

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



پروژه درس داده کاوی
خوشه‌بندی

استاد محترم درس
جناب آقای دکتر پاینده

دانشجو
آیدا اعلاییکی

در ابتدا بر خود واجب می‌دانم از زحمات استاد ارجمندم جناب آقای دکتر پاینده، کمال
تشکر را به جا آورم.

7	خوشه‌بندی.....
9	1-4- مقدمه.....
10	2-4- نقاط قوت روش خوشه‌بندی.....
11	3-4- نقاط ضعف روش خوشه‌بندی.....
11	4-4- تعیین تعداد خوشه.....
12	5-4- ارزیابی اعتبار در خوشه‌بندی.....
14	6-4- ماتریس تشابه و فاصله.....
14	1-6-4- تابع تشابه.....
15	2-6-4- تابع فاصله.....
15	3-6-4- ماتریس تشابه.....
16	4-6-4- ماتریس فاصله.....
16	7-4- روش‌های اصلی خوشه‌بندی.....
16	1-7-4- روش‌های افزایی.....
17	1-1-7-4- الگوریتم <i>k - means</i>
24	2-1-7-4- الگوریتم <i>k - medoids</i>
29	2-7-4- روش‌های سلسله مراتبی.....
31	1-2-7-4- الگوریتم <i>Brich</i>
33	2-2-7-4- الگوریتم <i>chameleon</i>

35	AGNES	الگوریتم	3-2-7-4
35	sin gel - link	خوشه‌بندی با روش	1-3-2-7-4
36	complete - link	خوشه‌بندی با روش	2-3-2-7-4
36	Average - link	خوشه‌بندی با روش	3-3-2-7-4
43		مقایسه خوشه‌بندی سلسله مراتبی و غیر سلسله مراتبی	4-2-7-4
44		روش‌های مبتنی بر چگال	3-7-4
45	DBSCAN	الگوریتم	1-3-7-4
46	DBSCAN	مزیتها و عیبهای الگوریتم	1-1-3-7-4
47	OPTICS	الگوریتم	2-3-7-4
50		روش مشبک‌مبنا	4-7-4
51	STING	الگوریتم	1-4-7-4
57	K۲	الگوریتم	5-3-6
59		یک پروژه نمونه	
59		مقدمه	
60		داده ها	
60		تحلیل اکتشافی داده ها	
70		مدل سازی	
70		یافتن تعداد بهینه خوشه ها	
73	K means		
74		نتیجه گیری	
76		پیوست	

خوشه بندی

4-1- مقدمه

خوشه‌بندی، گروه‌بندی نمونه‌های مشابه با هم در یک مجموعه داده می‌باشد. تجزیه و تحلیل خوشه‌ای، روشی برای گروه‌بندی داده‌ها یا مشاهدات با توجه به شباهت‌ها با درجه نزدیکی آن‌ها است. از طریق تجزیه و تحلیل خوشه‌ای، داده‌ها به دو رده‌ی همگن و ناهمگن (متمایز) تقسیم‌بندی می‌شوند. در روش خوشه‌بندی، هیچ رده‌ای از قبل وجود ندارد و در واقع متغیرها به صورت مستقل و وابسته تقسیم نمی‌شوند، بلکه به دنبال گروه‌هایی از داده‌ها هستیم که به هم شباهت دارند و با پیدا کردن این شباهت‌ها، توزیع و همبستگی بین داده‌ها را کشف می‌کنیم. در خوشه‌بندی داده‌هایی که در هر گروه قرار می‌گیرند، شباهت بسیار زیادی با یکدیگر دارند در حالیکه با داده‌های سایر گروه‌ها تفاوت چشم‌گیری دارند.

مسئله مهم در خوشه‌بندی عبارت است از: توزیع داده‌ها به k گروه مختلف که داده‌های هر گروه با یکدیگر مشابه بوده و داده‌های گروه‌های مختلف با یکدیگر نامتشابه باشند. در این تقسیم‌بندی، داده‌های گروه‌های مختلف باید با یکدیگر حداکثر تفاوت ممکن را با هم داشته باشند و داده‌های موجود در یک گروه باید بسیار به هم شبیه باشند. این تشابه یا عدم تشابه بر اساس معیارهای اندازه‌گیری فاصله تعریف می‌شود.

برخلاف رده‌بندی در خوشه‌بندی گروه‌ها از قبل مشخص نمی‌باشند و همچنین معلوم نیست که بر حسب کدام خصیصه‌ها گروه‌بندی صورت می‌گیرد. در نتیجه پس از انجام خوشه‌بندی باید یک فرد خبره خوشه‌های ایجاد شده را تفسیر کند و در بعضی مواقع لازم است که پس از بررسی خوشه‌ها بعضی پارامترهایی که در خوشه‌بندی در نظر گرفته شده ولی بی‌ربط بوده یا اهمیت چندانی ندارد حذف شده و جریان خوشه‌بندی از اول صورت گیرد. پس از این که داده‌ها به چند گروه منطقی و توجیه‌پذیر تقسیم شدند از این تقسیم‌بندی می‌توان برای کسب اطلاعات در مورد داده‌ها یا تقسیم داده‌های جدید استفاده کنیم. کیفیت نتیجه‌های خوشه‌بندی بستگی به روش اندازه‌گیری شباهت به کار رفته و همچنین پیاده‌سازی آن روش دارد. در روند کلی هر خوشه‌بندی، چهار مرحله اصلی وجود دارد:

1- انتخاب یا استخراج خصیصه‌ها

خوشه‌بندی ۱۰

2- طراحی یا انتخاب الگوریتم خوشه‌بندی

3- سنجش اعتبار خوشه‌بندی

4- تحلیل نتیجه‌ها.

خوشه‌بندی در موارد زیر مورد استفاده قرار می‌گیرد:

- تجزیه و تحلیل شباهت یا عدم شباهت: تجزیه و تحلیل این که کدام نقاط داده در یک نمونه به یکدیگر نزدیک تر می‌باشند.
- کاهش بُعد: اندازه و بُعد داده‌ها را می‌توان به وسیله خوشه‌بندی کاهش داد. این کاربرد بیش تر به عنوان پیش پردازش داده‌ها مورد استفاده قرار می‌گیرد.

خوشه‌بندی در زمینه‌های زیر کاربرد دارد:

- هوش تجاری
- علوم اقتصادی
- شناسایی متن
- تشخیص الگو
- جست و جوی وب
- پردازش تصویر
- بازاریابی و بیمه
- مطالعات زمین لرزه
- زیست‌شناسی
- علوم پزشکی و ژنتیک
- روان‌شناسی و جامعه‌شناسی

۴-۲- نقاط قوت روش خوشه‌بندی

- قدرت روش خوشه‌بندی به غیر مستقیم بودن آن است بدین معنی که روش را می‌توان حتی هنگامی که هیچ نوع اطلاعات قبلی از ساختار داخلی داده‌ها

خوشه‌بندی ۱۱

نداریم استفاده نمود. از این روش می‌توان برای کشف الگوهای پنهان و بهبود عملکرد روش‌های مستقیم نیز استفاده نمود.

- خوشه‌بندی را می‌توان برای داده‌های گوناگون استفاده نمود. با انتخاب درست اندازه، فاصله‌های گوناگون خوشه‌بندی را می‌توان برای بیش‌تر انواع داده‌ها استفاد نمود.
- استفاده از این روش آسان است. در این روش لازم نیست که بعضی از زمینه‌ها را به‌عنوان ورودی و بعضی دیگر را به‌عنوان خروجی در نظر بگیریم.

۳-۴- نقاط ضعف روش خوشه‌بندی

- انتخاب اندازه‌های دقیق فاصله‌ها و وزن‌ها کار آسانی نمی‌باشد. در این روش به پارامترهای اولیه نظیر تعداد خوشه‌ها، کمینه نزدیکی و خوشه‌های اولیه حساس است.
- تفسیر نتیجه‌های این روش می‌تواند مشکل باشد و به‌طور معمول نیاز به تحلیل افراد خبره است.

۴-۴- تعیین تعداد خوشه

استقلال یا وابستگی تمام متغیرها، در تعیین تعداد خوشه‌ها مؤثر است. اگر تمام متغیرها کاملاً مستقل باشند هیچ خوشه‌ای ایجاد نمی‌شود، برعکس اگر تمام متغیرها وابسته باشند آن‌گاه تمام داده‌ها تشکیل یک خوشه می‌دهند. در شرایط بین استقلال و وابستگی کامل، ما نمی‌دانیم که واقعاً چند خوشه وجود دارد. در انتخاب تعداد خوشه‌ها، تحلیل گر نقش به‌سزایی دارد. با توجه به کاربردهای متفاوت، ممکن است به تعداد خوشه‌ی بیش‌تر یا کم‌تر نیاز باشد. در بسیاری از مواقع با یک مقدار k خوشه‌بندی را انجام داده و نتیجه را ارزیابی می‌کنیم و به دنبال k دیگر می‌رویم. بعد از هر تکرار، ارزش نتیجه‌ها را به وسیله اندازه‌گیری میزان متوسط فاصله‌ها در داخل خوشه‌ها و میزان متوسط فاصله‌ها بین مراکز خوشه‌ها و یا روش‌های دیگر

خوشه‌بندی ۱۲

بررسی می‌کنیم. گاهی خوشه‌ها به وسیله قضاوت‌های ذهنی تحلیل‌گر مورد ارزیابی قرار می‌گیرند تا در کاربردهای خاصی استفاده شوند.

۴-۵- ارزیابی اعتبار در خوشه‌بندی

با توجه به این که نتیجه‌های حاصل از پیاده‌سازی الگوریتم‌های خوشه‌بندی بر داده‌ها با توجه به انتخاب پارامترها (شامل تعداد خوشه‌ها) می‌تواند بسیار متفاوت از یکدیگر باشد، لذا باید با در نظر گرفتن شاخص‌هایی، اعتبار و عملکرد خوشه‌ها مقایسه گردند تا خوشه‌هایی که بهترین تناسب را با داده‌ها دارند انتخاب شوند (آربلیتز، ۲۰۱۳).

دو معیار پایه اندازه‌گیری پیش‌نهاد شده برای ارزیابی و انتخاب خوشه‌های بهینه، شامل تراکم و جدایی است. در تراکم، داده‌های متعلق به یک خوشه باید تا حد ممکن به یکدیگر نزدیک باشند. معیار رایج برای تعیین میزان تراکم داده‌ها واریانس داده‌ها است. همچنین در معیار جدایی، خوشه‌ها باید به اندازه کافی از یکدیگر جدا باشند. در حالت کلی، جهت ارزیابی اعتبار در خوشه‌بندی از شاخص‌های درونی و بیرونی استفاده می‌گردد. شاخص‌های بیرونی از طریق ساختاری از پیش تعیین شده به ارزیابی خوشه‌ها می‌پردازد و برای محاسبه آن‌ها نیاز به داده‌های تاریخی است. در صورتی که شاخص‌های درونی، صرفاً از طریق اطلاعات درونی در داده‌های موجود اعتبار خوشه‌بندی را ارزیابی می‌کند. یکی از شاخص‌های درونی، شاخص سیلهوئت است که برای هر خوشه محاسبه می‌شود و در صورتی که به یک نزدیک باشد نشان می‌دهد که نمونه‌های درون خوشه بسیار شبیه به هم هستند و خیلی خوب خوشه‌بندی شده‌اند و در صورتی که نزدیک به صفر باشد، نشان می‌دهد تعداد نمونه‌های مرزی در خوشه زیاد است و مقدار منفی این معیار نشان می‌دهد که نمونه‌های درون خوشه شبیه به یکدیگر نبوده و عمل خوشه‌بندی نادرست صورت گرفته است.

• شاخص درونی

خوشه‌بندی ۱۳

شاخص سیلهوئت: این شاخص به صورت فرمول 1-4 می‌باشد (کافمن و روسو، 1990):

$$\begin{aligned} \text{شاخص سیلهوئت} &= \frac{\sum_{i=1}^n S(i)}{n} \\ S(i) &= \frac{b(i) - a(i)}{\max(a(i); b(i))} \quad (1-4) \\ a(i) &= \frac{\sum_{j \in C_r} d_{ij}}{n_r - 1} \\ b(i) &= \min_{s \neq r} d_{iC_s} \\ d_{iC_s} &= \frac{\sum_{j \in C_s} d_{ij}}{n_s} \end{aligned}$$

که $a(i)$ میانگین عدم تشابه i امین مشاهده نسبت به بقیه مشاهده‌ها در خوشه C_S است. C_S ، C_r نیز خوشه‌های r, s هستند. n_s, n_r نیز تعداد مشاهده‌های در خوشه‌های C_r, C_s هستند. شاخص مذکور بین دو مقدار ۱-، ۱+ قرار دارد و مقدار بزرگ‌تر برای این معیار نشان‌دهنده تعداد خوشه‌های مطلوب‌تر است.

• شاخص بیرونی

شاخص جاکارد: به منظور اندازه‌گیری ضریب شباهت بین دو مجموعه‌ی A و B از داده‌ها، ضریب شباهت جاکارد به صورت فرمول ۲-۴ تعریف می‌شود (جین و دویس، ۱۹۸۸).

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad (2-4)$$

ضریب مذکور بین دو مقدار ۰ و ۱ قرار دارد.

۴-۶- ماتریس تشابه و فاصله

به منظور رده‌بندی اشیا باید ملاکی برای اندازه‌گیری میزان تشابه یا فاصله بین آن‌ها در نظر گرفته شود. به این ترتیب اشیائی که به یکدیگر شبیه هستند می‌توانند در یک رده قرار بگیرند. برای نمایش تشابه چند شی، از ماتریس تشابه استفاده می‌شود. برای اندازه‌گیری عدم تشابه مابین دو سری داده، فاصله آن‌ها را محاسبه می‌کنند. می‌دانیم که ماتریس نمایشی منظم از اعداد و ارقام است که از تعدادی سطر و ستون تشکیل شده است. به این ترتیب ساختاری برای نمایش بهتر رابطه بین مجموعه‌ای از نقطه‌ها فراهم می‌شود. اگر تعداد سطرها و ستون‌های یک ماتریس برابر باشد، آن را ماتریس مربعی می‌گویند. اگر مقدارهای درون ماتریس فقط بر روی قطر اصلی قرار داشته باشند، ماتریس را قطری می‌نامند. همچنین اگر مقدارهای بالای قطر اصلی همگی صفر باشند، ماتریس را پایین مثلثی و به همین ترتیب مقدارهای پایین قطر اصلی صفر باشند، ماتریس را بالا مثلثی می‌گویند. برای معرفی ماتریس تشابه و فاصله ابتدا به معرفی تابع تشابه و فاصله می‌پردازیم.

۴-۶-۱- تابع تشابه

فرض کنید مجموعه داده از نقاط یک بعدی که آن‌ها را D می‌نامیم، داریم. اگر تابع S در شرایط زیر صدق کند، آن را تابع تشابه می‌نامیم.

$$۱) S(x, y) \geq 0 \quad x, y \in D$$

$$۲) S(x, y) = 1$$

$$۳) S(x, y) = S(y, x) \quad \forall x, y \in D$$

تابع تشابه تابعی نامنفی و برد آن فاصله بین ۰ تا ۱ است. همچنین رابطه آخر نشان می‌دهد که این تابع متقارن نیز هست. اگر رابطه آخر برای تابعی که به منظور اندازه‌گیری میزان تشابه به کار رفته، وجود نداشته باشد، آن را تابع نیمه تشابه می‌گویند.

۴-۶-۲- تابع فاصله

اگر d تابع با مقدارهای حقیقی و نامنفی باشد، در صورتی که خاصیت‌های زیر را برای آن صدق کند، می‌تواند به عنوان یک تابع فاصله به کار رود.

$$۱) d(x, y) \geq 0$$

$$۲) d(x, y) = 0 \leftrightarrow x = y$$

$$۳) d(x, y) = d(y, x)$$

$$۴) d(x, y) < d(x, y) + d(y, z) \quad \forall x, y, z \in D$$

۴-۶-۳- ماتریس تشابه

فرض کنید مجموعه D از نقاط x_1, \dots, x_n تشکیل شده است. بر اساس تابع تشابه S ، می‌توان برای هر زوج از اعضای D ، تشابه را اندازه‌گیری کرد و در ساختاری به شکل ماتریس قرار داد. بنا بر این اگر $S(x_i, x_j) = S_{ij}$ به معنی تشابه بین دو نقطه x_i, x_j باشد، این ماتریس به صورت زیر خواهد بود:

$$\begin{pmatrix} S_{11} & L & S_{1n} \\ \vdots & & \vdots \\ S_{n1} & L & S_{nn} \end{pmatrix}$$

همان‌طور که مشاهده می‌کنید این ماتریس یک ماتریس مربع است. اگر تابع S خاصیت‌های تابع تشابه را داشته باشد می‌توان آن را به فرم ساده‌تر نمایش داد.

$$\begin{pmatrix} S_{11} & L & S_{1n} \\ \vdots & & \vdots \\ S_{n1} & L & S_{nn} \end{pmatrix}$$

خوشه‌بندی ۱۶

از آنجایی که تابع تشابه متقارن است، در نتیجه ماتریس تشابه را می‌توان به صورت یک ماتریس بالا یا پایین مثلثی نمایش داد. توجه داشته باشید که ماتریس تشابه یک ماتریس بالا یا پایین مثلثی نیست بلکه برای نمایش آن از این حالت استفاده می‌کنیم.

۴-۶-۴- ماتریس فاصله

با توجه به مطالبی که در مورد ماتریس تشابه گفته شد، می‌توان ماتریس فاصله را نیز برای نقاط محاسبه کرد. کافی است که برای هر زوج از نقطه‌ها، تابع فاصله مورد نظر را به دست آورد و در یک ماتریس قرارا داد. در ماتریس فاصله لازم نیست که مقدارهای حقیقی غیرمنفی باشند، ممکن است عناصر اعداد صفر، منفی یا اعداد مختلط بگیرند. اگرچه در اکثر مواقع ماتریس فاصله در روی قطر اصلی دارای مقدار صفر است ولی می‌تواند غیر صفر را نیز روی قطر اصلی داشته باشند.

۴-۷- روش‌های اصلی خوشه‌بندی

روش‌های اصلی خوشه‌بندی عبارت‌اند از:

1. روش‌های افرازی
2. روش‌های سلسله مراتبی
3. روش‌های مبتنی بر چگال
4. روش مشبکی مبنا

۴-۷-۱- روش‌های افرازی

فرض کنید یک دادگان با n نمونه داریم. یک روش افرازی، k افراز از این داده‌های نمونه درست می‌کند به طوری که هر افراز یک خوشه را نشان می‌دهد و $k < n$. پس داده‌های نمونه در k گروه خوشه‌بندی شده و دارای دو شرط زیر می‌باشند:

- هر گروه حداقل یک نمونه دارد.

خوشه‌بندی ۱۷

- هر نمونه تنها به یک گروه تعلق دارد. این شرط تنها در روش‌های فازی می‌تواند قابل انعطاف باشد.

در روش افرازی برای k معلوم، یک افراز ابتدایی ایجاد می‌شود. سپس یک روش جابه‌جایی تکراری را به کار برده که تلاش به بهبود افرازبندی دارد. به این صورت که نمونه‌ها را از یک گروه به دیگر گروه‌ها می‌برد. یک معیار عمومی برای افرازبندی خوب این است که نمونه‌ها در یک خوشه به هم نزدیک یا به یکدیگر وابسته باشند و در مقابل نمونه‌ها در خوشه‌های مختلف، از یکدیگر دور یا تا حد امکان متفاوت باشند.

