

①

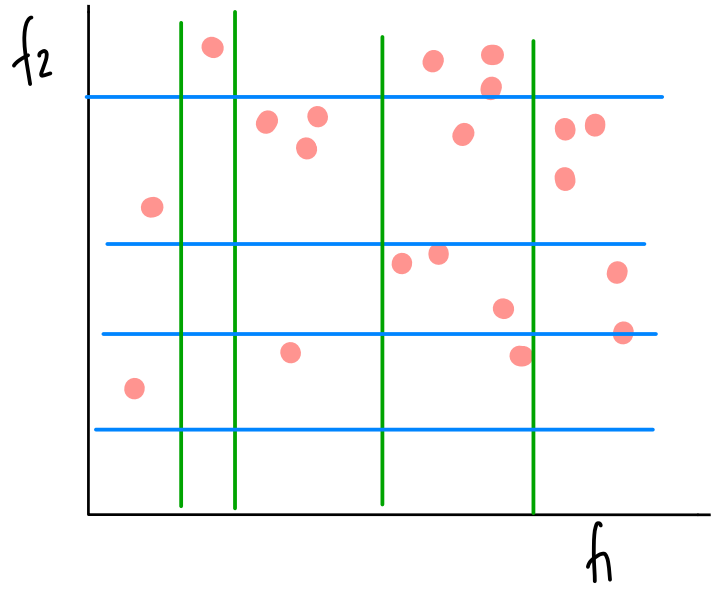
Isolation Forest \Rightarrow Based on DT

②

LOF = Local outlier Factor

Isolation Forest

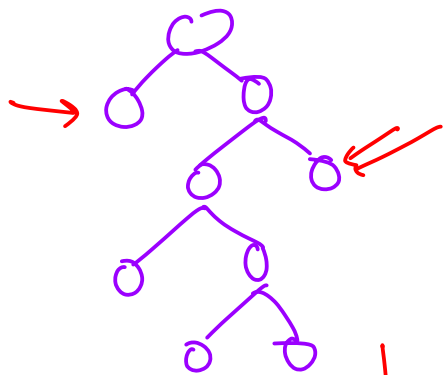
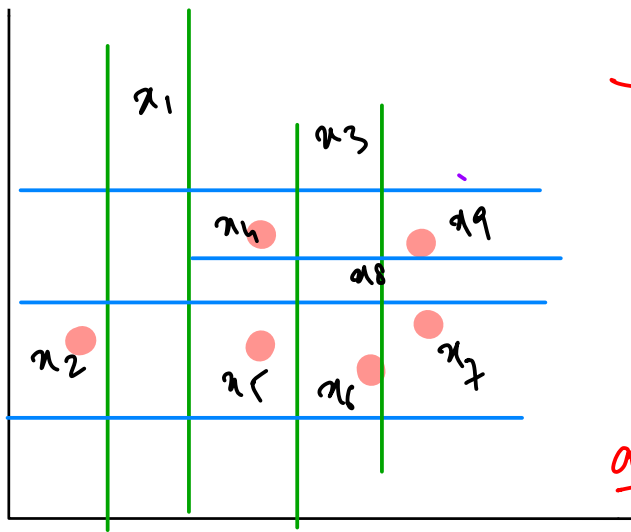
↳ Collection of trees
↳ Isolate points



Create multiple decision trees on f 's with thresholds and isolate point

for 2.

- ① Randomly pick a feature
- ② Randomly threshold the feature
- ③ Build the decision tree until each data point is a node itself.



