→ Class start at 9:05 pm

# Agenda

→ One - class SVM

→ Isolation Forest

→ LOF (Local Outlier factor)
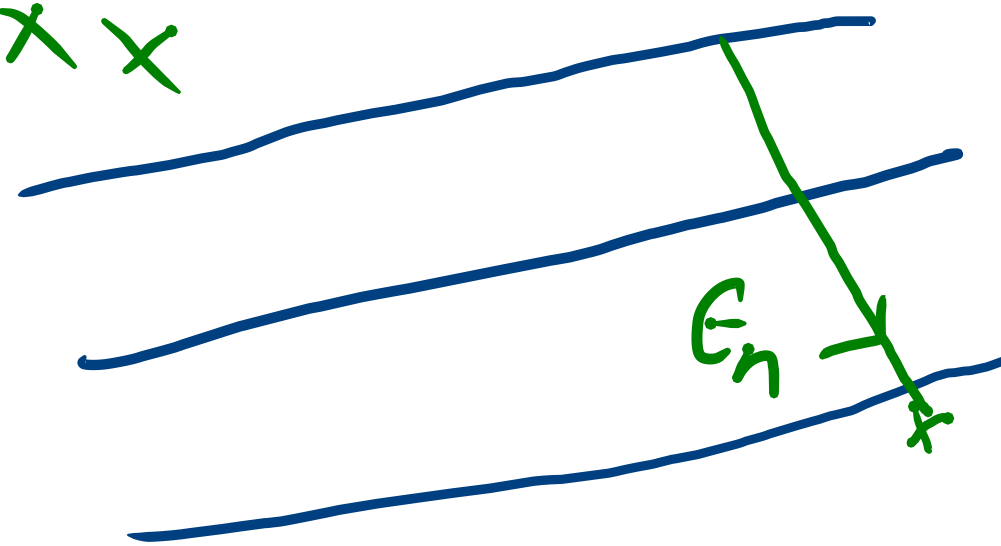
# SVM

$$w^T x + b = 1$$

$$w^T x + b = 0$$

$$w^T x + b = -1$$

$\varepsilon_{i} = 0$

x x
x x
x x

$\varepsilon_n$

$$\frac{1}{2} w^T w + C \sum_{i=1}^{N} \varepsilon_n$$

Subject to $y_n(w^T_n + b) \geq 1 - \varepsilon_n$

$$\varepsilon_n \geq 0$$

# Oneclass SVM

Minimize the hypersphere



$\rightarrow$ Find the centre $\mathcal{L}$ radius which contains most of the data point.