برای دستیابی به خوشه‌بندی بهینه در روش افرازی، به شمارش کامل همه افرازهای ممکن نیاز خواهد بود یعنی تمام حالات ممکن باید بررسی شوند که این روش برای دادگان‌های بزرگ ناممکن است، بنا بر این الگوریتم‌های زیر برای بررسی این گونه موارد استفاده می‌شوند:

1. الگوریتم k -means

2. الگوریتم k -medoids

این روش‌ها برای یافتن خوشه‌هایی به شکل کره در دادگان‌های کوچک تا متوسط به خوبی کار می‌کنند اما برای یافتن خوشه‌هایی با اشکال پیچیده و یا دارای مجموعه داده‌های بزرگ، باید توسعه داده شوند.

۴-۷-۱-۱- الگوریتم k -means

این روش، علی‌رغم سادگی یک روش پایه‌ای برای بسیاری از روش‌های خوشه‌بندی دیگر (مانند خوشه‌بندی فازی) محسوب شده و الگوریتمی بسیار ساده، قابل فهم و به‌طور منطقی قابل مقیاس‌بندی است. خوشه‌بندی k -means در سال ۱۹۶۵ توسط فوریجی ارائه شد که روش رده‌بندی بدون ناظر است. روش‌ها و الگوریتم‌های متعددی برای تبدیل نمونه‌ها به گروه‌های هم‌شکل یا مشابه وجود دارد. الگوریتم k -میانگین یکی از ساده‌ترین و

خوشه‌بندی ۱۸

محبوب‌ترین الگوریتم‌هایی است که در داده‌کاوی به‌خصوص در حوزه یادگیری بدون ناظر به کار می‌رود.

در k -means عملاً مجموعه داده‌ها به تعداد خوشه‌های از پیش تعیین شده تقسیم می‌شوند. ایده اصلی در این الگوریتم تعریف k مرکز برای هر یک از خوشه‌ها است. بهترین انتخاب برای مراکز خوشه‌ها در الگوریتم k -means قرار دادن آن‌ها در فاصله هر چه بیش‌تر از یکدیگر است. پس از آن هر نمونه در مجموعه داده به نزدیک‌ترین مرکز خوشه تخصیص می‌یابد. مرحله‌های الگوریتم k -means به‌صورت زیر است:

۱. ابتدا k مشاهده‌ی اولیه به‌عنوان هسته انتخاب می‌کنیم.
 ۲. متوسط مقدار خوشه را محاسبه می‌کنیم.
 ۳. در صورتی که فاصله مشاهده مورد نظر از میانگین خوشه خود زیاد و به خوشه دیگری نزدیک‌تر باشد این مشاهده به خوشه‌ای که نزدیک‌تر است اختصاص می‌یابد.
 ۴. این کار تا کمینه شدن تابع خطا، که به‌طور معمول مجموع فاصله‌های مشاهدات از مرکز خوشه خودش است، با تغییر نیافتن اعضای خوشه‌ها ادامه می‌یابد.
- این روش را هنگامی می‌توان استفاده کرد که مفهوم میانگین در یک مجموعه نمونه قابل تعریف باشد. برای رفع مشکلات الگوریتم k -means با تغییراتی روی آن این روش را توسعه داده و به جای استفاده از میانگین به‌عنوان مراکز خوشه از نماهای خوشه به‌عنوان نماینده یک خوشه استفاده می‌شود. این روش k -modes نام دارد. اگر به‌جای مرکز یا وسط یک خوشه، از میانه آن خوشه استفاده کنیم، آن‌گاه روش نسبت به داده‌های دور از مرکز حساس نمی‌شود زیرا میانه از مقدارهای بزرگ تأثیر نمی‌پذیرد.

عیب‌های الگوریتم k -means:

- معیاری برای تعیین تعداد و مراکز اولیه خوشه‌ها وجود ندارد.
- اگر در خوشه‌ای هیچ داده‌ای وجود نداشته باشد راهی برای تغییر و بهبود آن وجود ندارد.
- برای داده‌های دم‌سنگین مناسب نیست.

خوشه‌بندی ۱۹

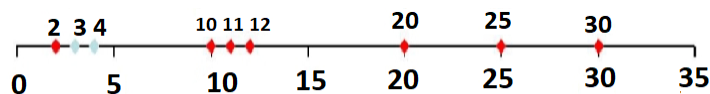
- جواب نهایی به انتخاب مراکز اولیه خوشه‌ها بستگی دارد.
- این روش برای کشف خوشه‌هایی با شکل‌های پیچیده مناسب نیست.
- امکان تولید خوشه‌های خالی توسط این روش وجود دارد.

مثال ۴-۱- فرض کنید سن ۹ نفر به صورت $\{2, 4, 10, 11, 12, 20, 25, 30, 35\}$ است. با استفاده

از روش k - میانگین، داده‌های به دست آمده را به دو خوشه تقسیم کنید؟

گام اول: ابتدا باید دو عدد تصادفی به عنوان مراکز خوشه‌ها تولید کنیم. فرض کنید اعداد ۳ و ۴ مراکز اولیه خوشه‌ها باشند.

حال هر یک از داده‌ها را با توجه به فاصله‌ی هر یک از آن‌ها از این مراکز به یکی از خوشه‌ها نسبت می‌دهیم. به عنوان مثال عدد ۲ به مرکز ۳ و عدد ۱۰ به مرکز ۴ نزدیک‌تر است.

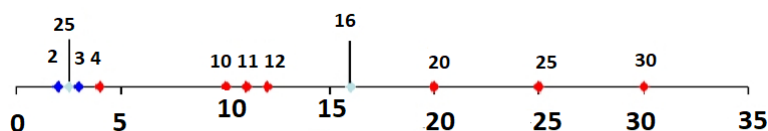


بنا بر این با ادامه‌ی روند بالا، دو خوشه‌ی زیر به دست می‌آیند.

$$C_1 = \{2, 3\}$$

$$C_2 = \{4, 10, 11, 12, 20, 25, 30\}$$

گام دوم: در این گام، مراکز جدید هر خوشه را محاسبه کرده و با توجه به مراکز جدید، دوباره عمل واگذاری داده‌ها را انجام می‌دهیم. مراکز جدید ۵/۲ و ۱۶ خواهند بود.



مشاهده می‌شود عدد ۴ که در مرحله قبل به خوشه دوم نسبت داده شد، در این مرحله به مرکز خوشه اول نزدیک‌تر است. در نتیجه باید عدد ۴ را در خوشه اول قرار دهیم. خوشه‌های جدید به شکل زیر خواهند بود:

خوشه‌بندی ۲۰

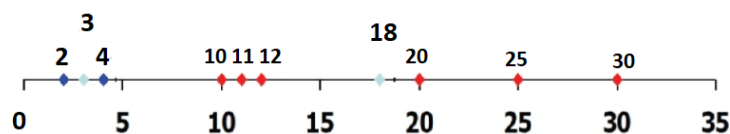
$$C_1 = \{2, 3, 4\}$$

$$C_2 = \{10, 12, 20, 30, 11, 25\}$$

روند بالا را تا زمانی که شرط توقف برقرار شود ادامه می‌دهیم.

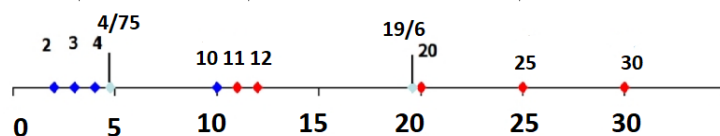
مرحله‌های بعدی به صورت زیر خواهند بود:

$$\mu_1 = 3, \mu_2 = 18 \rightarrow C_1 = \{2, 3, 4, 10\} \quad C_2 = \{12, 20, 30, 11, 25\}$$



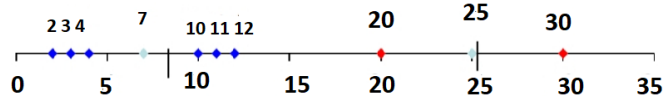
مرحله بعدی به صورت زیر خواهد بود:

$$\mu_1 = 4/75, \mu_2 = 19/6 \rightarrow C_1 = \{2, 3, 4, 10, 11, 12\} \quad C_2 = \{20, 30, 25\}$$



و مرحله نهایی به صورت زیر است که در آن مراکز ثابت شده است و خوشه‌ها تغییر نمی‌کنند.

$$\mu_1 = 7, \mu_2 = 25 \rightarrow C_1 = \{2, 3, 4, 10, 11, 12\} \quad C_2 = \{20, 30, 25\}$$

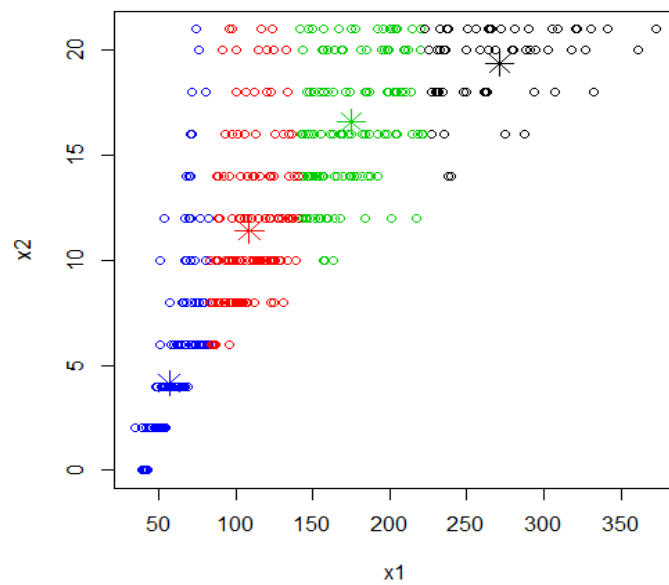


مثال ۲-۴- در مثال ۳-۳ به توضیح داده‌های وزن جوجه پرداختیم. همان‌طور که اشاره شد داده‌ها مربوط به ۵۰ جوجه و ۴ رژیم غذایی متفاوت است. حال روش *k-means* را برای داده‌های وزن جوجه با استفاده از بسته نرم‌افزاری *amap* پیاده‌سازی می‌کنیم. (کد ۱)

خوشه‌بندی ۲۱

```
> data(ChickWeight)
> ChickWeight
> install.packages("amap")
> library(amap)
> x1=ChickWeight$weight
> x2=ChickWeight$Time
> x=cbind(x1,x2)
> x
> cl<-Kmeans(x,4, iter.max =85, method = "euclidean")
> plot(x, col = cl$cluster)
> points(cl$centers, col = 1:4, pch = 8, cex=2)
```

خروجی دستورها را در شکل 4-1 مشاهده می‌کنیم.



شکل 4-1: خروجی الگوریتم *k-means*

در نمودار بالا تعداد خوشه‌ها را برابر 4 و حداکثر تعداد تکرار تا رسیدن به بهترین خوشه را 85 و معیار فاصله‌ی نقطه تا میانگین، فاصله اقلیدسی در نظر گرفته شده است.

مثال 4-3- داده‌های گل زنبق را در نظر بگیرید. همان‌طور که می‌دانیم این مجموعه داده دارای یک متغیر به نام *Species* است و گروه هر یک از گل‌ها را در سه سطح مشخص می‌کند که آن را از مجموعه داده‌ها حذف کرده‌ایم. پس تعداد خوشه‌ها را در تابع *kmeans()* برابر سه در نظر خواهیم گرفت. تعداد اعضای هر خوشه همان‌طور که مشاهده می‌کنیم برابر با 96,21,33 است.

میانگین متغیرها در هر خوشه حساب شده و خروجی *Clustering Vector* در حقیقت خوشه‌ی هر نمونه را تعیین می‌کند که این الگوریتم مجموع توان دوم درون خوشه‌ها یعنی *Within cluster sum of square by cluster* را برای هر خوشه کمینه می‌کند و مجموع توان دوم بیرون خوشه‌ای را بیشینه خواهد کرد. این مدل 79٪ از کل واریانس داده‌ها را بیان می‌کند و هر چه قدر تعداد خوشه‌ها بیش‌تر باشد این عدد نیز به 100٪ نزدیک خواهد شد. برای مقایسه‌ی خوشه‌های به دست آمده از این مدل و گروه‌بندی واقعی که توسط متغیر *Species* تعیین شده است از تابع *table()* استفاده کرده‌ایم. بنا بر این از گونه اول 17 مورد در خوشه دوم و 33 مورد در خوشه سوم قرار گرفته‌اند و از گونه دوم 46 مورد در خوشه اول و 4 مورد در خوشه دوم قرار گرفته و برای گونه سوم تمام موارد در خوشه اول قرار گرفته است. (کد 2)

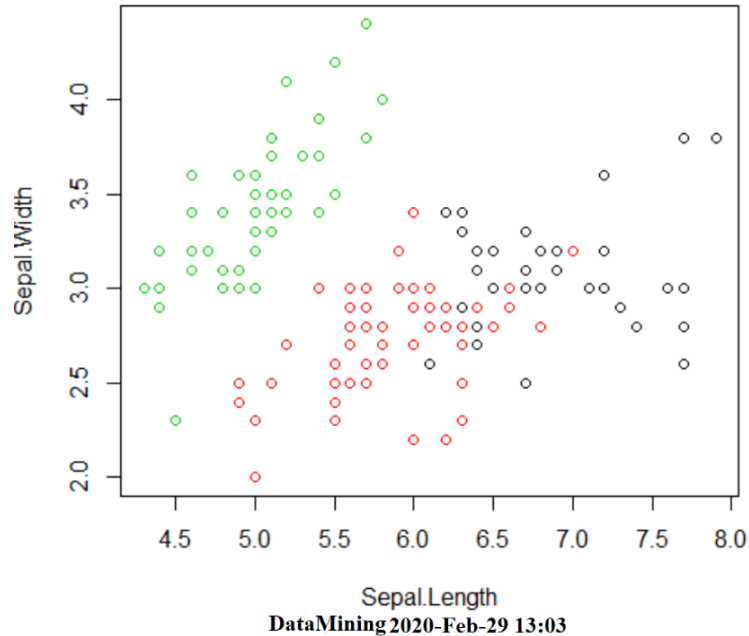
```

> data(iris)
> iris2=iris
> iris2$Species=NULL
> kmeans.result=kmeans(iris2,3)
> kmeans.result
K-means clustering with 3 clusters of sizes 96, 21, 33
Cluster means:
  Sepal.Length Sepal.Width Petal.Length Petal.Width
1    6.314583    2.895833    4.973958    1.7031250
2    4.738095    2.904762    1.790476    0.3523810
3    5.175758    3.624242    1.472727    0.2727273

Clustering vector:
 [1] 3 2 2 2 3 3 3 3 2 2 3 3 2 2 3 3 3 3 3 3 3 3 3 3
2 2 3 3 3 2 2 3 3 3 2 3 3 3 2 3 3
 [42] 2 2 3 3 2 3 2 3 3 1 1 1 1 1 1 1 2 1 1 2 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 [83] 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 2 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
[124] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1
Within cluster sum of squares by cluster:
[1] 118.651875 17.669524 6.432121
(between_SS / total_SS = 79.0 %)
Available components:
[1] "cluster" "centers" "totss" "withinss" "tot.withinss"
[6] "betweenss" "size" "iter" "ifault"
> table(iris$Species, kmeans.result$cluster)
      1  2  3
setosa  0 17 33
versicolor 46  4  0
virginica 50  0  0
> plot(iris2[,c("Sepal.Length", "Sepal.Width")], col=kmeans.result$cluster, sub=paste("DataMining", format(Sys.time(), "%Y-%b-%d %H:%S")))
> points(kmeans.result$centers[,c("Sepal.Length", "Sepal.Width")], col=1:3, pch=8, cex=2)

```

برای رسم نمودار دو بُعدی این خوشه‌بندی، از تابع `plot()` استفاده کرده و توسط تابع `points()` می‌توان مراکز هر خوشه را مشخص کرد.



شکل 4-2: نمودار دو بُعدی خوشه‌بندی با داده‌های گل زنبق

۴-۷-۱-۲- الگوریتم k -medoids

الگوریتم k -means به داده‌های دور افتاده حساس است و از آن‌جا که فاصله‌ی چنین داده‌هایی با اکثر داده‌ها بسیار زیاد است، هنگام قرار گرفتن در یک خوشه به‌صورت چشم‌گیری باعث انحراف و تغییر میانگین خوشه می‌شود. این موضوع به صورت غیر عمدی روی انتساب نمونه‌ها دیگر به خوشه‌ها تأثیر می‌گذارد. به جای در نظر گرفتن مقدار میانگین نمونه‌ها موجود در خوشه به عنوان یک نقطه ارتجاعی، می‌توان از نمونه‌های واقعی برای نمایش خوشه‌ها استفاده کرد. مدوئیدها نمونه‌هایی هستند که در مرکزی‌ترین محل یک خوشه قرار دارند. در این الگوریتم می‌توان به جای استفاده از مرکز یک خوشه به عنوان مرجع از مدوئیدها استفاده کرد. روش افراز به دنبال کمینه کردن مجموع اختلاف میان هر نمونه و نماینده نظیر آن است. گروه‌بندی n نمونه در k خوشه با کمینه کردن خطای مطلق، ایده اصلی روش k -مدوئید است.

یکی از الگوریتم‌های رایج در میان روش‌های خوشه‌بندی k -مدوئید، الگوریتم *PAM* است. در این الگوریتم همانند الگوریتم *k - means* در ابتدا برخی از نمونه‌ها به صورت اختیاری به عنوان نمایندگان اولیه انتخاب می‌شوند. پس از آن بررسی می‌شود که آیا جایگزین نمودن یک نمونه غیر نماینده یا یک نمونه نماینده، کیفیت خوشه‌بندی را بهبود می‌بخشد یا خیر. کلیه‌ی جایگزینی‌های ممکن آزمایش می‌شوند. این جایگزینی نمونه‌ها تا هنگامی که خوشه‌بندی بهبود یابد، ادامه می‌یابد. از آنجا که یک نماینده یا به عبارتی یک مدوئید نسبت به مقدار میانگین کم‌تر تحت تأثیر داده‌های دور افتاده قرار می‌گیرد، روش k -مدوئید در مواجهه با داده‌های دور افتاده و خطا مقاوم‌تر از روش *k - means* است.

