

# UNSUPERVISED LEARNING

ANOMALY DETECTION  
CLUSTERING



NO LABELS (OR JUST A few)

## ANOMALY DETECTION

### Applications

↳ if IN REAL-TIME  
→ QUICK AND DIRTY  
methods needed

↳ SEMI-SUPERVISED  
LEARNING

- \* Credit card fraud
  - fraud
  - No fraud (- class)
- \* Quality Control
  - OK
  - NOT OK (- class)

- class >>> + class

Pot stress  
on a classifier!

⇒ classifiers not ideal!

↳ Accuracy not good KPI  
USE ROC, AUC, F1 ...

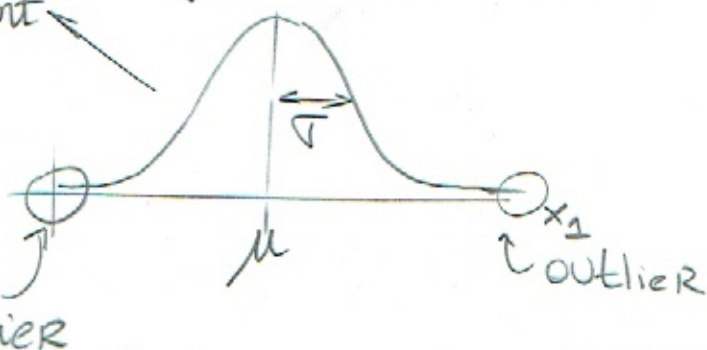
ANOMALY = OUTLIER

can change  
OVER TIME

↳ any deviation from NORMAL / Expectation

issue in anomaly detection

→ we don't  
know  $P_x$



pos neg

① simple stat model

$$P(x) < \epsilon \rightarrow \text{ANOMALY}$$

Z-SCORE →  $Z = \frac{x - \mu}{\sigma}$  → How many std. dev from mean is observation

② Density based anomaly → many variables ②

↳  $\frac{\text{mass} \rightarrow \# \text{ points}}{\text{volume} \rightarrow \text{space}}$

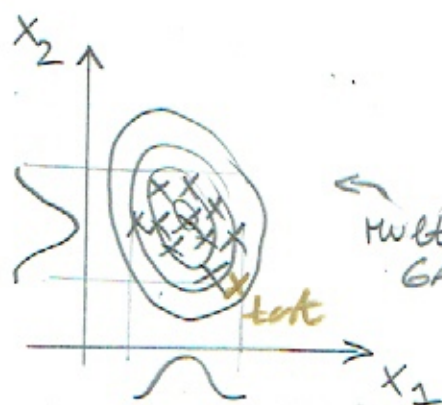
$D = \{x_1, x_2, \dots, x_d\}$   
 $d$  features  
 with unknown distribution

solution: try to estimate distribution through MLE parameters

A) Assume Gaussian!

$$x_1^{(i)} \sim N(\mu_1, \sigma_1^2)$$

$$x_d^{(i)} \sim N(\mu_d, \sigma_d^2)$$



Multi-variate Gaussian

$\mu_i, \sigma_i^2$  unknown

MLE

estimate  $\mu_i, \sigma_i^2$  so that the model best explains the observed data

B) Assume  $x_i$  ARE independent

(A) in (B) →  $p(x_{\text{test}}) = \prod_{j=1}^d p(x_j; \mu_j, \sigma_j^2)$

SO: 1) select features  $x_i$  that you think might be indicative of anomalous examples

2) fit parameters  $\mu_i, \sigma_i^2$

3) calculate  $p(x)$  for new point  $x_{\text{test}}$   
 if  $p(x_{\text{test}}) < \epsilon \rightarrow \text{ANOMALY}$

