

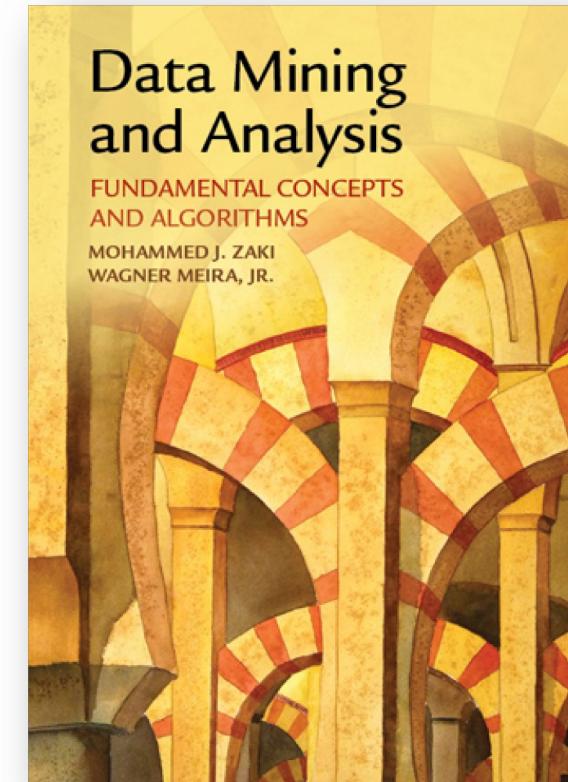
Density Based Clustering

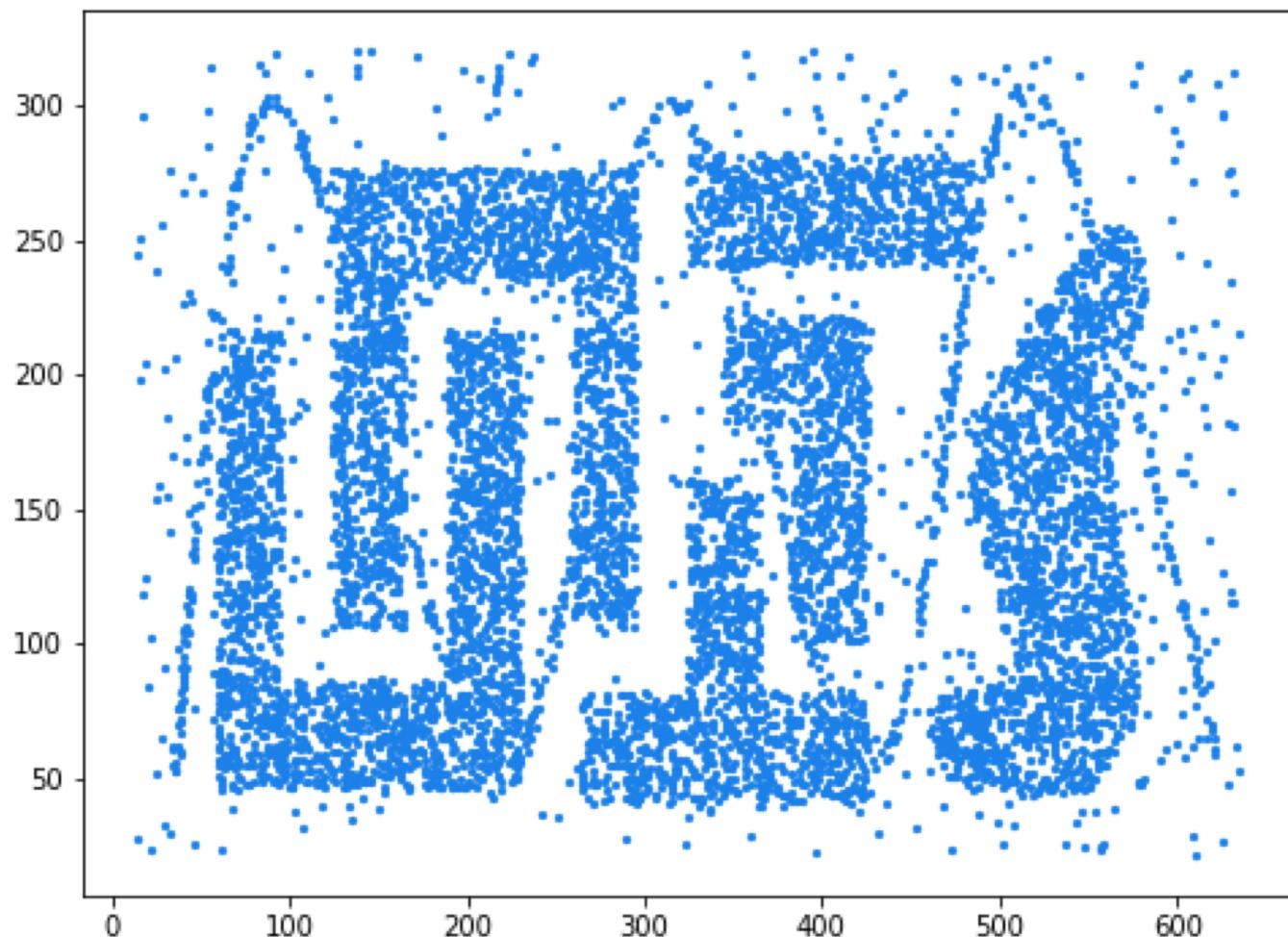
Data Science for Mobility