اما در استفاده از هر دو روش لازم است کاربر مقدار k یعنی تعداد خوشه‌ها را مشخص کند. الگوریتم *PAM* که در رده‌ی خوشه‌بندی مبتنی بر افراز می‌گنجد، عملکرد خوبی را در مواجهه با مجموعه داده‌های کوچک دارد. اما قابلیت اجرایی مقیاس‌پذیر را بر روی داده‌های بزرگ ندارد. جهت کار با داده‌های با اندازه بالا می‌توان از روشی مبتنی بر نمونه‌گیری با نام *CLARA* استفاده کرد. در این روش به جای این که از همه داده‌ها استفاده شود، تنها بخشی از آن‌ها با روش‌های نمونه‌گیری برای بررسی انتخاب می‌شوند. الگوریتم *CLARA* عمل خوشه‌بندی را چندین بار بر روی نمونه‌های مختلف اجرا می‌کند و بهترین خوشه‌بندی به عنوان خروجی برگردانده می‌شود. این الگوریتم به تعداد نمونه‌ها یا به عبارتی دیگر به اندازه‌ی نمونه‌گیری بستگی دارد. الگوریتم *PAM* در میان یک مجموعه داده‌ها به دنبال k نماینده است، درحالی که الگوریتم *CLARA* همین کار را بر روی نمونه‌های منتخبی از داده‌ها انجام می‌دهد. در الگوریتم *PAM* جهت یافتن نمایندگان بهتر، جانشینی و تعویض همه نمونه‌ها با نمایندگان جاری بررسی می‌شود، در حالی که در الگوریتم *CLARA* این بررسی تنها محدود به داده‌های نمونه‌گیری شده به صورت تصادفی می‌شود. پیچیدگی محاسبات *PAM* در هر تکرار متناظر با $O(k^2 + k(n - k))$ می‌باشد که k تعداد خوشه‌ها، S تعداد نمونه انتخاب شده و n کل نمونه‌ها می‌باشد. لذا پیچیدگی الگوریتم کاهش می‌یابد و از مرتبه تعداد نمونه انتخاب شده است.

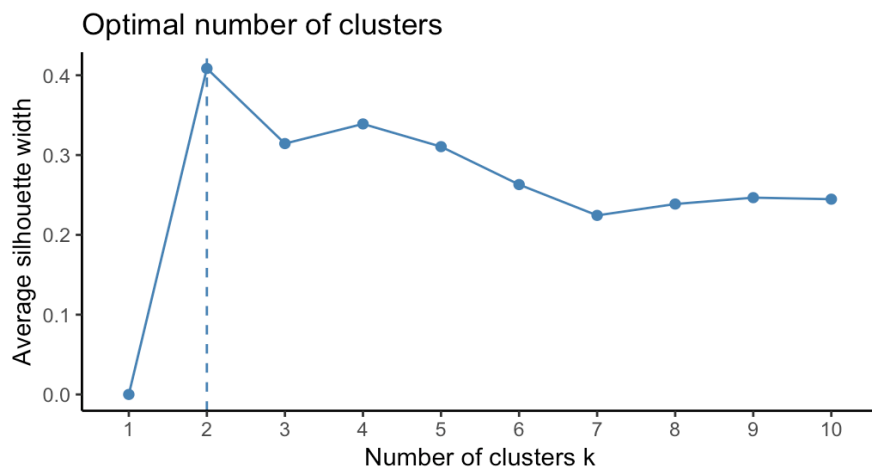
مثال ۴-۴- مجموعه داده‌های *USArrests* شامل اطلاعاتی مربوط به میزان قتل، اذیت و آزار، جمعیت شهری و میزان تجاوز در ایالات مختلف ایالات متحده‌ی آمریکا را نشان می‌دهد. روی این داده‌گان الگوریتم *PAM* را پیاده سازی می‌کنیم. (کد ۳)

```
> library(datasets)
> data(USArrests)
> df=scale(USArrests)
> head(df,n=3)
```

از طریق دستورهایی زیر در نرم‌افزار *R* می‌توان تعداد خوشه‌ها را برآورد کرد.

```
> install.packages(c("cluster","factoextra"))
> library(cluster)
> library(factoextra)
#the optimal number of clusters
> fviz_nbclust(data,pam.method="silhouette")
+geom_vline(xintercept=3)
```

خروجی این دستورها به شکل ۳-۴ نشان داده شده است:



شکل ۳-۴: تعداد خوشه‌ها

دستورهای زیر، الگوریتم *PAM* را به ازای $k = 2$ در نرم‌افزار *R* محاسبه می‌کند:

```
#pam(x,metric="euclidean",stand=FALSE)
> PAM=pam(df,2)
```

> Print(PAM)

خروجی اجرای این تابع در نرم افزار به صورت زیر می باشد:

```
Medoids:
      ID  Murder  Assault  UrbanPop  Rape
New Mexico 31  0.8292944  1.3708088  0.3081225  1.1603196
Nebraska  27 -0.8008247 -0.8250772 -0.2445636 -0.5052109
Clustering vector:
      Alabama  Alaska  Arizona  Arkansas  California  Colorado  Connecticut
      1         1         1         2         1         1         2
      Delaware  Florida  Georgia  Hawaii     Idaho      Illinois  Indiana
      2         1         1         2         2         1         2
      Iowa      Kansas   Kentucky Louisiana  Maine      Maryland  Massachusetts
      2         2         2         1         2         1         2
      Michigan  Minnesota Mississippi Missouri    Montana  Nebraska  Nevada
      1         2         1         1         2         2         1
      New Hampshire  New Jersey  New Mexico  New York  North Carolina  North Dakota  Ohio
      2         2         1         1         1         2         2
      Oklahoma      Oregon  Pennsylvania  Rhode Island  South Carolina  South Dakota  Tennessee
      2         2         2         2         1         2         1
      Texas         Utah     Vermont      Virginia    Washington  West Virginia  Wisconsin
      1         2         2         2         2         2         2
      Wyoming
      2
Objective function:
      build  swap
1.441358 1.368969
Available components:
[1] "medoids" "id.med"  "clustering" "objective" "isolation" "clusinfo" "silinfo" "diss"
[9] "call"   "data"
```

با توجه به خروجی نرم افزار، قسمت اول خروجی یعنی مدوئید نشان دهنده‌ی ماتریسی با سطرها‌ی مدوئید و ستون‌های مربوط به متغیرها است. قسمت دوم خروجی یعنی *Clustering Vector* یک بردار از اعداد حقیقی است که نشان دهنده‌ی اختصاص نقاط به هر یک از دو خوشه می باشد که به طور خلاصه از طریق دستورهای زیر در نرم افزار قابل ملاحظه است:

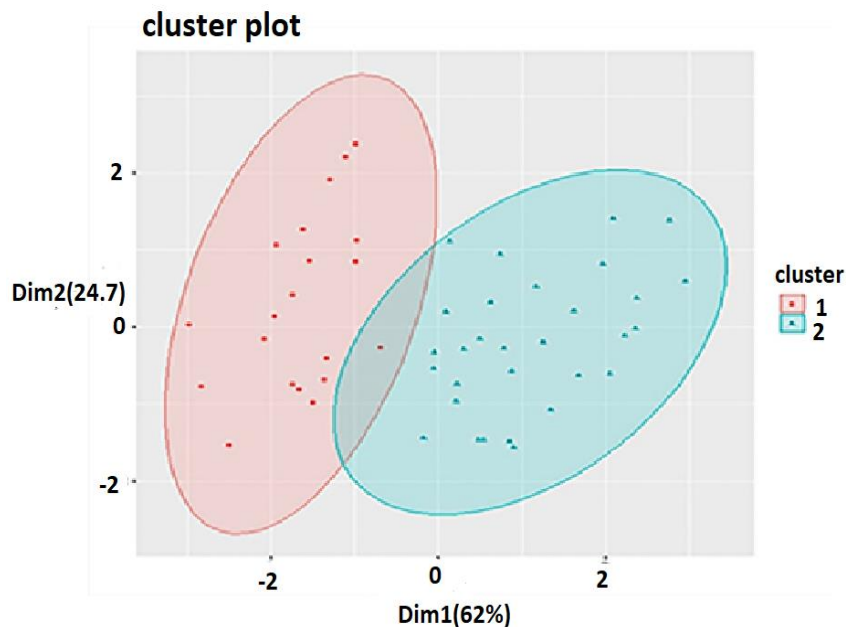
```
> d=cbind(USArrests,cluster=PAM$cluster)
> head(d,n=3)
```

	Murder	Assault	UrbanPop	Rape	cluster
Alabama	13.2	236	58	21.2	1
Alaska	10.0	263	48	44.5	1
Arizona	8.1	294	80	31.0	1

```
> |
```

خوشه‌های تشکیل شده از این داده‌ها را می‌توان با دستورهای زیر مشاهده کرد:

```
#visualizing PAM clusters
> library(factoextra)
> fviz_cluster(PAM,geom="point",ellipse.type="norm")
```



شکل ۴-۴ خوشه‌های تشکیل شده

مثال 4-5- داده‌های گل زنبق را در نظر بگیرید. در مثال قبل از تابع *pam()* برای خوشه‌بندی استفاده نمودیم. حال در این مثال از تابع *pamk()* استفاده می‌کنیم. ابتدا تعداد خوشه‌های بهینه را در بازه‌ی [۲,۱۰۰] به دست آورده‌ایم که برابر با ۲ شده است. (کد ۴)

```
> library(fpc)
> pamk.result=pamk(iris2,2:100)
> pamk.result$nc
[1] 2
```

سپس به مشخص کردن مراکز خوشه‌ها می‌پردازیم:

```
> pamk.result$pamobject$medoids
      Sepal.Length Sepal.Width Petal.Length Petal.Width
[1,]           5.0           3.4           1.5           0.2
[2,]           6.2           2.8           4.8           1.8
```

حال برای مقایسه‌ی خوشه‌های به دست آمده از این مدل و گروه‌بندی واقعی که توسط متغیر *species* تعیین شده است، می‌توانیم از تابع *table()* استفاده کنیم:

```
> table(pamk.result$pamobject$clustering,iris$Species)

      setosa versicolor virginica
1         50           1         0
2          0          49        50
```

همان‌طور که مشاهده می‌کنیم گونه‌های *verginica,versicolor* به یکدیگر شباهت دارند به طوری که در خوشه 2 قرار گرفته‌اند و گونه‌ی *setosa* با سایر گونه‌ها تفاوت دارد و در خوشه‌ی یک قرار گرفته است.

۴-۷-۲- روش‌های سلسله مراتبی

- خوشه‌بندی سلسله مراتبی، روشی است که در گروه‌بندی یا رده‌بندی داده‌ها به کار می‌رود. به عبارتی این روش، یک روش خوشه‌بندی می‌باشد که هدف آن ساخت یک سلسله مراتب از خوشه‌ها است. در روش خوشه‌بندی سلسله مراتبی، به خوشه‌های نهایی ساختاری سلسله مراتبی که به طور معمول به صورت درختی است،

نسبت داده می شود. به این درخت سلسله مراتبی درخت‌واره‌نگار می گویند. نتیجه - های یک خوشه‌بندی سلسله مراتبی عموماً به شکل یک درخت‌واره‌نگار نمایش داده می شوند.

روش‌های خوشه‌بندی سلسله مراتبی به دو دسته تقسیم می شوند: تجمعی و تجزیه‌ای.

- **تجمعی:** رویکرد این دسته پایین به بالا می باشد و با شروع از پایین، در هر مرحله دو خوشه که کم‌ترین تفاوت را دارند با یکدیگر تجمیع شده و یک خوشه جدید را شکل می دهند. خوشه‌های جدید در سطوح بالاتر قرار گرفته و این مرحله‌ها تا زمانی که تعداد خوشه‌ها به یک برسد تکرار می شود.
- **تجزیه‌ای:** رویکرد این دسته بالا به پایین می باشد و با شروع از بالا، در هر مرحله یک خوشه به خوشه‌های کوچک‌تری تجزیه می شود که در سطح پایین‌تر قرار می گیرند.

برای این که بدانیم کدام خوشه‌ها باید با هم تجمیع بشوند یا از یکدیگر تقسیم بشوند باید معیاری از تفاوت بین خوشه‌ها تعریف شود. در اکثر روش‌ها، این معیار به کمک تعریف یک تابع امتیاز و یک معیار پیوند حاصل می شود. تابع امتیاز فاصله‌ی بین دو تک مشاهده را تعیین کرده و معیار پیوند فاصله‌ی بین دو مجموعه مشاهده را توسط تابعی از فاصله دو به دو بین مشاهدات هر مجموعه تعریف می کند. انتخاب یک تابع امتیاز مناسب شکل خوشه‌ها را تحت تأثیر قرار می دهد، زیرا به ازای یک تابع امتیاز چند مشاهده می توانند به یکدیگر نزدیک باشند ولی به ازای تابع امتیاز دیگری فاصله‌ی آن‌ها از هم افزایش یابد. به عنوان مثال در یک فضای دو بُعدی فاصله‌ی بین نقاط $(0,0)$ و $(1,0)$ بنابر روش‌های معمول یک می باشد، اما فاصله‌ی بین همین دو نقطه با در نظر گرفتن فاصله منتهن مقدار 2، با در نظر گرفتن فاصله اقلیدسی جذر ۲ می باشد، و با در نظر گرفتن فاصله بیشینه مقدار 1 می باشد.

روش‌های خوشه‌بندی مبتنی بر سلسله مراتب به صورت زیر است:

- الگوریتم *Brich*

- الگوریتم *chameleon*
- الگوریتم *DIANA*
- الگوریتم *AGNES*

۴-۷-۲-۱- الگوریتم *Brich*

الگوریتم *Brich* با کنار هم قرار دادن خوشه‌بندی سلسله‌مراتبی و روش‌های دیگر خوشه‌بندی نظیر افراز تکرارشونده با هدف خوشه‌بندی مقدار بزرگی از داده‌های عددی، طراحی شده است.

در الگوریتم *Brich* برای تلخیص یک خوشه، از مفهوم ویژگی خوشه‌بندی و جهت نمایش سلسله‌مراتب خوشه از درختی موسوم به درخت ویژگی خوشه‌بندی استفاده می‌شود. با کمک این ساختارها، این روش سرعت مناسب و قابلیت مقیاس‌پذیری خوبی را در مواجهه با داده‌های بزرگ و حتی داده‌های جریانی از خود نشان می‌دهد. خوشه‌ای از n شیء یا نقطه‌ای d بُعدی را در نظر بگیرید. ویژگی خوشه‌بندی این خوشه یک بردار سه‌بعدی است که در آن اطلاعات خوشه‌ای نمونه‌ها خلاصه شده است. ویژگی خوشه‌بندی به صورت زیر معرفی می‌شود:

$$CF = (n, LS, SS)$$

که در آن LS مجموع خطی n نقطه (به عبارت دیگر $\hat{a}_{i=1}^n x_i$) و SS مجموع توان دوم نقاط (یعنی $\hat{a}_{i=1}^n x_i^2$) را نشان می‌دهند. ویژگی خوشه‌بندی خلاصه‌ای از آماره‌های یک خوشه را نشان می‌دهد. با کمک این ویژگی می‌توان به سادگی آماره‌های سودمندی در مورد یک خوشه را به دست آورد که این آماره‌ها و اطلاعات عبارت‌اند از: گرانیگاه خوشه، شعاع خوشه و قطر خوشه. تلخیص یک خوشه با کمک ویژگی خوشه‌بندی باعث می‌شود تا از ذخیره‌ی اطلاعات جزئی در مورد نمونه‌ها به صورت جداگانه جلوگیری شود. در عوض تنها به اندازه‌ی ثابتی از حافظه برای ذخیره‌ی این ویژگی‌ها نیاز است. کلید موفقیت الگوریتم *Brich* در مقدار حافظه‌ی مصرفی است. درخت ویژگی خوشه‌بندی یا همان *CF - tree*

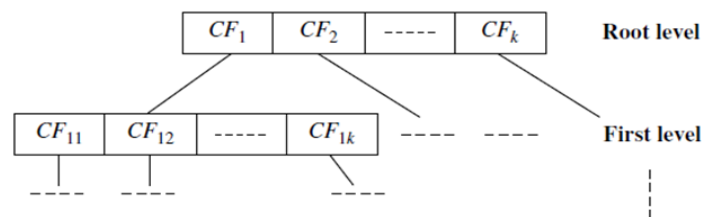
درخت متعادلی است که در خود، ویژگی‌های خوشه‌بندی را برای یک خوشه‌بندی سلسله‌مراتبی نگهداری می‌کند. هر $CF - tree$ دارای دو پارامتر با نام‌های فاکتور شاخه‌بندی B و حد آستانه‌ی T است. حداکثر تعداد فرزندان گره‌های غیر پایانی درخت، با فاکتور شاخه‌بندی B مشخص می‌شود و حد آستانه‌ی T حداکثر قطر (فاصله) خوشه‌های ذخیره شده در گره‌های پایانی یعنی برگ‌ها را تعیین می‌کند. با کمک این پارامترها می‌توان اندازه‌ی درخت حاصل را کنترل کرد.

یکی از مشخصات مهم و قابل توجه در الگوریتم $Brich$ کمینه کردن زمان لازم برای ورودی/خروجی است. الگوریتم $Brich$ از یک رویکرد خوشه‌بندی چند مرحله‌ای بهره می‌برد که در آن پس از پیمایش اول داده‌ها و خوشه‌بندی مناسب آن‌ها با پیمایش بعدی درصدد بهبود کیفیت و نتیجه‌های حاصل بر می‌آید. مرحله‌های ابتدایی این الگوریتم عبارت‌اند از:

1. الگوریتم با پیمایش دادگان‌ها، یک $CF - tree$ اولیه در حافظه می‌سازد. این درخت شکلی از فشردسازی چندسطحی داده‌ها است که سعی می‌کند تا ساختار خوشه‌بندی ذاتی موجود در داده‌ها را حفظ کند.
2. الگوریتم با انتخاب یک روش خوشه‌بندی، گره‌های برگ $CF - tree$ را خوشه‌بندی می‌کند. سپس در این مرحله خوشه‌های خلوت شناسایی شده و حذف می‌شوند و گروه‌هایی از خوشه‌ها که به یکدیگر بسیار نزدیک تر هستند، با یکدیگر ادغام می‌شوند تا خوشه‌ی بزرگ‌تری ایجاد شود.

در مرحله‌ی اول همان‌طور که نمونه‌ها درج می‌شوند، درخت $CF - tree$ به صورت پویا ساخته می‌شود. بنا بر این، این روش از نوع افزایشی است. اندازه‌ی $CF - tree$ را می‌توان با کمک مقدار حد آستانه تغییر داد. اگر اندازه‌ی حافظه‌ی مورد نیاز برای نگهداری $CF - tree$ بیش‌تر از اندازه‌ی حافظه‌ی اصلی باشد، می‌توان مقدار بزرگ‌تری را برای حد آستانه در نظر گرفت و درخت را بازسازی کرد.

