2008, Lew et al., 2006).

یکی از مسائل مطرح در حوزه بازیابی تصویر مبتنی بر محتوا، به‌کارگیری روش‌هایی برای سرعت‌بخشیدن به عملیات جستجو در پایگاه داده‌های تصویری است. بدین منظور، از ساختارهای شاخص‌گذاری چندبعدی<sup>۲</sup> برای سازماندهی بردارهای ویژگی چندبعدی تصویر استفاده می‌شود و سپس جستجوی پایگاه داده تصویری از طریق ساختار شاخص‌گذاری برای رسیدن به سرعت مناسب در فرآیند بازیابی انجام می‌شود (Gaede and Gunther, 1998, Markov et al., 2008, Chakrabarti and Mehrotra, 1999). تاکنون تلاش‌های گوناگونی به‌منظور توسعه ساختارهای شاخص‌گذاری چندبعدی انجام گرفته است (Bohm et al., 2001, Gaede and Gunther, 1998, Keyvanpour and Izadpanah, 2011(a), Keyvanpour and Izadpanah, 2011(b), Li et al., 2002)؛ اما با توجه به حجم وسیع تصاویر در پایگاه داده‌های تصویری و ابعاد بالای بردارهای ویژگی تصاویر، ارائه یک ساختار شاخص‌گذاری با کارایی مطلوب بسیار دشوار است.

افزایش ابعاد در حوزه ساختارهای شاخص‌گذاری چندبعدی بر پیچیدگی محاسباتی افزوده و هم‌گرایی روش جستجو را دچار اختلال می‌کند. این مسائل منجر به افزایش هزینه جستجو و در نتیجه کاهش کارایی این ساختارها در پایگاه داده‌های چندبعدی با ابعاد بالایی می‌شود. بنابراین ارائه ساختارهایی برای اداره داده‌های با ابعاد بالا ضروری است. (Keyvanpour and Izadpanah, 2011(b)) از ساختارهای شاخص‌گذاری چندبعدی در فضاهای با ابعاد بالا را ارائه داده است.

تاکنون الگوریتم‌های خوشه‌بندی فراوانی جهت خوشه‌بندی داده‌های چندبعدی ارائه شده‌اند؛ اما تلاش‌های اندکی به‌منظور ارائه روش‌های خوشه‌بندی برای به‌کارگیری آنها، جهت شاخص‌گذاری اطلاعات تصویری صورت گرفته است (Li et al., 2002). محققان الگوریتم‌های خوشه‌بندی را به‌طور کلی به چهار دسته الگوریتم‌های افراز<sup>۳</sup>، الگوریتم‌های سلسله‌مراتبی<sup>۴</sup>، الگوریتم‌های مبتنی بر توری<sup>۵</sup> و الگوریتم‌های

<sup>5</sup> Grid-based algorithms

<sup>6</sup> Density-based algorithms

<sup>7</sup> Divisive

<sup>8</sup> Agglomerative

<sup>9</sup> Hierarchical divisive clustering

<sup>10</sup> Data projection

<sup>11</sup> Negentropy

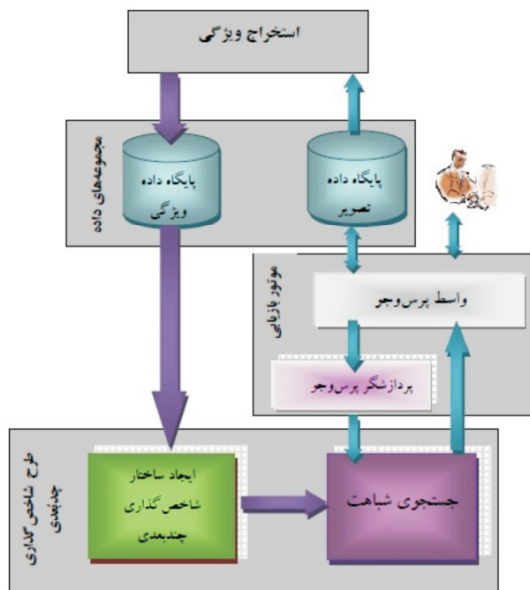
<sup>12</sup> No Overlapping-Non Gaussian Projection based-tree indexing

<sup>1</sup> Similarity measures

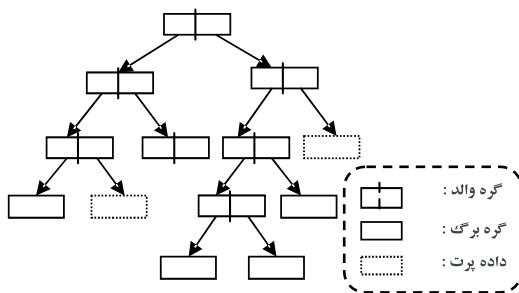
<sup>2</sup> Multi-dimensional indexing structures

<sup>3</sup> Partitioning algorithms

<sup>4</sup> Hierarchical algorithms



شکل (۱-۲): معماری پایه سامانه بازیابی تصویر و جایگاه ساختار شاخص‌گذاری پیشنهادی



شکل (۲-۲): ساختار شاخص‌گذاری پیشنهادی NO-NGP-tree

### ۳- ساختار شاخص‌گذاری چندبعدي پیشنهادی NO-NGP-tree

در سامانه بازیابی تصویر براساس ساختار شاخص‌گذاری پیشنهادی، هر تصویر با یک بردار ویژگی چندبعدي نمایش داده می‌شود و بردارهای ویژگی تصاویر در پایگاه‌داده ویژگی ذخیره می‌شوند. بنابراین پایگاه‌داده ویژگی را می‌توان به صورت ماتریس  $M_0(n, m)$  در نظر گرفت. ستون‌های ماتریس  $M_0$  بردارهای ویژگی می‌باشند.  $n$  برابر تعداد ابعاد بردارهای ویژگی و  $m$  برابر تعداد بردارهای ویژگی درون پایگاه داده ویژگی است. به عبارت دیگر، ماتریس  $M_0$  ماتریس

است. شکل (۱-۲) معماری پایه سامانه بازیابی تصویر و جایگاه ساختار شاخص‌گذاری پیشنهادی را نشان می‌دهد. بردارهای ویژگی تصویر، در بخش استخراج ویژگی، از تصاویر درون پایگاه‌داده تصویری استخراج و در پایگاه‌داده ویژگی ذخیره می‌شوند. کاربر پرس‌وجوی خود را از طریق واسط پرس‌وجو مطرح می‌کند، تصویر مورد جستجو توسط پردازشگر پرس‌وجو به بردار ویژگی چندبعدي تبدیل می‌شود. در قسمت طرح شاخص‌گذاری چندبعدي، در یک مرحله برون خط<sup>۱</sup> کلیه بردارهای ویژگی تصویر از پایگاه‌داده ویژگی گرفته می‌شوند و در ساختار شاخص‌گذاری مبتنی بر خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده NO-NGP-tree سازماندهی می‌شوند؛ سپس در مرحله جستجوی شباهت به صورت برخط<sup>۲</sup>، بردار ویژگی تصویر مورد جستجو گرفته می‌شود و از طریق یک الگوریتم جستجوی شباهت بازگشتی، ساختار شاخص‌گذاری NO-NGP-tree مورد جستجو قرار گرفته و k-نزدیک‌ترین بردارهای ویژگی به بردار پرس‌وجو برای نمایش تصاویر بازیابی شده به واسط پرس‌وجو فرستاده می‌شوند.

در این پژوهش، حاصل بخش شاخص‌گذاری مبتنی بر خوشه‌بندی یک ساختار درختی باینری نامتوازن است (شکل (۲-۲))، که در بالاترین سطح - یا به عبارتی ریشه - شامل یک گره است که کل بردارهای ویژگی تصویر چندبعدي را شامل می‌شود.

در ساختار پیشنهادی، گره‌هایی که دارای دو فرزند هستند گره‌های والد، و گره‌هایی که فرزندی ندارند، در صورتی که تعداد داده‌های درون آن از یک حد آستانه، که با مشخصه پیش‌نیاز Minpts تعیین می‌شود، کمتر باشند، "داده پرت"<sup>۳</sup> و در غیر این صورت، "گره برگ" می‌گویند. به عبارت دیگر، Minpts حداقل تعداد داده‌ها در گره‌های برگ را مشخص می‌نماید. در بخش جستجوی شباهت، از الگوریتم‌های جستجوی نزدیک‌ترین همسایه ارائه‌شده در (Taileb et al., 2007) برای جستجوی شباهت استفاده شده است.

<sup>1</sup>Out-line

<sup>2</sup>On-line

<sup>3</sup>Outlier data

خوشه انتخاب شده در تکرار  $i-1$  هستند و  $i-1$  outlier، زیرخوشه‌ای با تعداد داده کمتر از مقدار تعیین شده توسط پارامتر پیش‌نیاز Minpts در تکرار  $i-1$  است.  $(I_i)_j$  طبق رابطه زیر تعریف می‌شود.

(۲)

$$(I_i)_j = \{F_j, (a_i)_j, \text{IDX1}, \text{IDX2}, \text{CP1}, \text{CP2}\}$$

در رابطه (۲)،  $(I_i)_j$  اطلاعات استخراج شده از  $(M_i)_j$  است، که ماتریس  $(M_i)_j$ ، ماتریس مربوط به عضو  $i$ ام مجموعه  $M_i$  می‌باشد، و  $j$  شمارنده تعداد خوشه‌های موجود در مجموعه  $M_i$  است.  $F_j$  بردار افکنش ماتریس  $M_i$  مربوط به  $i$ امین برگ،  $(a_i)_j$  مؤلفه غیرگوسی معنادار به دست آمده از ماتریس مربوط به  $i$ امین برگ، در مجموعه برگ‌های تکرار  $i$  موجود در مجموعه  $\text{IDX1}, \text{IDX2}, M_i$  برابر زیرخوشه‌های افکنش  $\text{CP1}, \text{CP2}$  برابر مراکز زیرخوشه‌های افکنش هستند.

لازم به ذکر است که در اولین تکرار الگوریتم ( $i=1$ )، ماتریس  $M_0$  حاوی کل بردارهای ویژگی تصویر، به‌عنوان ورودی در نظر گرفته می‌شود. مرحله پیش‌افزاندی، عملیات لازم را برای به دست آوردن اطلاعات  $I_i$  مورد نیاز برای انتخاب و تجزیه خوشه بعدی به دو زیرخوشه در مرحله افزاندی فراهم می‌کند.  $I_i$  اطلاعات استخراج شده از مجموعه  $M_i$  را در برمی‌گیرد. سپس در مرحله افزاندی، خوشه بعدی برای افزاندی از مجموعه  $M_i$  انتخاب می‌شود و به دو زیرخوشه راست  $(MR_i)_s$  و زیرخوشه چپ  $(ML_i)_s$  تجزیه می‌شود.

از آنجا که افزاندی براساس افکنش داده‌ها روی مؤلفه غیرگوسی معنادار به دست آمده از ماتریس مربوط به خوشه انتخاب شده  $(a_i)_s$  صورت می‌گیرد، که در مرحله پیش‌افزاندی به دست آمده، مؤلفه  $(a_i)_s$  به مرحله مرزبندی زیرخوشه‌ها فرستاده می‌شود. در مرحله مرزبندی، زیرخوشه‌های راست و چپ توسط مستطیل‌های مرزبندی کمینه‌ای که در جهت مؤلفه  $(a_i)_s$  گسترش می‌یابند، مرزبندی می‌شوند. نتایج مرحله مرزبندی  $(BR_i)_s, (BL_i)_s$  هستند، که به ترتیب شامل اطلاعات مربوط به مرزبندی زیرخوشه راست و زیرخوشه چپ خوشه انتخاب شده در تکرار  $i$  می‌باشند؛ سپس در مرحله تشکیل ساختار درختی، دو زیرخوشه راست و چپ ایجاد شده در تکرار  $i$  بررسی شده و در صورتی که تعداد اعضای آن خوشه‌ها از حد آستانه تعیین شده کمتر باشد، به‌عنوان داده‌های پرت و در غیر این صورت به‌عنوان گره برگ علامت زده می‌شوند. در هر دو

<sup>2</sup>Matrix projection vector

حاوی بردارهای ویژگی چندبعدی استخراج شده از تصاویر (ماتریس مربوط به ریشه) است. برای ایجاد ساختار شاخص گذاری NO-NGP-tree، ماتریس  $M_0$  با استفاده از یک الگوریتم تکرار کننده<sup>۱</sup> به صورت بازگشتی افزاندی می‌شود. در هر تکرار یک ساختار باینری ناتمام تشکیل می‌شود. روال تکرار کننده تا زمانی که تعداد گره برگ مورد نظر ( $k$ ) ایجاد شود، ادامه پیدا می‌کند.