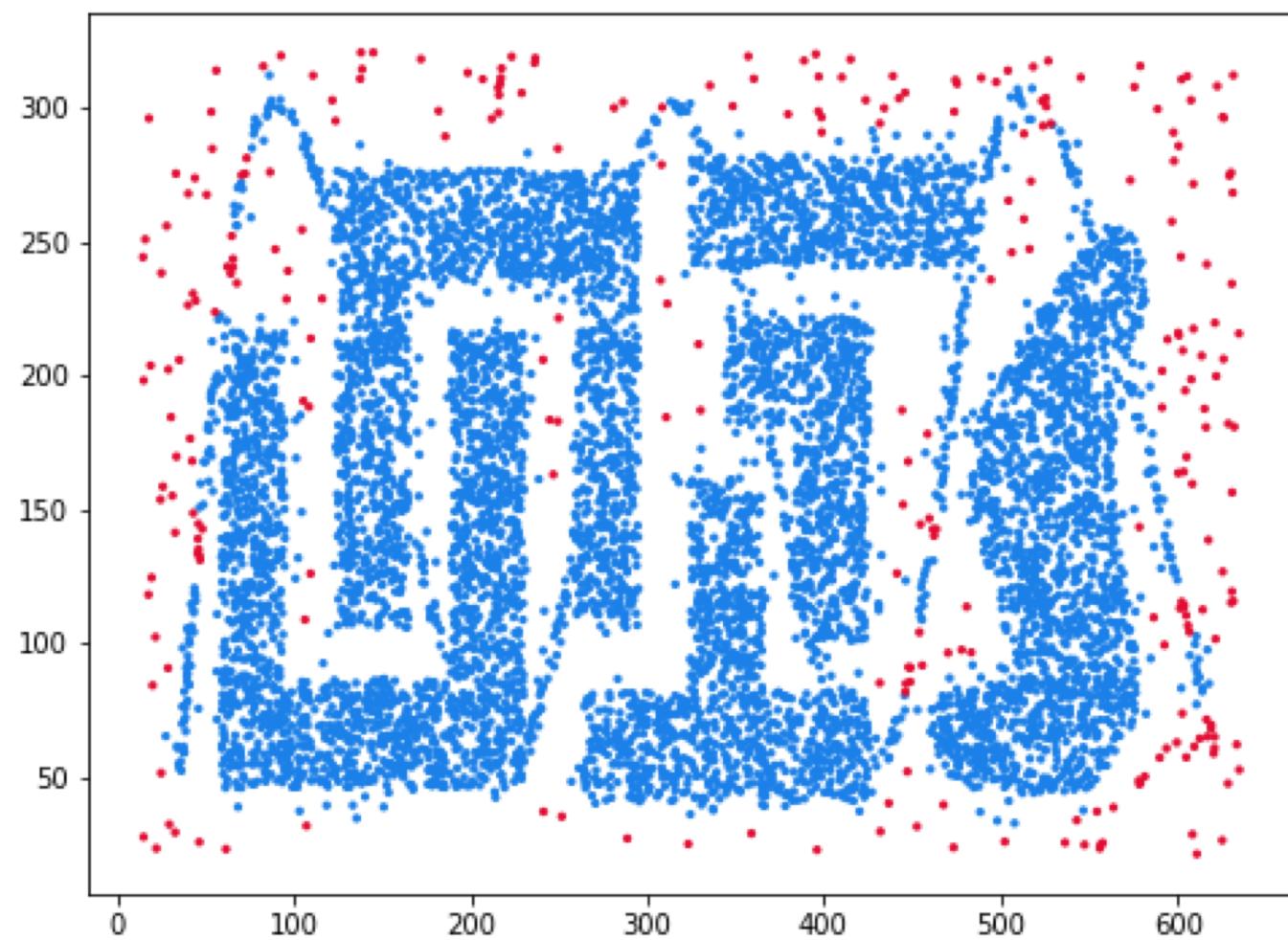


Readings

- “Data Mining and Analysis” by Zaki & Meira
 - Chapter 15
- <http://www.dataminingbook.info>