فرایند بازسازی، با ساخت یک درخت جدید از برگ‌های درخت قدیمی انجام می‌شود. بنا بر این برای اجرای این فرایند لازم است اطلاعات مربوط به نمونه‌ها یا نقاط دوباره خوانده شوند. لذا برای ساخت درخت، داده‌ها تنها یک بار خوانده می‌شوند. پیمایش‌های اضافی داده‌ها با استفاده از سلسله‌مراتب‌ها و روش‌های متعدد باعث می‌شود داده‌های دور افتاده شناخته شوند و کیفیت درخت‌های ویژگی خوشه‌بندی یعنی *CF-tree* افزایش یابد. پیچیدگی الگوریتم برابر با $O(n)$ است که در آن n به تعداد نمونه‌هایی که باید خوشه‌بندی شوند، اشاره می‌کند. چون الگوریتم *Brich* از مفهوم‌هایی چون شعاع و قطر برای کنترل کران‌های یک خوشه استفاده می‌کند، هنگامی که خوشه‌ها کروی شکل نباشند، این الگوریتم عملکرد مناسبی از خود نشان نمی‌دهد. مثالی از درخت ویژگی خوشه‌بندی یا همان *CF-tree* در شکل ۴-۵ نشان داده شده است:

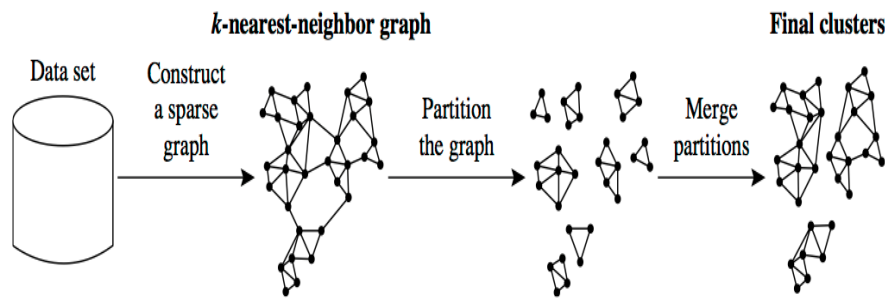


شکل ۴-۵: درخت ویژگی خوشه‌بندی

۴-۷-۲-۲- الگوریتم *chameleon*

الگوریتم *Chameleon* یک الگوریتم خوشه‌بندی سلسله‌مراتبی است که با کمک مدل‌سازی پویا، شباهت میان زوج خوشه‌ها را تعیین می‌کند. در این الگوریتم شباهت خوشه بر اساس چگونگی اتصال مناسب نمونه‌ها درون خوشه و مجاورت و نزدیکی خوشه‌ها ارزیابی می‌شود. بنا بر این دو خوشه‌ای ادغام خواهند شد که دارای اتصالات قوی در درون خود و به یکدیگر بسیار نزدیک باشند. بنا بر این الگوریتم وابسته به یک مدل ایستا که توسط کاربر تهیه شده است، نیست و به صورت خودکار و بر اساس خصیصه‌های داخلی خوشه‌های

ادغام شونده عمل می کند. در شکل ۴-۶ چگونگی کار الگوریتم *Chameleon* نشان داده شده است:



شکل ۴-۶: چگونگی کار کرد الگوریتم *Chameleon*

این الگوریتم از یک رویکرد خاص برای ساخت یک گراف استفاده می کند. رأس های این گراف، نمونه ها را نشان می دهند و وجود یک یال میان دو رأس نشان می دهد که یکی از نمونه ها نظیر رأس ها در میان k شیء مشابه به شیء دیگر قرار دارد. وزن های یال ها تشابه میان نمونه ها را نشان می دهند. الگوریتم *Chameleon* با کمک یک الگوریتم افراز گراف، گراف حاصل از مرحله ی قبل را به گونه ای به بخش های کوچک تر تقسیم می کند که برش یال ها کمینه باشند. بنا بر این خوشه ی C به نحوی به زیرخوشه های C_i و C_j تقسیم می شود که وزن یال های انتخابی برای برش کمینه باشند. الگوریتم *Chameleon* به هم پیوستگی مطلق میان خوشه های C_i و C_j را ارزیابی می کند و سپس این الگوریتم از یک روش خوشه بندی سلسله مراتبی تجمیعی استفاده می کند تا به صورت تکرار شونده زیرخوشه ها را بر اساس تشابه آن ها ادغام کند. جهت تعیین دو زیرخوشه ای که بیش ترین شباهت را با یکدیگر دارند، به هم پیوستگی و نزدیکی خوشه ها بررسی می شوند. الگوریتم *Chameleon* جهت تعیین تشابه میان خوشه های C_i و C_j از دو معیار به هم پیوستگی نسبی یعنی $RI(C_i, C_j)$ و نزدیکی نسبی یعنی $RC(C_i, C_j)$ استفاده می کند. الگوریتم *Chameleon* دارای توانایی بالاتری در کشف خوشه هایی با کیفیت مطلوب و شکل

دلخواه نسبت به الگوریتم‌های شناخته شده‌ای مانند *BIRCH* و الگوریتم چگالی مبنا *DBSCAN* می‌باشد.

۴-۷-۲-۳- الگوریتم *AGNES*

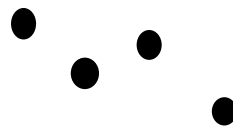
در الگوریتم *AGNES*، ابتدا هر کدام از نمونه‌ها در داخل یک خوشه قرار می‌گیرند. سپس خوشه‌ها مرحله به مرحله و بر اساس بعضی معیارها با یکدیگر ترکیب می‌شوند. اگر فاصله بین نمونه‌های هر خوشه با نمونه‌های خوشه‌های دیگر را حساب کنیم و دو نمونه متعلق به خوشه‌ها، کم‌ترین فاصله را داشته باشند آن دو خوشه با یکدیگر ترکیب می‌شوند. ترکیب خوشه‌ها آن‌قدر ادامه می‌یابد تا همه نمونه‌ها درون یک خوشه قرار گیرند. در این روش تحلیل‌گر از قبل می‌تواند تعداد خوشه‌ها را انتخاب کند و از آن برای پایان شرط الگوریتم استفاده کند. معیارهای متفاوتی برای فاصله بین خوشه‌ها وجود دارد که پیوند تکی، پیوند کامل، پیوند متوسط و پیوند مرکزی از جمله این معیارها هستند. از جمله روش‌های مهم الگوریتم *AGNES* می‌توان به *Singe link* و *Complete link* و *Average link* اشاره کرد که به توضیح هر یک از آن‌ها در ادامه می‌پردازیم.

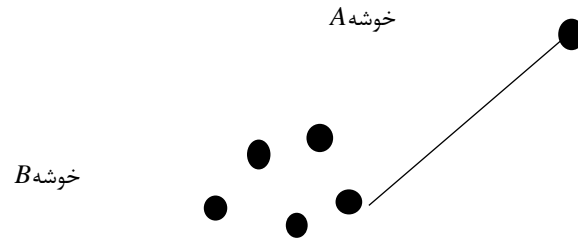
۴-۷-۲-۱- خوشه‌بندی با روش *sin gel-link*

این روش یکی از ساده‌ترین روش‌های خوشه‌بندی و جزء روش‌های خوشه‌بندی سلسله مراتبی محسوب می‌شود. به این روش خوشه‌بندی نزدیک‌ترین همسایه نیز گفته می‌شود. در این روش برای محاسبه شباهت بین دو خوشه *A* و *B* از معیار زیر استفاده می‌شود:

$$d_{AB} = \min d_{ij}, i \in A, j \in B$$

که *i* یک نمونه داده متعلق به خوشه *A* و *j* یک نمونه داده متعلق به خوشه *B* می‌باشد. در واقع در این روش بین دو خوشه، کم‌ترین فاصله بین یک عضو از یکی با یک عضو از دیگری است. در شکل ۴-۷ این مفهوم نشان داده شده‌است.





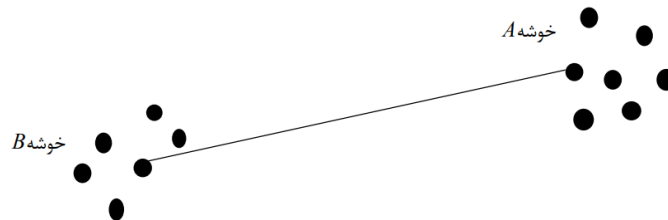
شکل ۴-۷: شباهت بین دو خوشه در روش *single-link* برابر است با کمترین فاصله بین داده‌های دو خوشه

۴-۷-۲-۳-۲ خوشه‌بندی با روش *complete-link*

این روش همانند *single-link* جزء روش‌های خوشه‌بندی سلسله مراتبی محسوب می‌شود. به این روش خوشه‌بندی دورترین همسایه نیز گفته می‌شود. در این روش برای محاسبه شباهت بین دو خوشه A و B از معیار زیر استفاده می‌شود:

$$d_{AB} = \max d_{ij}, i \in A, j \in B$$

که i یک نمونه داده متعلق به خوشه A و j یک نمونه داده متعلق به خوشه B می‌باشد. در واقع در این روش شباهت بین دو خوشه بیشترین فاصله بین یک عضو از یکی و یک عضو از دیگری است. در شکل ۴-۸ این مفهوم بهتر نشان داده شده است.



شکل ۴-۸: شباهت بین دو خوشه در روش *complete-link* برابر است با بیشترین فاصله بین دو داده‌های دو خوشه

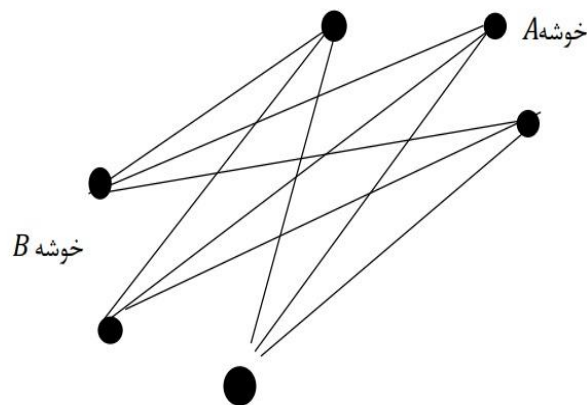
۴-۷-۲-۳-۳ خوشه‌بندی با روش *Average-link*

این روش همانند *single-link* جزء روش‌های خوشه‌بندی سلسله مراتبی است. از آنجا که هر دو روش *single-link* و *complete-link* به شدت به نطفه حساس می‌باشد، این

روش که محاسبات بیش تری دارد، پیش نهاد شد. در این روش برای محاسبه شباهت بین دو خوشه A و B از معیار زیر استفاده می شود:

$$d_{AB} = \frac{\sum_{i \in A, j \in B} d_{ij}}{N_A N_B}$$

تعداد اعضا خوشه A و N_B تعداد اعضاء خوشه B است. در واقع در این روش، شباهت بین دو خوشه میانگین فاصله بین تمام اعضا یکی با تمام اعضا دیگری است. در شکل ۴-۹ این مفهوم بهتر نشان داده شده است.



شکل ۴-۹: شباهت بین دو خوشه در روش *Average-link* برابر است با میانگین فاصله بین داده های دو خوشه

مثال ۴-۶- فرض کنید ۶ نمونه داده داشته باشیم و ماتریس فاصله بین آن ها در جدول ۴-۱ محاسبه شده باشد. با روش *single-Link* نحوه اعمال خوشه بندی را نشان داده و درخت وارهنگار مربوطه را رسم می کنیم.

جدول ۴-۱: ماتریس فاصله بین ۶ نمونه داده

	1	2	3	4	5	6
1	0	4	13	24	12	8
2		0	10	22	11	10
3			0	7	3	9
4				0	6	18
5					0	8·5
6						0

در این روش ابتدا هر داده به عنوان یک خوشه در نظر گرفته می شود و یافتن نزدیک ترین خوشه در واقع یافتن کمینه فاصله بین داده های بالا خواهد بود. با توجه به جدول ۴-۱ مشخص است که داده های ۳ و ۵ کم ترین فاصله را دارا هستند که در نتیجه دو خوشه ۳ و ۵ را با هم ترکیب می کنیم و خوشه جدیدی حاصل می شود که فاصله آن از سایر خوشه ها برابر است با کم ترین فاصله بین ۳ و ۵ از سایر خوشه ها. نتیجه در جدول ۴-۲ نشان داده شده است.

جدول ۴-۲: ماتریس فاصله بین ۵ خوشه حاصل از تکرار اول

	1	2	(3و5)	4	6
1	0	4	12	24	8
2		0	10	22	10
(3و5)			0	6	8·5
4				0	18
6					0

با توجه به جدول ۴-۲ خوشه های ۱ و ۲ کم ترین فاصله را دارا هستند پس آن ها را با هم ترکیب می کنیم و خوشه جدیدی حاصل می شود که فاصله آن از سایر خوشه ها برابر با کم ترین فاصله بین ۱ و ۲ از سایر خوشه ها است. نتیجه در جدول ۴-۳ نشان داده شده است.

جدول ۴-۳: ماتریس فاصله بین ۴ خوشه حاصل از تکرار دوم

	(1و2)	(3و5)	4	6
(1و2)	0	10	22	8

(3و5)		0	6	8.5
4			0	18
6				0

با توجه به جدول 3-4 مشخص است که خوشه‌های (3و5) و 4 کم‌ترین فاصله را دارا هستند. در نتیجه آن‌ها را با هم ترکیب کرده و خوشه جدیدی حاصل می‌شود که فاصله آن از سایر خوشه‌ها برابر است با کم‌ترین فاصله بین خوشه (3و5) و 4 از سایر خوشه‌ها. نتیجه در جدول 4-4 نشان داده شده است.

جدول 4-4: ماتریس فاصله 3 خوشه حاصل از تکرار سوم

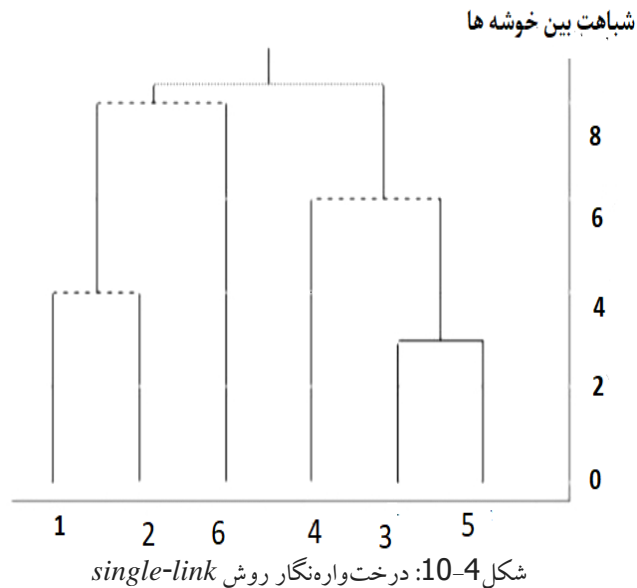
	(1و2)	(3و4و5)	6
(1و2)	0	10	8
(3و4و5)		0	8.5
6			0

با توجه به جدول 4-4 مشخص است که خوشه‌های (1و2) و 6 کم‌ترین فاصله را دارا هستند و در نتیجه آن‌ها را با هم ترکیب کرده و خوشه جدیدی حاصل می‌شود که فاصله آن از سایر خوشه‌ها برابر است با کم‌ترین فاصله بین (1و2) و 6 از سایر خوشه‌ها. نتیجه در جدول 5-4 نشان داده شده است.

جدول 5-4: ماتریس فاصله بین 2 خوشه حاصل از تکرار چهارم

	(1و2و6)	(3و4و5)
(1و2و6)	0	8.5
(3و4و5)		0

در نهایت این خوشه‌ها ترکیب می‌شوند و درخت‌واره‌نگار مربوطه به صورت شکل 4-10 حاصل می‌شود.



مثال 4-7- در این مثال به بررسی خوشه بندی سلسله مراتبی با داده های گل زنبق می-پردازیم که هدف ایجاد سه خوشه یعنی $k=3$ برای یک نمونه 150 تایی از داده ها می باشد. می دانیم داده های گل زنبق شامل سه نوع گل زنبق است پس به همین دلیل از هر دسته 50 تایی یک نمونه 5 تایی یعنی $n=5$ تهیه شده تا مشاهده نمودار و خروجی های خوشه بندی سلسله مراتبی بهتر دیده شود. در این مثال تابعی که برای اندازه گیری فاصله بین مشاهدات و یا خوشه ها در نظر گرفته شده، فاصله اقلیدسی است و معیار پیوند خوشه ها، روش میانگین است که محاسبات بر اساس آن صورت گرفته شده است. ماتریس فاصله با دستور *dist* تولید شده در متغیر *distance* و خروجی دستور *hclust* خوشه بندی سلسله مراتبی را می-دهد، به شکل 4-11 و 4-12 توجه کنید: (کد 5)


```

> library(stats)
> library(factoextra)
> n=5
> k=3
> data=iris[c(sample(x = 1:50,size = n),sample(x
=51:100,size=n),sample(x = 101:150,size= n)),1:4]
> distmethod=c('euclidean')
> linkagemethod=c("average")
> distance=dist(data,distmethod)
> hc=hclust(d = distance,method=linkagemethod)
> distance
29      41      20      2      35      98
54      94      75

41  0.2645751
20  0.4358899 0.3741657
2   0.5000000 0.5291503 0.8366600
35  0.4358899 0.4690416 0.7348469 0.1414214
98  3.2969683 3.4351128 3.2954514 3.3645208
3.2832910
54  3.0446675 3.1591138 3.1080541 2.9698485 2.9086079
0.9695360
94  2.3452079 2.4351591 2.4474477 2.1794495 2.1283797
1.7000000 0.9110434
75  3.3630343 3.5099858 3.3674916 3.4467376 3.3674916
0.2000000 1.1224972 1.8466185
74  3.6138622 3.7509999 3.6124784 3.6565011 3.5735137
0.4358899 1.0535654 1.8601075 0.5196152
115 4.3874822 4.4698993 4.3428102 4.4022721 4.3243497
1.4212670 1.6613248 2.4677925 1.4899664
119 6.4459289 6.5924199 6.4311741 6.5314623 6.4544558
3.1780497 3.7868192 4.6936127 3.0886890
133 4.8414874 4.9547957 4.8072861 4.8918299 4.8114447
1.5968719 2.1047565 2.9899833 1.5842980
117 4.6065171 4.7318073 4.5661800 4.6829478 4.5967380
1.3379088 1.9974984 2.8670542 1.3076697
142 4.5912961 4.7127487 4.5486262 4.7021272 4.6227697
1.4730920 2.1931712 3.0298515 1.3892444
74      115      119      133      117
41
20
2
35
98
54
94