$$\min_{c, r} \; r^2$$

$\epsilon_i = 0 \quad$ if $\text{dist}(x_i, c) \le r$

$\quad = d(x_i, c) - r \quad \text{dist}(x_i, c) > r$

$$\min_{c, r, \epsilon_i} \; r^2 + \lambda \sum_{i=1}^{n} \epsilon_i$$

$$\text{s.t} \quad \text{dist}(x_i, c) < r^2 + \epsilon_i^2 \quad \forall i$$

$$\epsilon_i > 0$$

# Disadvantage

1. As the no of data point $\uparrow$, time complexity $\uparrow$
2. Kernel selection
3. SVM complexity

Original data $\rightarrow$ 2D

RBF_Kernal

Hypersphere in in high dimension Space
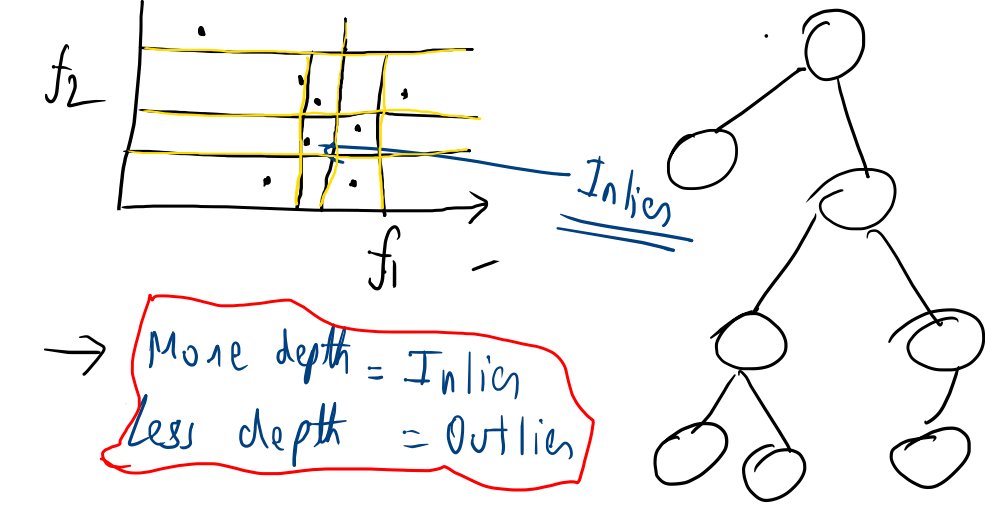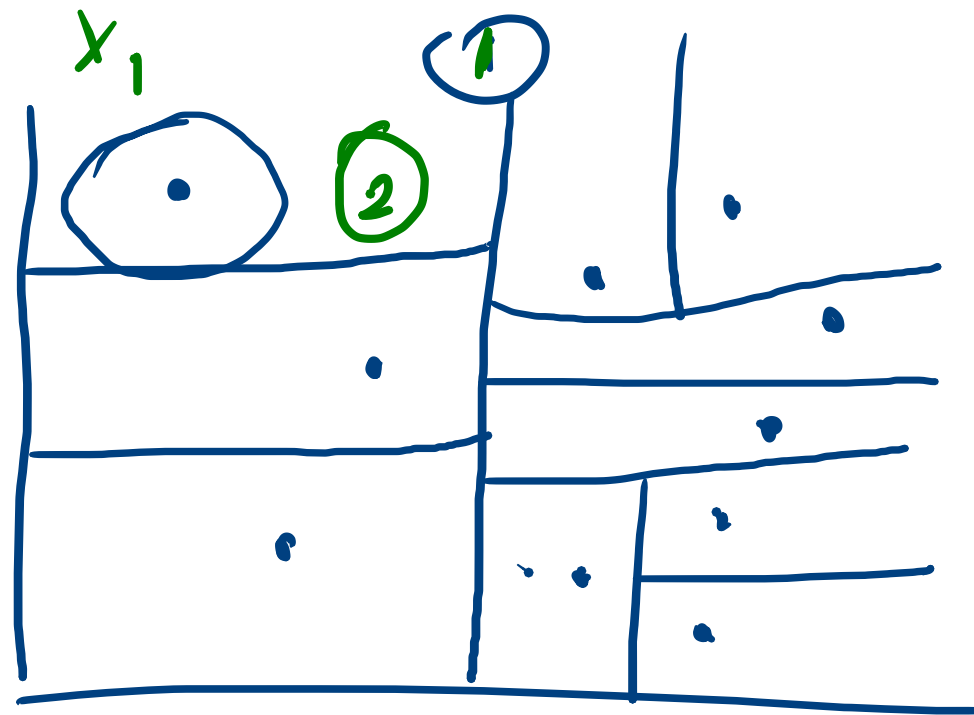
# Isolation Forest

Randomly sample datapoints?

