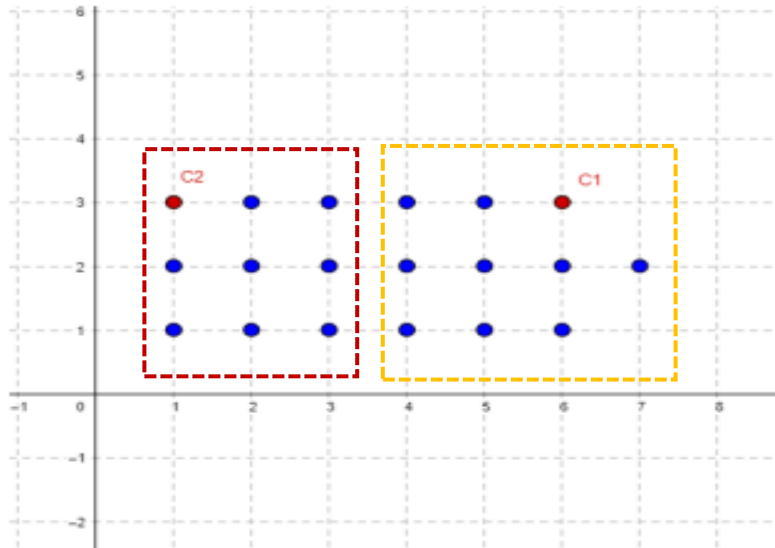


۱) ابتدا دو centroid داده شده است. حال باید فاصله تمامی نقاط را با دو نقطه c_1 و c_2 حساب کنیم و نسبت به هر کدام که نزدیک تر بود نقطه را به آن خوشه باید assign کنیم. خوشه بندی ها طبق شکل به صورت زیر انجام میشود:



نقطه های موجود در مستطیل قرمز رنگ در یک خوشه و نقاط موجود در مستطیل زرد رنگ در خوشه دیگر قرار میگیرند. پس خوشه های مربوط به هر نقطه مشخص شد حال باید مراکز خوشه هارا آپدیت کنیم، داریم:

خوشه ۲ :

$$X = (1 * 3 + 2 * 3 + 3 * 3) / 9 = 2$$

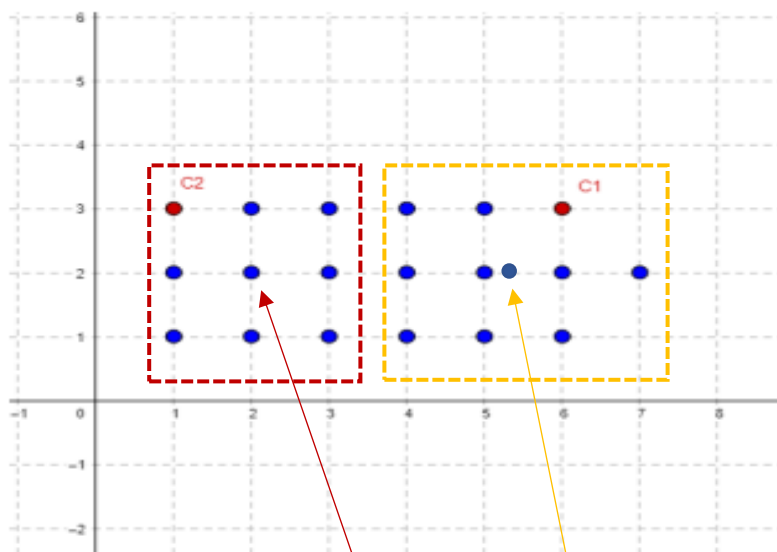
$$Y = (1 * 3 + 2 * 3 + 3 * 3) / 9 = 2$$

خوشه ۱ :

$$X = (4 * 3 + 5 * 3 + 6 * 3 + 7 * 1) / 10 = 5.2$$

$$Y = (1 * 3 + 2 * 4 + 3 * 3) / 10 = 2$$

حال مراکز آپدیت شدند. بار دیگر که فواصل نقاط با این مراکز را حساب بکنیم خوشه بندی ها عوض نمیشوند پس درنتیجه خوشه بندی نهایی به صورت زیر است:



$C2 (2, 2)$

$C1(5.2, 2)$

(۲)

برای دسته بندی داده های زیر از الگوریتم DBSCAN استفاده میکنیم. این الگوریتم ویژگی های زیر را دارد

- خوشه هایی که تشکیل میدهد میتوانند از نظر شکل و اندازه با یکدیگر متفاوت باشند
- تعداد خوشه ها لازم نیست از قبل مشخص شود، درواقع تعداد خوشه ها دیگر مثل الگوریتم Kmeans یک Hyperparameter نیست.
- داده های نویز و outlier ها را برخلاف الگوریتم Kmeans به خوبی هندل میکند و نمیگذارد این داده ها مثل الگوریتم Kmeans باعث ایجاد خوشه های نامناسب بشوند.
- مناطق چگال رو از مناطقی که نقاط چگالی کمتری دارند متمایز میکند.

طبق ویژگی های گفته شده برای داده های زیر از الگوریتم DBSCAN استفاده میکنیم زیرا طبق شکل میتوان متصور شد دو منطقه چگال داریم که الگوریتم DBSCAN به خوبی آنها را تشخیص میدهد.

الگوریتم DBSCAN دارای دو هایپرپارامتر با نام های ϵ و minPoints میباشد. حداقل تعداد نقاطی است که باید در همسایگی یک نقطه در شعاع همسایگی ϵ باشند تا بتوانیم یک خوشه ایجاد کنیم. متغیر ϵ نیز میزان شعاع همسایگی را مشخص میکند.

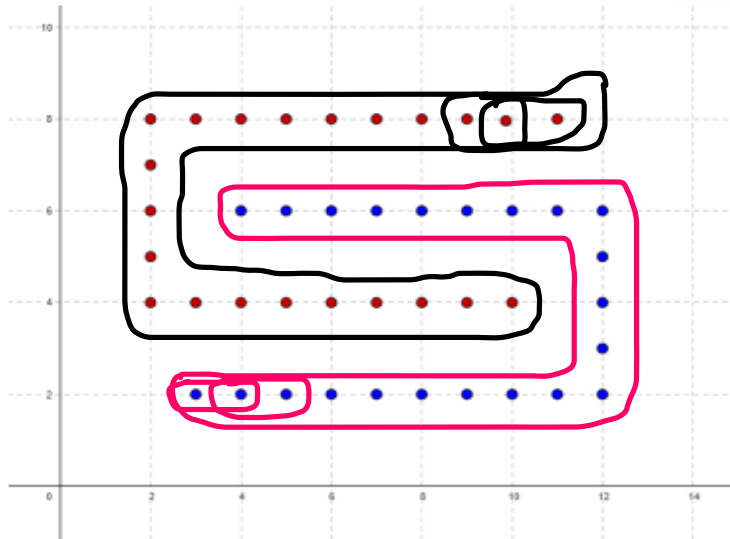
الگوریتم DBSCAN بدینصورت عمل میکند که ابتدا یک نقطه به صورت رندوم انتخاب میکنیم و تعداد همسایه های آن نقطه را بدست می آوریم. اگر تعداد همسایه ها از minPoints کمتر بود آنگاه آن نقطه را outlier تشخیص میدهیم. در غیر اینصورت آن نقطه هسته است و میتوانیم خوشه را ایجاد کنیم. بدین ترتیب تمام نقاط را بررسی میکنیم تا همه نقاط بررسی شوند. همچنین نقاطی که در داخل خوشه قرار میگیرند ولی خود هسته نیستند border نامیده میشوند و داده هایی که در داخل هیچ خوشه ای جای ندارند outlier نامیده میشوند.