شبه‌کد مربوط به فرآیند ایجاد ساختار شاخص گذاری چندبعدی پیشنهادی در ادامه، و فرآیند ایجاد ساختار شاخص گذاری چندبعدی پیشنهادی در شکل (۱-۳) نشان داده شده‌اند. همان‌طور که در شکل (۱-۳) مشخص است، فرآیند ایجاد ساختار شاخص گذاری چندبعدی پیشنهادی به صورت تکرار کننده از چهار مرحله اصلی شامل: پیش‌افزاندی، افزاندی، مرزبندی و ایجاد ساختار درختی تشکیل شده است. این مراحل در شبه‌کد زیر، به ترتیب با توابع  $\text{Partitioning}(), \text{Pre-Partitioning}(), \text{Bouding}(), \text{Build tree structure}()$  نشان داده شده‌اند. این مراحل تا زمانی که تعداد برگ تعیین شده توسط مشخصه پیش‌نیاز  $k$  ایجاد شود، به صورت بازگشتی تکرار می‌شوند.

```

NO-NGP-tree ( $M_i, K, \text{Minpts}$ )
Begin
 $i=1$ 
while  $LN_0 \leq k$ 
{
 $\{I_i\} = \text{Pre-partitioning}(M_i)$ 
 $\{(MR_i)_s, (ML_i)_s\} = \text{Partitioning}(M_i, I_i)$ 
 $\{(BR_i)_s, (BL_i)_s, (MR_i)_s, (ML_i)_s\} = \text{Bouding}((MR_i)_s, (ML_i)_s, (a_i)_s)$ 
 $\{TS_i\} = \text{Build tree structure}((BR_i)_s, (BL_i)_s, (MR_i)_s, (ML_i)_s)$ 
 $i=i+1$ 
}
Return (NO-NGP-tree)
End
    
```

ماتریس  $M_i$  مجموعه ماتریس‌های مربوط به برگ‌ها در تکرار  $i$ ام است. به عبارت دیگر، در هر تکرار، مجموعه ماتریس‌های مربوط به برگ‌ها در تکرار  $i$  ( $M_i$ ) دریافت می‌شود.  $M_i$  طبق رابطه زیر تعریف می‌شود.

(۱)

$$M_i = \{M_i - (M_{i-1}) \cup (MR_{i-1}) \cup (ML_{i-1}) - \text{outlier}_{i-1}\}$$

در رابطه (۱)،  $M_{i-1}$  حاوی مجموعه برگ‌ها در تکرار  $i-1$  ماتریس مربوط به برگ انتخاب شده در تکرار  $i-1$ ،  $(MR_{i-1})_s$  و  $(ML_{i-1})_s$  به ترتیب زیرخوشه‌های راست و چپ

<sup>1</sup>Iterative

به‌طوری که این مؤلفه در جهتی قرار گرفته باشد که بتواند به‌شکل مطلوبی ساختار داده‌ها و یا به‌عبارتی خوشه‌های طبیعی موجود در توزیع داده‌ها را نشان دهد.

افکنش: هدف از این مرحله، افکنش داده‌های ماتریس  $(M_i)_j$  روی مؤلفه به‌دست‌آمده در مرحله یافتن مؤلفه غیرگوسی معنادار از داده‌هاست.

خوشه‌بندی داده‌های افکنش‌شده: هدف این مرحله، به‌کارگیری الگوریتم خوشه‌بندی 2-means روی داده‌های افکنش‌شده در مرحله افکنش، برای به‌دست‌آوردن دو زیرخوشه به نام خوشه‌های افکنش و مراکز ثقل آنهاست.

```

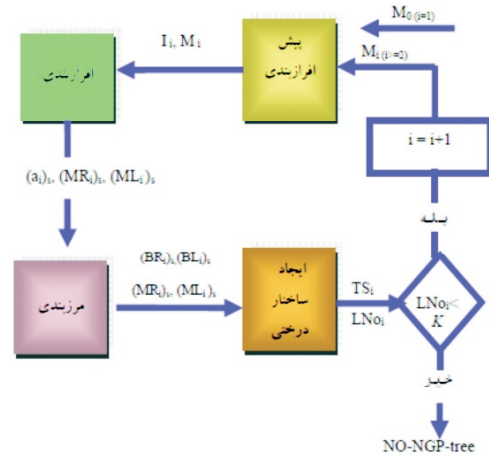
Pre-partitioning ( $M_i$ )
Begin
For  $j=1$  to  $LNO_i$ 
{
 $\{(a_j)_i\}$  = Find meaningful non-gaussian component ( $(M_i)_j$ )
 $\{F_j\}$  = Projection ( $\{(a_j)_i\}, (M_i)_j$ ):
 $\{Cp1, Cp2, IDX1, IDX2\}$  = 2-mean clustering ( $F_j$ )
}
Return ( $I_i$ )
End
    
```

### ۳-۱-۱- یافتن مؤلفه غیرگوسی معنادار

برای ایجاد ساختار پیشنهادی، در هر تکرار، ماتریس مربوط به برگ انتخاب‌شده (از برگ‌های درخت باینری تشکیل‌شده تا تکرار مربوطه) در مرحله افرازبندی، به دو زیرخوشه تجزیه می‌شود. بدین منظور، داده‌های ماتریس مربوط به برگ انتخاب‌شده روی یک مؤلفه از داده‌ها افکنش شده و سپس در مرحله افرازبندی از نقطه‌ای روی داده‌های افکنش‌شده به دو زیرخوشه تجزیه می‌شوند.

افکنش داده‌ها روی یک مؤلفه از داده‌ها ایده جدیدی نیست و در قبل در الگوریتم خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده PDDP، برای خوشه‌بندی داده‌های با ابعاد بالا استفاده شده است. در PDDP، داده‌ها روی اولین مؤلفه اصلی افکنش‌شده و بعد با ابرصفحه‌گذرنده از مرکز ثقل داده‌ها و عمودبر اولین مؤلفه اصلی داده‌ها به دو زیرخوشه تجزیه می‌شوند. این عملیات به‌صورت بازگشتی تکرار می‌شود تا تعداد خوشه موردنظر حاصل شود (Boley, 1998). درحقیقت در PDDP داده‌ها با استفاده از تبدیل خطی PCA<sup>۲</sup> (Jackson, 1991, Jolliffe, 2002) روی اولین مؤلفه اصلی افکنش می‌شوند. اولین مؤلفه اصلی، مؤلفه‌ای از داده‌هاست که امتداد آن منطبق با بیشترین پراکندگی قابل مشاهده در توزیع داده‌هاست. افکنش داده‌ها روی اولین

حالت  $(BR_i)_s$  و  $(BL_i)_s$  ذخیره می‌شوند؛ سپس،  $TS_i$  که نشان‌دهنده درخت باینری تشکیل‌شده در تکرار  $i$  است، به‌همراه  $LNO_i$  نتایج این مرحله را تشکیل می‌دهند.  $LNO_i$  تعداد برگ‌های ساختار درختی باینری تشکیل‌شده در تکرار  $i$  است.



شکل (۳-۲): فرآیند ایجاد ساختار شاخص‌گذاری چندبعدی پیشنهادی NO-NGP-tree

### ۳-۱-۲- پیش‌افرازبندی

در این مرحله، اطلاعات  $I_i$  شامل اطلاعات مورد نیاز برای انتخاب و تجزیه خوشه بعدی به دو زیرخوشه در مرحله افرازبندی فراهم می‌شود. مرحله پیش‌افرازبندی، مجموعه  $M_i$  شامل ماتریس‌های مربوط به برگ‌های موجود در تکرار  $i-1$  را می‌گیرد و برای هریک از اعضای مجموعه  $M_i$  طی سه مرحله، اطلاعات  $(I_i)_j$  به‌دست می‌آید و اطلاعات نهایی برای  $j=1$  تا  $j=LNO_i$  در مجموعه  $I_i$  قرار می‌گیرند و همراه مجموعه  $M_i$  به مرحله افرازبندی فرستاده می‌شوند (ژشمارنده تعداد برگ‌ها در  $M_i$  است، بنابراین  $(M_i)_j$  برابر زامین برگ در مجموعه  $M_i$  می‌باشد).

فرآیند پیش‌افرازبندی از سه زیرمرحله که برای  $j=1$  تا  $j=LNO_i$  تکرار می‌شوند، تشکیل شده است. در این شبه‌کد، زیرمراحل به ترتیب با توابع  $\text{Find meaningful non-gaussian component}$ ،  $\text{Projection}$ ،  $\text{gaussian component}$  و  $\text{2-mean clustering}$  نشان داده شده‌اند.

یافتن مؤلفه غیرگوسی معنادار: هدف این مرحله، پیدا کردن مؤلفه  $(a_i)_j$  از داده‌های ماتریس  $(M_i)_j$  است،

<sup>1</sup> First principal component  
<sup>2</sup> Principal component analysis (PCA)



نشان داده است که توزیع گوسین<sup>۴</sup> کمترین معنا را در توزیع داده‌ها نشان می‌دهد و جهت‌های مطلوب معنادار، جهت‌هایی هستند که از توزیع گوسین فاصله دارند.

بهینه‌سازی محلی تابع هدف برای به‌دست‌آوردن مؤلفه‌های غیرگوسی، به‌طور دقیق در تحلیل مؤلفه‌های مستقل<sup>۵</sup>، که یکی از روش‌های پیگرد افکنش محسوب می‌شود نیز انجام می‌شود (Hyvarinen and Oja, 2000). در (Hyvarinen and Oja, 1997) با استفاده از تبدیل قوانین یادگیری شبکه عصبی<sup>۶</sup> به تکرار ممیز ثابت<sup>۷</sup>، الگوریتمی به نام FastICA ارائه شده که در اصل برای پیدا کردن مؤلفه‌های مستقل پیشنهاد شده، اما درحقیقت یک روش پیگرد افکنش است که می‌توان از آن برای پیدا کردن مؤلفه‌های غیرگوسی استفاده کرد. این الگوریتم هر مؤلفه را تک‌به‌تک پیدا می‌کند، بنابراین می‌توان از آن برای پیدا کردن تنها یک مؤلفه استفاده کرد. تابع هدف در این الگوریتم می‌تواند هر یک از توابع هدفی که در تحقیقات برای به‌دست‌آوردن مؤلفه‌های پیگرد افکنش استفاده شده، باشد. در این پژوهش از آن‌جایی که هدف به‌دست‌آوردن مؤلفه‌ای است که بتواند ساختار خوشه‌ای داده‌ها را نشان دهد، از نسخه‌ای از الگوریتم FastICA که در آن تقریب آشفتگی منفی به‌عنوان تابع هدف در نظر گرفته می‌شود، استفاده شده است (Hyvarinen and Oja, 2000).

از آن‌جایی که متغیرهای گوسین بیشترین آشفتگی<sup>۸</sup> را در میان متغیرهای با واریانس یکسان دارند، بنابراین آشفتگی می‌تواند یک معیار اندازه‌گیری عدم‌گوسین بودن قرار داده شود. درحقیقت، این امر نشان می‌دهد که توزیع گوسین کمترین ساخت‌یافتگی را در میان انواع توزیع داده‌ها دارد. در تحقیقات برای به‌دست‌آوردن معیار مناسبی برای اندازه‌گیری میزان عدم‌گوسین از آشفتگی منفی استفاده می‌شود. آشفتگی منفی<sup>۹</sup> در رابطه زیر تعریف می‌شود (Hyvarinen and Oja, 2000).

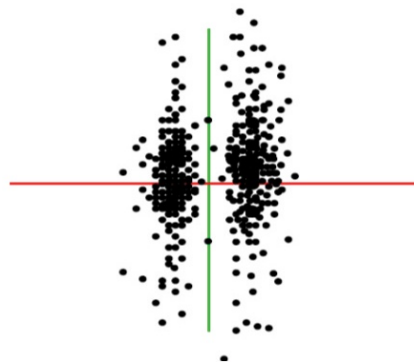
$$J(y) = H(y_{gauss}) - H(y) \quad (3)$$

در رابطه (۳)،  $y_{gauss}$  متغیر گوسین با واریانس یکسان و  $H(y)$  برابر آشفتگی متغیر گسسته  $y$  است. از آن‌جا که

مؤلفه اصلی منجر به تولید خوشه‌های متراکم می‌شود؛ اما در صورتی که توزیع داده‌ها به‌گونه‌ای باشد که خوشه‌های طبیعی آن در راستای اولین مؤلفه اصلی از هم جدا نباشند، اولین مؤلفه اصلی نمی‌تواند حاوی اطلاعات مفیدی برای تجزیه داده‌ها باشد.