0.02



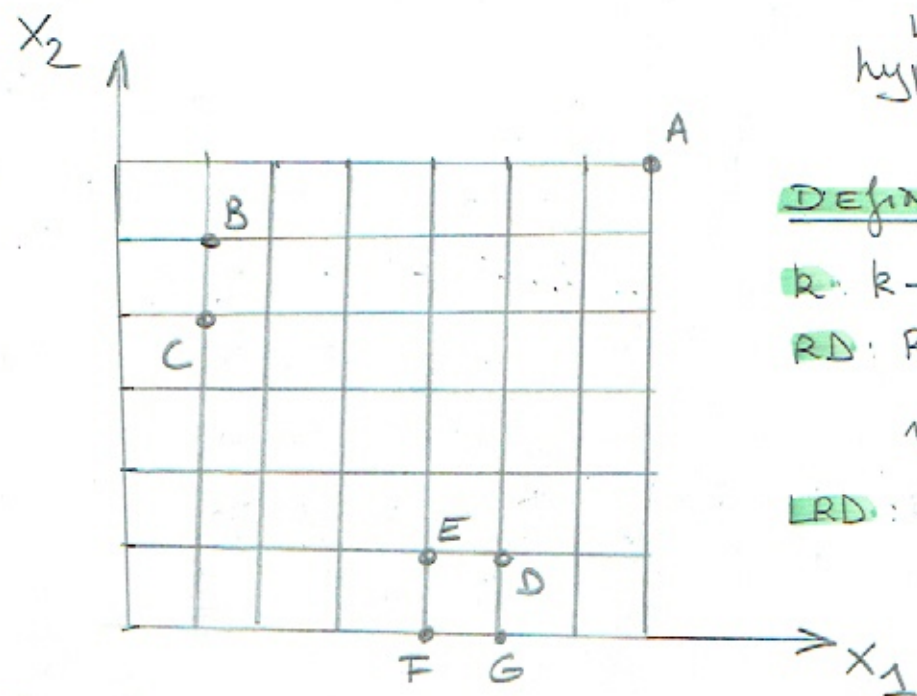
# example of density-based ANomaly Detection System<sup>(3)</sup>

↳ Local Outlier Factor: LOF (2000)

↳ compares local density of a point to the local density of its k-neighbors

INTUITION: outliers are isolated

↳ low density area  
hyper parameter!



## DEFINITIONS

k: k-th neighbor of a point

RD: Reachability distance of a point A with respect to its k neighbors

LRD: local Reachability Density

$$AV(RD_A) = \sum_k \max(k\text{-dist}_k, \text{dist}(A, k))$$

↳ Average distance from pt A to its 3 neighbors!

$$LRD_A = [AV(RD_A)]^{-1}$$

to go from distance to density!

$$RD_A \overset{k=3}{=} \max \left( \begin{array}{l} \text{3rd-dist B, distance (AB)} \quad 5.65 \quad (6.08) \\ \text{3rd-dist C, distance (AC)} \quad 5 \quad (6.32) \\ \text{3rd-dist D, distance (AD)} \quad 1.4 \quad (5.4) \end{array} \right)$$

↳ B, C, D don't see A as their k-NN

$$\Rightarrow AV(RD_A) = \frac{1}{3} (6.08 + 6.32 + 5.4) = 5.9 \quad \rightarrow LRD_A = 0.168$$

$$RD_B = \left[ \begin{array}{l} \text{MAX} \left( \overset{6.32}{\text{3rd dist A}}, \overset{6.08}{\text{dist(AB)}} \right) + \\ \text{MAX} \left( \overset{5}{\text{3rd dist C}}, \overset{1}{\text{dist(BC)}} \right) + \\ \text{MAX} \left( \underset{1.4}{\text{3rd dist D}}, \overset{5.65}{\text{dist(BD)}} \right) \end{array} \right]$$

$$\rightarrow AV(RD_B) = \frac{1}{3} (6.32 + 5 + 5.65) = 5.65$$

$$LRD_B = 0.176$$

$$RD_C = \left[ \begin{array}{l} \text{MAX} \left( \overset{6.32}{\text{3rd dist A}}, \overset{6.32}{\text{dist(AC)}} \right) + \\ \text{MAX} \left( \overset{5.65}{\text{3rd dist B}}, \overset{1}{\text{dist(BC)}} \right) + \\ \text{MAX} \left( \underset{1.4}{\text{3rd dist D}}, \overset{5}{\text{dist(CD)}} \right) \end{array} \right]$$

$$\rightarrow AV(RD_C) = \frac{1}{3} (6.32 + 5.65 + 5)$$

$$= 5.65 \rightarrow LRD_C = 0.176$$

$$RD_D = \left[ \begin{array}{l} \text{MAX} \left( \overset{6.32}{\text{3rd dist A}}, \overset{5.4}{\text{dist(AD)}} \right) + \\ \text{MAX} \left( \overset{5.65}{\text{3rd dist B}}, \overset{5.65}{\text{dist(BD)}} \right) + \\ \text{MAX} \left( \underset{5.2}{\text{3rd dist C}}, \overset{5}{\text{dist(CD)}} \right) \end{array} \right]$$

$$\rightarrow AV(RD_D) = \frac{1}{3} (6.32 + 5.65 + 5)$$

$$= 5.65 \rightarrow LRD_D = 0.176$$

$$LOF_A = \frac{LRD_B + LRD_C + LRD_D}{3} = \frac{0.176}{0.168} = \underline{1.05}$$

outlier?

