

Tutorial DBSCAN

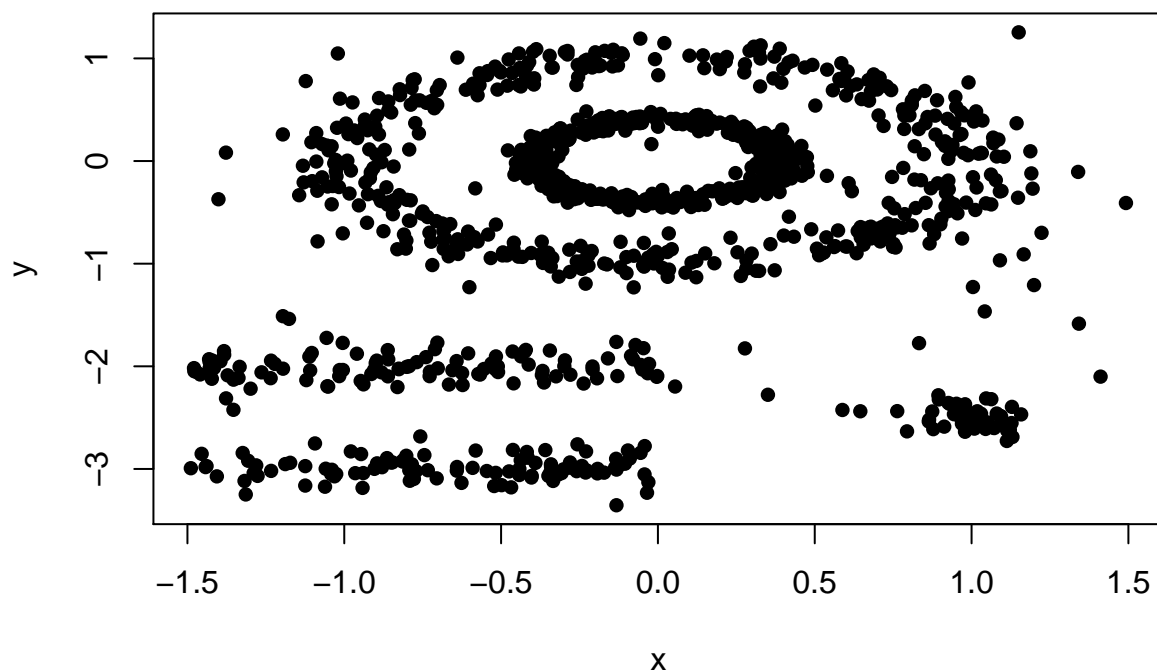
Tutorial ini adalah bagian dari mata kuliah Fundamen Sains Data, Informatika, UII.

Pada tutorial kali ini kita akan membanding hasil dari K-Means dan DBSCAN pada sebuah data set yang memiliki *odd shape* atau bentuk yang tidak biasa. Di sini kita akan menggunakan package **factoextra**. Kita akan menggunakan data set **multishapes**, yang berisi 1100 observations dengan dua variabel (kita mengambil dua kolom pertama pada baris ke-3 dari script di bawah ini) yang merepresentasikan berbagai macam bentuk.

```
1 library(factoextra)
2 data("multishapes")
3 df <- multishapes[, 1:2]
```