در این پژوهش از یک روش پیگرد افکنش<sup>۱</sup> (الگوریتم FastICA) با هدف پیدا کردن یک مؤلفه غیرگوسی معنادار از داده‌ها استفاده شده است. افکنش داده‌ها روی این مؤلفه غیرگوسی معنادار تقریب آشفتگی منفی را بیشینه می‌کند و در نتیجه افکنش داده‌ها روی آن برای تجزیه داده‌ها به دو زیرخوشه مفید خواهد بود. نمونه‌ای از مطلوب بودن مؤلفه‌هایی با بیشترین خاصیت غیرگوسی (مؤلفه حاصل از به‌کارگیری روش‌های پیگرد افکنش) نسبت به مؤلفه‌های اصلی، از نظر توانایی در نمایش توزیع معنادار داده‌ها در شکل (۲-۳) نشان داده شده است.

همان‌طور که در شکل (۲-۳) مشاهده می‌شود، اولین مؤلفه اصلی که با رنگ سبز (در جهت عمودی) نشان داده شده است، ساختار خوشه‌ای داده‌ها را نشان نمی‌دهد؛ اما مؤلفه غیرگوسی حاصل از به‌کارگیری روش‌های پیگرد افکنش که با رنگ قرمز (در جهت افقی) نشان داده شده است، جداسازی بین خوشه‌ها را نشان می‌دهد و بنابراین ساختار داده‌ها را مشخص می‌کند. گوسین بودن مؤلفه‌ها سادگی محاسباتی و کاهش میزان نیاز به اطلاعات اضافی را در بر دارد.



شکل (۲-۳): مطلوب بودن مؤلفه غیرگوسین نسبت به اولین مؤلفه اصلی از نظر توانایی در نمایش توزیع معنادار داده‌ها

ایده اصلی روش‌های پیگرد افکنش، استخراج مؤلفه‌های معنادار توسط بهینه‌سازی محلی تابع شاخص افکنش<sup>۲</sup> (تابع هدف<sup>۳</sup>) مورد نظر است. تحقیقات در این زمینه

<sup>1</sup> Projection pursuit  
<sup>2</sup> Projection index

<sup>3</sup> Objective function

<sup>4</sup> Gauss distribution

<sup>5</sup> Independent component analysis (ICA)

<sup>6</sup> Neural network

<sup>7</sup> Fixed point iteration

<sup>8</sup> Entropy



برآورد آشفتنگی منفی دشوار است، در مرجع (Hyvarinen and Oja, 2000) در عمل از تقریب آشفتنگی منفی به‌عنوان تابع هدف استفاده شده است. تقریب آشفتنگی منفی در رابطه زیر تعریف شده است.

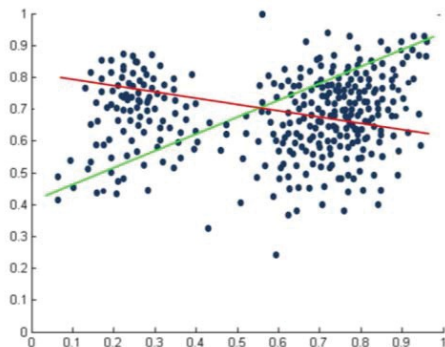
$$J(y) \cong [E\{G(y)\} - E\{G'(V)\}]^2 \quad (4)$$

در رابطه (4)، تابع E امید ریاضی و V متغیر گوسی با میانگین صفر و واریانس یک است. تابع G به صورت رابطه (5) تعریف شده و G' برابر مشتق تابع G می‌باشد.  $1 \leq c \leq 2$  ثابت است که در عمل برابر یک در نظر گرفته می‌شود.

$$G(u) = \tanh(cu) \quad (5)$$

مکان‌هایی که تراکم داده‌ها در آنها کمتر است، مناسبی برای تجزیه خوشه مربوطه است. مکان‌هایی که تراکم داده‌ها در آنها بیشتر است، نزدیک به مرکز ثقل زیرخوشه‌های اصلی فرضی است.

در شکل (3-3) مثالی از خوشه‌های نامتوازن غیرمجزا در فضای دوبعدی نشان داده شده است. در این شکل، ساختار خوشه‌ای داده‌ها که از خوشه‌های نامتوازن و غیرمجزا تشکیل شده مشاهده می‌شود. همچنین، اولین مؤلفه اصلی با رنگ سبز و مؤلفه غیرگوسی معنادار با رنگ قرمز نشان داده شده است.



شکل (3-3): مثالی از خوشه‌های نامتوازن و غیرمجزا در فضای دوبعدی

در این پژوهش از الگوریتم FastICA برای به‌دست آوردن مؤلفه غیرگوسی معنادار از داده‌ها استفاده شده است. این الگوریتم به صورتی تنظیم شده که تنها یک مؤلفه غیرگوسی را به‌دست می‌آورد. در الگوریتم FastICA اصلی از یک بردار وزن تصادفی<sup>1</sup> برای آغاز الگوریتم استفاده شده است، اما از آن‌جاکه در این پژوهش به دنبال تنها یک مؤلفه غیرگوسی می‌باشیم و می‌خواهیم آن مؤلفه بر جهت پراهمیتی از داده‌ها با تراکم داده‌ای بالا و توزیع با خاصیت عدم گوسین بالا قرار گیرد، از اولین مؤلفه اصلی، که در پیش‌پردازش الگوریتم FastICA استفاده می‌شود، به‌عنوان بردار وزن اولیه استفاده می‌کنیم.

### 3-1-2-2- افکنش

هدف از مرحله افکنش، افکنش داده‌های ماتریس  $(M_i)$  روی مؤلفه  $(a_i)$  به‌دست آمده در مرحله قبل است. هدف از افکنش داده‌ها روی چنین مؤلفه‌ای، تبدیل داده‌های چندبعدی با استفاده از افکنش آنها روی یک مؤلفه به داده‌های تک‌بعدی است؛ سپس داده‌ها از محلی بر داده‌های افکنش شده و توسط ابرصفحه‌ای عمود بر آن، به دو زیرخوشه تجزیه می‌شوند. افکنش داده‌ها روی مؤلفه  $(a_i)$  طبق رابطه زیر صورت می‌گیرد.

$$F_j(x_b) = (a_i)_j^T (x_b) \quad (6)$$

در رابطه (6)،  $x_b$  بردار ویژگی  $(b=1..nodesize)$  م در ماتریس  $(M_i)$ ، علامت ترانهاده و  $F_j(x_b)$  حاوی افکنش‌های تک‌بعدی بردارهای ویژگی چندبعدی  $x_b$  است.

### 3-1-3- خوشه‌بندی داده‌های افکنش شده

هدف از این مرحله به‌کارگیری یک الگوریتم خوشه‌بندی مانند k-means (MacQueen, 1967) روی داده‌های افکنش شده  $F_j$  به‌دست آمده در مرحله قبل، برای به‌دست آوردن دو زیرخوشه افکنش (IDX1, IDX2) و مراکز ثقل آنها (CP1, CP2) است.

افکنش داده‌ها روی مؤلفه  $(a_i)$  می‌تواند محلی با تراکم کمتر و محلی با تراکم بالا را نشان دهد. در این پژوهش برای این که محلی‌های مورد نظر حدس زده شود از

<sup>2</sup>Scalar

<sup>1</sup> Weight vector

مرحله است که در شبه‌کد زیر به‌ترتیب با توابع Selection() و Split() نشان داده شده‌اند:

انتخاب خوشه: هدف از این مرحله، انتخاب یکی از برگ‌ها است که داده‌های ماتریس مربوط به آن توزیع ساخت‌یافته‌تری را دارد.

تجزیه: هدف این مرحله، تجزیه خوشه انتخاب‌شده به دو زیرخوشه است.

```

Partitioning (Mi, Ii):
Begin
(Mi)s = Selection (Mi, Ii)
{(MRi)s, (MLi)s} = Split ((Mi)s, Ii)
Returns ((MRi)s, (MLi)s)
End
    
```

### ۳-۲-۱- انتخاب خوشه

در این مرحله، خوشه‌ای که باید در آن تکرار به دو زیرخوشه تجزیه شود، از میان برگ‌هایی که تا آن تکرار ایجاد شده‌اند، انتخاب می‌شود. در مرحله انتخاب خوشه، براساس اطلاعات مجموعه‌های  $M_i$  و  $I_i$  خوشه مربوط به یکی از برگ‌ها انتخاب و در مرحله تجزیه به دو زیرخوشه راست و چپ تجزیه می‌شود. در این پژوهش از اطلاعات حاصل از افکنش داده‌های مربوط به زیرخوشه‌ها در مرحله پیش‌افرازبندی استفاده شده و معیاری برای انتخاب برگ‌ها که حاوی داده‌هایی با بیشترین ساخت‌یافتگی باشد ارائه شده است.

انتخاب خوشه بعدی در هر تکرار یکی از چالش‌های عمده روش‌های سلسله‌مراتبی تقسیم‌کننده به حساب می‌آید؛ زیرا در این مرحله هدف پیدا کردن خوشه‌ای است که بیشترین ساخت‌یافتگی را داشته باشد. در اکثر تحقیقات پیشین مانند PDDP، خوشه‌ای که بیشترین پراکندگی از مرکز ثقل را داشته باشد به‌عنوان خوشه مورد تجزیه در هر تکرار انتخاب می‌شود. رابطه (۷)، پراکندگی<sup>۱</sup> از مرکز ثقل را نشان می‌دهد.

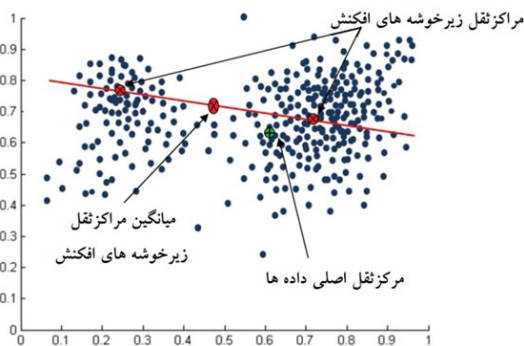
$$Scattvalue = \frac{1}{N} \sum_{i=1}^N \|x_i - w\|^2 \quad (7)$$

در این رابطه، Scattvalue مقدار پراکندگی از مرکز ثقل،  $N$  تعداد نقاط داده‌ای خوشه مربوطه،  $w$  مرکز ثقل خوشه و  $\| \cdot \|$  نرم اقلیدسی است.

اما همان‌طور که در مرجع (Savaresi et al., 2002) مطرح شده، پراکندگی از مرکز ثقل خوشه نمی‌تواند معیار خوبی برای نشان دادن ساختار خوشه‌بندی‌شده داده‌ها باشد.

<sup>۱</sup> Scatter

خوشه‌بندی k-means با در نظر گرفتن  $k=2$  لروی بردار  $F_j$  حاوی داده‌های تک‌بعدی افکنش، استفاده شده است. به این ترتیب که دو مرکز ثقل به‌دست‌آمده همان محل‌های متراکم داده‌ها و میانگین این مراکز ثقل به‌عنوان محلی که به‌طور معمول تراکم داده‌ای کمتری دارد در نظر گرفته می‌شود. اگر ساختار خوشه‌های طبیعی داده‌ها به‌گونه‌ای باشد که خوشه‌ها متوازن و به‌طور کامل مجزا باشند، محل میانگین زیرخوشه‌های افکنش به احتمال زیاد در محل مرکز ثقل داده‌ها یا نزدیک به آن خواهد بود؛ اما اگر خوشه‌ها نامتوازن و غیرمجزا باشند، محل میانگین مراکز ثقل خوشه‌های افکنش، محل کم‌تراکم‌تری را نسبت به محل مرکز ثقل اصلی داده‌ها نشان می‌دهند. در شکل (۳-۴)، مرکز ثقل اصلی داده‌ها و میانگین خوشه‌های افکنش به‌همراه مؤلفه غیرگوسی معنادار نشان داده شده است. همان‌طور که مشاهده می‌شود، محل میانگین مراکز ثقل خوشه‌های افکنش، که روی مؤلفه غیرگوسی معنادار قرار دارد، محل کم‌تراکم‌تری را نسبت به مرکز ثقل اصلی داده‌ها نشان می‌دهد.