100's ←
DT

avg(d) x1 [2, 2, 1, 3, 2, 2, 3, ...]
1.8 x2 [2, 1, 1, 1, 2, 3, 4, 2, 1]
2.9 x3 [3, 4, 3, 4, 5, 6, 3]
5 x4
2.9 x5
3.2 x6
x7
x8
x9

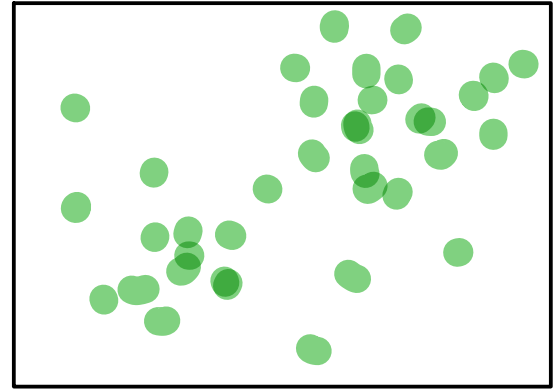
We use avg depth as a metric.

sort point based on avg (d)

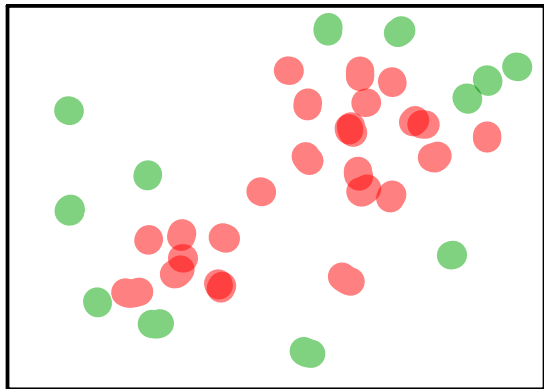
$$\overset{2\%}{x_1} < x_2 < x_3, \dots, \dots$$