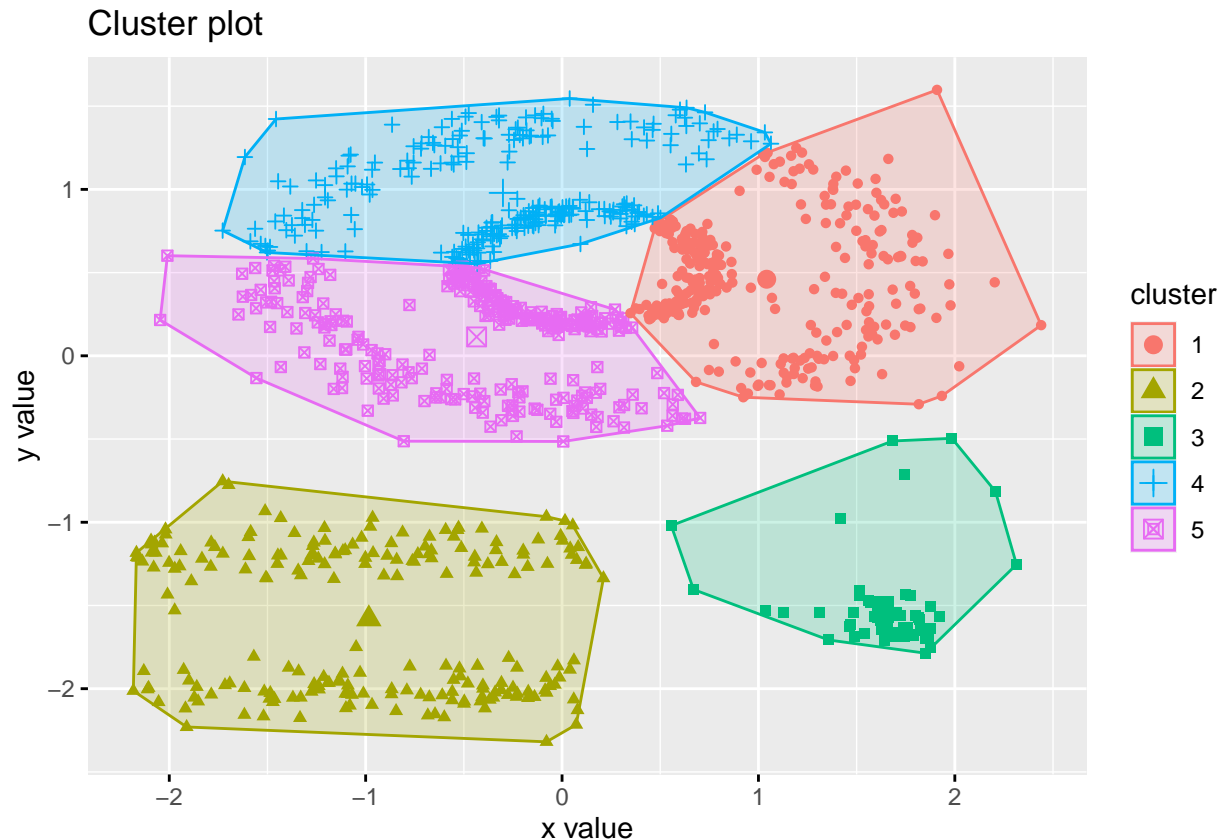
Selanjutnya kita plot data set tersebut, untuk memperoleh gambaran. Maksud dari **pch=16** adalah tipe simbol ke-16 untuk sebuah data point; di sini, kita memilih simbol sebuah lingkaran solid untuk sebuah data point.

```
plot(df, pch=16)
```



Berdasarkan plot data set di atas, secara intuisi kita melihat ada sekitar 5 cluster. Mari kita coba aplikasi K-Means dengan 5 cluster.

```
1 set.seed(123)
2 km.res <- kmeans(df, 5, nstart = 25)
3 fviz_cluster(km.res, df, frame = FALSE, geom = "point")
```