```

```

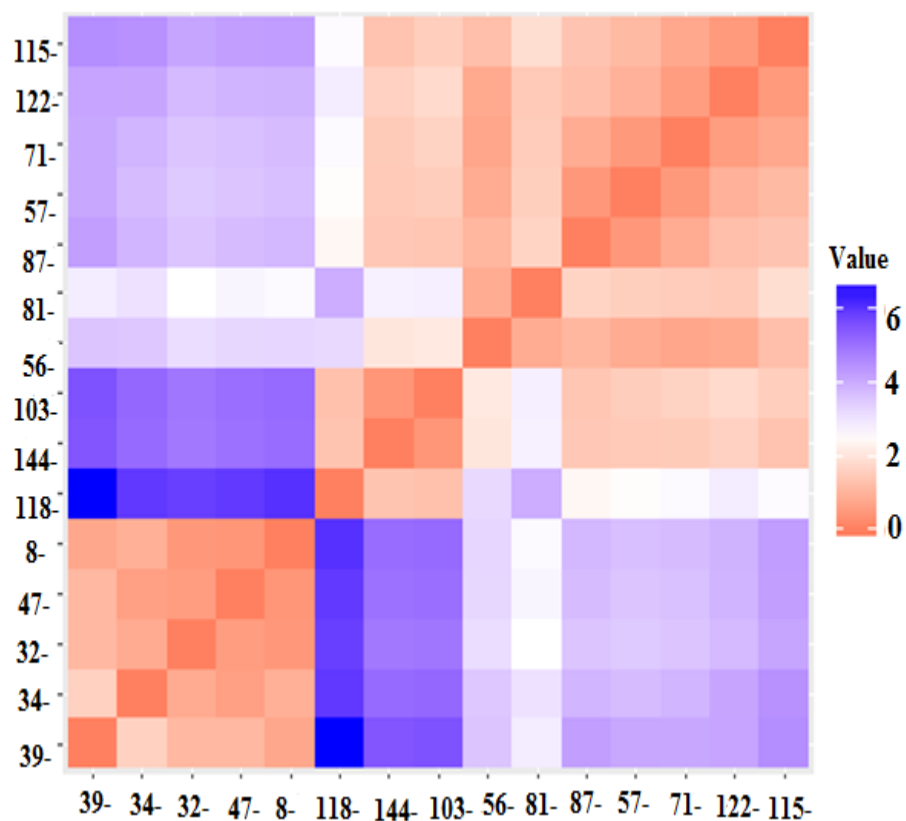
75
74
115 1.3000000
119 2.9410882 2.6267851
133 1.3784049 0.8062258 1.8520259
117 1.0954451 1.0246951 1.9519221 0.4690416
142 1.4491377 1.1445523 2.0322401 0.7745967 0.7615773

> fviz_dist(dist.obj = distance)
> hc

Call:
hclust(d = distance, method = linkagemethod)

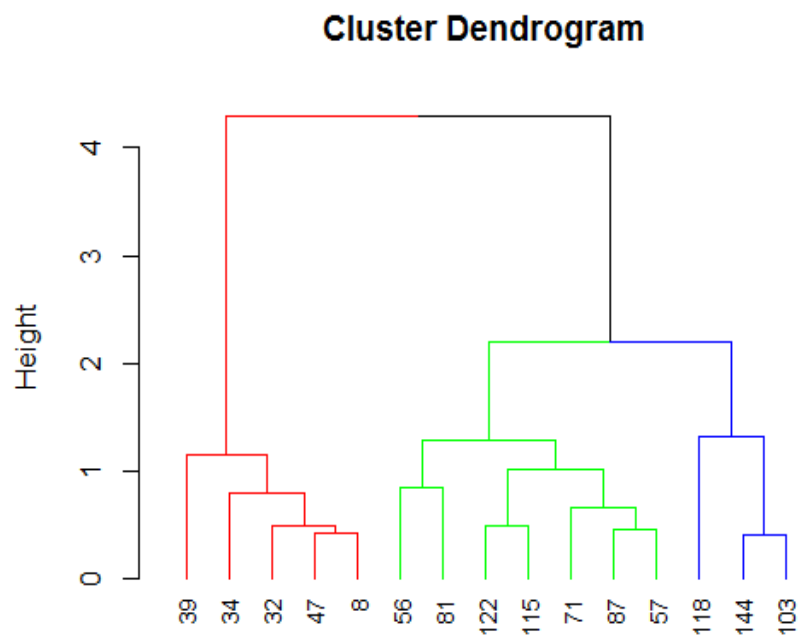
Cluster method      : average
Distance            : euclidean
Number of objects: 15
> fviz_dend(x = hc,k = k)

```



شکل 4-11: نمایش تصویری ماتریس فاصله

شکل 4-11 نمایشی از ماتریس فاصله را نشان می‌دهد که در آن هر شیء با خودش کم‌ترین فاصله را دارد و با رنگ خاکستری نشان داده شده است و نمونه‌هایی که بیش‌ترین فاصله را با هم دارند با رنگ خاکستری نشان داده شده است.



شکل 4-12: نمودار سلسله مراتبی داده‌های گل زنبق

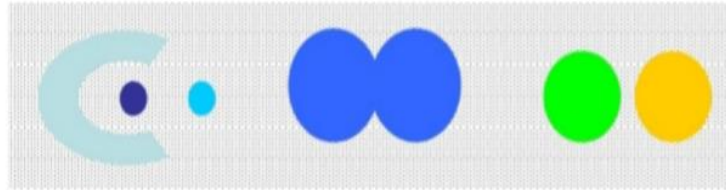
۴-۷-۲-۴- مقایسه خوشه‌بندی سلسله مراتبی و غیر سلسله مراتبی

روش‌های خوشه‌بندی غیر سلسله مراتبی نیاز به تصمیم‌گیری‌هایی از سوی کاربر دارد به‌طور مثال انتخاب تعداد خوشه‌ها. در این روش به‌طور معمول خوشه‌های اولیه ایجاد می‌شود و سپس در ادامه بهبود می‌یابد. با توجه به این نکته که تعداد خوشه‌های اولیه و انتخاب آن مهم است و دارای حساسیت می‌باشد، گاهی اوقات به کمک روش سلسله مراتبی، تعداد خوشه‌ها برآورد می‌شود. ازجمله ایرادات روش سلسله مراتبی این است که در یک مرحله قابلیت تغییر در مرحله‌های بعد را ندارد که این موجب تصمیمات نادرست می‌شود.

۴-۷-۳- روش‌های مبتنی بر چگال

در تعیین تعداد خوشه‌ها استقلال یا وابستگی کامل داده‌ها نقش مهمی دارد. اگر تمام متغیرها وابسته باشند، آن‌گاه تمام داده‌ها تشکیل یک خوشه می‌دهند، یا به عکس اگر تمام متغیرها به‌طور کامل مستقل باشند، هیچ خوشه‌ای ایجاد نمی‌شود (تمام فضا به‌صورت تصادفی با نقاط داده پر می‌شود). در شرایط استقلال یا وابستگی کامل نمی‌دانیم که واقعاً چند خوشه وجود دارد. به‌طور معمول در انتخاب تعداد خوشه‌ها نقش تحلیل‌گر بیش‌تر از رایانه است، با توجه به کاربردهای متفاوت خوشه‌بندی، تعداد خوشه‌ها ممکن است بیش‌تر یا کم‌تر باشد. در بسیاری از مواقع با یک مقدار k خوشه‌بندی را انجام داده و نتیجه‌ها را بررسی می‌کنند و دوباره یک مقدار k دیگر را امتحان می‌کنند. بعد از هر تکرار، ارزش نتیجه‌ها به‌وسیله میزان متوسط فاصله‌ها در داخل خوشه‌ها و یا میزان متوسط فاصله‌ها بین مراکز خوشه‌ها، اندازه‌گیری و با یکدیگر مقایسه می‌شوند. این موضوع را هم باید مد نظر قرار داد که گاهی خوشه‌ها به‌وسیله قضاوت‌های ذهنی تحلیل‌گر مورد ارزیابی قرار می‌گیرد تا ارزش آن‌ها در کاربردهای خاصی مشخص شوند. معیارهایی برای ارزیابی رده‌های تشکیل شده و همچنین تعیین k مناسب وجود دارد. اکثر روش‌های افراز بر اساس فاصله‌ی میان نمونه آن‌ها را خوشه‌بندی می‌کنند. این روش‌ها می‌توانند تنها خوشه‌های کروی شکل را پیدا کنند و برای کشف خوشه‌هایی با شکل دلخواه با مشکل مواجه می‌شوند. ایده کلی در روش‌های مبتنی بر چگال به این صورت است که مادامی که چگالی (تعداد نمونه یا داده‌ها) در همسایگی یک خوشه، از مقدارهایی بیش‌تر باشد، خوشه به رشد خود ادامه می‌دهد. چنین روشی را می‌توان برای یافتن داده‌های دور افتاده یا کشف خوشه‌هایی با شکل دلخواه استفاده کرد. روش‌های چگالی مبنای می‌توانند مجموعه‌ای از نمونه را در چندین خوشه‌ی انحصاری یا یک سلسله‌مراتبی از خوشه‌ها تقسیم کنند. به‌طور معمول در این روش‌ها خوشه‌ها انحصاری هستند، به این معنی که هر شیء تنها به یک خوشه تعلق دارد. همچنین با تعمیم این روش‌ها می‌توان از آن‌ها برای خوشه‌بندی زیرفضا نیز استفاده کرد. یکی از رایج‌ترین الگوریتم‌های

خوشه بندی مبتنی بر چگال الگوریتم *DBSCAN* است. همچنین می توان الگوریتم *optics* را به عنوان یک روش خوشه بندی مبتنی بر چگال اشاره کرد.



شکل ۴-۱۳: یک مثال برای خوشه های چگالی مبنا

۴-۷-۳-۱- الگوریتم *DBSCAN*

یکی از معروف ترین الگوریتم ها در حوزه خوشه بندی الگوریتم *DBSCAN* است. در این الگوریتم نیاز به تعیین تعداد خوشه ها توسط کاربر نیست و خود الگوریتم می تواند خوشه ها را مبتنی بر غلظت آن ها شناسایی کند، یا داده ها را بر اساس تراکم آن ها خوشه بندی کند. الگوریتم *DBSCAN* به دو پارامتر ورودی نیاز دارد: شعاع هر خوشه و حداقل نقطه در درون هر خوشه (جانسون، ۲۰۱۴). در ابتدا یک نمونه انتخاب می شود و با توجه به شعاع، به دنبال همسایه برای این نقطه در فضا (مجموعه داده ها) می گردد. اگر در شعاع مشخص، الگوریتم بتواند حداقل نقاط موجود در داده ها را پیدا کند، آن گاه همه ی آن داده ها با هم به یک خوشه تعلق می گیرند. سپس به دنبال نقاط هم جوار نقطه فعلی می رود تا با شعاع *Epsilon* در آن نقطه به دنبال همسایه دیگر بگردد، و اگر تعداد نقاط همسایه جدید باز هم پیدا نشوند، این الگوریتم دوباره همه ی آن نقاط جدید را با نقاط قبلی به یک خوشه متعلق می کند و اگر نقطه ی جدیدی پیدا نشود این خوشه تمام شده است و برای پیدا کردن خوشه های دیگر در نقاط دیگر، به صورت تصادفی یک نقطه دیگر را انتخاب می کند و شروع به یافتن همسایه و تشکیل خوشه جدید برای آن نقطه می کند. این کار آن قدر ادامه می یابد تا تمامی نقاط بررسی شوند.

حال با توجه به توضیحات بالا می توان نتیجه گرفت که یک خوشه چگالی مینا، مجموعه ای از نقاط متصل به یکدیگر از نظر چگالی است که هر داده را که خارج از این خوشه ها باشند، به عنوان داده دور افتاده در نظر می گیرد. یکی از مشکلات الگوریتم *DBSCAN* مشخص نبودن تعداد *Epsilon* و حداقل نقاط است که باید توسط کاربر مشخص شود.

۴-۷-۳-۱-۱- مزیت ها و عیب های الگوریتم *DBSCAN*

مزیت ها:

- عدم محدودیت به شکل خوشه ها
- سریع برای داده های با بُعد کم
- یافتن خوشه ها برای اشکال نامنظم
- تشخیص نوفه

عیب ها:

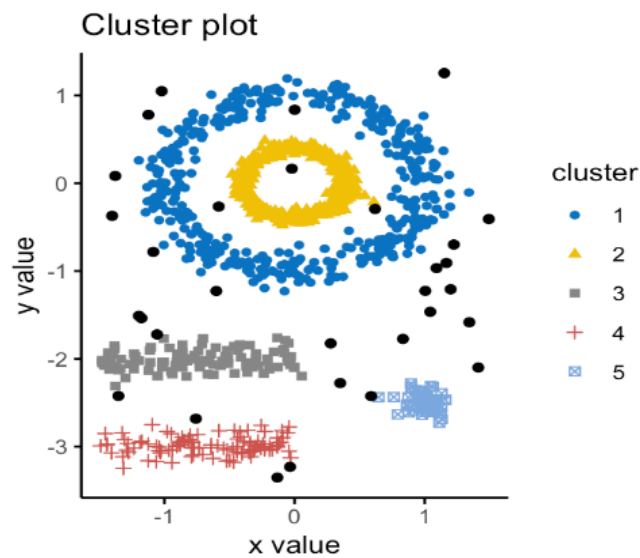
- نقاط مرزی که می توانند در دو خوشه نیز باشند، ممکن است به هر یک از خوشه ها تعلق گیرند
- مشکل بودن تعیین مقدار دقیق پارامترهای ورودی
- عدم تشخیص خوشه های با چگالی متفاوت

مثال ۴-۸- مجموعه داده های *Multishapes* مجموعه ای در داده گان های نرم افزار *R* است که مناسب برای روش های خوشه بندی چگالی مینا می باشد، این مجموعه داده با ۱۱۰۰ مشاهده می باشد که دارای سه متغیر x ، y و $shape$ (یک بردار عددی مربوط به شماره خوشه ها) است. برای پیاده سازی الگوریتم *DBSCAN* روی این مجموعه داده ها به صورت زیر عمل می شود: (کد ۶)

```
> data("multishapes", package="factoextra")
> df=multishapes[1:2,]
> library("fpc")
> set.seed(123)
> db=fpc::dbscan(df, eps=0.15, minpts=5)
> print(db)
```

```
> library("factoextra")
> fviz_cluster(db, data=df, stand=FALSE, ellipse=FALSE,
show.clust.cent=FALSE
> geom="point", palette="jco", ggtheme=theme_classic())
```

خروجی دستورهای بالا شکل ۴-۱۴ خواهد بود:



شکل ۴-۱۴: نشان می‌دهد که این مجموعه به ۵ خوشه تقسیم شده

در این شکل نقاط دایره‌ای شکل مربوط به داده‌های نوفه می‌باشد. برای تغییر پیکربندی خوشه‌ها نیز می‌توان مقدارهای *eps* و *MinPts* را در تابع *fpc::dbscan()* تغییر داد.

۴-۷-۳-۲- الگوریتم OPTICS

اگرچه الگوریتم *DBSCAN* با کمک پارامترهای *Epsilon* (حداکثر شعاع یک همسایه) و حداقل نقاط لازم در همسایگی یک شیء قادر است نمونه را خوشه‌بندی کند، اما کاربرها می‌بایست بهترین مقدارها را برای این پارامترها مشخص کنند تا الگوریتم خوشه‌های قابل قبولی را استخراج کند. این مشکل در بسیاری از الگوریتم‌های خوشه‌بندی دیگر نیز وجود دارد. به‌طور معمول تنظیم چنین پارامترهایی به‌صورت تجربی انجام می‌شود و تعیین آن‌ها

به خصوص برای داده‌های دنیای واقعی با تعداد ابعاد بالا، بسیار دشوار است. اکثر الگوریتم‌ها به مقدارهای این پارامترها حساس هستند، به نحوی که تغییر ناچیز آن‌ها ممکن است نتیجه - های بسیار متفاوتی در خوشه‌بندی را به همراه داشته باشد. اغلب در دنیای واقعی، داده‌ها دارای ابعاد بالایی هستند و همچنین توزیع‌های آن‌ها با چولگی زیاد است. بنا بر این ساختار طبیعی خوشه‌های موجود در این گونه داده‌ها را نمی‌توان تنها با کمک چند پارامتر چگالی سراسری توصیف نمود.

جهت غلبه بر مشکلاتی که استفاده از پارامترهای سراسری در تحلیل خوشه ایجاد می‌کنند، یک روش تحلیل خوشه به نام *OPTICS* پیش‌نهاد شد. این الگوریتم به صورت آشکار عمل خوشه‌بندی داده‌ها را انجام نمی‌دهد و در عوض یک ترتیب خوشه را در خروجی خود قرار می‌دهد. این خروجی یک فهرست خطی از کلیه‌ی نمونه تحت بررسی و تحلیل است و ساختار خوشه‌بندی چگالی مبنا داده‌ها را نمایش می‌دهد. در ترتیب خوشه، نمونه‌هایی که در یک خوشه‌ی متراکم‌تر هستند در این فهرست به یکدیگر نزدیک‌تر خواهند بود. این نظم معادل خوشه‌بندی چگالی مبنا است که با کمک طیف وسیعی از تنظیمات پارامترها به دست آمده است. بنا بر این در *OPTICS* لازم نیست کاربر حد آستانه‌ی چگالی مشخصی را تعیین کند. ترتیب خوشه را می‌توان برای استخراج اطلاعات پایه‌ای خوشه‌بندی (برای مثال مراکز خوشه‌ها یا خوشه‌هایی با شکل‌های دلخواه)، کشف ساختار طبیعی و موجود خوشه‌ها و همچنین مصورسازی خوشه‌بندی استفاده کرد.

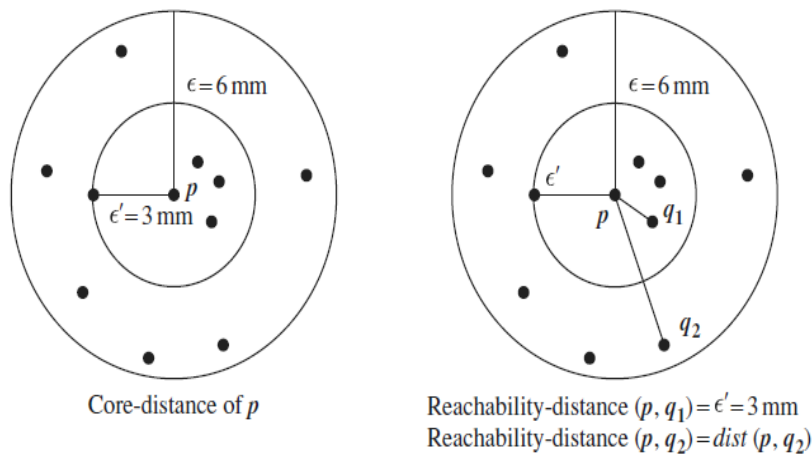
برای ساخت خوشه‌بندی‌های مختلف به‌طور هم‌زمان، نمونه به یک نظم خاصی پردازش می‌شوند. در این نظم شی‌ای که در رابطه با کم‌ترین مقدار ϵ قابل دسترس است، انتخاب می‌شود. بدین ترتیب خوشه‌هایی با تراکم بالاتر (ϵ کوچک‌تر) زودتر به پایان خواهند رسید. بر اساس همین ایده الگوریتم *OPTICS* برای هر شی به دو فقره اطلاعات مهم نیاز دارد:

1. فاصله‌ی هسته‌ی یک شی p به کوچک‌ترین مقداری از ϵ' اطلاق می‌شود که در

شعاع همسایگی ϵ' از شی p بتوان حداقل تعداد $MinPts$ شی پیدا کرد.

2. فاصله‌ی قابل دسترس به شیء p از q به مقدار حداقل شعاعی اطلاق می‌شود که باعث می‌شود p از q قابل دسترس باشد. طبق تعریف قابلیت دسترسی، شیء q باید یک شیء هسته باشد و p نیز باید در همسایگی q قرار گرفته باشد. چنانچه با توجه به مقدارهای ϵ و $MinPts$ شیء q یک شیء هسته نباشد، فاصله‌ی قابل دسترس شیء p از q نیز تعریف نشده‌است.