شکل (۳-۴): مقایسه مرکز ثقل اصلی داده‌ها و میانگین مراکز ثقل خوشه‌های افکنش در مثال شکل (۳-۳)

### ۳-۲-۲- افرازبندی

در مرحله افرازبندی، که شبه‌کد آن در ادامه آمده است، اطلاعات  $I_i$  به‌دست‌آمده از مرحله پیش‌افرازبندی و مجموعه  $M_i$  دریافت می‌شود.  $M_i$  حاوی برگ‌های تشکیل‌شده تا تکرار است. در مرحله افرازبندی ابتدا با استفاده از اطلاعاتی که در مرحله پیش‌افرازبندی به‌دست‌آمده یک معیار انتخاب برای گزینش یکی از برگ‌های تولیدشده تا تکرار  $i$  ارائه، سپس یکی از برگ‌ها براساس معیار ارائه‌شده انتخاب و خوشه مربوط به آن براساس اطلاعات به‌دست‌آمده از افکنش داده‌ها به دو زیرخوشه تجزیه می‌شود. مرحله افرازبندی شامل دو



محل میانگین مراکز ثقل زیرخوشه‌های افکنش، توسط ابرصفحه‌ای عمود بر  $(a_i)_s$  به دو زیرخوشه راست و چپ تجزیه می‌شود و  $(a_i)_s$  برای هر دو زیرخوشه راست و چپ ذخیره می‌شود.

دلیل تجزیه خوشه از محل میانگین، در بخش ۳-۱-۳ بیان شد. شکل (۳-۳) ابرصفحه جداکننده<sup>۱</sup> را در مورد شکل (۳-۳) نشان می‌دهد. در مورد داده‌هایی که توزیع داده‌ها در آنها به‌گونه‌ای است که از خوشه‌های غیرمتوازن و غیرمجزا تشکیل شده‌اند، انتخاب مرکز ثقل داده‌ها، به‌عنوان محلی که ابرصفحه جداکننده از آن بگذرد، انتخاب جالبی نیست؛ زیرا خوشه‌های طبیعی را قطع می‌کند (بخش الف) شکل (۳-۳). اما میانگین مراکز ثقل زیرخوشه‌های افکنش، محل مناسب‌تری است. تجزیه داده‌ها توسط ابرصفحه عمود بر مؤلفه غیرگوسی معنادار و گذرنده از میانگین مراکز ثقل زیرخوشه‌های افکنش، داده‌ها را به دو زیرخوشه فشرده تجزیه می‌کند که در نهایت پس از مرزبندی، مستطیل‌های کمینه فشرده‌تری را منجر می‌شوند (بخش ب) شکل (۳-۳).

تجزیه خوشه توسط ابرصفحه جداکننده عمود بر مؤلفه غیرگوسی معنادار و گذرنده از میانگین مرکز ثقل زیرخوشه‌های افکنش، نقاط داده‌ای درون خوشه انتخاب شده برای تجزیه، که افکنش آنها در راست محل میانگین مرکز ثقل‌های افکنش باشد به زیرخوشه راست نسبت داده شده و در ماتریس  $(MR_i)_s$  قرار داده می‌شوند. نقاط داده‌ای که افکنش متناظر با آنها در چپ قرار دارد به زیرخوشه چپ نسبت داده شده و در  $(ML_i)_s$  قرار داده می‌شوند. درحقیقت، برای تمام  $x_b$ ها، رابطه (۱۰) اعمال و در صورتی که مقدار به‌دست‌آمده مثبت باشد،  $x_b$  در  $(MR_i)_s$  و در غیر این صورت در  $(ML_i)_s$  قرار داده می‌شود.

$$F_s(x_b) - (a_i)_s^T c_{mean} \quad (10)$$

در رابطه (۱۰)،  $x_b$  بردار داده  $b$ ام،  $F_s(x_b)$  افکنش داده  $b$ ام در خوشه انتخاب شده و  $c_{mean}$  میانگین مراکز ثقل زیرخوشه‌های افکنش است.

بنابراین، در این پژوهش از اطلاعات  $(I_i)_z$  به‌دست‌آمده از افکنش داده‌های مربوط به هر برگ در تکرار  $\alpha$ م، روی مؤلفه غیرگوسی معنادار و زیرخوشه‌ها و مراکز ثقل حاصل از به‌کارگیری خوشه‌بندی 2-means روی افکنش‌ها استفاده شده و معیاری برای انتخاب خوشه (تعریف شده در رابطه (۸)) ارائه شده است. هرچه مقدار به‌دست‌آمده از معیار ارائه شده برای خوشه‌ای بزرگ‌تر باشد، ساختار داده‌های مربوط به آن خوشه ساختار یافته‌تر است. بنابراین در تکرار  $\alpha$ م روش پیشنهادی، در زیرمرحله انتخاب خوشه، خوشه‌ای که این معیار را بیشینه نماید از میان برگ‌های موجود در مجموعه  $M_i$  انتخاب می‌شود.

$$selvalue = \left| \frac{CPI - CP2}{\max_{c=1,2}(diameter(IDX_c))} \right| \quad (8)$$

در رابطه (۸)  $selvalue$  مقدار معیار انتخاب،  $CPI, CP2$  مراکز ثقل زیرخوشه‌های افکنش مربوط به برگ مورد بررسی و  $diameter(IDX_c)$  قطر خوشه افکنش  $c$ ام ( $c=1,2$ ) را نشان می‌دهد. قطر یک خوشه افکنش از رابطه زیر به‌دست می‌آید:

$$diameter(IDX) = \max(F_p) - \min(F_p) \quad (9)$$

در رابطه (۹)،  $IDX$  خوشه افکنش و  $F_p$  مقدار داده افکنش  $p=1..psize$  را در زیرخوشه افکنش نشان می‌دهد، طوری که  $psize$  برابر تعداد افکنش‌هاست.

ایده معیار انتخاب ارائه شده این است که ما اعتقاد داریم اگر زیرخوشه‌های افکنش که از خوشه‌بندی داده‌های افکنش شده روی مؤلفه غیرگوسی معنادار به‌دست می‌آیند استقلال بیشتری داشته باشند، یا به عبارتی فاصله مرکز ثقل‌های آنها از هم بیشتر و قطر آنها کوچک‌تر باشد، زیرخوشه‌های طبیعی در ساختار داده نیز استقلال بیشتری دارند؛ بنابراین خوشه مورد نظر ساختار خوشه‌بندی شده‌تری را نشان می‌دهد. شکل (۳-۵) مثالی را در فضای دوبعدی نشان می‌دهد که در آن نقاط داده‌ای، افکنش‌ها و مراکز ثقل زیرخوشه‌های افکنش نشان داده شده است.

### ۳-۲-۲- تجزیه

در مرحله تجزیه، ماتریس مربوط به برگ انتخاب شده در تکرار  $\alpha$  یا  $(M_i)_s$  به همراه مؤلفه غیرگوسی معنادار  $(a_i)_s$  و همچنین بردار حاوی داده‌های افکنش شده برگ انتخاب شده از  $(F_s)$  گرفته می‌شود و زیرخوشه مربوط به برگ انتخاب شده از

<sup>1</sup> Separating hyper-plane

### ۳-۳- مرزبندی

همان‌طور که در شکل (۳-۱) فرآیند ایجاد ساختار شاخص‌گذاری چندبعدی پیشنهادی نشان داده شد، در هر تکرار، بعد از مرحله افزایشی، زیرخوشه‌های راست و چپ تولید شده در زیرمرحله افزایشی، همراه  $(a_i)_s$  به مرحله مرزبندی زیرخوشه‌ها فرستاده می‌شوند تا در این مرحله، توسط مستطیل‌های مرزبندی کمینه گسترده شده در امتداد  $(a_i)_s$  مرزبندی شوند. شبه‌کد مربوط به مرزبندی در زیر آمده است.

```

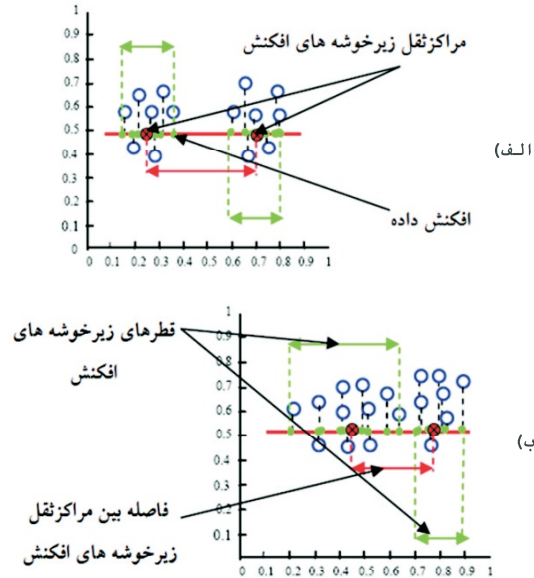
Bouding ((MRi)s, (MLi)s, (ai)s)
Begin
{(ai)s, (NRi)s, (NLi)s} = Change-reference mark ((MRi)s,
(MLi)s, (ai)s)
{(BRi)s, (BLi)s} = Minimum bounding rectangle
construction ((NRi)s, (NLi)s)
Returns((BRi)s, (BLi)s, (MRi)s, (MLi)s)
End
    
```

مرحله مرزبندی شامل دو مرحله است که در شبه‌کد به ترتیب با تابع  $\text{Change-reference mark}()$  و  $\text{Minimum bounding rectangle construction}()$  مشخص می‌شوند: تغییر مبدأ مختصات: هدف از این مرحله، تغییر مبدأ مختصات در جهت  $(a_i)_s$  برای جلوگیری از هم‌پوشانی مستطیل‌های مرزبندی کمینه است. تشکیل مستطیل‌های مرزبندی کمینه: هدف مرزبندی زیرخوشه‌های راست و چپ در مبدأ مختصات جدید است.

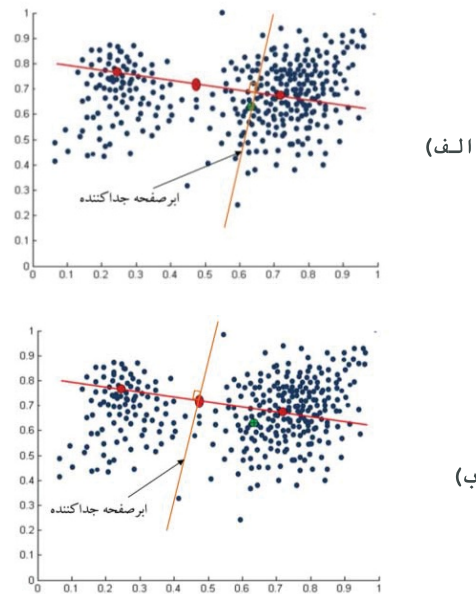
### ۳-۳-۱- تغییر مبدأ مختصات

همان‌طور که در قبل توضیح داده شد، برای جلوگیری از هم‌پوشانی بین گره‌ها در ساختار شاخص‌گذاری پیشنهادی، در این مرحله بعد از تولید زیرخوشه‌های راست و چپ در هر تکرار، ماتریس داده‌های مربوط به هریک از زیرخوشه‌های راست و چپ در مختصات جدید، که در امتداد  $(a_i)_s$  قرار داده شده است، به دست می‌آیند.