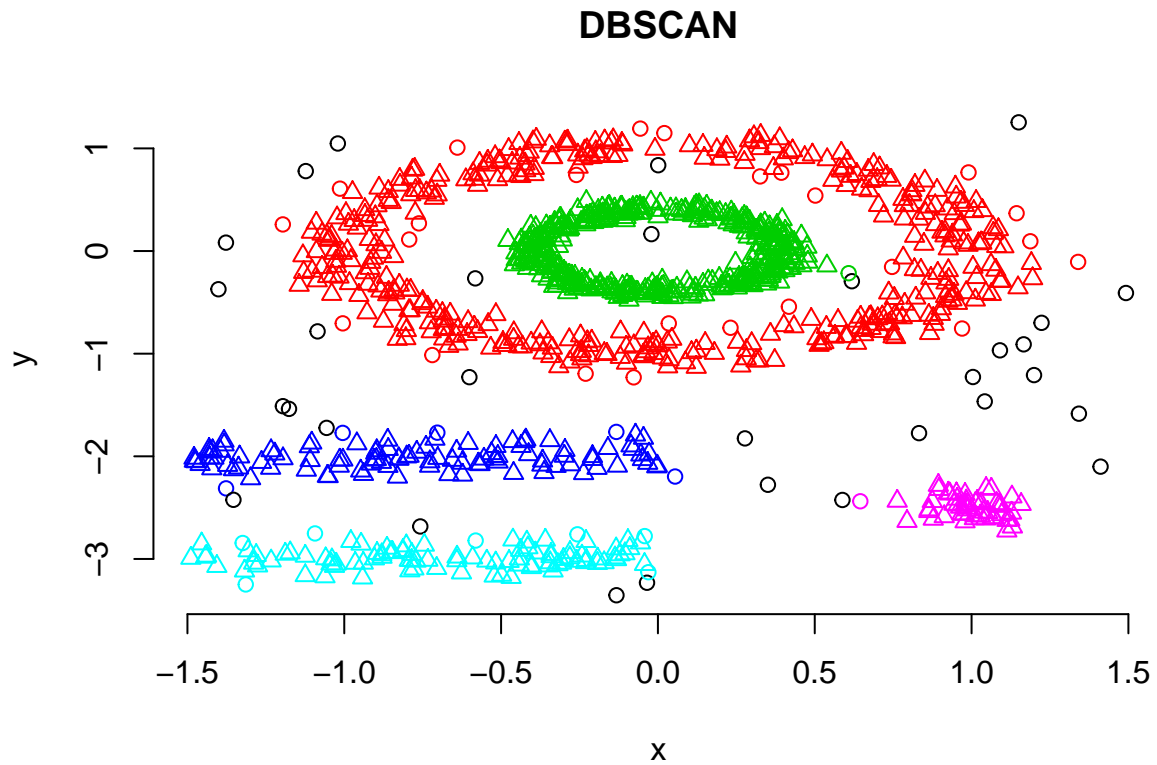


Dapat kita lihat, K-Means memberikan 5 cluster yang berbeda dengan apa yang kita bayangkan. Di sini, adalah sebuah contoh di mana sebaiknya kita tidak menggunakan K-Means pada data yang memiliki bentuk sebaran yang tidak lazim.

Selanjutnya, mari kita aplikasikan DBSCAN pada data yang sama. Ingat pada DBSCAN, kita tidak perlu menentukan jumlah cluster di awal.

Aplikasi DBSCAN Di sini kita akan menggunakan fungsi `dbscan()` dari package `fpc`, dengan parameter $\epsilon = 5$ dan `minPts = 5`.

```
1 library("fpc")
2 set.seed(123)
3 db <- dbscan(df, eps = 0.15, MinPts = 5)
4 plot(db, df, main = "DBSCAN", frame = FALSE)
```



Dari hasil di atas, dapat kita lihat jika DBSCAN mampu menemukan 5 cluster yang kita bayangkan secara intuisi beserta data point yang sekiranya termasuk noise (simbol lingkaran tanpa warna).

Sebagai alternatif untuk plot, kita juga dapat menggunakan sebuah fungsi untuk plot cluster dari **factoextra**.

```
library("factoextra")  
fviz_cluster(db, df, stand = FALSE, frame = FALSE, geom = "point")
```


Dari informasi di atas, kita dapat lihat setiap data point (dari 1 hingga ke 1100) masuk ke dalam cluster berapa; jika 0, dianggap sebagai noise.