- Clustering based on density (local cluster criterion), such as density-connected points
- Major features:
 - Discover clusters of arbitrary shape
 - Handle noise
 - One scan
 - Need density parameters as termination condition
- Several interesting studies:
 - DBSCAN: Ester, et al. (KDD'96)
 - OPTICS: Ankerst, et al (SIGMOD'99).
 - DENCLUE: Hinneburg & D. Keim (KDD'98)
 - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)

- The neighborhood within a radius ϵ of a given object is called the ϵ -neighborhood of the object

$$N_\epsilon(x) = \{y | \delta(x, y) \leq \epsilon\}$$

- **Core Object**
 - If the ϵ -neighborhood of an object contains at least **minpts** objects, then the object is a core object
- **Directly density reachable**
 - An object x is directly density-reachable from object y if x is within the ϵ -neighborhood of y and y is a core object

- **Density Reachable**
 - An object x is density-reachable from object y if there is a chain of objects x_1, \dots, x_n where $x_1=x$ and $x_n=y$ such that x_{i+1} is directly density reachable from x_i
- **Density Connected**
 - An object p is density-connected to q with respect to ϵ and MinPts if there is an object o such that both p and q are density reachable from o
- **Density-Based Cluster**
 - A density-based cluster is defined as a maximal set of density connected points.

- Density corresponds to have at least minpts points within a specified radius ϵ
- A border point has fewer than minpt within ϵ , but is in the neighborhood of a core point
- A noise point is any point that is not a core point nor a border point