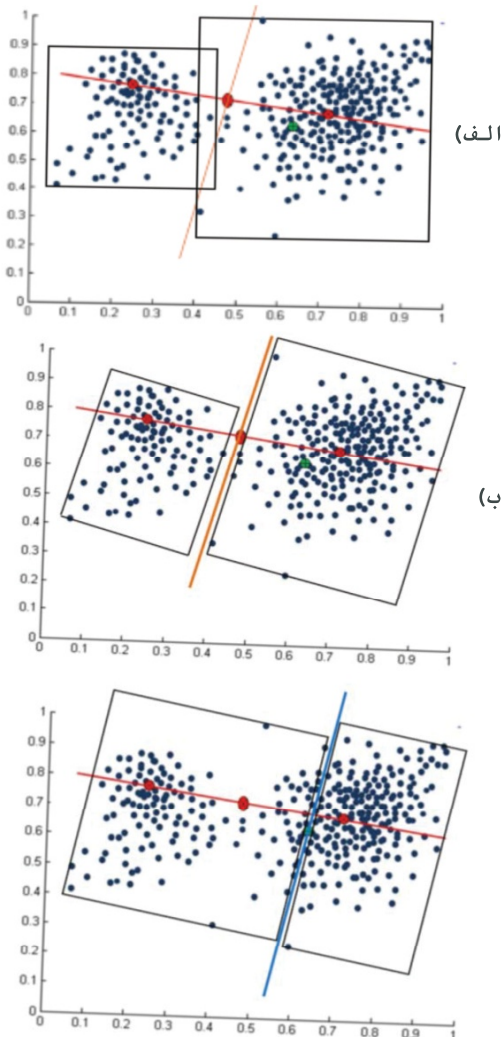
مسئله تغییر جهت مبدأ مختصات برای جلوگیری از هم‌پوشانی گره‌ها در ساختارهای شاخص‌گذاری ایده جدیدی نیست. پیش‌تر در NOHIS-tree برای جلوگیری از هم‌پوشانی مستطیل‌های مرزبندی کمینه، از تغییر جهت مبدأ مختصات و سپس تشکیل مستطیل‌های مرزبندی کمینه در آن جهت استفاده شده است (Taileb et al., 2007).



شکل (۳-۵): مثالی از داده‌های افکنش شده روی مؤلفه غیرگوسی معنادار در فضای دوبعدی الف) خوشه‌ای که در آن داده‌ها دارای توزیع غیرساخت یافته می‌باشند، مقدار selvalue=2.8 خوشه‌ای که در آن داده‌ها دارای توزیع ساخت یافته می‌باشند، مقدار selvalue=1.5



شکل (۳-۶): تجزیه توسط ابرفسحه جداکننده. الف) ابرفسحه جداکننده عمود بر مؤلفه غیرگوسی معنادار و گذرنده از مرکز ثقل اصلی داده‌ها. ب) ابرفسحه جداکننده عمود بر مؤلفه غیرگوسی معنادار و گذرنده از میانگین مرکز ثقل زیرخوشه‌های افکنش



(شکل ۳-۷): تشکیل مستطیل‌های مرزبندی کمینه. الف)

تشکیل مستطیل‌های مرزبندی کمینه در راستای مبدأ اصلی. ب)  
تشکیل مستطیل‌های مرزبندی کمینه در راستای مبدأ جدید. پ)  
تجزیه داده‌ها از محل مرکز ثقل اصلی داده‌ها و تشکیل  
مستطیل‌های مرزبندی کمینه در راستای مبدأ جدید

### ۳-۴ ایجاد ساختار درختی

در این مرحله، در هر تکرار از فرآیند ایجاد ساختار شاخص‌گذاری پیشنهادی، زیرخوشه‌های راست و چپ  $(MR_i)_s$ ،  $(ML_i)_s$  به همراه ماتریس حاوی اطلاعات مرزبندی زیرخوشه‌ها  $(BR_i)$  و  $(BL_i)$  به ساختار درختی  $TS_i$  اضافه می‌شوند. شبه‌کد مربوطه در ادامه آمده است.

### ۳-۳-۲- تشکیل مستطیل‌های مرزبندی کمینه

در این مرحله، زیرخوشه‌های راست و چپ در مختصات جدید  $(NR_i)_s$  و  $(NL_i)_s$  گرفته می‌شوند تا توسط مستطیل‌های مرزبندی کمینه (MBRs) مرزبندی شوند. از آن‌جاکه  $(NR_i)_s$  و  $(NL_i)_s$  شامل ماتریس داده‌های مربوط به زیرخوشه‌های راست و چپ در مختصات جدید می‌باشند، مستطیل‌های مرزبندی کمینه در جهت مؤلفه غیرگوسی معنادار مربوط به گره والد، یا به عبارتی مؤلفه غیرگوسی معنادار خوشه انتخاب شده در تکرار  $i$   $(a_i)_s$  تشکیل می‌شوند. سپس اطلاعات به دست آمده مربوط به مستطیل‌های مرزبندی کمینه زیرخوشه‌های راست و چپ  $(BR_i)_s$ ،  $(BL_i)_s$  ذخیره می‌شود. برای مرزبندی روش استفاده شده در (Taileb et al., 2007) به کار گرفته شده است.

در شکل (۳-۷) نشان داده شده که چگونه تغییر مبدأ مختصات و تشکیل مستطیل‌های مرزبندی کمینه در راستای مبدأ جدید از هم‌پوشانی بین گره‌ها جلوگیری به عمل می‌آورد. همان‌طور که در شکل (۳-۷) مشاهده می‌شود، تشکیل مستطیل‌های مرزبندی کمینه در راستای مبدأ جدید از هم‌پوشانی گره‌ها جلوگیری می‌کند. در شکل (۳-۷) الف) و ب) تجزیه داده‌ها توسط ابرصفحه عمود بر مؤلفه غیرگوسی معنادار و گذرنده از میانگین مراکز ثقل زیرخوشه‌های افکنش که با رنگ نارنجی نشان داده شده، صورت گرفته است. در شکل (۳-۷) پ) تجزیه داده‌ها توسط ابرصفحه عمود بر مؤلفه غیرگوسی معنادار و گذرنده از مرکز ثقل اصلی داده‌ها که با رنگ آبی نشان داده شده، صورت گرفته است. با این‌که هم‌چنان عدم هم‌پوشانی بین مستطیل‌های مرزبندی کمینه در راستای  $(a_i)_s$  تضمین شده، اما از آن‌جاکه تجزیه داده‌ها در این حالت از محل مرکز ثقل اصلی داده‌ها صورت گرفته است، تجزیه داده‌ها منجر به خوشه‌های مجزا نمی‌شود و بنابراین، در مجموع، حجم مستطیل‌های مرزبندی کمینه تشکیل شده در شکل (۳-۷) پ) نسبت به شکل (۳-۷) ب) بیشتر است؛ زیرا مستطیل‌های مرزبندی شکل (۳-۷) ب) فشرده‌تر می‌باشند.

از آن جاکه هدف اصلی ساختارهای شاخص‌گذاری تصویر، بالابردن سرعت بازیابی تصاویر است، بازیابی داده‌ها از طریق ساختار شاخص‌گذاری پیشنهادی در این پژوهش، با معیار زمان پاسخ اندازه‌گیری شده‌اند. این معیار، درحقیقت زمانی است که صرف جستجوی پایگاه‌داده از طریق ساختار شاخص‌گذاری برای پیدا کردن  $k$  نزدیک‌ترین داده‌های درون پایگاه‌داده نسبت به داده پرس‌وجو می‌شود، و طبق رابطه (۱۲) تعریف می‌شود (Li et al., 2002, Taieb et al., 2007, Yu and Zhang, 2003).

$$R = r_1 - r_2 \quad (12)$$

در رابطه (۱۲)،  $R$  زمان پاسخ،  $r_1$  زمان خروج نتایج نهایی از بخش جستجوی شباهت و  $r_2$  زمان ورود پرس‌وجو به بخش جستجوی شباهت است.

#### ۴-۱-۲- روش آزمون

همه آزمون‌ها با استفاده از روش Cross-validation انجام گرفته‌اند. بدین منظور هر آزمون ده‌بار انجام و هر بار، بیست بردار ویژگی از پایگاه‌داده ویژگی استفاده‌شده برای آزمون موردنظر، در مجموعه آزمون قرار داده شده و برای هر یک، کل پایگاه‌داده، با جستجوی ترتیبی<sup>۲</sup> و معیار فاصله اقلیدسی، پیمایش شده و ۲۰- نزدیک‌ترین بردارهای مربوط به هر بردار موجود در مجموعه آزمون پیدا و سپس در مجموعه داده‌های مرتبط (مجموعه  $\alpha$ ) قرار داده شده و در مرحله بعد ساختارهای شاخص‌گذاری مورد ارزیابی با تصاویر باقی‌مانده ساخته شده است. از داده‌های مجموعه آزمون به‌عنوان بردارهای پرس‌وجو و به‌منظور پرس‌وجو از پایگاه‌داده استفاده و میانگین نتایج حاصل برای بیست‌بردار پرس‌وجو به‌دست‌آمده و در نهایت نتایج میانگین ده آزمون گزارش شده است (Savaresi et al., 2002, Yu and Zhang, 2003).

#### ۴-۳-۱- مجموعه داده‌ها

برای ارزیابی‌ها و مقایسه روش پیشنهادی با دیگر روش‌های مطرح در این حوزه، از مجموعه داده به‌کاررفته در آزمون‌های مطرح در (Taieb et al., 2007)، که پایگاه‌داده حاوی ۱۱۹۳۶۴۷ بردار ویژگی استخراج‌شده از تصاویر رنگی می‌باشد، استفاده شده است. بردارهای ویژگی تصاویر موجود در این پایگاه‌داده، ویژگی‌های محلی براساس نقاط مطلوب استخراج‌شده از ۴۹۹۶ تصویر مختلف رنگی می‌باشند که در

```
Build tree structure ((BRi)s, (BLi)s, (MRi)s, (MLi)s)
Begin
  If number of rows of ((MRi)s/(MLi)s) < Minpts
    let ((MRi)s/(MLi)s) be outlier
  Else let it be leaf
  Add (MRi)s, (MLi)s to the binary tree structure and save
  (BRi)s, (BLi)s, LNOi = LNOi+1
  Returns (TSi, LNOi)
End
```

همان‌طور که در شبه‌کد بالا ملاحظه می‌شود، در فرآیند مرحله ساخت ساختار درختی، در صورتی که تعداد اعضای داده‌های هر یک از آنها کوچک‌تر از مقدار تنظیم‌شده برای پارامتر Minpts باشد، آن زیرخوشه به‌عنوان داده پرت علامت‌زده می‌شود و در غیر این صورت به‌عنوان برگ علامت‌زده می‌شوند و در تکرار بعدی الگوریتم در مجموعه  $M_i$  قرار می‌گیرد.

## ۴- پیاده‌سازی و آزمون

در این بخش نتایج پیاده‌سازی و آزمون روش پیشنهادی ارائه شده است.

### ۴-۱- پیاده‌سازی

در این بخش به جزئیات پیاده‌سازی، شامل معیارهای ارزیابی، روش آزمون، مجموعه داده‌ها و محیط پیاده‌سازی پرداخته شده است.

#### ۴-۱-۱- معیارهای ارزیابی

در این پژوهش، دو معیار Recall و زمان پاسخ<sup>۱</sup> به‌منظور ارزیابی روش پیشنهادی در نظر گرفته شده‌اند. کیفیت خوشه‌های نهایی که مورد جستجو قرار می‌گیرند با استفاده از معیار Recall اندازه‌گیری شده است. معیار Recall، نسبت تعداد داده‌های مرتبط با بردار پرس‌وجو را که بازیابی شده‌اند به کل داده‌های مرتبط با بردار پرس‌وجوی مطرح‌شده در پایگاه‌داده، اندازه‌گیری می‌کند، و طبق رابطه زیر تعریف می‌شود (Li et al., 2002, Yu and Zhang, 2003, Mitchell, 1997).

$$\text{Recall} = \frac{|\alpha \cap \beta|}{|\alpha|} \quad (11)$$

در رابطه (۱۱)،  $\alpha$  مجموعه داده‌های مرتبط بازیابی‌شده،  $\beta$  مجموعه داده‌های بازیابی‌شده و  $|\alpha|$ ، نشان‌دهنده تعداد است.

<sup>2</sup>Sequential-scan

<sup>1</sup> Response time



برای تنظیم دو پارامتر پیش‌نیاز  $k$  و  $Minpts$ ، آزمون‌هایی انجام، و متوسط زمان پاسخ برحسب ثانیه گزارش شده است. این آزمون‌ها با تغییر مقادیر پارامتر  $Minpts$  و مقدار ثابت  $k$  روی هر چهار پایگاه‌داده ۶۰،۴۰،۲۵ و ۸۰ بعدی انجام و این آزمون برای مقادیر مختلف  $k$  تکرار شده است. جدول (۴-۱) نتایج متوسط زمان پاسخ برحسب ثانیه (S) به‌دست‌آمده از آزمون روش پیشنهادی را، با در نظر گرفتن مقادیر مختلف  $Minpts$  و  $k$ ، روی چهار پایگاه‌داده ۶۰،۴۰،۲۵ و ۸۰ بعدی نشان می‌دهد.