ممکن است شیء‌ای مانند p به‌طور مستقیم از چندین شیء هسته قابل دسترس باشد. بدین ترتیب شیء p با توجه به نمونه هسته‌ی مختلف دارای چندین فاصله‌ی قابل دسترس است. از میان آن‌ها کوچک‌ترین مقدار برای ما اهمیت دارد، چون با کمک این مقدار، کوتاه‌ترین مسیری که شیء p به یک خوشه‌ی متراکم متصل می‌شود، مشخص می‌گردد. شکل ۴-۱۵ دو مفهوم فاصله‌ی هسته و فاصله‌ی قابل دسترس را به‌صورت تصویری نشان می‌دهد:



شکل ۴-۱۵: فاصله‌ی هسته و فاصله قابل دسترس

الگوریتم *OPTICS* نظم و ترتیب را برای کلیه‌ی نمونه موجود در داده‌گان‌ها محاسبه می‌کند و مقدارهای فاصله‌ی هسته و یک فاصله‌ی قابل دسترس مناسبی برای هر یک از نمونه را نیز ذخیره می‌کند. در این الگوریتم برای تولید نظم خروجی، از فهرستی به نام

OrderSeeds استفاده می شود. نمونه در این فهرست بر اساس فاصله ی قابل دسترس از نزدیک ترین نمونه هسته ی مربوط به خود، به صورت صعودی مرتب شده اند. الگوریتم *OPTICS* کار خود را با انتخاب یک شیء مانند p به صورت اختیاری از دادگان ها آغاز می کند. همسایگان p (به شعاع ϵ) بازیابی می شوند، فاصله ی هسته نیز برای p تعیین می گردد. از آن جا که فاصله ی قابل دسترس این شیء یعنی p مشخص نیست، مقدار فاصله ی قابل دسترس آن با برچسب تعریف نشده علامت گذاری می شود. شیء p به عنوان یک شیء پردازش شده در خروجی نوشته می شود. اگر شیء p یک شیء هسته نباشد، الگوریتم به سراغ شیء بعدی در فهرست *OrderSeeds* می رود و یا اگر فهرست *OrderSeeds* خالی باشد، از دادگان شیء ای انتخاب می شود. نظم به دست آمده از الگوریتم را می توان به صورت گرافیکی نمایش داد. این کار به مصورسازی و فهم ساختار خوشه بندی موجود در مجموعه - داده ها کمک می کند.

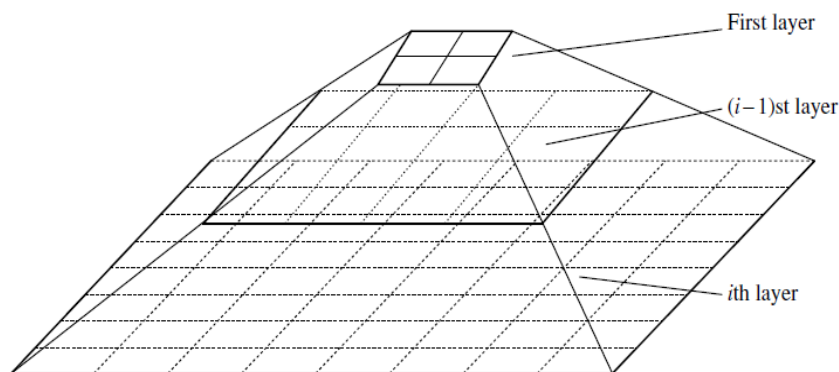
۴-۷-۴- روش مشبکی مبنا

روش های خوشه بندی که تا این جا بحث شد، داده محور هستند. آن ها با توجه به توزیع داده ها مجموعه ای از نمونه را بخش بندی می کنند. یک روش خوشه بندی مشبکه ای از یک رویکرد فضا محور بهره می گیرد و این کار را با افراز فضای تعبیه شده به سلول ها و مستقل از توزیع نمونه ورودی انجام می دهد. روش مشبکی مبنا به سلول های مختلف، امکان کار بر روی اطلاعات با درجه وضوح و شفافیت های متفاوت را فراهم می کند. در این رویکرد ابتدا فضا به سلول هایی تقسیم شده و سپس کلیه عملیات خوشه بندی بر روی این فضا انجام می گیرد. یکی از مزیت های این روش سرعت بالای پردازش آن است، زیرا پیچیدگی وابسته به تعداد سلول ها است. از دیگر مزیت های این روش کارا بودن است. ابتدایی ترین روش در این دسته روش شبکه اطلاعات آماری یا *STING* است.

۴-۷-۴-۱- الگوریتم *STING*

روش *STING* یکی از روش‌های مشبکی مبنا است که در این روش فضا به سلول‌های ابتدایی تقسیم می‌شود، این تقسیم‌بندی می‌تواند به روش سلسله مراتبی یا بازگشتی باشد. با افزایش هر سلول در سطوح بالا، تعداد سلول در سطوح پایین شکل می‌گیرد. به طور مثال از ترکیب هر ۴ سلول یک سلول در لایه‌ای بالاتر با درجه تفکیک کم‌تر شکل می‌گیرد. سپس برای هر سلول پارامترهایی مانند میانگین، میانه، کمینه و بیشینه محاسبه و به عنوان پارامترهای آماری ذخیره می‌شوند.

شکل ۴-۱۶ یک ساختار سلسله‌مراتبی را برای خوشه‌بندی *STING* نشان می‌دهد. پارامترهای آماری سلول‌های سطوح بالاتر را می‌توان با کمک پارامترهای سلول‌های سطوح پایین‌تر محاسبه کرد. این پارامترها می‌توانند گونه‌های مختلفی باشند: پارامتر مستقل از صفت خاصه مانند شمارش مقدارها، پارامترهای وابسته به صفت خاصه نظیر میانگین، انحراف استاندارد، بیشینه و کمینه و گونه‌ی سوم نوع توزیعی است که مقدارهای صفت خاصه‌ی موجود در سلول از آن پیروی می‌کنند مانند توزیع نرمال، یکنواخت، نمایی و یا حتی توزیع ناشناخته (هان و کامبر، ۲۰۰۶).



شکل ۴-۱۶: خوشه‌بندی به روش *STING*

در این جا پارامتر محاسبه‌شده، یک معیار منتخب برای تحلیل است و هنگام بارگذاری داده‌ها درون داده‌گان‌ها، پارامترهای شمارش، میانگین، انحراف استاندارد، بیشینه و کمینه سلول‌های

سطح پایین را می توان به طور مستقیم از روی داده ها محاسبه نمود. چنانچه از قبل نوع توزیع مشخص باشد، کاربر آن را برای سلول معین می کند و یا در صورت عدم آگاهی کاربر از توزیع، می توان آن را با کمک آزمون هایی نظیر آزمون خی دو به دست آورد. نوع توزیع سلول های سطوح بالاتر با رأی اکثریت نوع توزیع در سلول های متناظر سطوح پایین تر محاسبه می شوند. برای این کار می توان از حد آستانه ای نیز استفاده کرد. چنانچه توزیع های سلول های سطح پایین تر بر اساس انتخاب یک توزیع مناسب برای سلول سطح بالاتر به توافق نرسند و آزمون حد آستانه نیز این مشکل را برطرف نسازد، آن گاه نوع توزیع سطح بالا با برچسب ناشناخته تنظیم و علامت گذاری می شود.

چگونه این اطلاعات آماری برای پاسخ به پرسش سودمند است؟ از پارامترهای آماری می توان در یک روش مشبکی مبنا و یک روش بالا به پایین استفاده کرد. در ابتدا لایه ای از ساختار سلسله مراتب تعیین می شود، لایه ای که فرایند پاسخ به پرسش از آن آغاز می شود. به طور معمول این لایه شامل تعداد کمی از سلول ها است و برای هریک از سلول های سطح جاری، فاصله اطمینان محاسبه می شود تا ارتباط هر یک از آن ها با پرسش مورد نظر مشخص شود. بر روی سلول های بی ربط بررسی های بیش تری انجام نمی شود و آن ها از این فرایند حذف می شوند. جهت پردازش لایه ی بعدی تنها سلول های مرتبط و باقی مانده در نظر گرفته می شوند. این فرایند تا سطح پایین ساختار ادامه می یابد. چنانچه در این زمان مشخصات پرسش دیده شود، مناطق از سلول های مرتبط که پرسش را در بر می گیرد، برگردانده می شوند. در غیر این صورت داده هایی که درون سلول های مرتبط هستند، بازیابی می شوند و در ادامه تا رسیدن به نیازمندی های پرسش پردازش های بیش تر انجام می شود.

یکی از خصیصه های جالب توجه روش *STING* این است که چنانچه خوشه بندی به صفر نزدیک شود (به عبارت دیگر به سمت داده های سطح پایین حرکت کند)، نتیجه های خوشه بندی در این الگوریتم به نتیجه های حاصل از الگوریتم *DBSCAN* نزدیک می شود. به عبارت دیگر با کمک الگوریتم *STING* و همچنین استفاده از اطلاعات مربوط به

اندازه‌ی سلول و شمارش آن‌ها می‌توان خوشه‌های متراکم را شناسایی نمود. بدین ترتیب می‌توان الگوریتم *STING* را یک روش خوشه‌بندی چگالی مبنا نیز تصور نمود.

از آن‌جا که الگوریتم *STING* از یک رویکرد با چندین سطح وضوح و درستی جهت تحلیل خوشه استفاده می‌کند، کیفیت خوشه‌بندی در آن وابسته به دانه‌بندی پایین‌ترین سطح از ساختار شبکه می‌باشد. چنانچه دانه‌بندی بسیار ریز انتخاب شود، هزینه‌ی پردازش به‌صورت قابل ملاحظه‌ای افزایش خواهد یافت و از طرفی اگر سطح پایین ساختار شبکه از دانه‌بندی درشت استفاده کند، کیفیت تحلیل خوشه کاهش خواهد یافت.

مثال 4-9: در این مثال به بررسی تعداد بهینه‌ی خوشه‌ها برای دادگان گل زنبق می‌پردازیم. از تابع *fviz_nbclust()* در بسته نرم‌افزاری *factoextra* و تابع *NbClust()* در بسته نرم‌افزاری *NbClust* برای تعیین تعداد بهینه‌ی خوشه‌ها می‌توان استفاده کرد. در این مثال به بررسی تعداد بهینه‌ی خوشه‌ها با استفاده از تابع *NbClust()* می‌پردازیم. (کد 7)

تابع *NbClust()* در بسته‌ی نرم‌افزاری *NbClust*، 26 شاخص متفاوت برای تعیین بهترین تعداد خوشه ارائه می‌کند که همگی آن‌ها لزوماً با یکدیگر سازگار نمی‌باشند، اما می‌توان از نتایج این تابع به‌عنوان راهنمایی برای انتخاب مقادارهای ممکن k (تعداد خوشه) استفاده کرد.

تابع *fviz_nbclust()* در بسته نرم‌افزاری *factoextra*: از این تابع می‌توان در محاسبه‌ی متدهای *silhouette* و *elbow* استفاده کرد. همچنین این تابع قابلیت به‌کارگیری در انواع مختلف تابع‌های خوشه‌بندی مانند *k-means*، *k-medoids* و *CLARA* را دارد. تابع *NbClust()* در بسته نرم‌افزاری *NbClust*: ورودی‌های تابع *NbClust()*، دادگان، معیار فاصله، روش خوشه‌بندی، کم‌ترین تعداد خوشه و بیش‌ترین تعداد خوشه است. این تابع تعداد شاخص‌های هر خوشه همراه با بهترین تعداد خوشه‌ی پیشنهادی را به‌عنوان خروجی نمایش می‌دهد.

کاربر با تغییر متغیرهایی چون تعداد خوشه‌ها، معیار محاسبه فاصله بین داده‌ها و روش خوشه‌بندی مورد استفاده، بهترین الگوی خوشه‌بندی را به‌عنوان خروجی تابع خواهد داد. تنها با

۵۴ خوشه بندی

یکبار اجرای این تابع، مقدارهای بهینه‌ی این شاخص‌ها محاسبه شده و تعداد بهینه خوشه‌ها تعیین می‌شود.

```

> data(iris)
> iris2=iris
> iris2$Species=NULL
> library(NbClust)
> nc=NbClust(iris2,min.nc=2,max.nc=15,method=
"kmeans")

*** : The Hubert index is a graphical method
of determining the number of clusters.
      In the plot of Hubert index,
we seek a significant knee that corresponds t
o a
      significant increase of the v
alue of the measure i.e the significant peak
in Hubert
      index second differences plot.
*** : The D index is a graphical method of de
termining the number of clusters.
In the plot of D index, we seek a significant
knee (the significant peak in Dindex second d
ifferences plot) that corresponds to a signif
icant increase of the value of the measure.
*****
*****
* Among all indices:
* 11 proposed 2 as the best number of cluster
s
* 11 proposed 3 as the best number of cluster
s
* 1 proposed 8 as the best number of clusters
* 1 proposed 12 as the best number of cluster
s
      ***** Conclusion *****
* According to the majority rule,the best num
ber of clusters is 2
*****
*****
> BAR=table(nc$Best.nc[1,])
> BAR

```

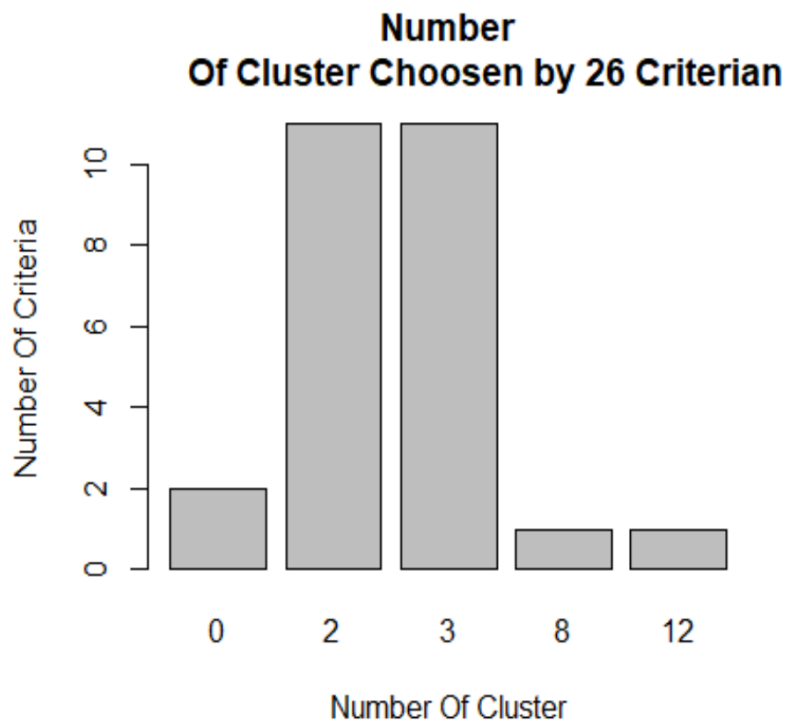
```

0  2  3  8 12
2 11 11  1  1
> barplot(BAR,xlab="Number Of Cluster",ylab="
Number Of Criteria", main = "Number
+Of Cluster Chooosen by 26 Criterion")

```

که همان‌طور که پیدا است تعداد خوشه بهینه 2 و 3 در نظر گرفته شده است. یعنی اگر $k = 2$ در نظر گرفته شود تعداد 11 شاخص تأیید می‌شود و اگر $k = 8$ در نظر گرفته شود تعداد 1 شاخص تأیید می‌شود.

نمودار میله‌ای تعداد خوشه‌ها در مقابل تعداد شاخص‌ها به صورت شکل است.



شکل 4-17: نمودار میله‌ای تعداد خوشه‌ها در مقابل تعداد شاخص‌ها

۶-۳-۵- الگوریتم K_2

الگوریتم K_2 یکی از روش‌های مبتنی بر جست‌وجو و رتبه‌بندی است که از جست‌وجو حریصانه استفاده می‌کند. ایده اصلی الگوریتم K_2 بیشینه کردن احتمال ساختار بر اساس داده‌های موجود است. الگوریتم K_2 احتمالاً یکی از بهترین الگوریتم‌هایی است که تاکنون برای یادگیری ساختار شبکه‌های یزی ابداع شده است. این الگوریتم از یک معیار امتیازدهی یزی استفاده می‌کند. هدف اصلی این الگوریتم جست‌وجوی پایگاه داده D و بیشینه کردن $P(B_s, D)$ برای دستیابی به بهینه B_s^* (شبکه یزی) است.

الگوریتم کار خود را با فرض این که تمام گره‌ها فاقد والدین هستند، شروع می‌کند؛ سپس در هر مرحله به تدریج برای هر گره، متغیرهایی که منجر به بیش‌ترین افزایش در احتمال ساختار نهایی می‌شوند، به عنوان والدین گره انتخاب می‌کند (اسکندری و همکاران، 2018).

• ورودی و خروجی الگوریتم

ورودی: مجموعه‌ای از n گره یا متغیر تصادفی گسسته، یک ترتیب روی گره‌ها، یک حد بالای u روی تعداد والدین هر گره و یک دادگان شامل m نمونه مستقل.

خروجی: مجموعه والدین مربوط به هر گره.

به طور کلی مرحله‌های اجرای الگوریتم به صورت زیر است:

شروع الگوریتم K2

برای i از یک تا n انجام شود:

مجموعه والدین متغیر i را صفر در نظر بگیرید، یعنی: $Pa(x_i) = \emptyset$

$P_{old} = g(x_i, Pa(x_i))$ (تابع g تابع امتیاز K2 شبکه بیزی است)

تا زمانی که متغیری برای ادامه دادن وجود دارد و تعداد والدین گره x_i از u کمتر است، ادامه بدهید.

شروع

تا زمانی که گره Z از مجموعه متغیرهایی که x_i می تواند به عنوان والدین خود انتخاب نماید و هنوز انتخاب نشده است، وجود دارد، آن را طوری انتخاب کنید که عبارت $g(x_i, Pa(x_i) \cup \{Z\})$ را ماکسیمم کند.

$$P_{New} = g(x_i, Pa(x_i) \cup \{Z\})$$

اگر $P_{New} > P_{old}$

شروع

$$P_{old} := P_{New}$$

$$Pa(x_i) = Pa(x_i) \cup \{Z\}$$

پایان

اگر متغیری که به عنوان والدین متغیر x_i انتخاب شود، در ترتیب وجود نداشت،

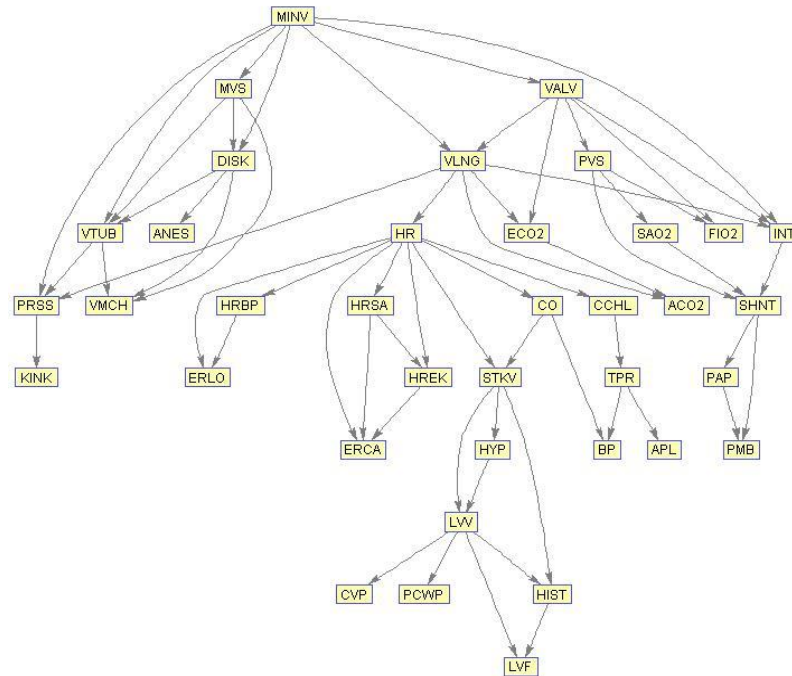
پایان

والدین مربوط به متغیر x_i را چاپ کنید.

