

# Machine Learning in High Dimension

## IA317

### Anomaly Detection

Thomas Bonald

2023 – 2024



# Context

## Anomaly

Deviation from the **expected** behavior

Various sources of anomaly:

- ▶ Errors
- ▶ Frauds
- ▶ Failures
- ▶ Attacks
- ▶ Specific events

How to detect **anomalies** / **outliers** without **supervision**?

# Outline

## ► **Algorithms**

Isolation metrics

Isolation tables

Isolation forests

## ► **Metrics**

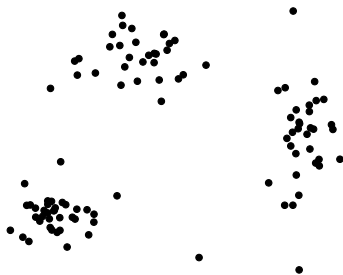
Nearest neighbors

Locally sensitive hashing

Ensemble methods

# Isolation metrics

**Idea:** Isolated samples are likely outliers

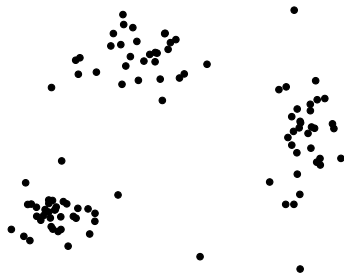


## Local outlier factor

### Definition

$$\text{LOF}(i) = \frac{\text{Local density of } k \text{ NN of } i}{\text{Local density of } i}$$

$\text{LOF} \gg 1 \Leftrightarrow$  potential outlier

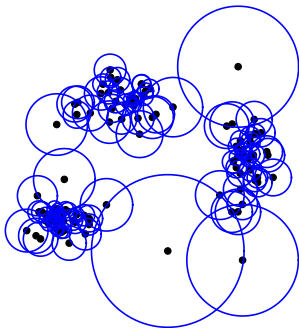


## Ball radius

The **local density** around sample  $i$  can be estimated as:

$$\frac{1}{r(i)}$$

where  $r(i)$  is the distance to the  $k$ -th nearest neighbor (ball radius)



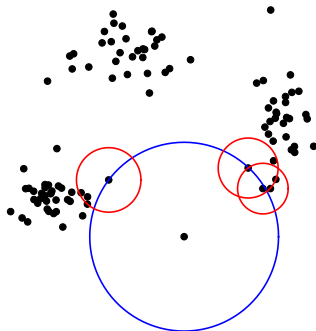
**Example:** Balls formed by the nearest neighbors ( $k = 3$ )

## Simple local outlier factor

Let  $N(i)$  be the  $k$  nearest neighbors of sample  $i$  (excluding itself)

### Definition

$$\text{Simple-LOF}(i) = \frac{\text{Local density of } k \text{ NN of } i}{\text{Local density of } i} = \frac{1}{k} \sum_{j \in N(i)} \frac{r(i)}{r(j)}$$



## Example: MNIST

$$X \in \{0, \dots, 255\}^{n \times d}$$

$n = 10,000$  samples

$$d = 28 \times 28 = 784$$

Simple-LOF ( $k = 10$ , cosine similarity)



1	6	9	3	4	9	4	5	1	8
0	9	2	0	1	5	4	7	3	2
0	9	4	3	8	9	6	2	3	7
3	9	9	8	1	5	3	8	3	1
0	5	3	3	7	1	6	4	8	1
2	3	1	9	2	7	1	9	9	1
1	8	2	3	6	7	8	5	1	7
2	1	6	2	8	9	7	3	4	9
6	9	3	4	4	2	9	6	6	3
9	2	0	4	3	8	1	6	1	0

Random samples



## Example: MNIST

$$X \in \{0, \dots, 255\}^{n \times d}$$

$n = 10,000$  samples

$$d = 28 \times 28 = 784$$

Simple-LOF ( $k = 10$ , cosine similarity)



Random samples



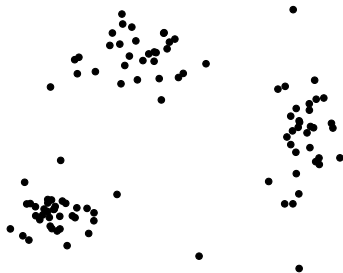
Outliers

## Another measure of local density

The **local density** around sample  $i$  can be estimated as:

$$\frac{1}{R(i)}$$

where  $R(i)$  is the **average reachability distance** to the  $k$  nearest neighbors

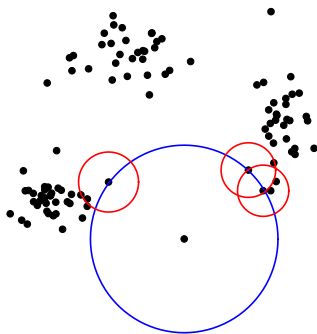


## Local outlier factor

Let  $N(i)$  be the  $k$  nearest neighbors of sample  $i$  (excluding itself)

### Definition

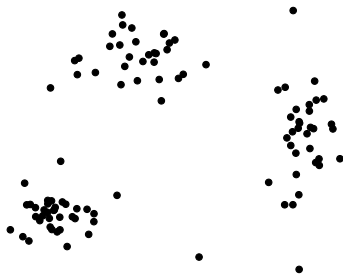
$$\text{LOF}(i) = \frac{\text{Local density of } k \text{ NN of } i}{\text{Local density of } i} = \frac{1}{k} \sum_{j \in N(i)} \frac{R(i)}{R(j)}$$



