

- هر image بعد  $70 \times 70$  دارد ← هر sample 5600 تا جمع دارد (از همان پیکسل ها تصویر را به این اندازه استفاده میکنیم)  
(مقادیر gray scale اند)

: K Means

- از کتابخانه sklearn استفاده میکنیم.

- انتخاب  $seed$  مناسب،  $k=41$ ، یکبار اجرا کنیم  
Rand Index = 0.97

: Agglomerative

: Average link

- از کتابخانه sklearn استفاده میکنیم.

- انتخاب  $rand-index = 0.95$  ( $k=41$ )

: Single link

- انتخاب  $rand-index = 0.63$  ( $k=60$ )

(با انتخاب  $k=41$  خوشه ها بیش از حد مطلوب برود  
Over clustering  
قابل تحمل است)

over clustering

: Complete link

- انتخاب  $rand-index = 0.97$  ( $k=41$ )

: DBSCAN

- از کتابخانه sklearn استفاده میکنیم. (از همان مقادیر استفاده میکنیم)

- برای انتخاب  $\epsilon$  از آنتریج  $k$ -nearest neighbors استفاده میکنیم. (توصیحات استفاده از KNN در لینک که)

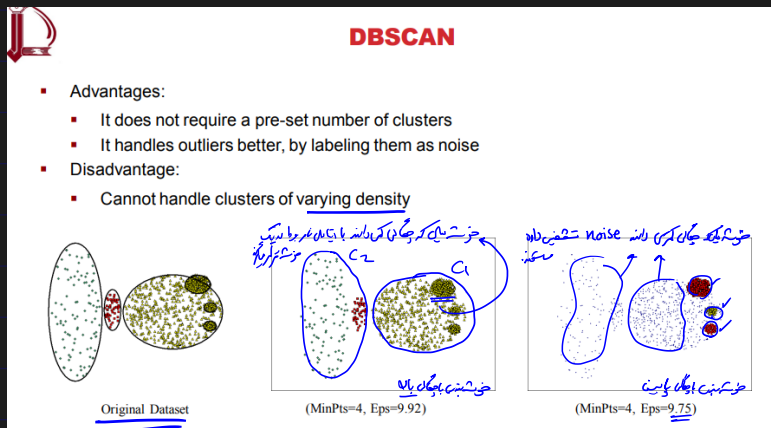
در Notebook در قسمت DBSCAN  
تغییر داد و تست میکنیم.

- دسته ها را بر DBSCAN داده مارا میزنیم (تا انتخاب صحیح را متوجه شویم)

- انتخاب  $rand-index = 0.82$  ( $\epsilon=55$ ,  $min\_clusters=4$ )

داده های نویزی noise تشخیص داده میشوند، راه حل این مشکل که ناسازگار است  
چک کنیم ضربه ما  
در قسمت به توضیح داده نموده است

- از مشکلات اصلی اگرستج dbscan این است که خوشه‌ها را با چگالی‌ها متفاوت را به یکدیگر پیوند می‌دهد.



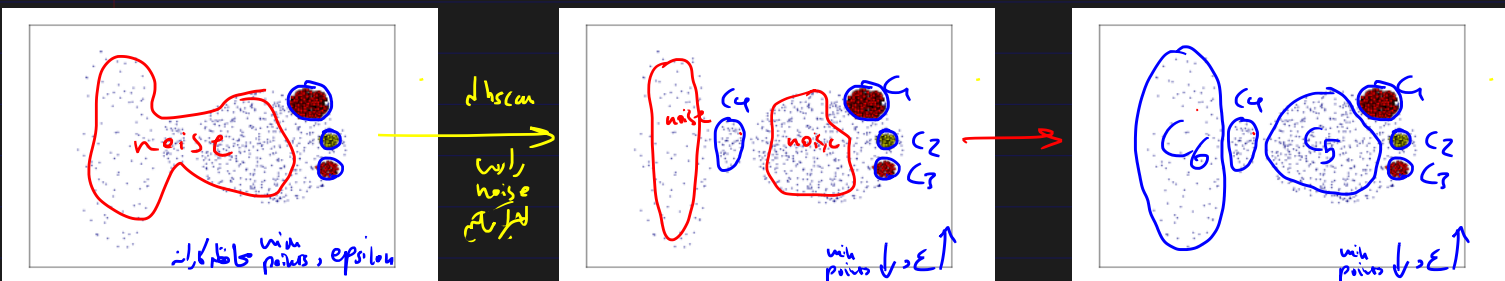
یک راه حل برای حل این مشکل این است که اگرستج DBSCAN را بر اساس  $\epsilon$  و  $\text{min points}$  مختلف اجرا کنیم.

- این مشکل در دیتاست  $\epsilon$  و  $\text{min points}$  هم دیده می‌شود (تغییر  $\epsilon$  و  $\text{min points}$  با این  $\epsilon$  به دنبال  $\text{noise}$  شایسته می‌شوند)

- به این رفع این مشکل در دیتاست داده شده، از این روش استفاده می‌کنیم که به صورت حفاظت‌گزارانه  $\epsilon$  و  $\text{min points}$  را انتخاب می‌کنیم (به کم و بیش)

و در مرحله بعدی در دیتاست که در مرحله قبل  $\text{noise}$  تشخیص داده شده دوباره DBSCAN را با  $\epsilon$  و  $\text{min points}$  را افزایش می‌دهیم (در مرحله  $\text{min points}$  را کم می‌کنیم)

با اینکار خوشه‌ها را با چگالی مختلف به هم می‌زنیم. (از آنجایی که ما  $\text{min points}$  را کم می‌کنیم)



- این روش در کتاب Enhanced DBSCAN به یاد می‌آید.

به این صورت که در هر مرحله از اجرای DBSCAN، خوشه‌ها را به هم می‌زنیم و  $\text{min points}$  را افزایش می‌دهیم و  $\epsilon$  را کم می‌کنیم.

انتخاب  $\epsilon$  و خوشه‌ها را به هم می‌زنیم و  $\text{min points}$  را افزایش می‌دهیم و  $\epsilon$  را کم می‌کنیم.

همین روش  $\epsilon$  و  $\text{min points}$  را برای دیتاست  $\text{noise}$  انتخاب می‌کنیم و خوشه‌ها را به هم می‌زنیم و  $\text{min points}$  را افزایش می‌دهیم و  $\epsilon$  را کم می‌کنیم.

با انتخاب  $\epsilon = 5.5$  و  $\text{min points} = 4$  در مرحله اول و  $\epsilon = 6.0$  و  $\text{min points} = 3$  در مرحله دوم است.  $\text{rand index} = 0.97$  را می‌بینیم.

و مهم‌تر از همه  $\text{noise}$  کمتر می‌شود.

## \* Rand Index

$TP+FP =$  تعداد اجتناب شده در یک cluster

Some cluster

$$= \sum_{i=1}^k \binom{N_i}{2}$$

$N_i =$  population of cluster  $i$ ,  
 $k =$  number of clusters

$TP =$  Similar Pairs in a cluster (Count)

$$FP = (TP + FP) - TP$$

$$FN = FN_{i_1} + FN_{i_2} + \dots + FN_{i_t}$$

$$= \left( \binom{N_1}{2} - TP_1 \right) + \left( \binom{N_2}{2} - TP_2 \right) + \dots + \left( \binom{N_t}{2} - TP_t \right)$$

$$\underline{\underline{1}} = \text{false negative} = \text{اجتناب شده (similar)} - \text{true positive}$$

similar in diffrent

similar in same cluster

$$TN = N - (FN + TP + FP)$$