100

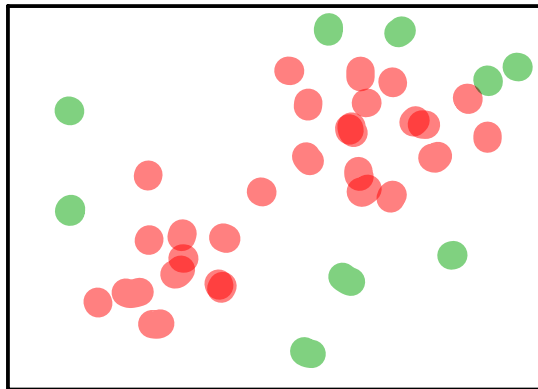
Contamination: 2%



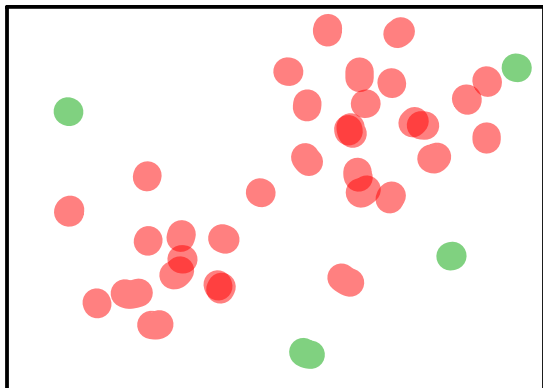
Contamination: 20



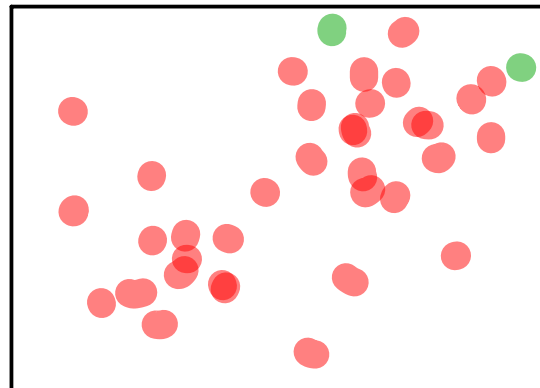
Contamination = 10



contamination: 5

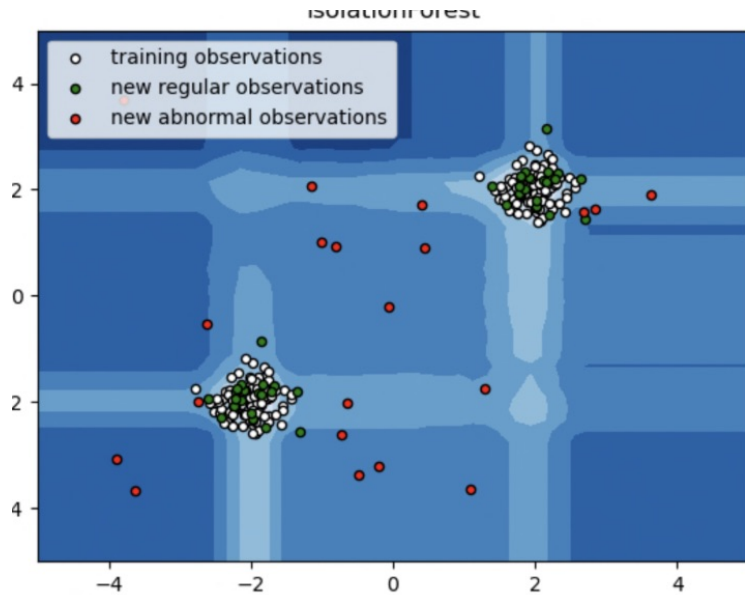


contamination = 1



Disadvantages

- * Slow if the size of dataset $\uparrow\uparrow\uparrow$ huge.
- * Isolation Forests are biased towards axis parallel split.
- * In the diagram below, shades of color represent likelihood of a data point being inlier/outlier.



Advantage.

* Effective.

* Find out outliers b/w 2

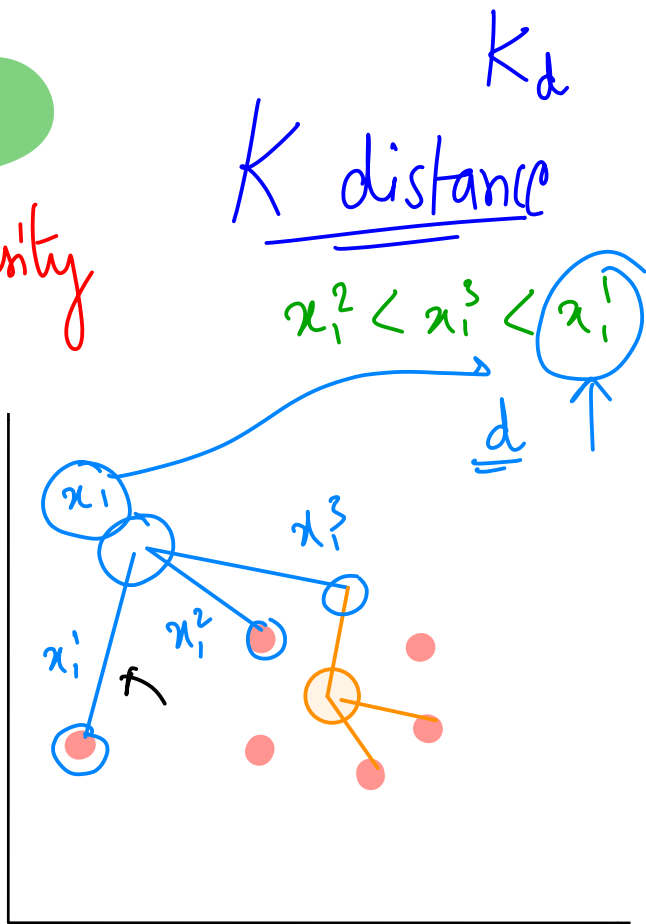
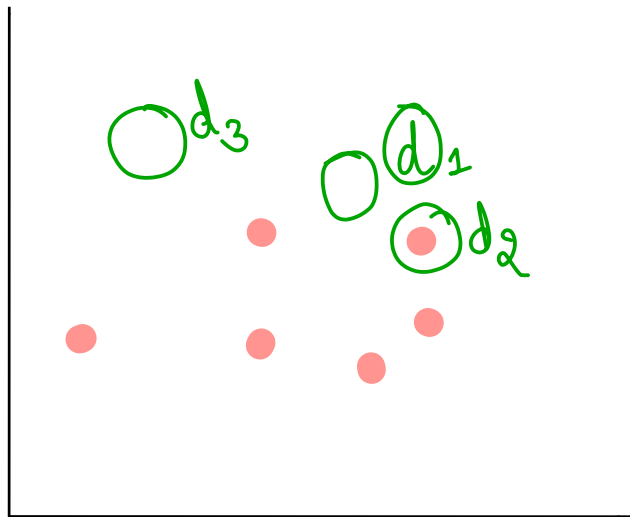
Cluster.