جدول (۴-۱): متوسط زمان پاسخ برای مقادیر مختلف  $Minpts$  و  $k$  روی پایگاه‌داده‌ها با ابعاد مختلف

۶۵	۴۵	۳۵	۲۵	۱۵	۵	Minpts
متوسط زمان پاسخ (S) روی پایگاه داده‌ها با ابعاد مختلف، $k=200$						
۰.۱۴۶	۰.۱۲۵	۰.۱۱۲	۰.۱۰۴	۰.۱۲۵	۰.۱۴۶	پایگاه داده ۲۵ بعدی
۰.۴۲۷	۰.۳۸۶	۰.۳۷۵	۰.۳۵۲	۰.۳۶۸	۰.۴۶۲	پایگاه داده ۴۰ بعدی
۰.۶۱۳	۰.۵۹۳	۰.۵۵۸	۰.۵۳۵	۰.۵۷۸	۰.۶۶۳	پایگاه داده ۶۰ بعدی
۰.۷۳۲	۰.۷۰۳	۰.۶۷۹	۰.۶۶۴	۰.۶۸۴	۰.۷۷	پایگاه داده ۸۰ بعدی
متوسط زمان پاسخ (S) روی پایگاه داده‌ها با ابعاد مختلف، $k=600$						
۰.۰۷۴	۰.۰۶۸	۰.۰۶۳	۰.۰۶۲	۰.۰۶۹	۰.۰۸۴	پایگاه داده ۲۵ بعدی
۰.۳۱۶	۰.۲۵۷	۰.۲۳	۰.۲۲۸	۰.۲۴۳	۰.۳۶	پایگاه داده ۴۰ بعدی
۰.۴۳۸	۰.۴۱	۰.۳۲۲	۰.۳۴۵	۰.۳۶	۰.۴۷۵	پایگاه داده ۶۰ بعدی
۰.۶۰۷	۰.۵۶۱	۰.۵۲۱	۰.۵۱۳	۰.۵۴	۰.۶۳۲	پایگاه داده ۸۰ بعدی
متوسط زمان پاسخ (S) روی پایگاه داده‌ها با ابعاد مختلف، $k=800$						
۰.۱۲۴	۰.۱۱	۰.۱۰۳	۰.۰۹۸	۰.۱۰۹	۰.۱۷۳	پایگاه داده ۲۵ بعدی
۰.۴۰۱	۰.۳۵	۰.۳۲۱	۰.۳۰۱	۰.۳۴۳	۰.۴۴۷	پایگاه داده ۴۰ بعدی
۰.۴۶	۰.۴۵۱	۰.۳۸۷	۰.۳۶	۰.۴۰۶	۰.۵۲۴	پایگاه داده ۶۰ بعدی
۰.۶۱۷	۰.۵۸	۰.۵۳۲	۰.۵۲	۰.۵۶۷	۰.۶۵	پایگاه داده ۸۰ بعدی
متوسط زمان پاسخ (S) روی پایگاه داده‌ها با ابعاد مختلف، $k=1000$						
۰.۲۰۱	۰.۱۵۵	۰.۱۴۸	۰.۱۴	۰.۱۴۵	۰.۲۴۶	پایگاه داده ۲۵ بعدی
۰.۵۶۳	۰.۵۲۶	۰.۴۹	۰.۴۵۳	۰.۵۱۸	۰.۶۱	پایگاه داده ۴۰ بعدی
۰.۷۶۳	۰.۷۲۴	۰.۶۸۳	۰.۶۵۲	۰.۷۰۵	۰.۷۹۲	پایگاه داده ۶۰ بعدی
۰.۹۷	۰.۹۴۶	۰.۹	۰.۸۷۳	۰.۹۲۵	۱.۰۲۶	پایگاه داده ۸۰ بعدی

بررسی نتایج جدول (۴-۱)، نشان می‌دهد که با بالا رفتن ابعاد پایگاه‌داده، متوسط زمان پاسخ برای هر چهار مقدار  $k$  افزایش یافته و با تغییر مقادیر  $Minpts$ ، نتایج مختلفی حاصل شده است. بنابراین برای تحلیل دقیق‌تر نتایج، در شکل (۴-۱) از ابعاد صرف‌نظر شده و نتایج روی پایگاه‌داده ۲۵ بعدی به‌صورت هیستوگرام متوسط زمان پاسخ، با تغییر مقادیر مختلف  $Minpts$ ، برای مقادیر مختلف  $k$ ، نشان داده شده است. در شکل (۴-۱)، مشاهده می‌شود که با تنظیمات انجام‌شده روی پارامترهای پیش‌نیاز، بهترین نتایج

تحقیقات پیشین توسط محققین، به‌منظور ارزیابی روش‌ها، مورد استفاده قرار گرفته‌اند. در آزمون‌های انجام‌شده از ۴ پایگاه‌داده حاوی ۵۰،۰۰۰ بردار ویژگی در ابعاد مختلف استفاده شده است. پایگاه‌داده اول شامل ۵۰،۰۰۰ بردار ویژگی ۲۵ بعدی، پایگاه‌داده دوم شامل ۵۰،۰۰۰ بردار ویژگی ۴۰ بعدی، پایگاه‌داده سوم شامل ۵۰،۰۰۰ بردار ویژگی ۶۰ بعدی و پایگاه‌داده چهارم شامل ۵۰،۰۰۰ بردار ویژگی ۸۰ بعدی است.

#### ۴-۱-۴- محیط پیاده‌سازی

برای پیاده‌سازی و آزمون روش پیشنهادی و سایر روش‌های مورد مقایسه، از محیط MATLAB 7.8.0.347 در یک دستگاه رایانه شخصی با امکانات CPU.4GHZ و RAM4GB استفاده شده است.

#### ۴-۲- آزمون

در این بخش، جزئیات مربوط به آزمون روش پیشنهادی مطرح شده است.

#### ۴-۲-۱- تنظیم پارامترهای پیش‌نیاز

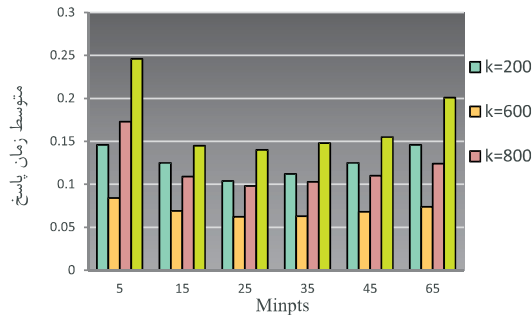
در این بخش، به‌منظور تنظیم پارامترهای پیش‌نیاز طرح شاخص‌گذاری پیشنهادی، آزمون‌هایی انجام شده و نتایج آن مورد تحلیل قرار گرفته است.

اگر به شبهه‌کدهای بخش ۳ رجوع شود، مشاهده می‌شود که در طرح پیشنهادی دو پارامتر پیش‌نیاز  $k$  و  $Minpts$  در نظر گرفته شده‌اند. پارامتر  $k$ ، تعداد خوشه‌های نهایی (برگ‌ها و داده‌های پرت نهایی) است. پارامتر  $k$  در حقیقت تعداد برگ‌ها به‌علاوه تعداد داده‌های پرت نهایی، پس از ایجاد ساختار شاخص‌گذاری چندبعدی پیشنهادی  $NO-NGP-tree$  را مشخص می‌کند. مقدار حداقل در نظر گرفته شده برای تعداد داده‌های برگ با پارامتر  $Minpts$  مشخص می‌شود. در آزمون‌های انجام‌شده، پارامتر  $Minpts$  به‌صورت درصد متوسط تعداد اعضای داده خوشه‌های نهایی در نظر گرفته شده است. به‌عنوان مثال، مقدار ۱۵، به‌معنای ۱۵ درصد متوسط تعداد اعضای داده خوشه‌های نهایی است. این مقدار متوسط طبق رابطه (۱۳) به‌دست می‌آید.

$$(13)$$

$k$ /تعداد کل داده‌ها=متوسط تعداد اعضای داده خوشه‌های نهایی





شکل (۴-۱): نتایج روی پایگاه داده ۲۵ بعدی به صورت هیستوگرام متوسط زمان پاسخ با تغییر مقادیر مختلف Minpts برای مقادیر مختلف k

NGP-tree: این ساختار شاخص گذاری با حذف مرحله تغییر مبدأ مختصات از فرآیند ایجاد ساختار شاخص گذاری پیشنهادی No-NGP-tree، و با هدف مقایسه اثر عدم هم پوشانی گره ها در ساختار شاخص گذاری پیشنهادی ایجاد شده است.

PDDP-tree: با اضافه کردن MBRها به روش خوشه بندی سلسله مراتبی تقسیم کننده PDDP ایجاد شده است.

NOHIS-tree: ساختار شاخص گذاری مبتنی بر خوشه بندی سلسله مراتبی تقسیم کننده PDDP است.

لازم به ذکر است در چهار روش بررسی شده، در بخش جستجوی شباهت، از الگوریتم جستجوی شباهت پیشنهاد شده در (Taieb et al., 2007) استفاده شده است. آزمون های این بخش با تنظیم پارامتر Minpts برابر ۲۵ درصد، و با در نظر گرفتن مقادیر مختلف برای پارامتر k، برای روش های NO-NGP-tree و NGP-tree انجام شده است. هدف از انجام این ارزیابی، با در نظر گرفتن مقادیر مختلف پارامتر k، نشان دادن تأثیر تنظیم این پارامتر و مقایسه روش پیشنهادی در این حالات با روش های دیگر است.

برای نشان دادن بهتر نتایج، نمودارهای متوسط Recall نسبت به تعداد خوشه های نهایی جستجو شده روی پایگاه داده های ۲۵ و ۸۰ بعدی رسم شده است. شکل (۴-۲)، نمودارهای متوسط Recall نسبت به تعداد خوشه های نهایی جستجو شده روی پایگاه داده ۲۵ و ۸۰ بعدی، در حالی که پارامتر k برابر ۶۰۰، ۸۰۰، ۱۰۰۰، و پارامتر Minpts برابر ۲۵ درصد را قرار داده شده است نشان می دهد.

در مورد NO-NGP-tree، جستجو پس از جستجوی به طور متوسط چهارده خوشه نهایی متوقف می شود، و

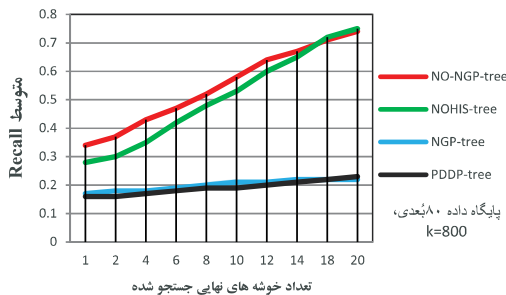
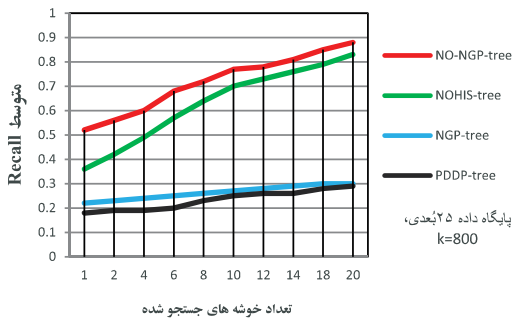
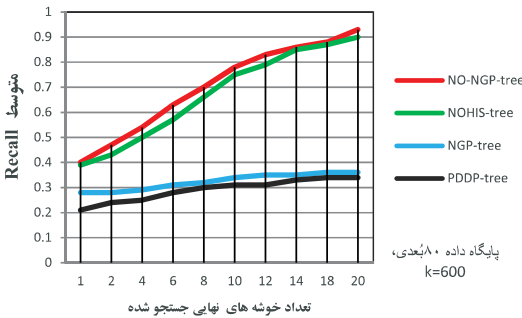
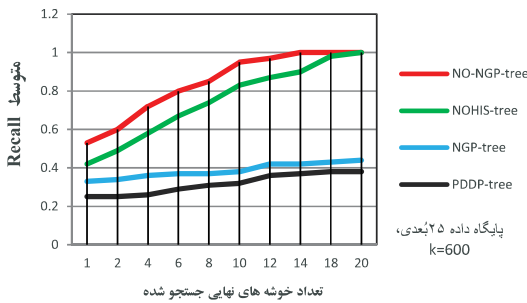
زمانی حاصل می شود که Minpts برابر ۲۵ درصد متوسط تعداد اعضای داده خوشه های نهایی، و k نیز برابر ششصد قرارداد شده است. همچنین وقتی Minpts روی مقادیر بزرگ (۶۵ درصد) و یا مقادیر کوچک (۵ درصد) تنظیم شد، متوسط زمان پاسخ بالا رفته، و هنگامی که Minpts روی مقادیر ۱۵ تا ۲۵ درصد تنظیم شده، نتایج بهتری حاصل شده است.

