

Clustering

- Top Down, K-means
- Bottom Up, Hierarchical

Distance

- Euclidean

$$\sqrt{\Delta x_1^2 + \dots + \Delta x_n^2}$$

Minkowski $(\Delta x_1^m + \dots + \Delta x_n^m)^{1/m}$

$L_{m=1}$, Manhattan distance

- Differences in scale across attributes
 - ↳ large magnitudes drown out smaller ones

Scale all attributes to a common range

- Non-numeric attributes
 - Categorical values (e.g. colour)
 - Binary values

Mixed attributes

- Some numeric, some not

Problem specific

Assume all attributes are non-numeric

- Simplest distance - fraction of agreement

x_1	x_2	x_3	x_4	0.75
1	1	1	1	Similar
y_1	y_2	y_3	y_4	on a 0-1 scale

Distance = 1 - Similarity

Example

DNA sequences

Special case - binary

	0	1
0	a	b
1	c	d

Each $x_i | y_i$ is in one of four categories

Similarity

$$\frac{\#a + \#d}{\#a + \#b + \#c + \#d}$$

Asymmetric case

Boolean document model

Vocabulary $V = \{w_1, \dots, w_N\}$

Document d is a vector $\{0, 1\}^N$

$d_i = 1$ iff w_i occurs in d

Compare documents d & d' Words absent in d & d'

$$\text{Similarity} = \frac{\#a + \#d - \text{words in } d \text{ & } d'}{\#a + \#L + \#c + \#d} \quad N$$

Most words are absent in most documents

$$N \approx 10^5$$

$$\text{Doc length} \approx 2500 \text{ (10 pages)}$$

At least 97,500 of 100,000 words do not appear in d

$$\#a \gg \#d$$

Similarly ≈ 1 for any pair d, d'