* Use probability / shade of color
in determining effectively
how sure are we in saying that
a datapoint is outlier / inlier.

LOF

Local Outlier Factor

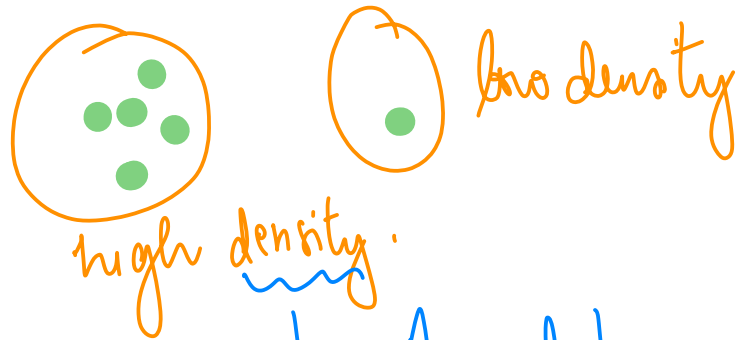
↳ KNN + Density



$K_d \Rightarrow$ We will find the distance of farthest point among K Nearest Neighbour.

Idea:

Outliers exist in low density region



\hookrightarrow Avg distance b/w point.

$$\text{Density} \propto \frac{1}{\text{avg}(d)} \Rightarrow \text{Density} = \frac{1}{\text{avg}(d)}$$

① Find k nearest neighbours. of point (A) $k=3$.

② Calculate avg distance. μ .

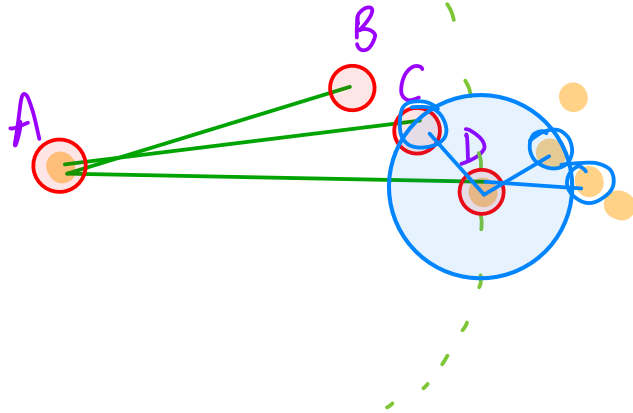
$$\text{Density}(A) = \frac{1}{d_1 + d_2 + d_3}$$

Variation → Instead of averages, we can also take maximum out of d_1, d_2 & d_3 → K_{dist} [max value]

$$\text{density}_i = \frac{1}{K_d(A)}$$

$\max(k \text{ nearest neighbours}) = K_d = \text{distance from sorted } K^{\text{th}} \text{ neighbour.}$

