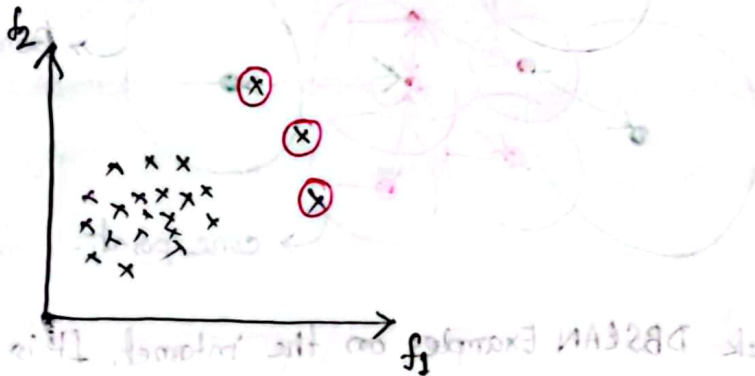


Anomaly detection Isolation Forest: (To detect outliers)

Suppose, we have a disease dataset, where the normal points shows the healthy people and the outliers shows the people who are in diseases.

So, not always we have to remove outliers.

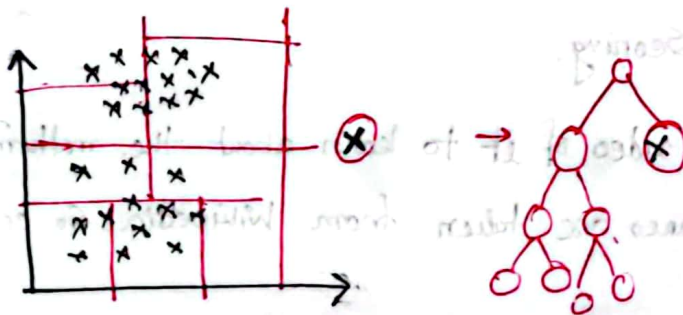
In this type of case, to detect outliers, we use Anomaly Detection algorithm.



There are two types of Anomaly Detection:

① Isolation Forest (Decision Tree)

We will run decision Tree on data points and will find the leaf nodes



The data points who are isolated, their leaf node can be found easily.

Like in our case we can find it in depth 2.

So, the isolated node, how quickly it is found, that becomes the ^{high} probability to ~~not~~ be an outlier. For that, we will calculate **Anomaly score**.

Because with many features, many decision trees will be created and many isolated nodes will be found. To measure the probability of a leaf node to become an outlier, we will use Anomaly score.

Mathematical score of anomaly score:

Anomaly score for a new score: $m = \text{no. of data points}$

$$S(x, m) = 2 \frac{-E(h(x))}{c(m)}$$

$E(h(x)) = \text{Average search depth}$

of x from the isolated tree

$$E(h(x)) \ll c(m) = S(x, m) \approx 1 \Rightarrow \text{Anomaly score} \rightarrow \text{outlier}$$

We will define the threshold, suppose 0.5

$c(m) = \text{Average depth of all the datapoints.}$

> 0.5 will be outliers.

$$E(h(x)) \gg c(m) = S(x, m) \approx 0.5 \Rightarrow \text{Normal data}$$

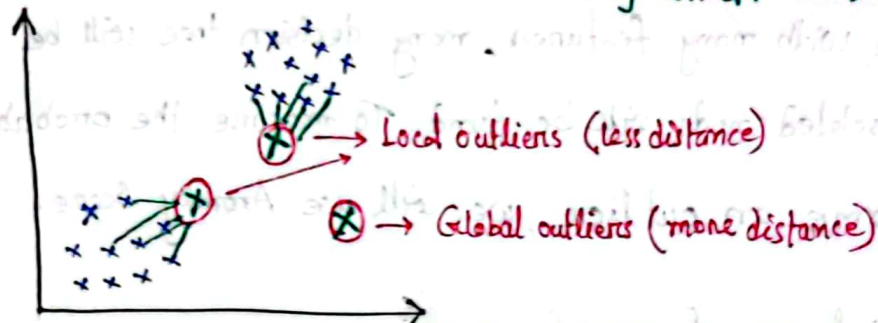
DBSCAN Clustering For Anomaly Detection:

(Already the basics of it is discussed)

It can also be used to detect outliers.

Local outlier factor anomaly detection (LOF)

Key concept: Measuring local density



It follows the technique of K Nearest Neighbour. It calculates the distance of point x with its K Nearest Neighbour (Suppose $K=5$).

So it calculates the distance of x with its 5 nearest neighbour and then calculates each neighbour's distance with its other 5 nearest neighbours. Then it compares the density of itself and the density of the neighbours with their neighbours. If the density of x is less than its neighbours then it is an outlier.

Distance of point x with its 5 nearest neighbour

Distance of each neighbour with its 5 nearest neighbour