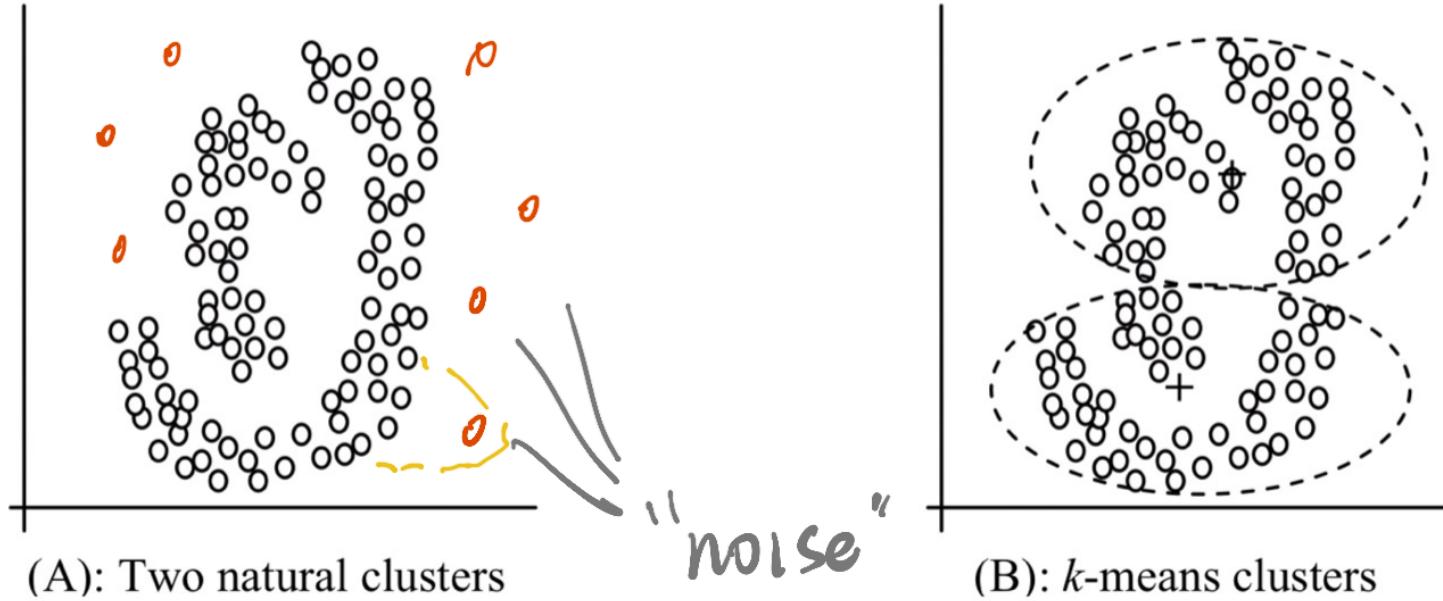
In such cases, discard #a

Focus on words that occur in $d \cup d'$

$$\frac{\#d}{\#b + \#c + \#d}$$

Jaccard Distance

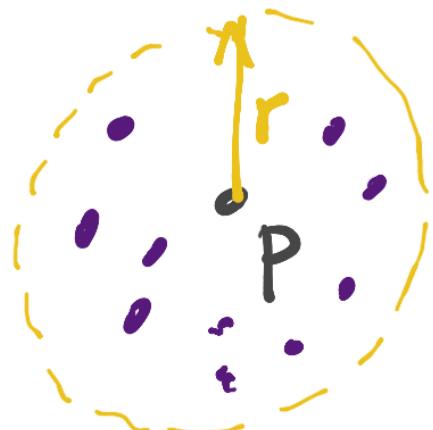
Finding strangely shaped clusters



Use "density" to define a cluster

Formalize

Fix a radius & count the points

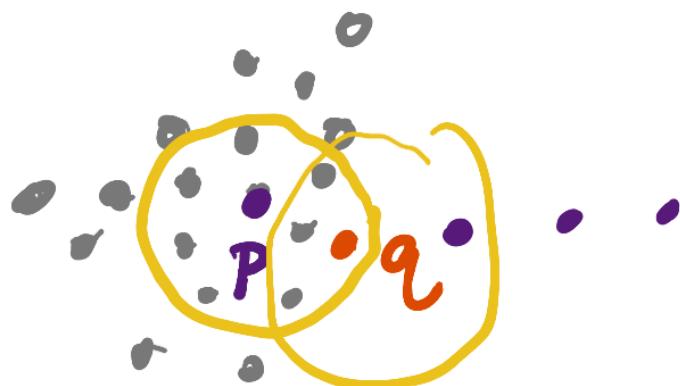


fixed radius

Fix p , $\text{nbdr}_r(p)$

$q \in \text{nbdr}_r(p)$

added to p 's
cluster



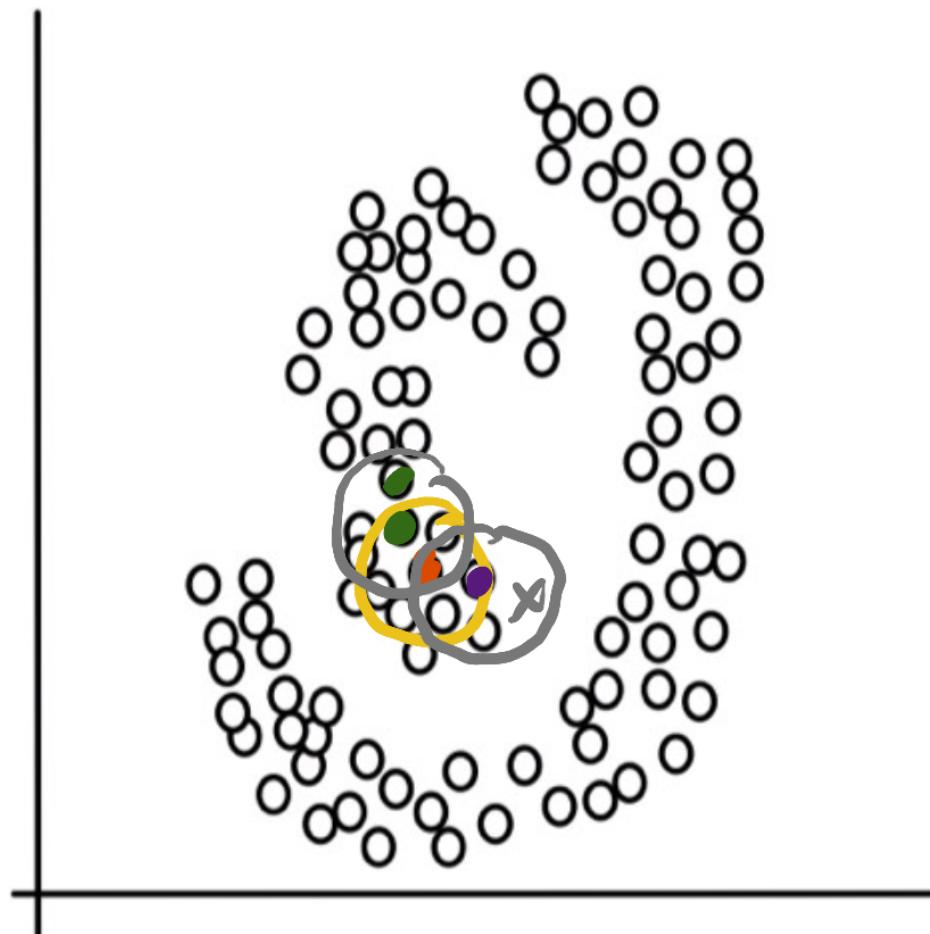
Should we continue the process with q ?