1M ↳ 1 lakh sample

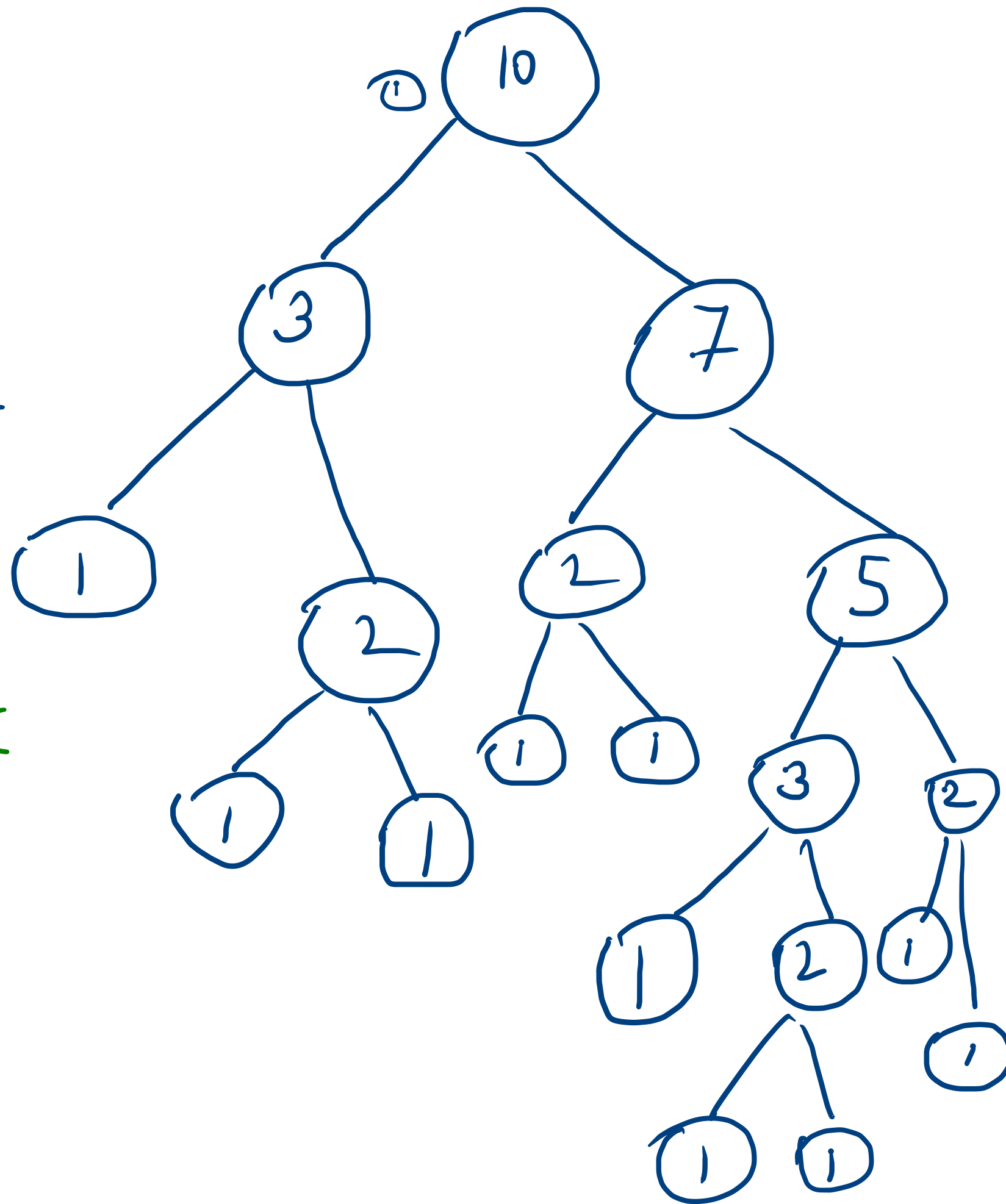→ Building many trees ( like Random Forest)

