

As a graph search:

- verts are core points
- edges when 2 core points are w/in ϵ of each other

Find connected components

- each component is a cluster

Finally - add in border pts w/in ϵ

Computational complexity

High level

1. $\forall \vec{x} \in D$

- compute $N_\epsilon(\vec{x})$
- check if \vec{x} is a core point
- init \vec{x} as not in a cluster

2. \forall core point

find all densely connected points

for each part

time proportional to

$O(n)$ in high dim
 $O(1)$
 $O(1)$

$O(\log n)$ in low dim
hidden d's

$O(n)$ not going to dominate

Total time:

- $O(n^2)$ in high dim
- or $O(n \log n)$ in low dim

Things to consider

- user supplied params minpts, ϵ
- only 1 ϵ value is allowed (can't handle varying densities)

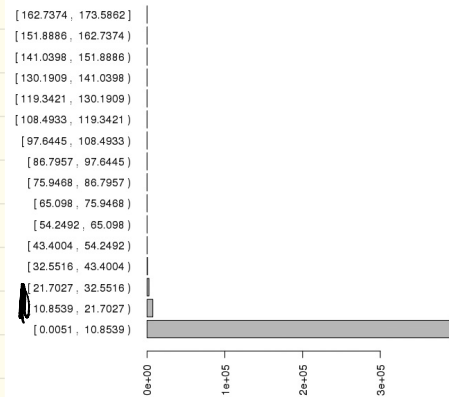
Soln 1: OPTICS - improvement of DB Scan



Soln 2: in practice

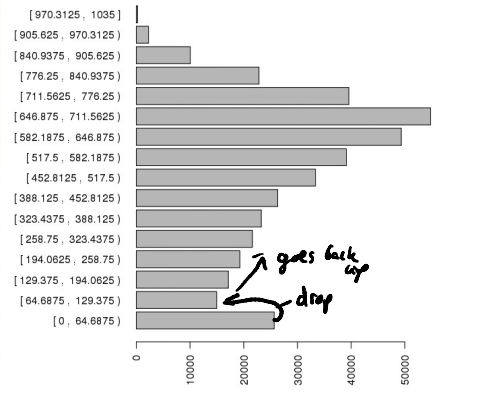
https://github.com/alitouka/spark_dbscan/wiki/Choosing-parameters-of-DBSCAN-algorithm

dist to nearest neighbor



pick $\epsilon=22$

fix $\epsilon=22$ look at # of core pts for varying minpts



choose min pts = 129