Impose a density threshold m (minimum)

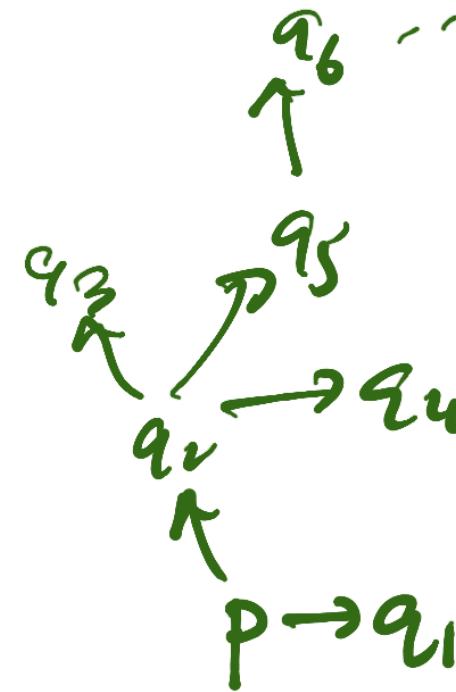
Explore $\text{nbdr}_r(p)$ only if $\text{nbdr}_r(p)$ has at least m points

Core point p has $\geq m$ neighbours in $\text{nbdr}_r(p)$

$p \rightarrow q$ if p is a core point & $q \in \text{nbdr}_r(p)$



(A): Two natural clusters



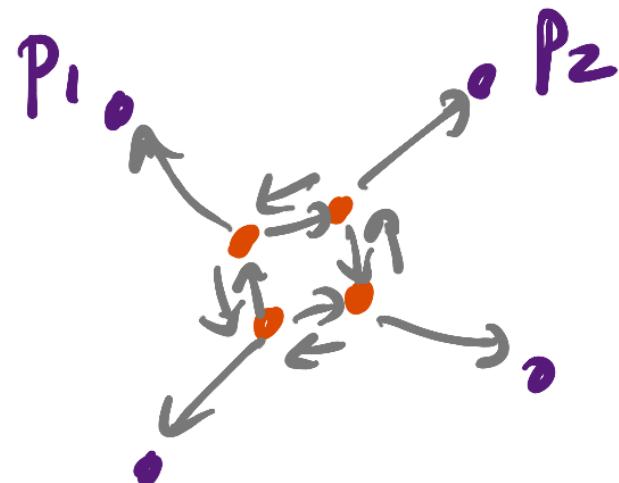
$p \rightarrow q$
is not
symmetric

Defines a directed graph on the points

p & p' are both core

and $p \rightarrow p' \Rightarrow p' \rightarrow p$

Connected component form clusters



Directed
reachability
 \rightarrow Undirected
Connectivity

Strategy

- Pick a point p in data set
- Explore all reachable points if p is a core point
 - Identifies cluster containing p