برای این مسئله هایپرپارامتر ها را به صورت زیر در نظر میگیریم:

$$\text{minPoints} = 2$$

$$\epsilon = 1$$

پس الگوریتم به صورت زیر ایجاد میشود:



مطابق شکل مشاهده میکنیم تعداد خوشه ها پس از اجرای الگوریتم DBSCAN برابر ۲ میشود و خوشه به درستی بر اساس تراکم نقاط تعیین شدند درحالی که اگر از الگوریتم Kmeans استفاده میکردیم خوشه بندی ها با هر نقاط ابتدایی به این صورت نمیشدند و خوشه بندی با الگوریتم Kmeans کیفیت کمتری داشت.

(۳

Closed frequent itemset ها زیر مجموعه frequent itemset ها هستند. همچنین مجموعه ای بسته است و مقدار support آن بزرگتر و مساوی minsup میباشد. همچنین itemset ای بسته است که هیچ superset ای از این مجموعه وجود نداشته باشد که مقدار support یکسانی با آن داشته باشد. در اقع رابطه $\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$ برقرار است.

همچنین کنار هر نود میزان support آن را مشخص میکنیم و بدین صورت closed frequent itemset ها را بدست می آوریم. طبق رابطه ذکر شده حال با داشتن closed frequent itemset ها میتوانیم مقدار support را برای frequent itemset ها بدست آوریم.

(۴)

$$\text{support} = 33\%$$

$$\text{confidence} = 6\%$$

1- itemsets:

count	support	
سب	۴	$\frac{4}{12} \geq \frac{1}{3}$ ✓
پرتقال	۳	$\frac{3}{12} \geq \frac{1}{3}$ ✓
هوز	۳	$\frac{3}{12} \geq \frac{1}{3}$ ✓
انار	۳	$\frac{3}{12} \geq \frac{1}{3}$ ✓
نارنگی	۲	$\frac{2}{12} \geq \frac{1}{3}$ ✓

دو itemset ای هوز و سب سرور →

2- itemsets:

count	support	
سب و پرتقال	۲	$\frac{2}{12} \geq \frac{1}{3}$ ✓
سب و هوز	۲	$\frac{2}{12} \geq \frac{1}{3}$ ✓
سب و انار	۱	$\frac{1}{12} < \frac{1}{3}$ ✗
سب و نارنگی	۲	$\frac{2}{12} \geq \frac{1}{3}$ ✓
پرتقال و هوز	۲	$\frac{2}{12} \geq \frac{1}{3}$ ✓
پرتقال و انار	۱	$\frac{1}{12} < \frac{1}{3}$ ✗
پرتقال و نارنگی	۰	$\frac{0}{12} < \frac{1}{3}$ ✗
هوز و انار	۱	$\frac{1}{12} < \frac{1}{3}$ ✗
هوز و نارنگی	۰	$\frac{0}{12} < \frac{1}{3}$ ✗
انار و نارنگی	۱	$\frac{1}{12} < \frac{1}{3}$ ✗

→

نسب و برتقال و صوز	count ۲	✓
نسب و برتقال و نارنگی	۰	X
نسب و صوز و نارنگی	۰	X
برتقال و صوز و نارنگی	۰	X

نسب و برتقال و صوز { نسبت و برتقال و صوز } به صورت زیر نشان داده می شود:

{ نسبت } → { برتقال و صوز }	$confidence = \frac{6 \text{ (نسب و برتقال و صوز)}}{6 \text{ (نسب)}} = \frac{2}{3} < 0.6 \times$
{ برتقال } → { نسبت و صوز }	$confidence = \frac{3}{3} > 0.6 \checkmark$
{ نسبت و برتقال } → { صوز }	$confidence = \frac{2}{3} > 0.6 \checkmark$
{ نسبت و برتقال } → { صوز }	$confidence = \frac{2}{3} > 0.6 \checkmark$
{ نسبت و صوز } → { برتقال }	$confidence = \frac{2}{3} > 0.6 \checkmark$
{ صوز و برتقال } → { نسبت }	$confidence = \frac{2}{3} > 0.6 \checkmark$

نسب روابط به صورت زیر هستند:

- ① { نسبت } → { صوز و برتقال } $c = 1$
- ② { برتقال } → { نسبت و صوز } $c = 1$
- ③ { صوز } → { نسبت و برتقال } $c = 1$
- ④ { نسبت و برتقال } → { صوز } $c = 0.66$
- ⑤ { نسبت و صوز } → { برتقال } $c = 0.66$

(۵۲)

Single link و Single link باید کتبی قواعد را در جدول انتخاب کنیم و در جدول ماتریس قواعد را بروز کنیم. (قطر اصلی شکل است که آثار را در نظر بگیریم. داریم %

① کتبی قواعد در ماتریس $P_1 P_2 = 0, 1$ است
لیست داریم و انتخاب اولی و دومی بصورت $P_1 P_2$.

② حال باید ماتریس را بروز کنیم. داریم
در روز شنبه باید روزی که $P_1 P_2$ با نقاط دیگر قواعد است
با آن نقاط است.

③ حال کمترین فاصله، $P_1 P_2$ با P_5 یعنی ۳ است. (۳)
پس نویسیم در آنجا که باید انتخاب شود.

④ حال ماتریس را بروز می کنیم

⑤ حال کمترین مقدار ۴ است پس داریم
و در آنجا انتخاب می شود.

⑥ ماتریس را بروز می کنیم

⑦ در نهایت P_4 جویندگی می شود داریم

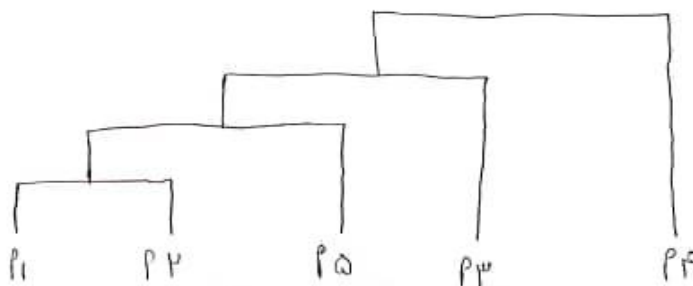
$P_1 P_2 P_5 P_3 P_4$

	$P_1 P_2$	P_3	P_4	P_5
$P_1 P_2$	0	0, 1	0, 4	0, 3
P_3	0, 4	0	0, 4	0, 1
P_4	0, 4	0, 4	0	0, 7
P_5	0, 3	0, 1	0, 7	0

	$P_1 P_2 P_5$	P_3	P_4
$P_1 P_2 P_5$	0	0, 4	0, 4
P_3	0, 4	0	0, 4
P_4	0, 4	0, 4	0

	$P_1 P_2 P_5 P_3$	P_4
$P_1 P_2 P_5 P_3$	0	0, 4
P_4	0, 4	0

در درخت آماری به صورت زیر است:



complete link و برگشت single link هنگام آپدیت کردن ماتریس به جای min از max استفاده کنیم.

① کمترین فاصله برای P_1P_2 است با مقدار ۱۰۰

	P_1P_2	P_3	P_4	P_5
P_1P_2	۰	۰/۴۴	۰/۵۵	۰/۹۸
P_3	۰/۴۴	۰	۰/۴۴	۰/۱۸۵
P_4	۰/۵۵	۰/۴۴	۰	۰/۷۴
P_5	۰/۹۸	۰/۱۸۵	۰/۷۴	۰

② بروز کردن ماتریس :

③ کمترین مقدار P_3P_4 با مقدار ۰ است

	P_1P_2	P_3P_4	P_5
P_1P_2	۰	۰/۴۴	۰/۹۸
P_3P_4	۰/۴۴	۰	۰/۱۸۵
P_5	۰/۹۸	۰/۱۸۵	۰

④ بروز کردن ماتریس :

⑤ کمترین فاصله برای $P_1P_2-P_3P_4$ است با ۰/۴۴

	$P_1P_2P_3P_4$	P_5
$P_1P_2P_3P_4$	۰	۰/۹۸
P_5	۰/۹۸	۰

⑥ بروز کردن ماتریس :

⑦ در آخر P_5 انتخابی شود داریم
 $P_1P_2P_3P_4P_5$