{ → Each tree → randomly pick a feature

↳ randomly threshold it

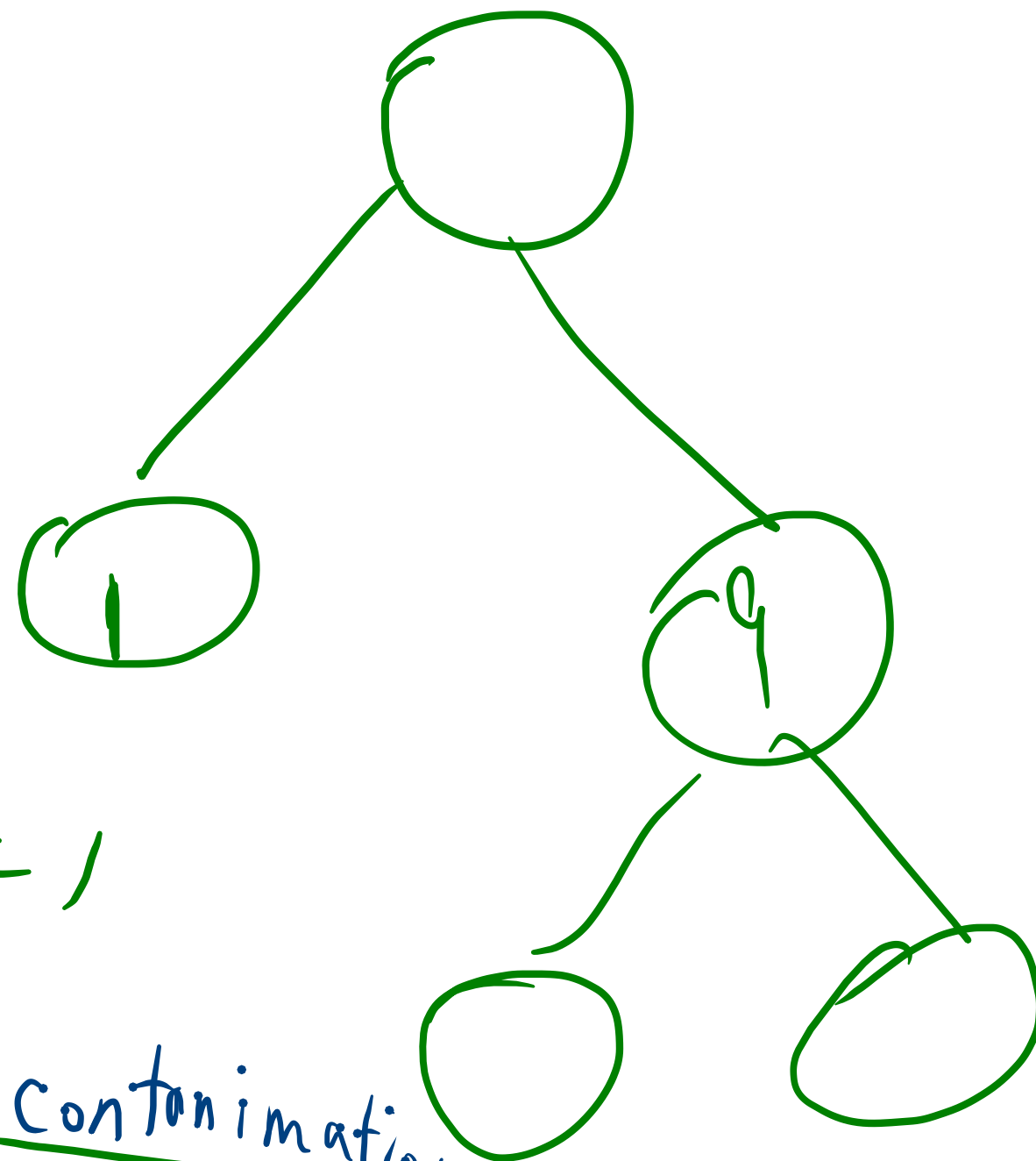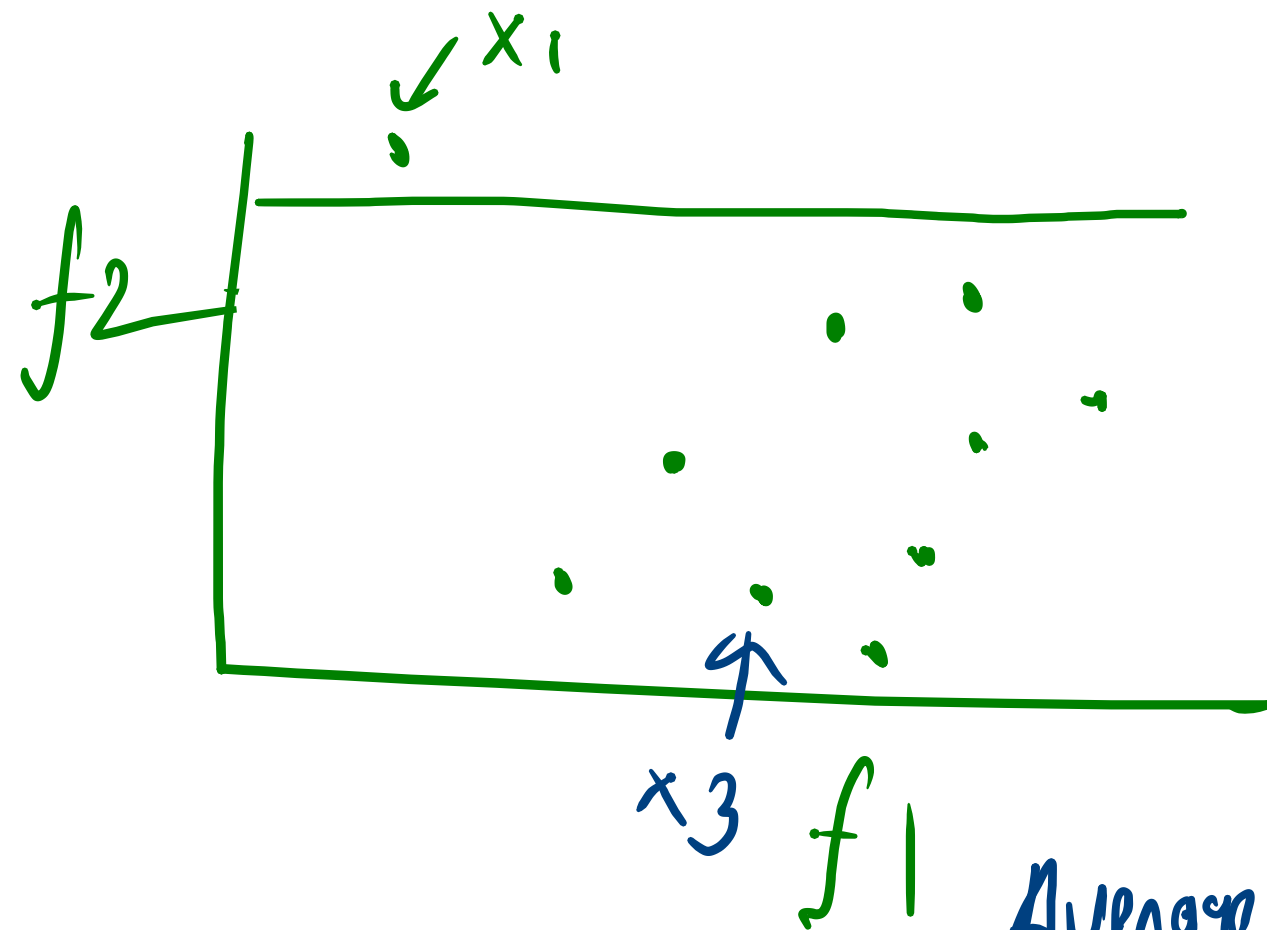→ Build each tree till the leaf node contain one datapoint