پایان

پایان الگوریتم K2

مثال ۶-۸- شبکه بیزی ساخته شده حاصل از الگوریتم K2 روی مجموعه داده آلازم در شکل ۶-۱۴ نشان داده شده است.



شکل شماره ۶-۱۴: ساختار شبکه بیزی مجموعه داده آلام با استفاده از الگوریتم K_2

لازم به ذکر است که ترتیب انتخاب متغیرها بر اساس مقایسه بی‌نظمی (آنتروپی) بین متغیرها تعیین شده است.

یک پروژه نمونه

مقدمه

هدف این پروژه، یافتن گروه‌هایی از مشتریان عمده فروشی بر اساس کالای خریداری شده است. مجموعه داده استفاده شده، مجموعه داده مشتریان عمده فروشی از دامنه عمومی خواهد بود.

این مجموعه داده به عنوان خوشه بندی بخش مشتری انتخاب شده است که یکی از موضوعات رایج است. خوشه بندی متغیرها/ویژگی ها با استفاده از یادگیری بدون نظارت اغلب امکان کشف شباهت زیاد بین طبقاتی و عدم تشابه زیاد درون کلاسی را فراهم می کند. این خوشه بندی به راحتی در زمینه های دیگر در کنار بخش مشتری قابل استفاده است. حوزه های دیگری مانند بیمه خودرو، کاربران کارت اعتباری، ریزش مشتریان مخاطراتی و خریداران تجارت الکترونیک از حوزه هایی هستند که خوشه بندی برای آن ها قابل استفاده است.

از آنجایی که این مدل سازی ها براساس داده های تاریخی با فرضیات معین است. این پروژه به سایر عوامل خارجی مانند نحوه خرید (کارت اعتباری / نقدی)، تحویل مستقیم یا خود جمع آوری و زمان تراکنش خرید توجهی نکرده است.

داده ها

FRESH: هزینه سالانه (m.u.) برای محصولات تازه (پیوسته)

MILK: هزینه سالانه (m.u.) برای محصولات شیر (پیوسته)

GROCERY: هزینه سالانه (m.u.) برای محصولات خواربارفروشی (پیوسته)

FROZEN: هزینه سالانه (m.u.) برای محصولات منجمد (پیوسته)

DETERGENTS_PAPER: هزینه سالانه (m.u.) برای مواد شوینده و محصولات کاغذی

(پیوسته)

DELICATESSEN: هزینه سالانه (m.u.) برای محصولات و اغذیه فروشی ها (پیوسته)

CHANNEL: کانال مشتریان – Horeca (هتل/رستوران/کافه) یا کانال خرده فروشی (اسمی)

REGION: منطقه مشتریان – Oporto، Lisnon یا دیگر (اسمی)

تحلیل اکتشافی داده ها

```
library(stats)
library(ggplot2)
library(plyr)
library(dplyr)
```

```
library(data.table)
library(corrplot)
library(factoextra)
```

نگاهی اجمالی به داده ها به شرح زیر است.

```
library(NbClust)
library(cluster)
df.original<-read.csv('Wholesale customers data.csv')
set.seed(1) #Ensure reproducible code
df<- sample_frac(df.original,0.7) #split into test and train data by 7:3 ratio
df.index<- as.numeric(rownames(df))
df.test<- df.original[-df.index,]
head(df)
```

##	Channel	Region	Fresh	Milk	Grocery	Frozen	Deterge
nts_Paper	Delicassen						
## 117	1	3	11173	2521	3355	1517	
310	222						
## 164	2	3	5531	15726	26870	2367	
13726	446						
## 251	1	1	3191	1993	1799	1730	
234	710						
## 397	2	3	4515	11991	9345	2644	
3378	2213						
## 88	1	3	43265	5025	8117	6312	
1579	14351						
## 391	1	3	3352	1181	1328	5502	
311	1000						

برای مقایسه با توضیحات ارائه شده درباره داده ها، به رکوردهای نمونه در مجموعه داده نگاه کنید.

۶۲ خوشه بندی

۶ دسته محصول با هزینه سالانه توسط گروه مشتریان عمده فروشی توسط کانال و منطقه وجود دارد.

summary(df)					
##	Channel	Region	Fresh		
##	Milk				
##	Min. : 1.000	Min. : 1.000	Min. : 18	Min	
.	: 55				
##	1st Qu.: 1.000	1st Qu.: 2.000	1st Qu.: 3094	1st	
	Qu.: 1610			Qu.	
##	Median : 1.000	Median : 3.000	Median : 8130	Med	
	ian : 3684			ian	
##	Mean : 1.331	Mean : 2.513	Mean : 12059	Mea	
	n : 5716			n	
##	3rd Qu.: 2.000	3rd Qu.: 3.000	3rd Qu.: 16851	3rd	
	Qu.: 7119			Qu.	
##	Max. : 2.000	Max. : 3.000	Max. : 112151	Max	
.	: 73498			.	
##	Grocery	Frozen	Detergents_Paper		
##	Delicassen				
##	Min. : 3	Min. : 33.0	Min. : 3.0	Min.	
	Min. : 3.0			Min.	
##	1st Qu.: 2156	1st Qu.: 805.8	1st Qu.: 282.0	1st	
	1st Qu.: 395.0			Qu.	
##	Median : 4904	Median : 1619.0	Median : 824.5	Median	
	Median : 944.5			Median	
##	Mean : 7729	Mean : 3089.2	Mean : 2764.6	Mean	
	Mean : 1558.2			Mean	
##	3rd Qu.: 10550	3rd Qu.: 3532.5	3rd Qu.: 4003.2	3rd	
	3rd Qu.: 1820.2			Qu.	
##	Max. : 59598	Max. : 60869.0	Max. : 26701.0	Max.	
	Max. : 47943.0			Max.	

یک اطلاعات خلاصه سریع در مورد متغیرها را می توان به شرح زیر نتیجه گرفت:
فقط دسته های غذایی از نظر آماره های توصیفی معنی دارند.

۶۳ خوشه بندی

نیازی به مقیاس نیست زیرا واحدهای اندازه گیری یکسان است. ما از دسته بندی محصولات برای خوشه بندی استفاده خواهیم کرد.

```
summary(is.na(df))
```

##	Channel	Region	Fresh	Milk
##	Mode :logical	Mode :logical	Mode :logical	Mode :logical
##	FALSE:308	FALSE:308	FALSE:308	FALSE:308
##	Grocery	Frozen	Detergents_Paper	Delicassen
##	Mode :logical	Mode :logical	Mode :logical	Mode :logical
##	FALSE:308	FALSE:308	FALSE:308	FALSE:308

در این بخش به بررسی متغیرهای Null پرداخته شده است. می توان دید که مجموعاً ۴۴۰ مشاهده وجود دارد و هیچ کدام از متغیرها دارای مقدار گم شده نیستند.

```
table(df$Channel)
```

##	
##	1 2
##	206 102

در این بخش فراوانی متغیرهای رسته ای مشاهده می شود. می توان دید که ۲۹۸ مشاهده وجود دارد که از دسته بندی Horeca (هتل، رستوران، کافه) خرید کرده اند و ۱۴۲ مشاهده از خرده فروشی خرید کرده اند.

```
table(df$Region)
```

##	
##	1 2 3
##	57 36 215

در رابطه با متغیر منطقه می توان ترکیب بالا را مشاهده کرد.

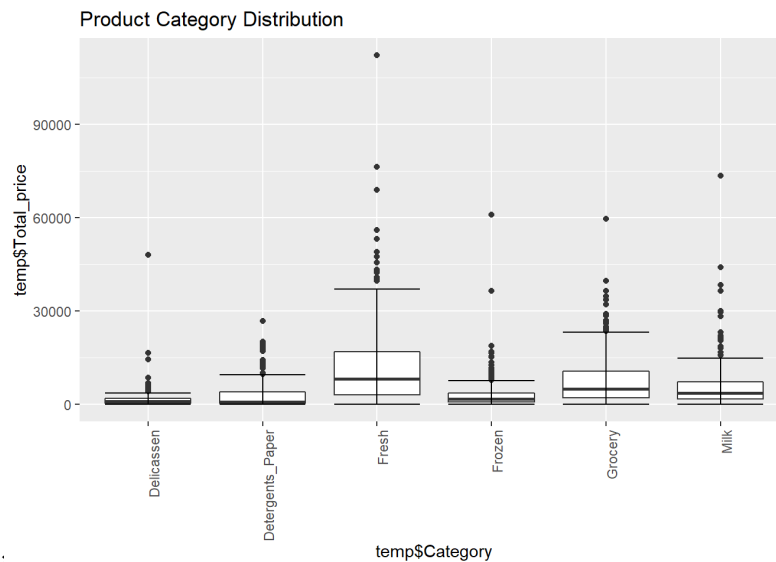
```
df %>%
  group_by(Channel,Region) %>%
  # multiple group columns
  summarise(total_fresh = sum(df$Fresh), total_Milk = sum(df$Milk), total_Grocery= sum(df$Grocery),total_Frozen=
sum(df$Frozen),total_Detergents_Paper=sum(df$Detergents_Paper), total_Delicassen= sum(df$Delicassen)) # multiple summary columns
```

```
## # A tibble: 6 x 8
## # Groups:   Channel [?]
##   Channel Region total_fresh total_Milk total_Grocery
total_Frozen
##   <int> <int> <int> <int> <int>
<int>
## 1      1      1 3714211 1760547 2380616
951463
## 2      1      2 3714211 1760547 2380616
951463
## 3      1      3 3714211 1760547 2380616
951463
## 4      2      1 3714211 1760547 2380616
951463
## 5      2      2 3714211 1760547 2380616
951463
## 6      2      3 3714211 1760547 2380616
951463
## # ... with 2 more variables: total_Detergents_Paper <
int>,
```

در بخش بالا مجموع هزینه برای دسته های مختلف محصول را نشان می دهد. دسته های تازه و مواد غذایی پرفروش ترین ها هستند.


```
temp <- reshape(df, direction="long", varying=c("Fresh",
"Milk", "Grocery", "Frozen", "Detergents_Paper", "Delicassen"),
v.names= "Total_price", timevar="Category", time=c(
"Fresh", "Milk", "Grocery", "Frozen", "Detergents_Paper", "
Delicassen"))
```

```
ggplot(temp, aes(x=temp$Category, y =temp$Total_price))
+geom_boxplot() +stat_boxplot(geom ='errorbar') + theme(
axis.text.x= element_text(angle=90,hjust=1))+ ggtitle("P
roduct Category Distribution")
```



نمودار فوق،

یافته های زیر را ارائه می کند:

1. تعداد نقاط دور افتاده کمی برای برخی از گروه ها وجود دارد. در نتیجه دو راه حل برای خوشه

بندی می توان ارائه داد.

1.1. نقاط دور افتاده را حذف کنید و سپس الگوریتم K mean را اجرا کنید.

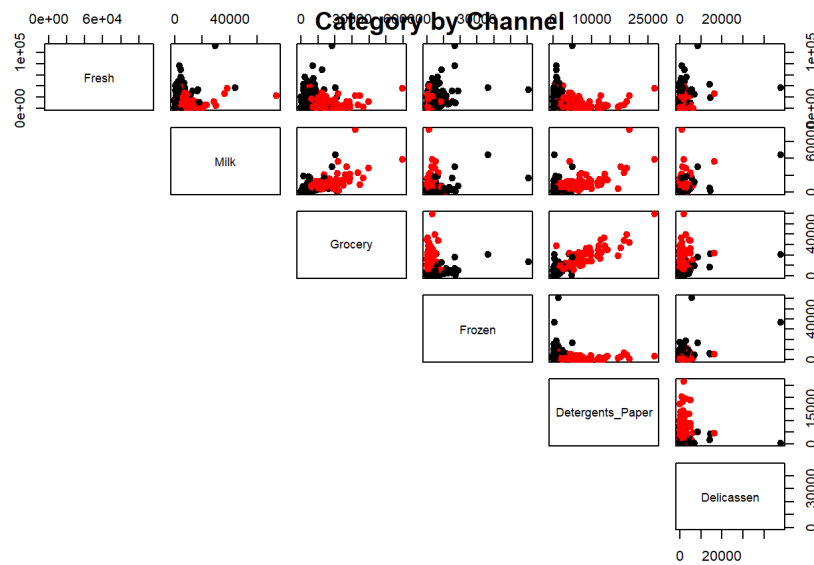
1.2. بدون حذف موارد پرت، الگوریتم (PAM) k medoids را با فاصله منتهن اجرا کنید.

2. واریانس بالایی برای کالاهای دسته مواد غذایی تازه و غلات وجود دارد.

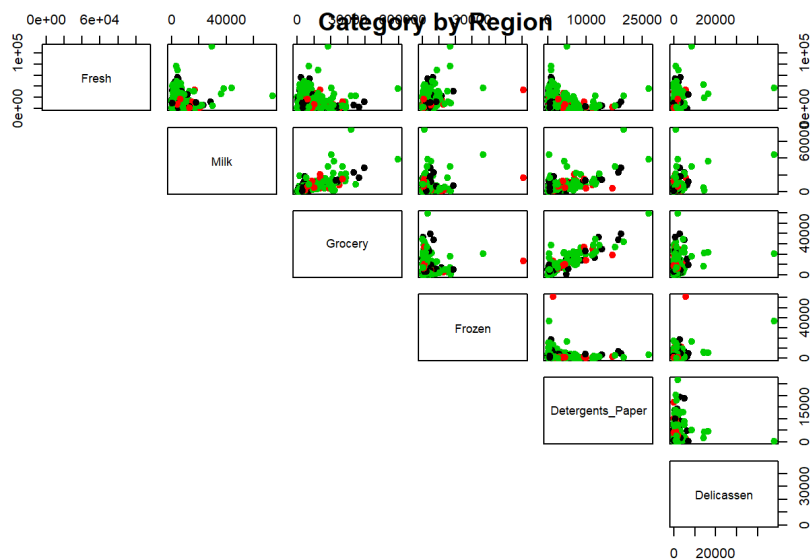
در ادامه، ما می خواهیم نموداری را به صورت سالیانه ارائه کنیم تا دیدی جامع از فروش داشته باشیم.

```
cor.result<- cor(df)

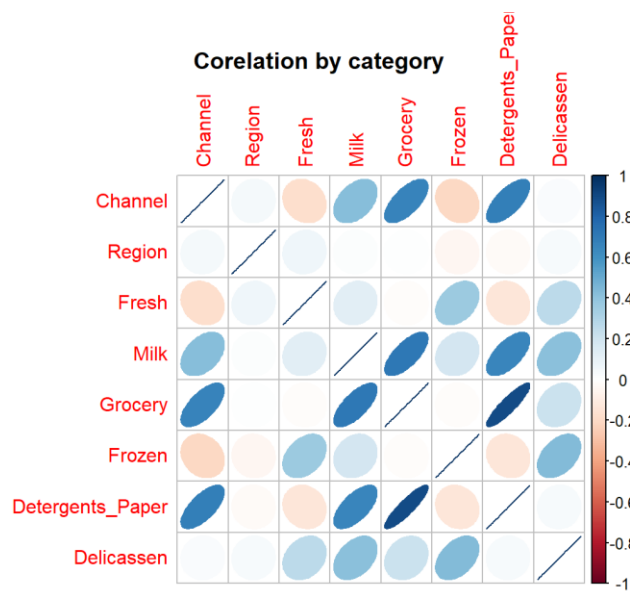
pairs(df[, -c(1:2)], col=df$Channel, pch=19, lower.panel
= NULL) +title(main = "Category by Channel")
```



```
plot(df[, -c(1:2)], col=df$Region, pch=19, lower.panel =
NULL) +title(main = "Category by Region")
```



```
corrplot(cor.result, method="ellipse") +title(main = "Co
relation by category")
```



نمودارهای فوق، یافته های زیر را ارایه می کنند:

۶۸ خوشه بندی

مواد شوینده، کاغذ و مواد غذایی بیشترین همبستگی را دارند.
در ادامه بررسی می کنیم که کدام متغیرها دارای نقاط پرت هستند.

```
#Remove Region & Channel Columns as they are not
sensible

apply(X= df[, -c(1:2)], MARGIN=2, FUN = function(x) length(
boxplot.stats(x)$out))
```

##	Fresh	Milk	Grocery
Frozen			
##	14	20	17
29			
##	Detergents_Paper	Delicassen	
##	21	18	

در بخش بالا، عدد نشان داده شده، تعداد نقاط پرت در هر دسته است. برای حذف نقاط پرت، از تکنیک Winsorizing استفاده خواهد شد. به طور خلاصه، نقاط پرت با صدک خاصی از داده ها، معمولاً ۹۰ یا ۹۵ جایگزین می شوند. ابتدا مقدار هر دسته را مرتب می کنیم.

```
sort(boxplot.stats(df$Grocery)$out)
```

##	[1]	23596	23998	24708	24773	25957	26839	26866	26870
		28540	28921	28986					
##	[12]	32114	33586	34792	36486	39694	59598		

صدک ها به شکل زیر حاصل می شود.

```
quantile(df$Grocery, probs=seq(from =0.9, to=1,by=0.025)
)
```

##	90%	92.5%	95%	97.5%	100%
##	19258.40	21198.98	23857.30	28663.83	59598.00

از بخش بالا، صدک ۹۵٪ انتخاب می شود. در مرحله بعد، صدک ۹۵ جایگزین نقطه پرت باقی مانده خواهد شد.

```
grocery.max <- as.numeric(quantile(df$Grocery,probs=0.95
))

df$Grocery[df$Grocery > grocery.max] <- grocery.max
```

```
sort(boxplot.stats(df$Detergents_Paper)$out)

quantile(df$Detergents_Paper, probs=seq(from =0.9, to=1,
by=0.025))

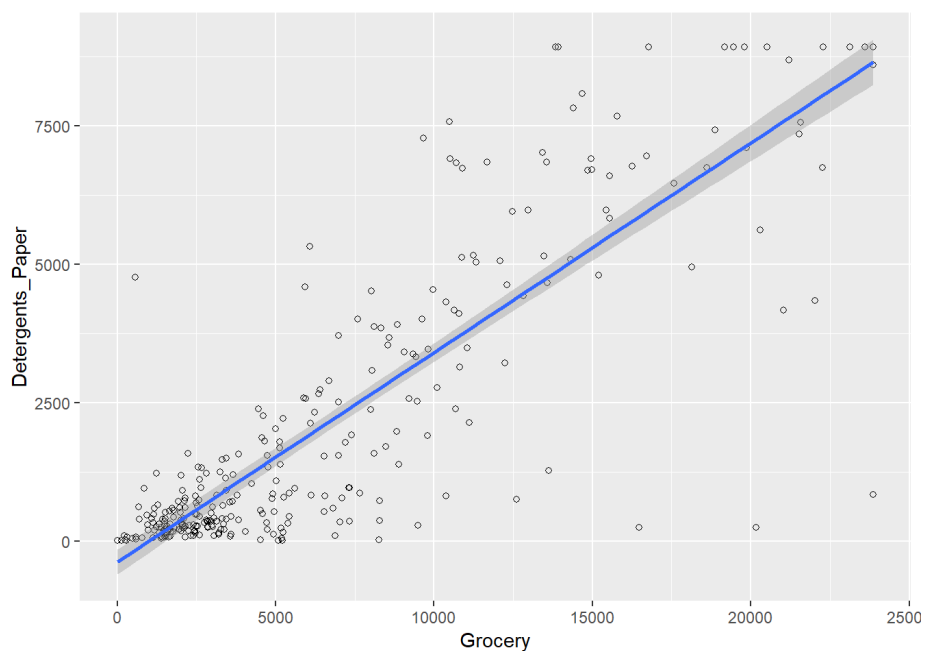
##          90%          92.5%          95%          97.5%          100%
## 7464.900  8926.725 11710.900 13629.475 26701.000

grocery.max <- as.numeric(quantile(df$Detergents_Paper,p
robs=0.925))

df$Detergents_Paper[df$Detergents_Paper > grocery.max] <-
- grocery.max
```