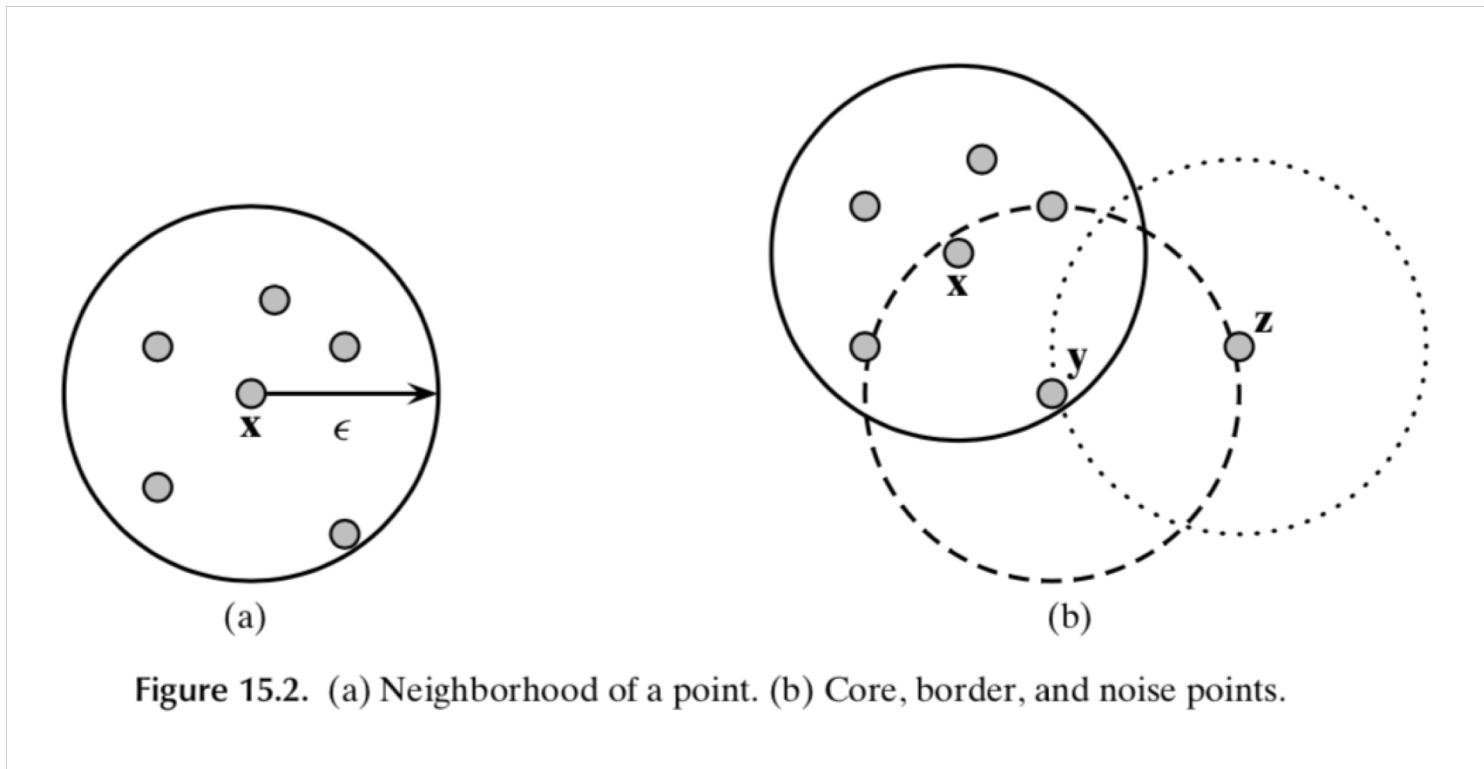


Figure 15.2. (a) Neighborhood of a point. (b) Core, border, and noise points.

Core, border and noise points when minpts is 6

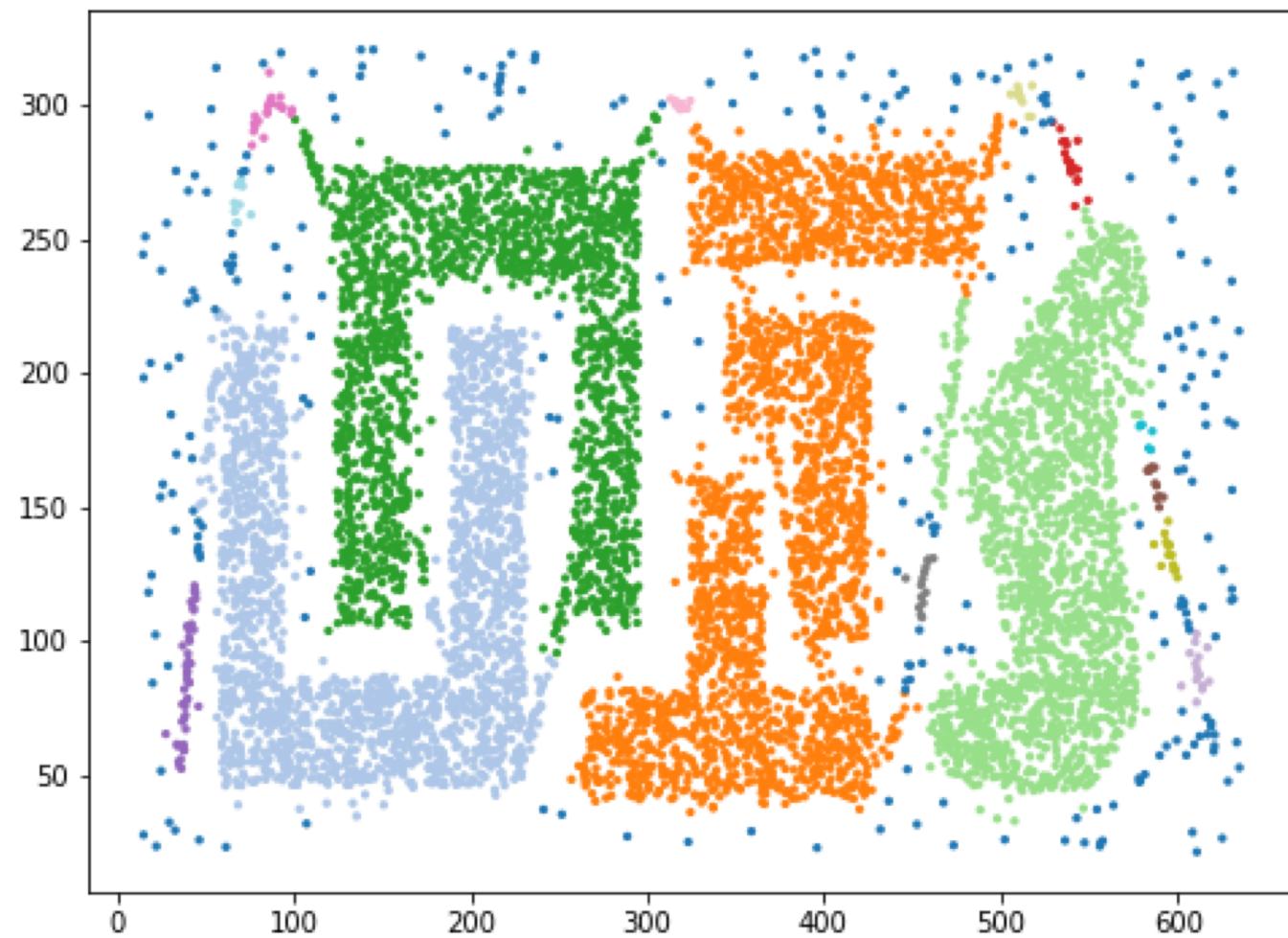
ALGORITHM 15.1. Density-based Clustering Algorithm