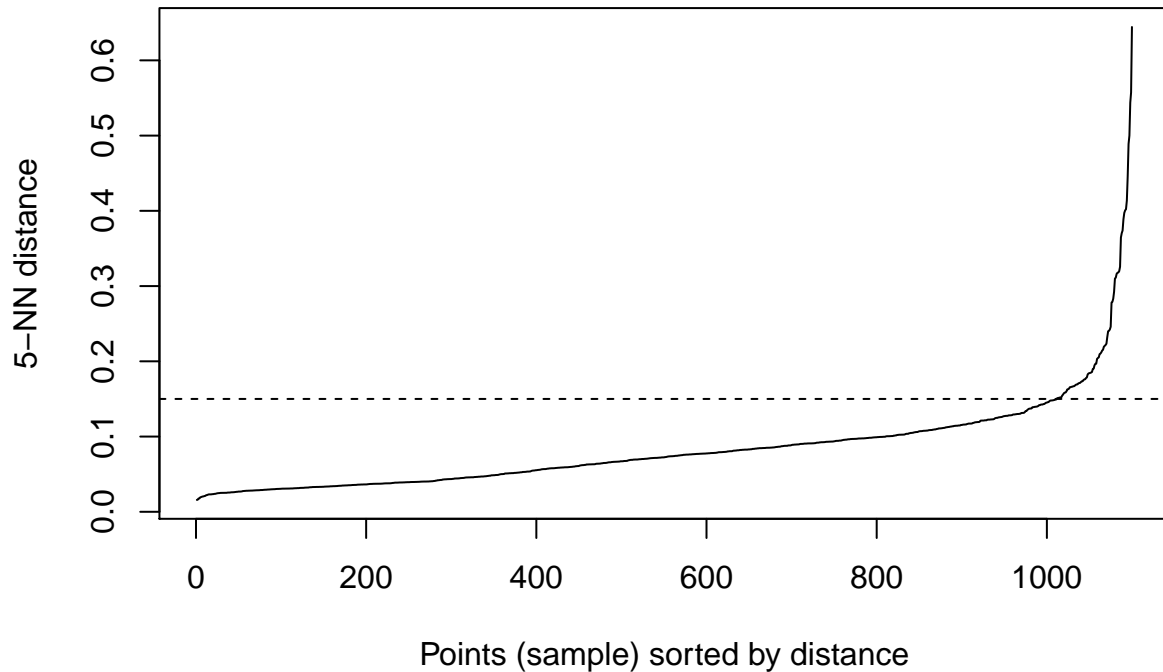
Dalam menentukan ϵ , kita dapat menggunakan k-distance graph. Metode ini melihat jarak setiap data point dari k-th nearest neighbor (data point terdekat). Jarak-jarak yang telah dihitung, diurutkan dari besar ke kecil; kemudian kita tentukan threshold yang menjadi nilai ϵ pada bagian “lembah” atau bagian dari kurva di mana kurva mengalami pembelokan (*lembah*). Di sini kita akan menggunakan fungsi `kNNdistplot` dari package `dbscan`.

Masih menggunakan data set yang sama (`df`) dengan contoh di atas, kita ingin melihat k-distance graph dengan `k=5`.

```
##
## Attaching package: 'dbscan'

## The following object is masked from 'package:fpc':
##
##      dbscan
```

```
kNNdistplot(df, k = 5)
abline(h = 0.15, lty = 2)
```



Dari plot yang di dapat, kita dapat lihat bahwa lembah pada kurva adalah ketika jarak = 0.15. Nilai inilah yang kita jadikan sebagai nilai ϵ .

Alternatif lain dalam menentukan *minPts* adalah dengan menghitung rata-rata dari nilai k-distance dari semua data point. Nilai rata-rata tersebut

Latihan Lakukan komputasi DBSCAN yang sama seperti di atas atau pada data lain yang Anda pilih sendiri, kemudian lihat perbedaan ketika

1. kita mengubah nilai ϵ ketika nilai *minPts* = 5 tetap
2. kita mengubah nilai *minPts* ketika nilai ϵ = 0.15 tetap

Jelaskan secara singkat, apa yang terjadi jika kita set *minPts* semakin besar atau semakin kecil, dan yang terjadi ketika ϵ semakin besar atau kecil.

Referensi http://www.sthda.com/english/wiki/wiki.php?id_contents=7940