این مسأله ناشی از این حقیقت است که اگر حداقل در نظر گرفته شده برای تعداد اعضای خوشه هایی که آنها را برگ می نامیم در هر تکرار بسیار بزرگ باشد، از آن جاکه معیار انتخاب خوشه بعدی ارائه شده در روش پیشنهادی تمایل به انتخاب خوشه های کوچک تر دارد، به طور مداوم خوشه های کوچک تجزیه می شوند و این مسأله منجر به افزایش ارتفاع ساختار شاخص گذاری به طور نامتوازن تر شده و در نتیجه سرعت جستجوی پایگاه داده از طریق ساختار شاخص گذاری پیشنهادی کاهش می یابد. همچنین، تنظیم مقادیر بسیار کوچک برای Minpts باعث می شود داده های پرت به درستی تشخیص داده نشوند و فقط خوشه هایی با تعداد داده های بسیار کم به عنوان داده های پرت علامت زده شوند. در نتیجه منجر به افزایش پرتی بی جا و کاهش کیفیت خوشه های نهایی می شود، که در نهایت افزایش زمان پاسخ را در بر خواهد داشت. در شکل (۴-۱)، مشاهده می شود که زمان پاسخ برای مقادیر مختلف k و مقدار ثابت Minpts متفاوت است. این مسأله می تواند به این علت باشد که اگر مقدار k نزدیک به تعداد خوشه های طبیعی موجود در توزیع داده ها تنظیم شود، کیفیت خوشه های نهایی بهتر خواهد بود و در نتیجه از آن جاکه جستجو در خوشه های کمتری را موجب می شود، جستجو از طریق ساختار شاخص گذاری پیشنهادی زمان کمتری خواهد برد.

#### ۴-۲-۲- نتایج حاصل از ارزیابی کیفیت خوشه های نهایی جستجو شده

در این بخش نتایج حاصل از ارزیابی کیفیت خوشه های نهایی جستجو شده، با استفاده از معیار Recall و در مقایسه با سه روش زیر نشان داده شده است:

مؤلفه غیر گوسی معنادار از داده‌ها در پایگاه داده‌های با ابعاد بالاتر است.



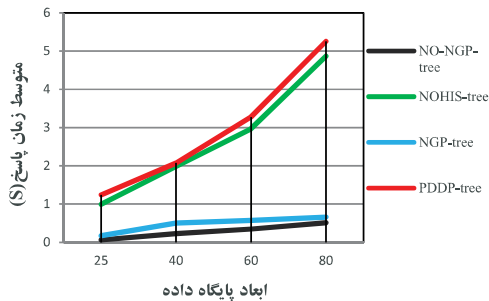
در مورد NOHIS-tree جستجو بعد از جستجوی به‌طور متوسط بیست خوشه نهایی متوقف می‌شود. همان‌طور که در شکل (۲-۴) مشاهده می‌شود متوسط Recall با افزایش تعداد خوشه‌های جستجو شده افزایش می‌یابد و Recall بعد از جستجوی به‌طور تقریبی چهارده و بیست خوشه به ترتیب در روش NO-NGP-tree و NOHIS-tree یک شده است. همچنین در مورد روش‌های PDDP-tree و NGP-tree که در ساختار آنها گره‌ها هم‌پوشانی دارند، متوسط Recall به‌کندی افزایش یافته و در نهایت در NGP-tree پس از جستجوی به‌طور تقریبی ۲۵ خوشه، و در PDDP-tree پس از جستجوی به‌طور تقریبی ۲۷ خوشه، جستجو متوقف می‌شود. دلیل افزایش کندتر متوسط Recall، در مورد دو روش که در ساختار آنها گره‌ها هم‌پوشانی دارند؛ را می‌توان عدم دقت محاسبه فاصله MINDIST (فاصله بین مستطیل مرزبندی کمینه و بردار پرس‌وجو (Taieb et al., 2007))، زمانی که گره‌های میانی هم‌پوشانی دارند، توجیه کرد، همچنین هر اندازه مستطیل‌های مرزبندی کمینه فشرده‌تر باشند، دقت محاسبه فاصله MINDIST بیشتر خواهد بود. همان‌طور که در شکل (۲-۴) مشاهده می‌شود، روش NO-NGP-tree و NOHIS-tree نتایج نسبتاً بهتری را نسبت به NOHIS-tree تولید کرده است. دلیل این نتایج بهتر را می‌توان با توجه به افزایش‌های بهینه‌ای که در طول ساخت ساختار شاخص‌گذاری دارد و منجر به ایجاد مستطیل‌های مرزبندی کمینه به‌نسب فشرده‌ای می‌شود، توجیه کرد.

همان‌طور که در شکل (۲-۴) مشاهده می‌شود، متوسط Recall نسبت به نتایج حاصله روی پایگاه داده با ابعاد ۲۵، در هر چهار روش کندتر رشد کرده است. همچنین، مشاهده می‌شود، با بالا رفتن مقدار پارامتر پیش‌نیاز  $k$ ، ساختار پیشنهادی NO-NGP-tree و NGP-tree افزایش متوسط Recall کمتری را نشان می‌دهند. دلیل کاهش متوسط Recall این است که هرچه مقدار  $k$  از تعداد خوشه‌های واقعی موجود در توزیع داده‌ها فاصله می‌گیرد، کیفیت خوشه‌های نهایی به دست آمده پایین می‌آید، همچنین این کاهش کارایی در بعد ۸۰ بیشتر خود را نشان می‌دهد، به‌طوری که بعد از جستجوی دوازده خوشه، متوسط Recall برای NOHIS-tree مقدار بالاتری را نسبت به متوسط Recall برای NO-NGP-tree نشان می‌دهد. دلیل این مسأله، کاهش کارایی الگوریتم استفاده‌شده در مرحله پیش‌افرازبندی برای یافتن یک



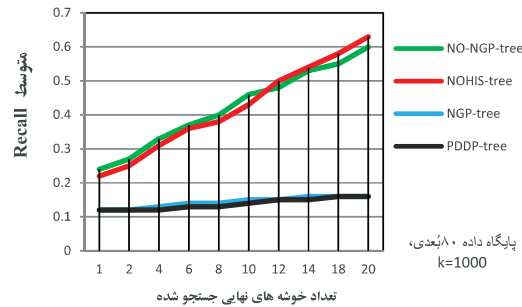
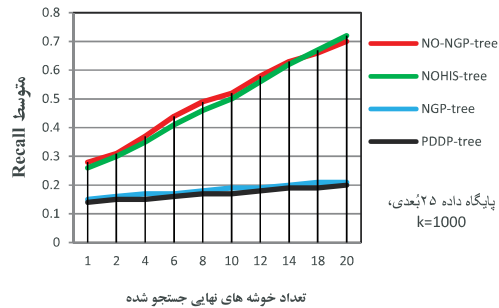
جدول (۲-۴): متوسط زمان پاسخ روی پایگاه داده‌های با ابعاد مختلف

طرح شاخص گذاری	متوسط زمان پاسخ (S)			
	پایگاه داده ۸۰ بعدی	پایگاه داده ۶۰ بعدی	پایگاه داده ۴۰ بعدی	پایگاه داده ۲۵ بعدی
NO-NGP-tree	۰.۵۱۳	۰.۳۴۵	۰.۲۲۸	۰.۶۲
NGP-tree	۴.۸۷	۲.۸۶۹	۱.۹۸۶	۰.۹۹۲
NOHIS-tree	۰.۶۶۰	۰.۵۷۲	۰.۵۰۴	۰.۱۷۲
PDDP-tree	۵.۲۵۲	۳.۲۷۶	۲.۵۰۳	۱.۲۴۳
Sequential-scan	۶.۶۵۲	۴.۹۲۸	۳.۱۰۲	۱.۷۲۳



شکل (۳-۴): متوسط زمان پاسخ روی پایگاه داده‌های ۲۵، ۴۰، ۶۰ و ۸۰ بعدی برای مقایسه طرح‌های شاخص گذاری NO-NGP-tree، PDDP-tree و NOHIS-tree و NGP-tree

در شکل (۳-۴)، طرح شاخص گذاری پیشنهادی-NO در NGP-tree که با رنگ مشکی نشان داده شده است، متوسط زمان پاسخ کمتری را نسبت به NOHIS-tree، NO-NGP-tree و PDDP-tree در همه ابعاد دارد. همان‌طور که مشاهده می‌شود زمان پاسخ برای هر چهار روش با افزایش ابعاد افزایش یافته است. در این روش‌ها افزایش ابعاد پایگاه داده، منجر به افزایش حجم مستطیل‌های مرزبندی کمینه گره‌ها می‌شود. بدین ترتیب، متوسط فاصله بردار پرس‌وجو از مستطیل‌های مرزبندی کمینه افزایش و در نتیجه دقت محاسبه فاصله MINDIST در گره‌ها کاهش می‌یابد. کاهش دقت محاسبه فاصله منجر به افزایش تعداد پیمایش‌های ساختار شاخص گذاری، افزایش تعداد خوشه‌های نهایی جستجو شده تا پیش از اتمام جستجو، و بنابراین موجب افزایش زمان پاسخ می‌شود. همچنین در شکل (۳-۴) مشاهده می‌شود دو روش NO-NGP-tree و PDDP-tree که در ساختارشان از هم‌پوشانی گره‌ها جلوگیری نشده است، نسبت به روش‌های NOHIS-tree و NO-NGP-tree متوسط زمان پاسخ بیشتری را دارند. دلیل بیشتر بودن زمان پاسخ این روش‌ها وجود هم‌پوشانی بین گره‌هاست که موجب افزایش پیمایش‌های ساختار شاخص گذاری هنگام جستجو



شکل (۲-۴): متوسط Recall نسبت به تعداد خوشه‌های نهایی جستجو شده، Minpts=25

### ۳-۴- نتایج حاصل از ارزیابی سرعت جستجو از طریق طرح شاخص گذاری پیشنهادی

در این بخش نتایج حاصل از ارزیابی سرعت جستجو از طریق طرح شاخص گذاری پیشنهادی، با استفاده از معیار زمان پاسخ و در مقایسه با طرح‌های شاخص گذاری NOHIS-tree، NO-NGP-tree و PDDP-tree، جستجوی ترتیبی نشان داده شده است.

آزمون‌های این بخش با تنظیم پارامترهای پیش‌نیاز در ساختارهای NO-NGP-tree و NGP-tree، روی بهترین حالت به دست آمده در بخش تنظیم پارامترهای موجود انجام شده است. بنابراین در تمام آزمون‌های این بخش Minpts برابر ۲۵ درصد و k برابر ۶۰۰ تنظیم شده است. همچنین در همه روش‌های بررسی شده، در بخش جستجوی شباهت، از الگوریتم جستجوی پیشنهاد شده در (Taieb et al., 2007) استفاده شده است. جدول (۲-۴)، متوسط زمان پاسخ بر حسب ثانیه را برای پایگاه داده ۴۰، ۶۰، ۲۵، ۸۰ بعدی نشان می‌دهد. شکل (۳-۴) نمودارهای متوسط زمان پاسخ را برای مقایسه طرح‌های شاخص گذاری NO-NGP-tree، NOHIS-tree، NO-NGP-tree و PDDP-tree، روی پایگاه داده‌های ۴۰، ۶۰ و ۸۰ بعدی نشان می‌دهد.

پژوهش با بهره‌گیری از الگوی روش‌های خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده، ساختار شاخص‌گذاری چندبعدی جهت مدیریت پایگاه‌داده‌های حاوی بردارهای ویژگی تصویر در ابعاد بالا ارائه شد. در بخش چهارم، پارامترهای پیش‌نیاز ساختار شاخص‌گذاری چندبعدی پیشنهادی NO-NGP-tree با مقادیر مختلف تنظیم شده، ساختار شاخص‌گذاری چندبعدی پیشنهادی موردارزیابی قرار گرفت و نتایج به‌دست‌آمده گزارش شد. همچنین، کیفیت خوشه‌های نهایی حاصل از روش خوشه‌بندی سلسله‌مراتبی تقسیم‌کننده پیشنهادی، که ساختار پیشنهادی براساس آن NO-NGP-tree شکل گرفته است، موردارزیابی قرار داده شد. در نهایت با تنظیم پارامترهای پیش‌نیاز در بهترین حالت، زمان بازیابی از طریق ساختار شاخص‌گذاری چندبعدی پیشنهادی موردارزیابی قرار گرفته و نتایج مطلوبی حاصل شد.