$f_2$

Inlier

$f_1$

More depth = Inlier
Less depth = Outlier

$f_2$    $x_1$

$x_3$   $f_1$

$x_1$    2, 1

$x_2$

$x_3$    6   7

     2, 3

$x_7$   5, 8

$x_{10}$

Average depth $x_1$

1.5

Depth = 1

6.5 —

2.5

6.5

contamination

10% of the data

$x_1$, $x_2$

→ LOF

→ Break until 10:10 pm

DBSCAN : Noise point ( Neither core point
                       nor boundary point)

Elliptic Envelope : Probability is very low
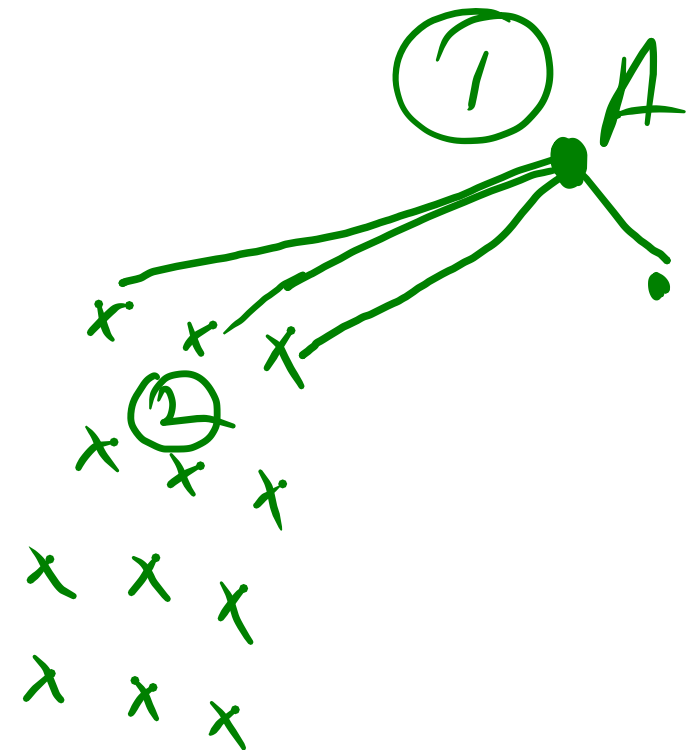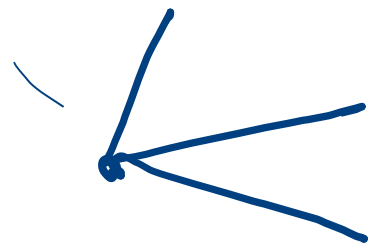            as points are far from centre $\mu_x, \mu_y$