DBSCAN ($\mathbf{D}, \epsilon, minpts$):

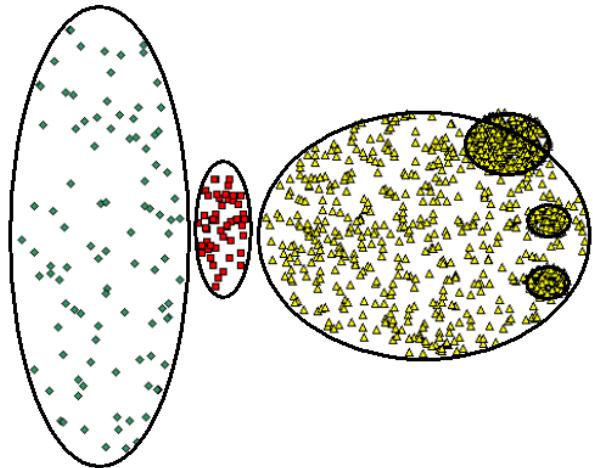
```
1 Core  $\leftarrow \emptyset$ 
2 foreach  $\mathbf{x}_i \in \mathbf{D}$  do // Find the core points
3   Compute  $N_\epsilon(\mathbf{x}_i)$ 
4    $id(\mathbf{x}_i) \leftarrow \emptyset$  // cluster id for  $\mathbf{x}_i$ 
5   if  $N_\epsilon(\mathbf{x}_i) \geq minpts$  then  $Core \leftarrow Core \cup \{\mathbf{x}_i\}$ 
6    $k \leftarrow 0$  // cluster id
7   foreach  $\mathbf{x}_i \in Core$ , such that  $id(\mathbf{x}_i) = \emptyset$  do
8      $k \leftarrow k + 1$ 
9      $id(\mathbf{x}_i) \leftarrow k$  // assign  $\mathbf{x}_i$  to cluster id  $k$ 
10    DENSITYCONNECTED ( $\mathbf{x}_i, k$ )
11   $\mathcal{C} \leftarrow \{C_i\}_{i=1}^k$ , where  $C_i \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = i\}$ 
12   $Noise \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = \emptyset\}$ 
13   $Border \leftarrow \mathbf{D} \setminus (Core \cup Noise)$ 
14 return  $\mathcal{C}, Core, Border, Noise$ 
```

DENSITYCONNECTED (\mathbf{x}, k):

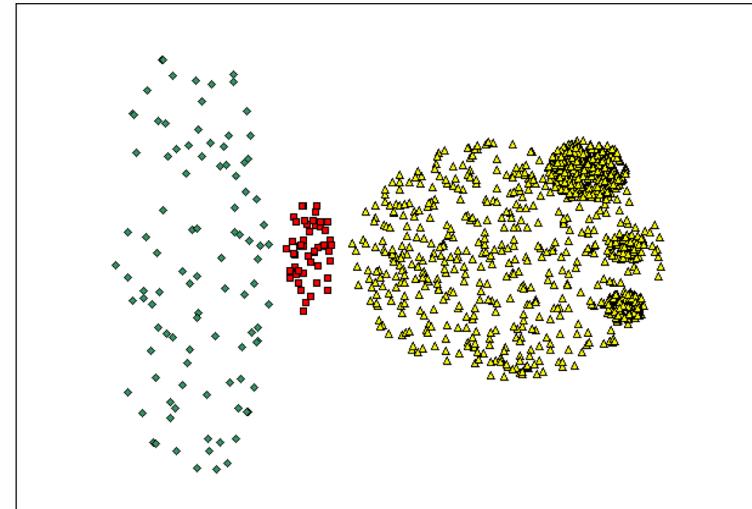
```
15 foreach  $\mathbf{y} \in N_\epsilon(\mathbf{x})$  do
16    $id(\mathbf{y}) \leftarrow k$  // assign  $\mathbf{y}$  to cluster id  $k$ 
17   if  $\mathbf{y} \in Core$  then DENSITYCONNECTED ( $\mathbf{y}, k$ )
```