k=3



low density region

A has a large radius

D has lower valued radius
high density region

LOF

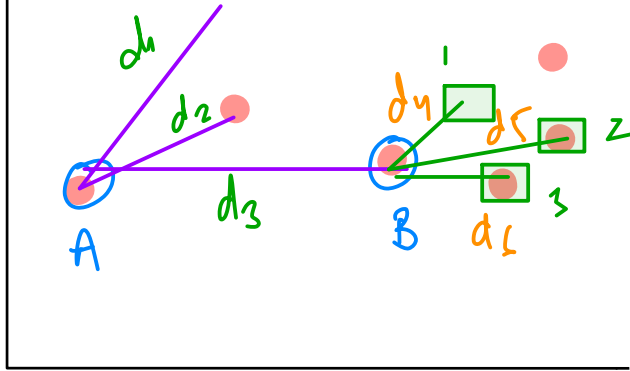
$$\text{Local Outlier Factor} = \frac{\text{Avg distance of } k \text{ neighbours}}{\text{Density of } A}$$

Local Reachability Density

reachability distance $(A, B) = rd(A, B)$

$$rd(A, B) = \max \{ \text{dist}(A, B), K \cdot \text{dist}(B) \}$$

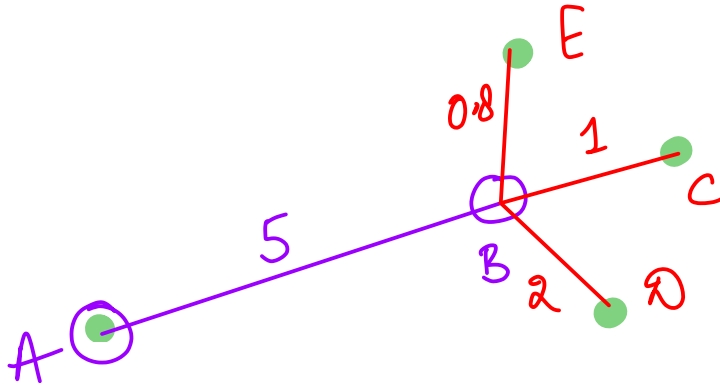
$$\text{Density}(A) = \frac{1}{\frac{1}{K} \sum RD(A, \text{neighbor})}$$



$$rd(A, B) = \max \{ d_3, d_5 \}$$

R_d (me, you)

if you are in my 'top k' but I am not in
your 'top k' then let's use actual distance
b/w us, but if we both are in 'top k' of
each other then let's use your k_{dist} .



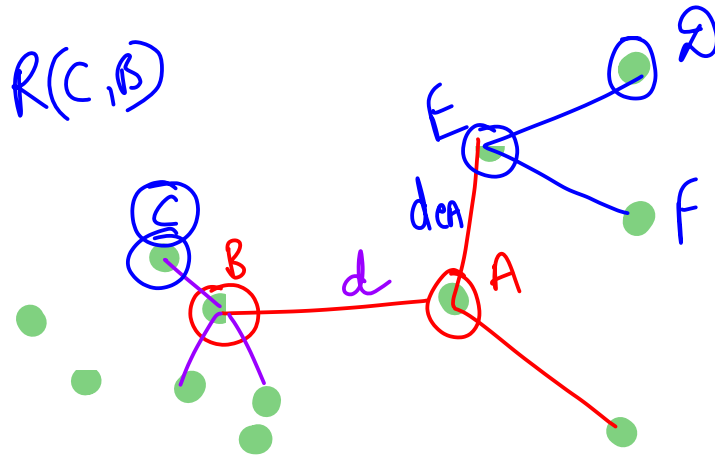
$$R_d(A, B) = 5$$

$$R_d(E, B) = 2$$

$$R_d(C, B) = 2$$



This is less intuitive bcoz some points



$$R(D, E) = denA$$

$$R(F, E) = denA$$

- ① Code for all
- ② Lox Interpretation + Adv + disadv. } next class .

2,00,000 ✓

Doubt

$$n = 160$$

$$y_{+ve} = 60$$

$$y_{-ve} = 40$$

$k = 100$ Impact of outlier ↓↓

$k = 1$ Impact of outlier ↑↑↑