Isolation Forest : Outlier are at less tree
            depth

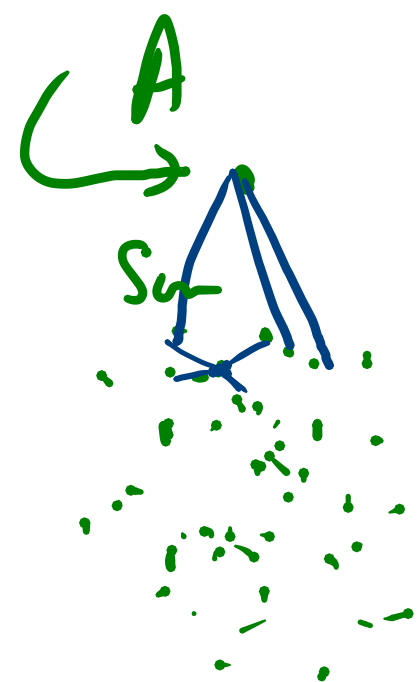LOF    : Outlier density is lower w.r.t
density of it neighbors

LOF

$\longrightarrow$ Concept of Density & K-NN



① A

Density $\alpha \frac{1}{dist}$

②

$$LOF = \frac{Avg\ density\ of\ Neighbours\ of\ A}{Density\ of\ A}$$

$A$

Sur

Sum $B$

Which one is more probable to be outlier $A$ on $B$?

LOF $=$ $\dfrac{density \uparrow}{density \downarrow}$

density $\downarrow$

$= \uparrow$

$\dfrac{density \downarrow}{density \downarrow} \approx 1$

LOF $\gg 1 \longrightarrow$ outlier $\longrightarrow \underline{A}$

① k-distance

$K = 3 \rightarrow$ dist(A,D)

dist of point A to it $k^{th}$ nearest neighbor

② 

Reachability dist

Reach(A,B) = max( dist(A,B) , K-dist(B))
= dist(A,B)

K = 3



Reach(C,B) = max(dist(C,B), K-dist(B))
= K-dist(B)

Is Reach(A,B) == Reach(B,A) Symetric ??

$\hookrightarrow$ max( dist(A,B) / k-dist(B) )

max( dist(B,A) , k-dist(A) )
||
dist(A,B)

# Local Reachability density

$$Lnd(A) = \cfrac{1}{\cfrac{\displaystyle\sum_{B \in N_k} \text{Reach dist}(A,B)}{N_K}}$$

# Local oullier facton

$$\text{Lof}_A = \frac{\text{Avg density (lnd) of neigh of A}}{\text{Density (lnd) of A}} \checkmark$$

$$= \frac{\sum\limits_{B \in N_k} \text{lnd}_k(B)}{|N(A)| \, \text{lnd}_k(A)}$$

For oullier LoF

Ⓐ $= 1$     Ⓑ $< 1$     Ⓒ $\gg 1$

$1/4$  $1/5$ $x_2$  $1/10$

$x_1$

$1/5$  $1/5$

$1/6$  $1/4$

$1/7$

$$LOf_{x_1} = \frac{\frac{1}{3}\left(\frac{1}{5} + \frac{1}{5} + \frac{1}{4}\right)}{\frac{1}{10}}$$

$$\approx 1.5$$

$$LOf_{x_2} = \frac{\frac{1}{3}\left(\frac{1}{4} + \frac{1}{5} + \frac{1}{5}\right)}{\frac{1}{5}} = \frac{2.5}{3} \lesssim 1$$

# Disadvantage of LOF

→ Find optimal K (K-NN)

→ High dimension calculting K-NN is expensive

→ Dependent on contamination value

LOf

$X_1 \rightarrow 10$

$X_2 \rightarrow 5$

$X_3 \rightarrow 1$

$X_{10} \rightarrow$

$$LOF >> \quad 1$$

$\searrow$

Oulliers

Contam = 0.1

$\downarrow$

10%