Suppose p & p' are core points in same cluster

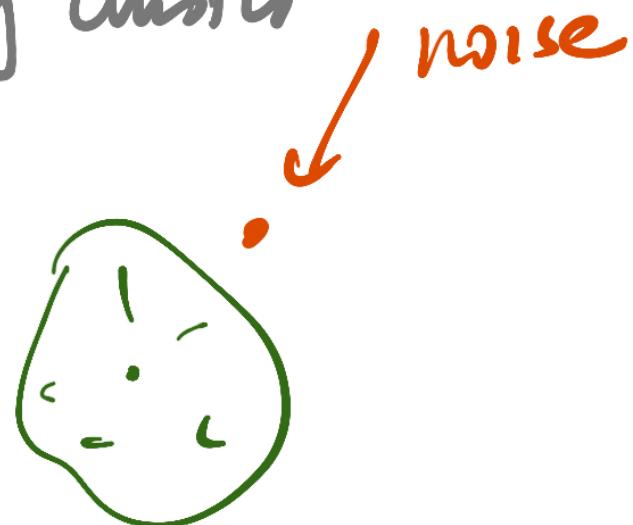


For all core nodes p

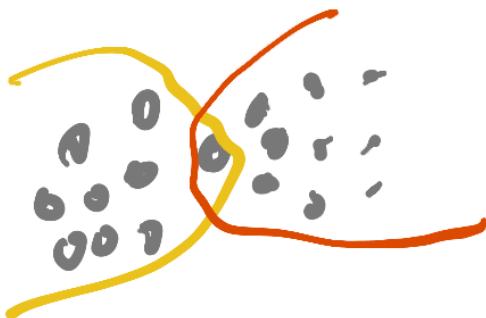
if $\text{cluster}(p)$ not defined

explore from p & label $\text{cluster}(p)$

Some non core points may not be
in any cluster



Two natural clusters that overlap at boundary



Not clear how to deal with overlapping boundary points

Preprocessing + efficient data structure
to extract underlying graph from
raw data