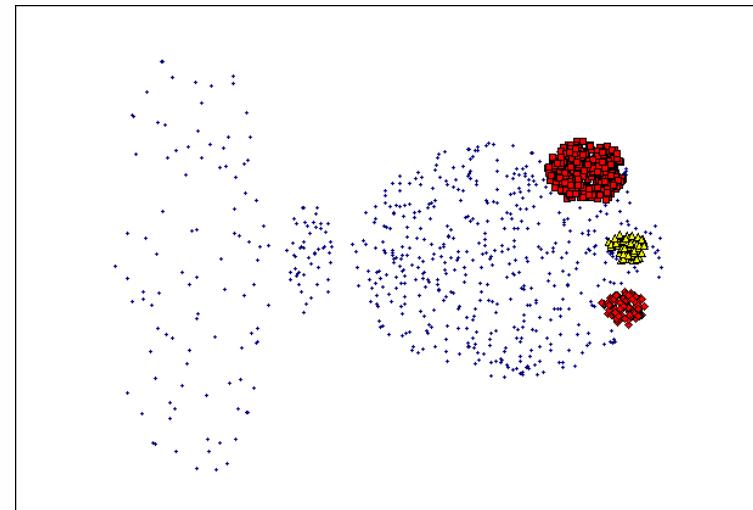
- Varying densities
- High-dimensional data



Original Points



$(\text{MinPts}=4, \text{Eps}=9.75)$.



$(\text{MinPts}=4, \text{Eps}=9.92)$

Examples using R

```
library(fpc)

set.seed(665544)
n <- 600
x <- cbind(runif(10, 0, 10)+rnorm(n, sd=0.2), runif(10, 0,
10)+rnorm(n, sd=0.2))

par(bg="grey40")
ds <- dbSCAN(x, 0.2, showplot=1)
```