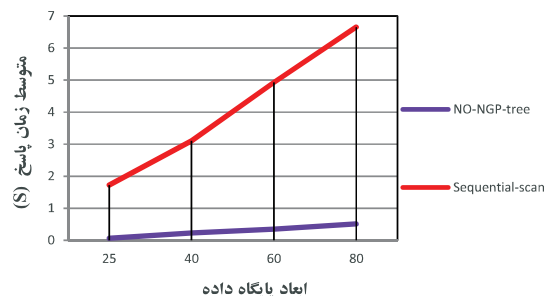
برخی از فعالیت‌هایی که به‌منظور بهبود و توسعه ساختار پیشنهادی NO-NGP-tree می‌توانند در آینده در نظر گرفته شوند عبارتند از: ۱- به‌کارگیری توابع هدف دیگر به‌جای تقریب آشفته‌گی منفی در فرآیند پیدا کردن مؤلفه غیرگوسی معنادار. مقایسه نتایج به‌دست‌آمده می‌تواند در فعالیت‌های آینده در نظر گرفته شود؛ ۲- ساختار شاخص‌گذاری چندبعدی پیشنهادی از یک ساختار درختی باینری نامتوازن تبعیت می‌کند، بنابراین تجزیه داده‌ها به چندین زیرخوشه به‌جای تنها دو زیرخوشه در هر تکرار الگوریتم می‌تواند تا حد مطلوبی کارایی ساختار شاخص‌گذاری پیشنهادی را بهبود دهد؛ ۳- ارائه روشی برای تخمین پارامتر پیش‌نیاز  $k$ ، با استفاده از روش‌های انتخاب مدل؛ ۴- تحلیل و ارزیابی عملکرد ساختار پیشنهادی در مجموعه داده‌هایی با بردارهای ویژگی در ابعاد بسیار بالاتر.

## ۶- منابع

Asbahi sh., Keyvanpour M.R., Amiri A., "Learning-Based Approach for Semantic Image Retrieval by Using a Dynamic Semantic Network", DEXA, 23rd International Workshop on Database and Expert Systems Applications, pp. 107-111, 2008.

Bohm C., Berchtold S., Keim D. A., "Searching in high-dimensional spaces-Index structure for Improving the performance of Multimedia Databases", ACM Computing surveys, Vol. 33, No. 3, pp. 322-373, 2001.

می‌شود. در این پژوهش ادعا شده که دلیل کاهش زمان پاسخ در طرح شاخص‌گذاری پیشنهادی NO-NGP-tree نسبت به طرح شاخص‌گذاری NOHIS-tree، افزایش بازیابی‌های بهینه و تولید زیرخوشه‌های متراکم‌تر و در نتیجه کاهش حجم مستطیل‌های مرزبندی کمینه در ساختار سلسله‌مراتبی ارائه شده NO-NGP-tree است. کاهش حجم مستطیل‌های مرزبندی کمینه موجب افزایش دقت محاسبه MINDIST، و در نتیجه باعث کاهش زمان پاسخ می‌شود. به همین دلیل زمان پاسخ برای طرح شاخص‌گذاری PDDP- از NO-NGP-tree کمتر شده است. شکل (۴-۴) نمودار زمان پاسخ برای مقایسه طرح شاخص‌گذاری NO-NGP-tree و Sequential-scan، در پایگاه‌داده‌های ۲۵،۴۰، ۶۰ و ۸۰ بعدی را نشان می‌دهد.



(شکل ۴-۴): متوسط زمان پاسخ روی پایگاه‌داده‌های ۲۵،۴۰،

۶۰ و ۸۰ بعدی برای مقایسه طرح‌های شاخص‌گذاری NO-NGP-tree و Sequential-scan

دلیل مقایسه طرح شاخص‌گذاری پیشنهادی با جستجوی ترتیبی این است که اکثر طرح‌های شاخص‌گذاری سنتی در فضاهای با ابعاد بالا کارایی خود را به‌گونه‌ای از دست می‌دهند که گاهی جستجوی ترتیبی داده‌ها بهتر عمل می‌کند (Markov et al., 2008, Chakrabarti and Mehrotra, 1999). همان‌طور که در شکل (۴-۴) مشاهده می‌شود، طرح شاخص‌گذاری پیشنهادی NO-NGP-tree با فاصله زیادی توانسته است زمان پاسخ را نسبت به جستجوی ترتیبی بهبود دهد.

## ۵- نتیجه‌گیری

در اکثر کاربردهای بازیابی تصویر، تصاویر توسط بردارهای ویژگی با ابعاد بالا مورد استفاده قرار می‌گیرند. نتایج تحقیقات نشان داده که ساختارهای شاخص‌گذاری چندبعدی رایج با افزایش ابعاد داده‌ها کارایی خود را از دست می‌دهند. در این

- Keyvanpour M.R. and Tavoli R., "Feature Weighting for Improving Document Image Retrieval System Performance", *International Journal of Computer Science Issues*, Vol 9, Issue 3, No 3, pp. 125-130, 2012.
- Lew M. S., Sebe N., Djeraba C. and Jain R., "Content-based Multimedia Information Retrieval: State of the Art and Challenges", *ACM Transactions on Multimedia Computing, Communications and Applications*, 2006.
- Li C., Chang E., Garcia-Molina H. and Wiederhold G., "Clustering for approximate similarity search in high-dimensional spaces", *IEEE Transactions on Knowledge and Data Engineering*, Vol. 14, No. 4, pp. 792-808, 2002.
- MacQueen J., "Some Methods for Classification and Analysis of Multivariate Observations", In *Proceeding of 5th Berkeley Symp. Mat. Statist. Prob.1*, pp. 281-297, 1967.
- Markov K., Ivanova K., Mitov I. and Karastanev S., "Advance of the Access Methods", *International Journal of Information Technologies and Knowledge*, Vol. 2, 2008.
- Mitchell T. M., "Machine Learning", McGraw-Hill, 1997.
- Rui Y. and Huang T. S., "Image Retrieval: current techniques, promising, Directions and Open Issues", *Journal of Visual communication and Image representation*, Vol. 10, Issue. 1, pp. 39-62, 1999.
- Savaresi S. M., Boley D., Bittanti S. and Gazzaniga G., "Cluster Selection in Divisive Clustering Algorithms", In *Proceedings of SDM*, 2002.
- Srinivasa Rao Ch., Srinivas Kumar S. and Chandra Mohan B., "Content based Image Retrieval Using Exact Legendre Moments And Support Vector Machine", *The international journal of Multimedia & Its applications*, Vol. 2, No. 2, 2010.
- Taileb M., Lamrous S. and Touati S., "Non-overlapping Hierarchical Index Structure for similarity search", *International Journal of Computer Science*, Vol. 3, No. 1, 2007.
- Xu H., Xu D. and Lin E., "An Applicable Hierarchical Clustering Algorithm for Content-Based Image Retrieval", Springer-Verlag Berlin Heidelberg, 2007.
- Yu D. and Zhang A., "ClusterTree: Integration of Cluster Representation and Nearest-Neighbor Search for Large Data Sets with High Dimensions", *IEEE transaction on Knowledge and data Engineering*, VOL. 15, NO. 3, MAY/JUNE 2003.
- Boley D., "Principal Direction Divisive Partitioning", *Data Mining and Knowledge Discovery* 2, 325-344, 1998.
- Chakrabarti K. and Mehrotra Sh., "The hybrid tree: a structure for high dimensional feature spaces", *IEEE international conference on data engineering*, 1999.
- Datta R., Joshi Jiali D. and Wang J. Z., "Image Retrieval: Ideas, Influences and Trends of the New Age", *The Pennsylvania State University*, *ACM Transactions on Computing Surveys*, 2008.
- Gaede V. and Günther O., "Multi-dimensional Access Methods", *ACM Computing Surveys*, Vol. 30, No. 2, 1998.
- Guha S., Rastogi R. and Shim K., "CURE: an efficient clustering algorithm for large databases", *Proc. of the ACM SIGMOD Int'l Conf. on Management of Data*, ACM Press, pp. 73-84, 1998.
- Hyvärinen A. and Oja E., "A fast fixed-point algorithm for independent component analysis", *Neural Computation*, Vol. 9, NO.7, pp. 1483-1492, 1997.
- Hyvärinen A. and Oja E., "Independent Component Analysis: Algorithms and Applications", *Neural Networks*, Vol. 13, No. (4-5), pp. 411-430, 2000.
- Jackson J. E., "A User's Guide to Principal Components", New York: John Wiley and Sons, Inc, 1991.
- Jolliffe I. T., "Principal Component Analysis", Second edition, New York: Springer-Verlag New York, Inc, 2002.
- Keyvanpour M.R. and Asbaghi sh., "A new approach for interactive semantic image retrieval using the high level semantics", *Proceedings of the ACM Symposium on Applied Computing (SAC)*, pp. 1175-1179, 2008.
- Keyvanpour M.R. and Izadpanah N., "Analytical Classification of Multimedia Index Structures by Using a Partitioning Method-Based Framework", *The International Journal of Multimedia & Its Applications*, Vol. 3, No. 1, February 2011.
- Keyvanpour M.R. and Izadpanah N., "Classification and Evaluation of High-dimensional Image Indexing Structures", *IEEE 3<sup>rd</sup> international conference on Machin Learning and Computing (ICMLC)*, 2011.
- Keyvanpour M.R. and Tavoli R., "Document Image Retrieval: Algorithms, Analysis and Promising Directions", *International Journal of Software Engineering and Its Applications*, Vol 7, No. 1, pp. 93-106, 2013.



Zhang T., Ramakrishnan R. and Livny M., "BIRCH: An Efficient Data Clustering Method for VeryLarge Databases", In Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data, Montreal, Canada, pp. 103-114, 1996.



**محمدرضا کیوان‌پور:** دوره

کارشناسی خود را در سال ۱۳۷۶ در رشته مهندسی کامپیوتر (نرم‌افزار) در دانشگاه علم و صنعت به‌پایان رساند.

همچنین مدرک کارشناسی ارشد خود

را در رشته مهندسی کامپیوتر (نرم‌افزار) در سال ۱۳۷۹ و مدرک دکترای خود را در سال ۱۳۸۶ در رشته مهندسی کامپیوتر (نرم‌افزار) از دانشگاه تربیت مدرس اخذ کرد. وی هم‌اکنون به‌عنوان عضو هیأت علمی دانشکده فنی و مهندسی دانشگاه الزهراء مشغول به فعالیت است. زمینه‌های علمی موردعلاقه وی پردازش تصویر و داده‌کاوی است.

نشانی رایانامه ایشان عبارت است از:

[keyvanpour@Alzahra.ac.ir](mailto:keyvanpour@Alzahra.ac.ir)



**سعیده رنجبران:** دوره کارشناسی

خود را در سال ۱۳۸۵ در رشته

مهندسی کامپیوتر (نرم‌افزار) در

دانشگاه آزاد اسلامی قزوین به‌پایان

رساند. وی اکنون در حال گذراندن

دوره کارشناسی ارشد در رشته مهندسی کامپیوتر (نرم‌افزار) در همان دانشگاه است. زمینه‌های علمی موردعلاقه وی رمزنگاری و داده‌کاوی است.

نشانی رایانامه ایشان عبارت است از:

[se.ranjbaran@gmail.com](mailto:se.ranjbaran@gmail.com)

**نجوا ایزدپناه:** دوره کارشناسی خود را در سال ۱۳۸۴ در

رشته مهندسی کامپیوتر (نرم‌افزار) در دانشگاه علم و

فرهنگ‌پایان رساند. همچنین مدرک کارشناسی ارشد خود

را در رشته مهندسی کامپیوتر (نرم‌افزار) در سال ۱۳۹۰ از

دانشگاه آزاد اسلامی قزوین اخذ کرد. زمینه‌های علمی

موردعلاقه وی شاخص‌گذاری تصویر و داده‌کاوی است.

نشانی رایانامه ایشان عبارت است از:

[n.izadpanah@qiau.ac.ir](mailto:n.izadpanah@qiau.ac.ir)