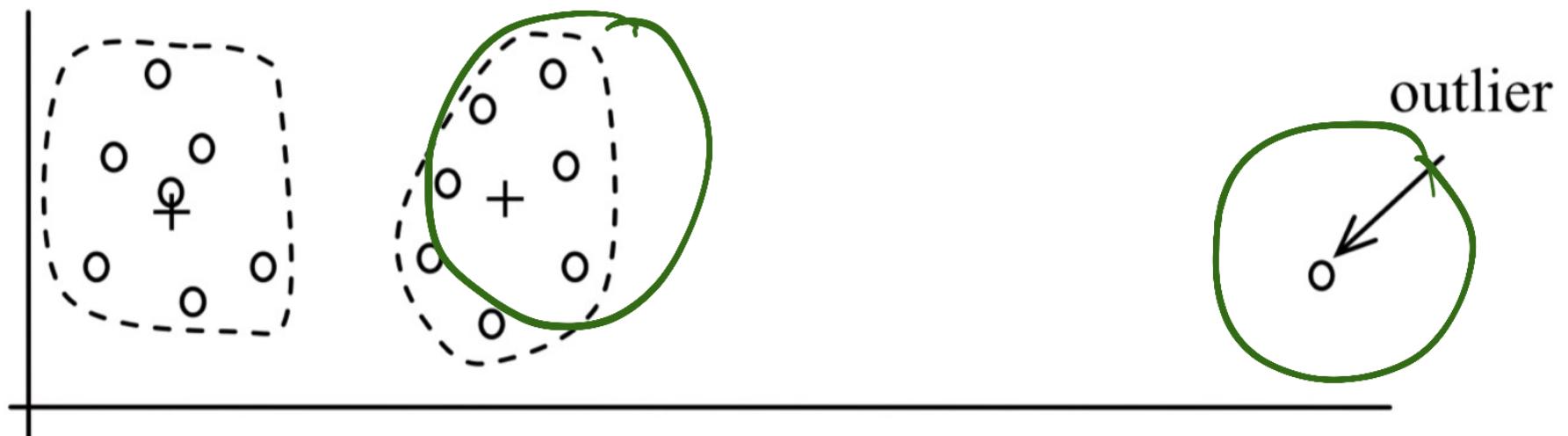
Outliers

Points "far away" from others

Simple numeric value (single),
define in terms of distance from
mean, etc

Multidimensional case

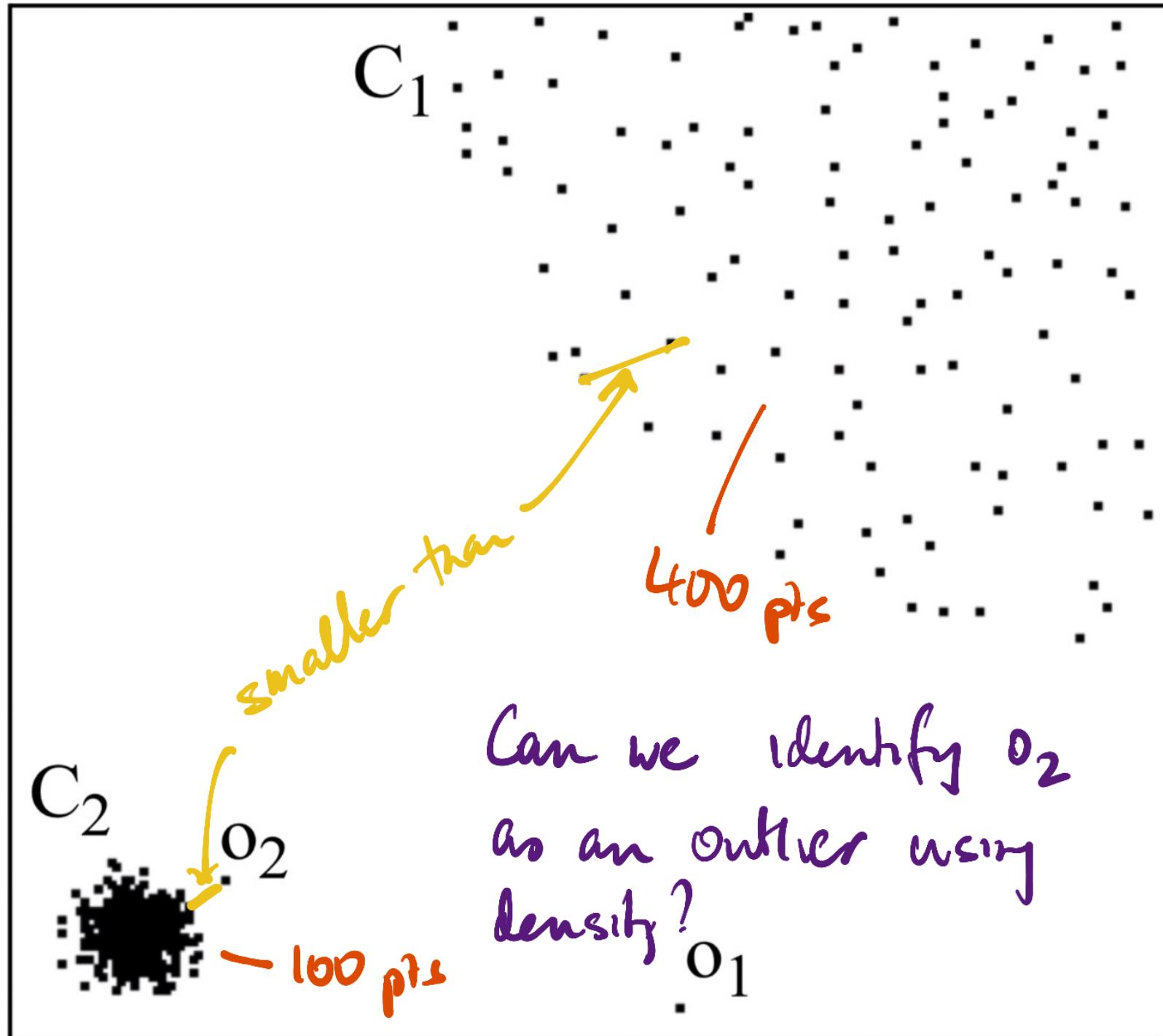
Use density



(B): Ideal clusters

An outlier will have a very sparse neighbourhood

- Fix radius r , threshold m
- Problem: r & m are global!

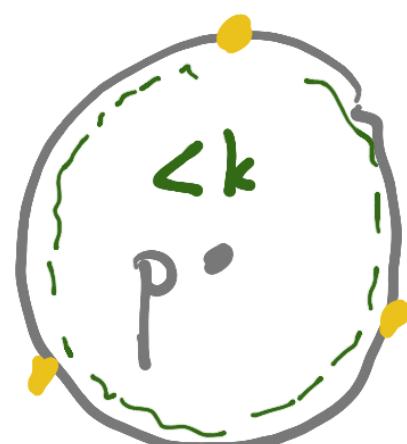


No good uniform values of $m \times r$ to classify O_2 accurately

Local notion of density

- Earler: fix r , count points in $\text{nbdr}(p)$
- Instead: fix k , find r s.t. $\text{nbdr}_r(p)$ has k points

Smallest r s.t. $\text{nbdr}(p)$ has $\geq k$ points, for any $r' < r$, $\text{nbdr}_{r'}(p) \leq k-1$ pts



Given P , $\text{local-radius}_k(P) = lr_k(P)$

For all q in $lr_k(P) \rightarrow lr_k(q)$

Take average: $\frac{1}{\text{size}(k\text{nsd}(P))} \sum lr_k(q)$

$\frac{lr_k(P)}{\text{Arg } lr_k(q) \text{ of neighbours}} = \begin{array}{l} \text{Local Outlier} \\ \text{Factor of} \\ P \end{array}$

Fix a threshold δ (empirically)

Define p to be an outlier if $\text{LOF}(p) \geq \delta$