در ادامه دو متغیر Detergents_Paper و Grocery برای خوشه بندی انتخاب می شوند.

```
ggplot(data=df, aes(x=Grocery, y =Detergents_Paper)) + g
eom_point(shape=1) +geom_smooth(method="lm")
```



```
df.subset1<-as.data.frame(df[,c("Grocery", "Detergents_Pa
per")])

summary(df.subset1)
```

۷۰ خوشه بندی

```
##      Grocery      Detergents_Paper
## Min.      :    3  Min.      :    3.0
## 1st Qu.: 2156  1st Qu.: 282.0
## Median : 4904  Median : 824.5
## Mean      : 7336  Mean      :2398.7
## 3rd Qu.:10550  3rd Qu.:4003.2
## Max.      :23857  Max.      :8926.7
```

برای متغیرهای انتخاب شده نیاز به نرمال سازی وجود دارد زیرا تفاوت زیادی میان آماره های نمایش داده شده وجود دارد.

```
df.subset1<- as.data.frame(scale(df.subset1))
summary(df.subset1)

##      Grocery      Detergents_Paper
## Min.      : -1.0922  Min.      : -0.8308
## 1st Qu.: -0.7715  1st Qu.: -0.7341
## Median : -0.3624  Median : -0.5459
## Mean      :  0.0000  Mean      :  0.0000
## 3rd Qu.:  0.4786  3rd Qu.:  0.5565
## Max.      :  2.4606  Max.      :  2.2640
```

مدل سازی

یافتن تعداد بهینه خوشه ها

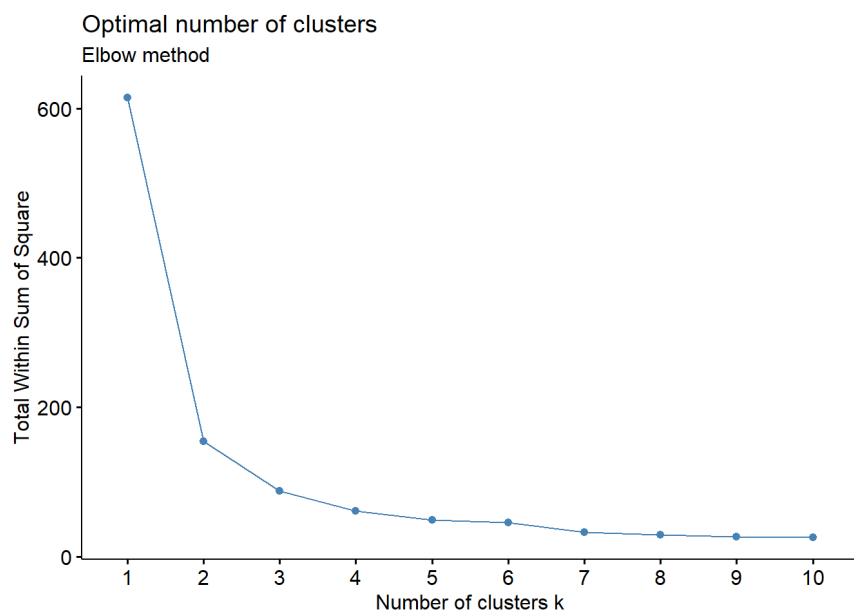
پس از انجام مقیاس بندی، روش های زیر برای یافتن تعداد بهینه خوشه ها اجرا می شود:

1. روش زانویی
2. روش میانگین سیلوئت
3. روش Gap Statistic

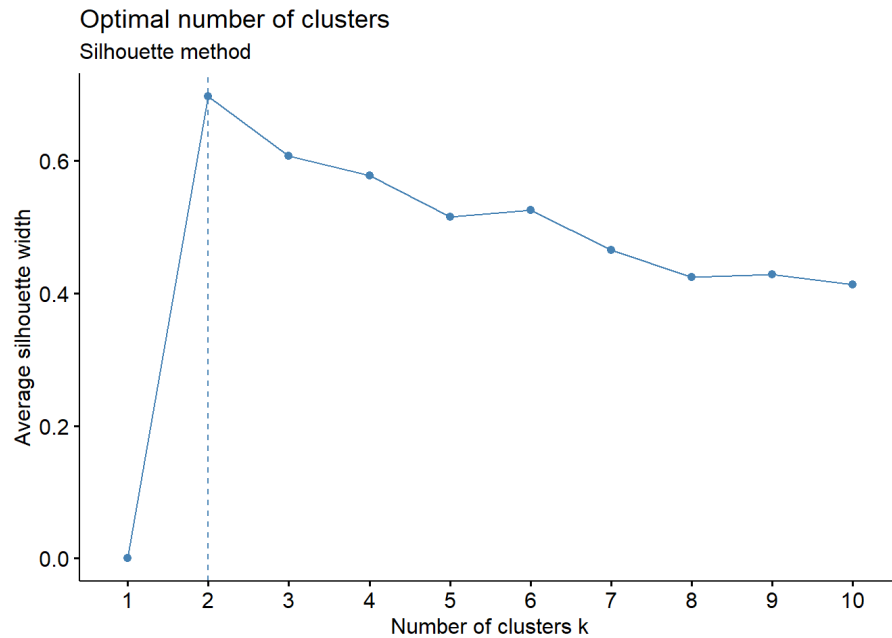
برای نمایش بهتر نتایج، ۳ تست اجرا خواهد شد.

```
set.seed(102)

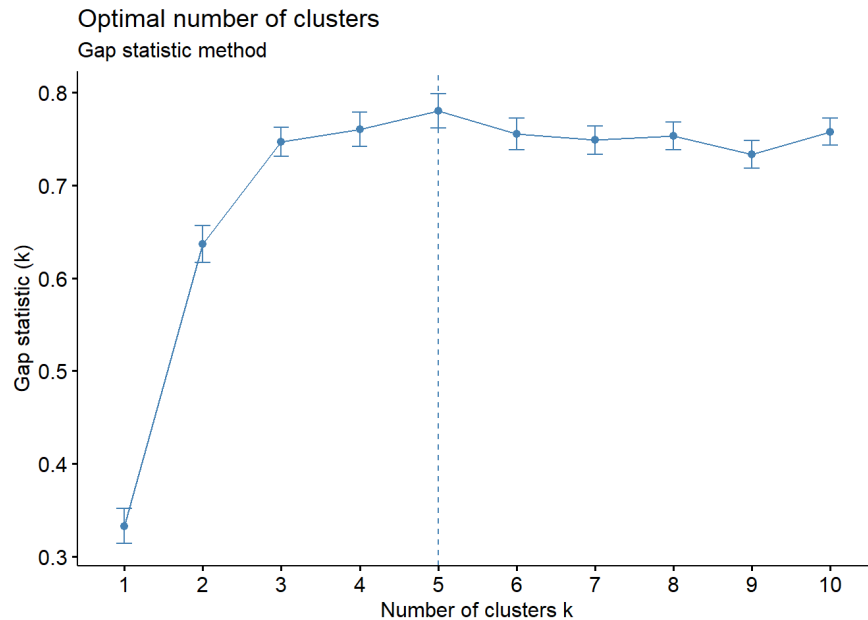
# Elbow method
fviz_nbclust(df.subset1, kmeans, method = "wss") +
  labs(subtitle = "Elbow method")
```



```
# Silhouette method
fviz_nbclust(df.subset1, kmeans, method = "silhouette") +
  labs(subtitle = "Silhouette method")
```



```
# Gap statistic
# nboot = 50 to keep the function speedy.
# recommended value: nboot= 500 for your analysis.
# Use verbose = FALSE to hide computing progression.
set.seed(123)
fviz_nbclust(df.subset1, kmeans, nstart = 25, method =
"gap_stat", nboot = 50)+
  labs(subtitle = "Gap statistic method")
```

۲ خوشه بر اساس تست های (۲ از ۳ تست) بالا انتخاب شده است. در ادامه، مرکز خوشه ها را با علامت مثلث به عنوان مرکز ترسیم می کنیم.

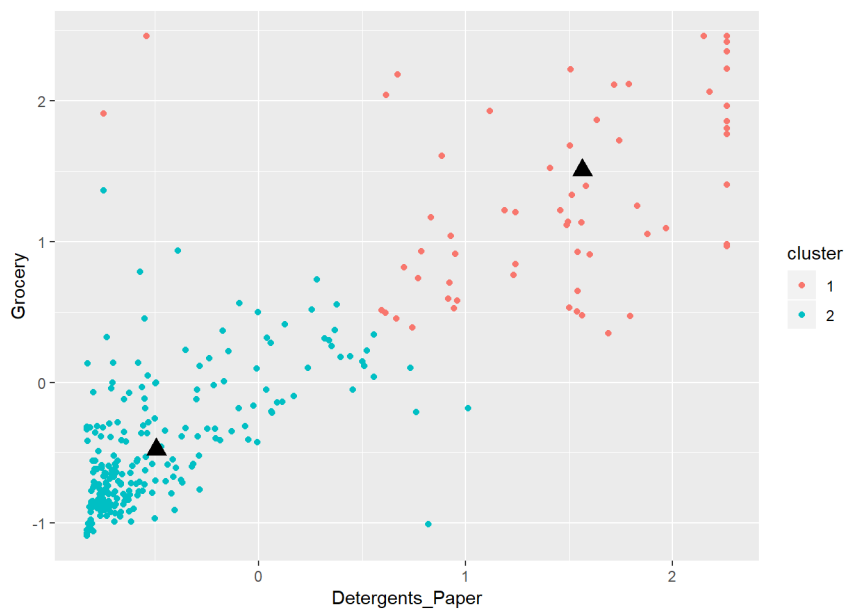
K means

```
set.seed(111)
kmean2.simple <- kmeans(df.subset1,centers=2, iter.max =
25, nstart=100)
df.subset1$cluster <- factor(kmean2.simple$cluster)
summary(df.subset1)
```

##	Grocery	Detergents_Paper	cluster
##	Min. : -1.0922	Min. : -0.8308	1: 74
##	1st Qu.: -0.7715	1st Qu.: -0.7341	2: 234
##	Median : -0.3624	Median : -0.5459	
##	Mean : 0.0000	Mean : 0.0000	
##	3rd Qu.: 0.4786	3rd Qu.: 0.5565	

```
## Max. : 2.4606 Max. : 2.2640
```

```
ggplot(data=df.subset1, aes(x=Detergents_Paper, y=Grocery, colour=cluster))+geom_point()+geom_point(data=as.data.frame(kmean2.simple$centers), color="black", size=4, shape=17)
```



نمودار بالا ۲ خوشه را بر اساس ۲ متغیر، *Grocery* و *Detergents_paper* بر اساس الگوریتم *k* میانگین نشان می دهد. می توان دید که این دو متغیر داده های ما را به خوبی از یکدیگر جدا کرده است و خوشه های مطلوبی ایجاد کرده است.

نتیجه گیری

همان طور که مشاهده کردیم نقاط دورافتاده تاثیر بسزایی در عملیات خوشه بندی دارند. در این مطالعه از دو متغیر *Grocery* و *Detergents_paper* برای خوشه بندی داده ها استفاده کردیم. این دو متغیر هم بستگی بالایی دارند. به عبارتی افزایش میزان *Grocery* می تواند نشان دهنده ی افزایش میزان *Detergents_paper* باشد. در خوشه بندی انجام شده دو خوشه یافت شد. یک خوشه که به صورت میانگین هزینه پایینی برای مواد شوینده، کاغذ و خواربار میپردازند و خوشه ای که هزینه بالایی برای این موارد پرداخت می کنند.

۷۵ خوشه بندی

یک تیم مارکتینگ قدرتمند می تواند از خروجی این خوشه بندی برای برنامه ریزی کمپین های مارکتینگ خود و ارایه پیشنهادات جذاب به مشتریان با توجه به خوشه ای که در آن قرار میگیرد

```
#kmeans
data(ChickWeight)
ChickWeight
install.packages("amap")
library(amap)
x1=ChickWeight$weight
x2=ChickWeight$Time
x=cbind(x1,x2)
x
cl<-Kmeans(x,4, iter.max =85, method = "euclidean")
plot(x, col = cl$cluster)
points(cl$centers, col = 1:4, pch = 8, cex=2)
data(iris)
iris2=iris
iris2$Species=NULL
kmeans.result=kmeans(iris2,3)
kmeans.result
table(iris$Species, kmeans.result$cluster)
plot(iris2[c("Sepal.Length", "Sepal.Width")],

col=kmeans.result$cluster,sub=paste("DataMining",format(Sys.time()
,"%Y-%b-%d %H:%S")))
points(kmeans.result$centers[,c("Sepal.Length", "Sepal.Width")],
col=1:3, pch=8, cex=2)
#####
#pam1
library(datasets)
data(USArrests)
df=scale(USArrests)
head(df,n=3)
install.packages(c("cluster","factoextra"))
library(cluster)
```

```

library(factoextra)
#the optimal number of clusters
library(cluster)
fviz_nbclust(df,cluster::pam,method="silhouette")
pam.res=pam(df,k=2,metric="euclidean" )
pam.res
d=cbind(USArrests,cluster=pam.res$cluster)
head(d,n=3)
library(factoextra)
fviz_cluster(pam.res,geom="point",ellipse.type="norm")
#pam2
library(fpc)
data(iris)
iris2=iris
iris2$Species=NULL
pamk.result=pamk(iris2,2:100)
pamk.result$nc
pamk.result$pamobject$medoids
table(pamk.result$pamobject$clustering,iris$Species)
#####
#سلسله مراتبی
library(stats)
library(factoextra)
n=5
k=3
data=iris[c(sample(x = 1:50,size = n),sample(x
=51:100,size=n),sample(x = 101:150,size= n)),1:4]
distmethod=c('euclidean')
linkagemethod=c("average")
distance=dist(data,distmethod)
hc=hclust(d = distance,method=linkagemethod)
distance
fviz_dist(dist.obj = distance)
hc
fviz_dend(x = hc,k = k)
#####

```

```

#DBSCAN
data("multishapes")
df=multishapes[1:2,]
library("fpc")
set.seed(123)
db=fpc::dbscan(df,eps=0.15,minpts=5)
print(db)
library("factoextra")
fviz_cluster(db,data=df,stand=FALSE,ellipse=FALSE,
show.clust.cent=FALSE
+geom="point",palette="jco",ggtheme=theme_classic())
#OPTICS
data(iris)
iris2=iris
iris2$Species=NULL
library(NbClust)
nc=NbClust(iris2,min.nc=2,max.nc=15,method="kmeans")
BAR=table(nc$Best.nc[1,])
BAR
barplot(BAR,xlab="Number Of Cluster",ylab="Number Of Criteria",
main = "Number
+Of Cluster Chosen by 26 Criterian")
#####3
library(stats)
library(ggplot2)
library(plyr)
library(dplyr)
library(data.table)
library(corrplot)
library(NbClust)
library(cluster)
df.original<-read.csv("C:/Users/ASUS/Desktop/Wholesale customers
data.csv")
df<- sample_frac(df.original,0.7)
df.index<- as.numeric(rownames(df))
df.test<- df.original[,-df.index,]

```

```

head(df)
summary(df)
summary(is.na(df))
table(df$Channel)
table(df$Region)
df %>% group_by(Channel,Region) %>%
  summarise(total_fresh = sum(df$Fresh), total_Milk = sum(df$Milk),
    total_Grocery= sum(df$Grocery),

    total_Frozen=sum(df$Frozen),total_Detergents_Paper=sum(df$Deter
    gents_Paper),
      total_Delicassen= sum(df$Delicassen))

temp <- reshape(df, direction="long",
  varying=c("Fresh","Milk","Grocery","Frozen","Detergents_Paper",
    "Delicassen"), v.names= "Total_price", timevar="Category",
  time=c("Fresh", "Milk","Grocery","Frozen","Detergents_Paper",
    "Delicassen"))

ggplot(temp, aes(x=temp$Category, y =temp$Total_price))
+geom_boxplot() +stat_boxplot(geom ='errorbar') +
theme(axis.text.x= element_text(angle=90,hjust=1))+ ggtitle("Product
Category Distribution")
cor.result<- cor(df)
pairs(df[, -c(1:2)], col=df$Channel, pch=19, lower.panel = NULL)
+title(main = "Category by Channel")
plot(df[, -c(1:2)], col=df$Region, pch=19, lower.panel = NULL)
+title(main = "Category by Region")

corrplot(cor.result, method="ellipse") +title(main = "Corelation by
category")
apply(X= df[, -c(1:2)], MARGIN=2, FUN =
function(x)length(boxplot.stats(x)$out))
sort(boxplot.stats(df$Grocery)$out)
quantile(df$Grocery, probs=seq(from =0.9, to=1,by=0.025))
grocery.max <- as.numeric(quantile(df$Grocery,probs=0.95))

```

```

df$Grocery[df$Grocery > grocery.max] <- grocery.max
sort(boxplot.stats(df$Detergents_Paper)$out)
quantile(df$Detergents_Paper, probs=seq(from =0.9, to=1,by=0.025))
grocery.max <-
as.numeric(quantile(df$Detergents_Paper,probs=0.925))
df$Detergents_Paper[df$Detergents_Paper > grocery.max] <-
grocery.max

ggplot(data=df, aes(x=Grocery, y =Detergents_Paper)) +
  geom_point(shape=1) +geom_smooth(method="lm")

df.subset1<-as.data.frame(df[,c("Grocery","Detergents_Paper")])
summary(df.subset1)

df.subset1<- as.data.frame(scale(df.subset1))
summary(df.subset1)

set.seed(102)
# Elbow method
fviz_nbclust(df.subset1, kmeans, method = "wss") +
  labs(subtitle = "Elbow method")

fviz_nbclust(df.subset1, kmeans, method = "silhouette")+
  labs(subtitle = "Silhouette method")

set.seed(123)
fviz_nbclust(df.subset1, kmeans, nstart = 25, method = "gap_stat",
  nboot = 50)+
  labs(subtitle = "Gap statistic method")

set.seed(111)
kmean2.simple <- kmeans(df.subset1,centers=2, iter.max = 25,
  nstart=100)
df.subset1$cluster <- factor(kmean2.simple$cluster)
summary(df.subset1)

```



```
ggplot(data=df.subset1, aes(x=Detergents_Paper, y=Grocery,
colour=cluster))+geom_point()+geom_point(data=as.data.frame(kme
an2.simple$centers), color ="black", size=4, shape =17)
```