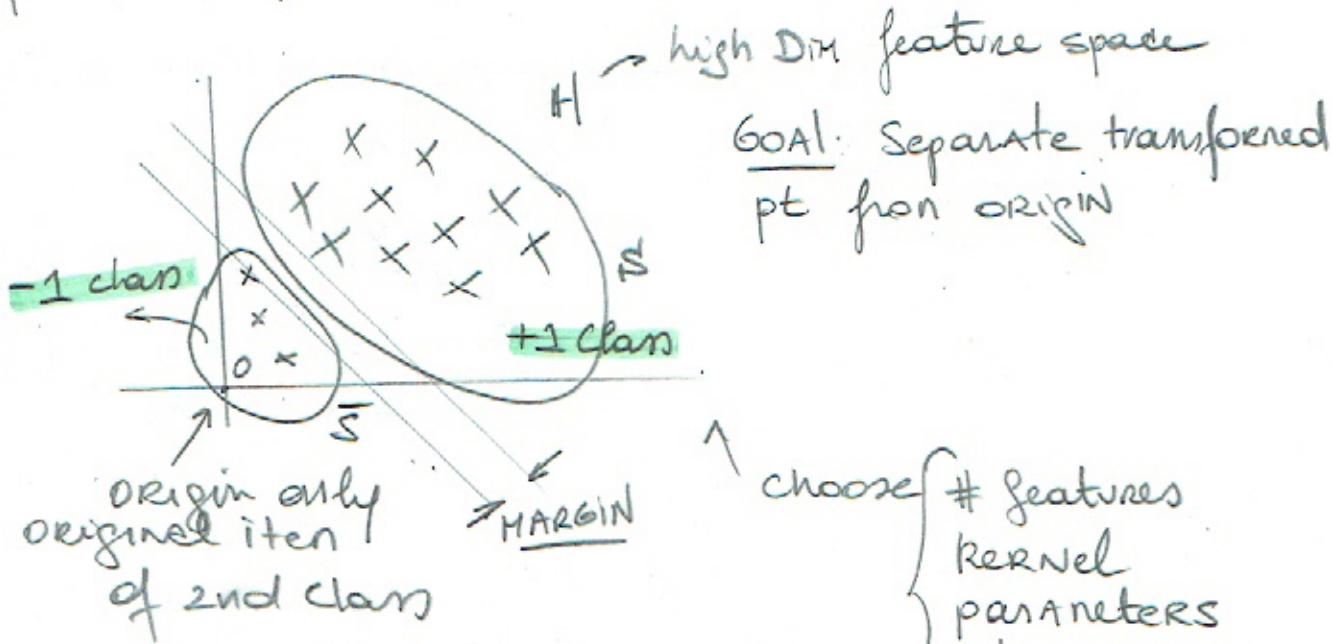
density A is smaller than its neighbors

$$LRD_A = 0.168$$



- ③ One Class SVM → UNSUPERVISED ⑤  
 Schölkopf 1999  
 $x_1, x_2, \dots, x_m \rightarrow$  training examples belonging to 1 class