Density-Based Clustering in R

```
library(fpc)

set.seed(665544)

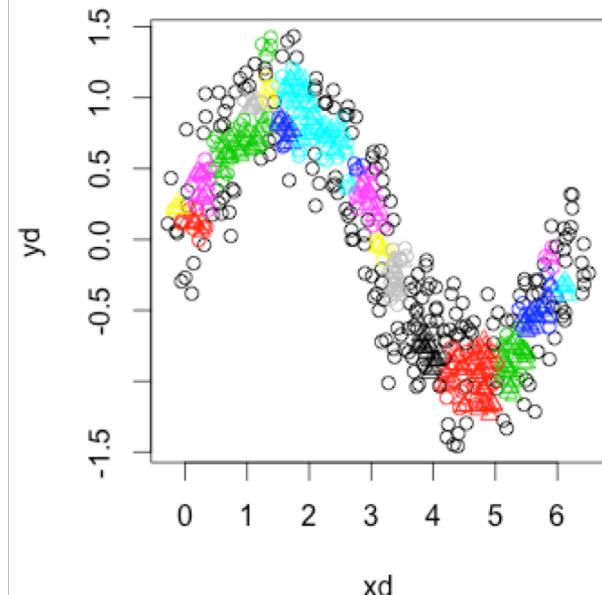
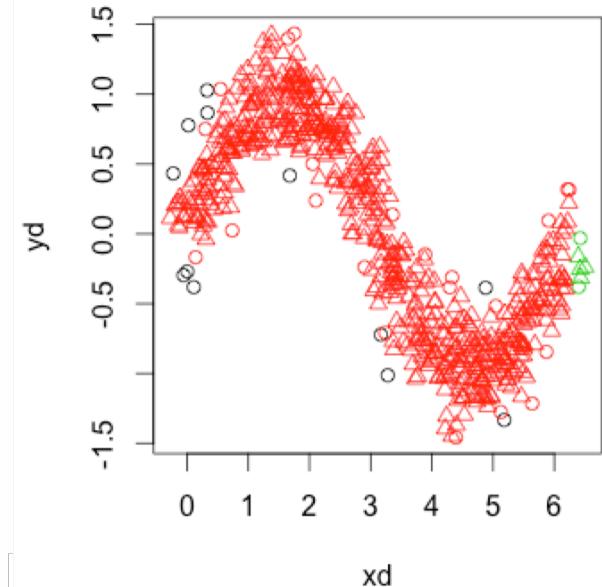
x <- seq(0, 6.28, 0.1)
y <- sin(x)

xd <- x+rnorm(630, sd=0.2)
yd <- y+rnorm(630, sd=0.2)
plot(xd, yd)

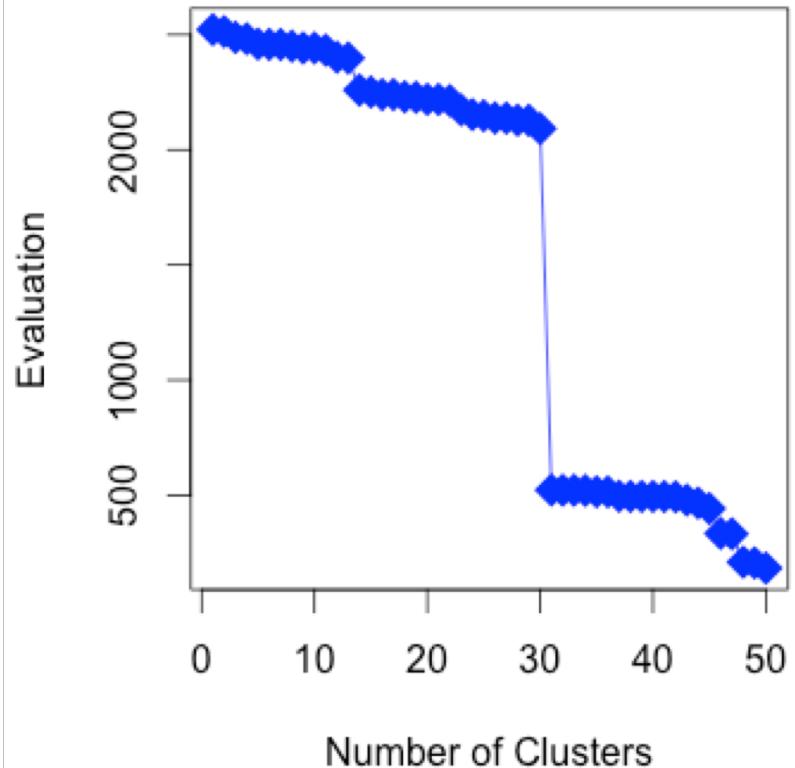
par(bg="grey40")
d <- cbind(xd, yd)

# this works nicely since the epsilon is
# the same size of the standard deviation (0.2)
# used to generate the data
ds <- dbSCAN(d, 0.2, showplot=1)

# this does not work so nicely
ds <- dbSCAN(d, 0.1, showplot=1)
```