$\phi: X \rightarrow H$  via kernel

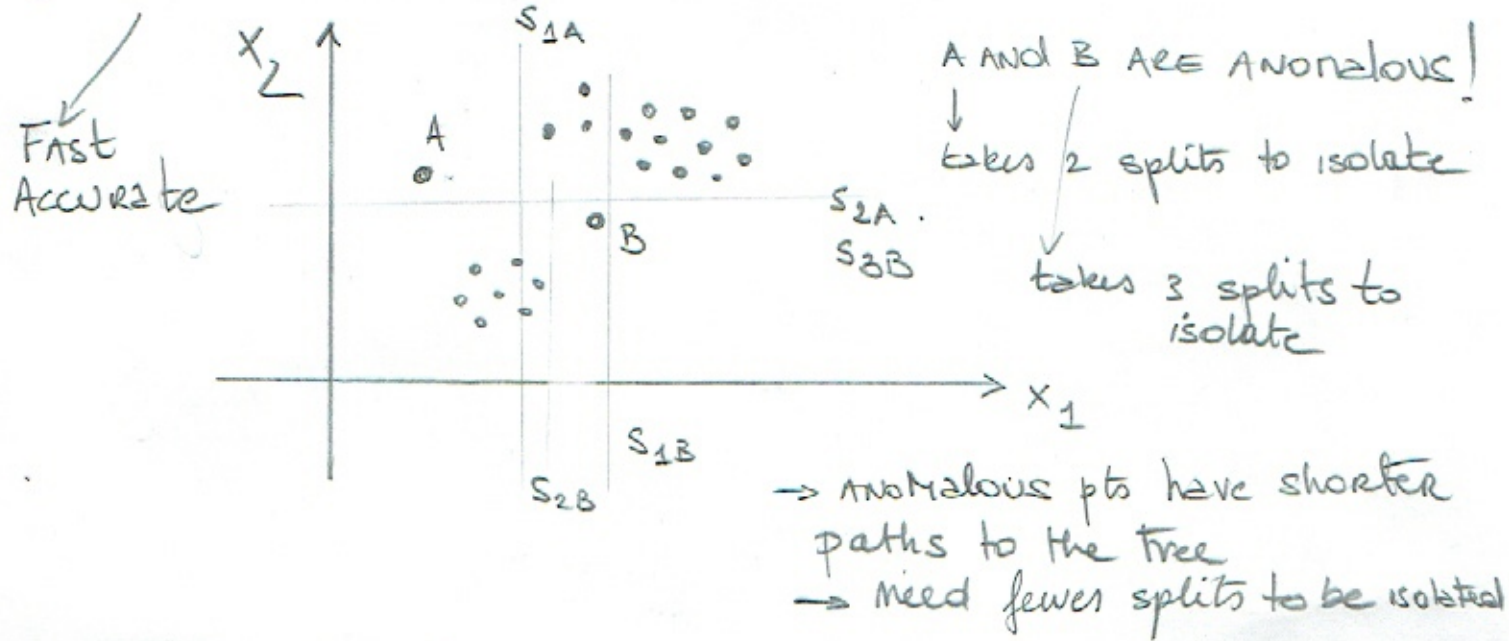


Find subspace  $S$  of  $H$  such that

$\rightarrow P(x_{\text{test}} \notin S) < \nu \rightarrow$  rest are outliers!

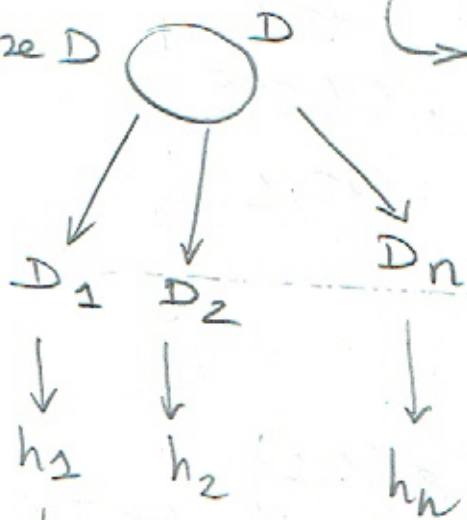
$$f(x) = \begin{cases} +1 & \text{if } x \in S \\ -1 & \text{if } x \in \bar{S} \end{cases} \text{ WITH } \bar{S} \text{ the complement of } S$$

#### ④ Isolation Forest



# Remember Random Forests

Size  $D_i = \text{Size } D$

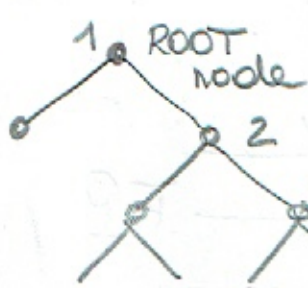


- ① split until you have 1 pt per leaf
- ② split on  $k < d$  features

$$\bar{h} = \frac{1}{M} \sum_{i=1}^M h_i$$

calculate path length for each point in each tree

→ Average path length for each pt is a measure for its anomaly score (shorter path means higher anomaly score)



→ 1 feature = 1 dimension  
 $X = \{x_1, \dots, x_m\}$  m data pts

→  $m-1$  INT Node → 2 children  
ext Node → no kids  
 → fully grown trees have  $m$  external nodes.

each split → binary!  
 → BINARY Trees!

→ memory requirement  $\sim n$

→ limit depth of tree because anomalies are isolated much earlier compared to normal points

anomaly score → SWAPPING → FALSE POSITIVE  
 MASKING → FALSE NEGATIVE

iForests do not need to calculate distances  
 → lowers comp costs  
 → avoid curse!

iForests work well with large DATASETS!