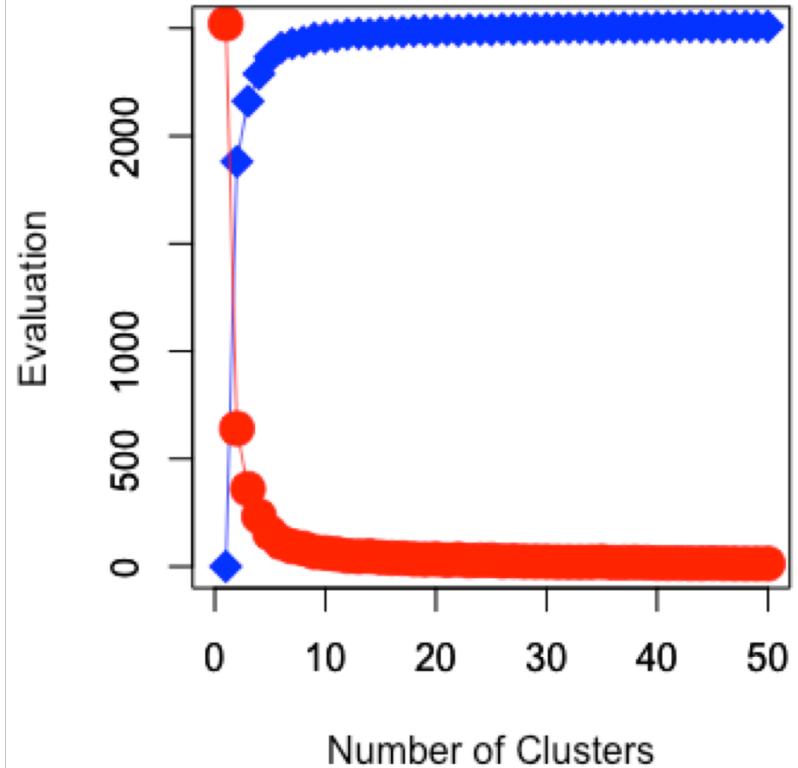


Between Cluster Sum-of-square



hierarchical clustering

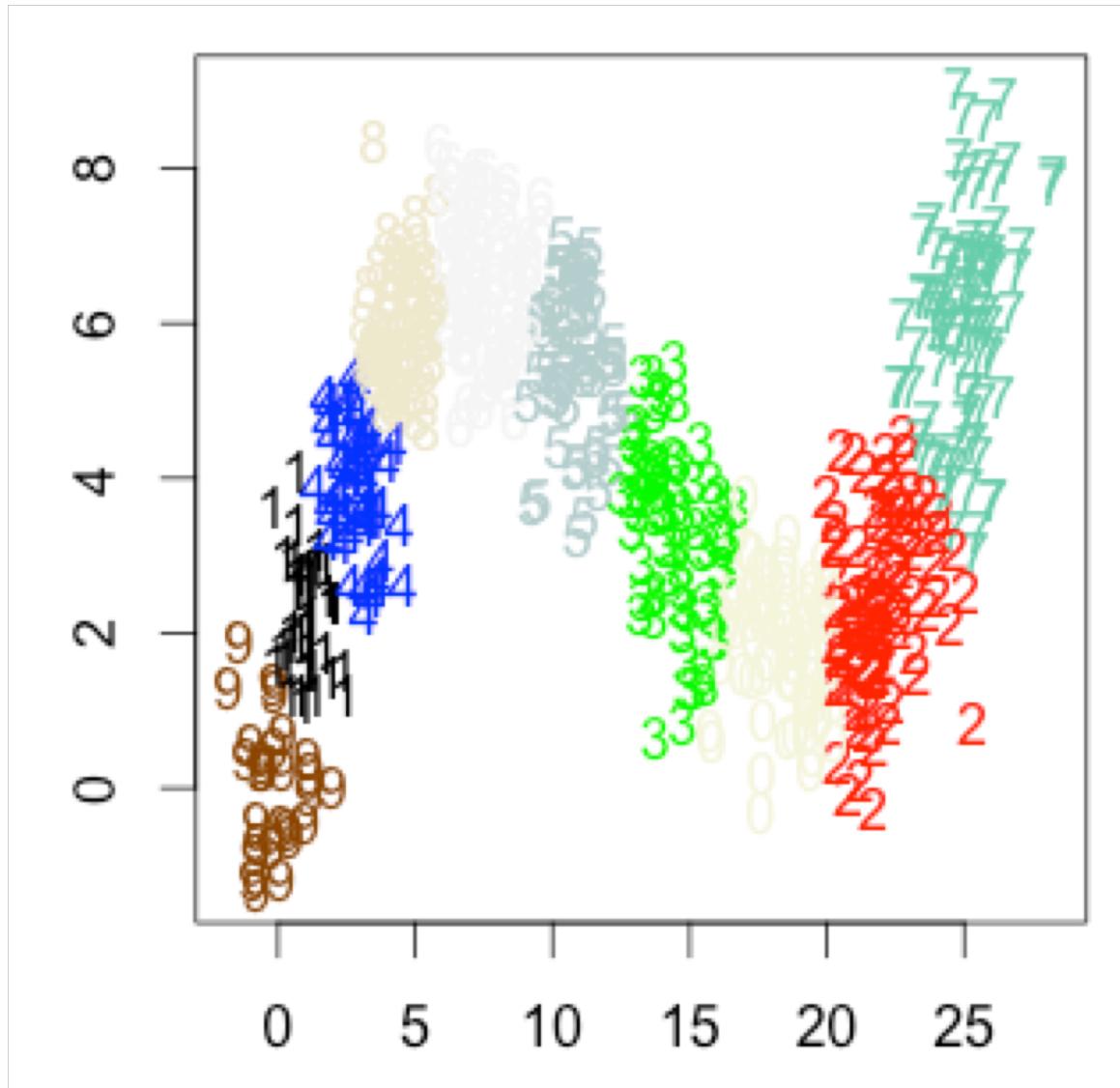
Within/Between Cluster Sum-of-square



kmeans clustering

Clustering Comparisons on Sin Data (k-means with 10 clusters)

17



Run the python notebooks and the
KNIME workflows for this lecture