## Reachability distance

The **reachability distance** of  $i$  from  $j$  is defined by:

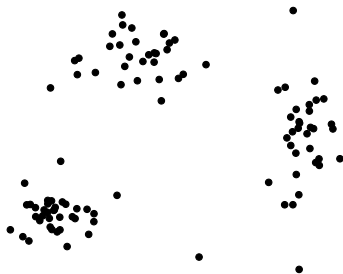
$$r(i,j) = \max(d(i,j), r(j))$$



## Reachability distance

The **reachability distance** of  $i$  from  $j$  is defined by:

$$r(i, j) = \max(d(i, j), r(j))$$



**Note:** Each  $k$ -NN of  $j$  has the reachability distance  $r(j)$  from  $j$

**Warning:** The reachability distance is **not** symmetric!

## Example: MNIST

$$X \in \{0, \dots, 255\}^{n \times d}$$

$n = 10,000$  samples

$$d = 28 \times 28 = 784$$

Outliers ( $k = 10$ , cosine similarity)



Simple LOF



LOF

# Outline

- ▶ **Algorithms**

  - Isolation metrics

  - Isolation tables**

  - Isolation forests

- ▶ **Metrics**

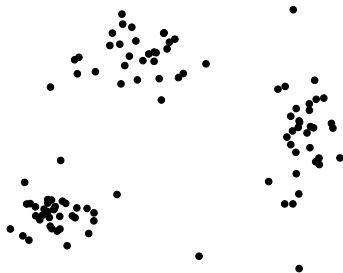
  - Nearest neighbors

  - Locally sensitive hashing

  - Ensemble methods

# Isolation tables

**Idea:** Isolated samples tend to have **specific signatures** by locally sensitive hashing





## Hashing isolation score

Data  $x_1, \dots, x_n \in \mathbb{R}^d$

Given some **locally sensitive hash** functions  $h_1, \dots, h_L$  chosen uniformly at random, we build  $L$  hash tables:

$H_1 : s \rightarrow$  set of samples  $i$  such that  $h_1(x_i) = s$

$H_2 : s \rightarrow$  set of samples  $i$  such that  $h_2(x_i) = s$

$\vdots$

$H_L : s \rightarrow$  set of samples  $i$  such that  $h_L(x_i) = s$

### Definition

$$\text{HIS}(i) = \frac{1}{L} \sum_{l=1}^L \frac{1}{|\underbrace{\{j : h_l(x_j) = h_l(x_i)\}}_{\text{bucket of sample } i}}| \in [0, 1]$$

## Example: MNIST

$$X \in \{0, \dots, 255\}^{n \times d}$$

$n = 10,000$  samples

$$d = 28 \times 28 = 784$$

Outliers



Simple LOF



HIS (random projection)

# Outline

- ▶ **Algorithms**

  - Isolation metrics

  - Isolation tables

  - Isolation forests**

- ▶ **Metrics**

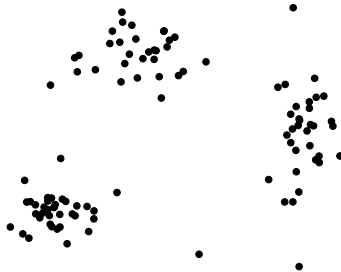
  - Nearest neighbors

  - Locally sensitive hashing

  - Ensemble methods

# Isolation forests

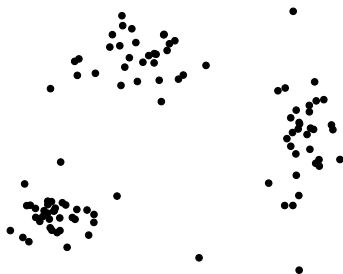
**Idea:** Isolated samples are **easily separable** by random splits



# Extra Random Tree

## Algorithm

Recursively split each set of samples at random (random feature, random threshold) until each sample is **isolated** or has the **same value** as all other samples



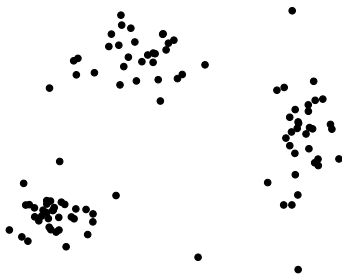
# Isolation forest

**Parameters:** Number of trees  $N$ , sampling size  $s < n$

## Training

For  $t = 1, \dots, N$ ,

- ▶ sample  $s$  data points without replacement
- ▶ build an Extra Random Tree using these  $s$  samples



# Isolation forests

## Evaluation

Evaluate the **average depth** of each sample  $i$  over the  $N$  trees:

$$D(i) = \frac{1}{N} \sum_{t=1}^N \text{depth of } i \text{ in tree } t$$

Low depth  $\Leftrightarrow$  potential outlier

# Anomaly score

## Definition

For each sample  $i$ ,

$$S(i) = 2^{-\frac{D(i)}{D}} \in [0, 1]$$

where  $D$  is the average depth of unsuccessful searches in a **random binary tree** with  $s$  values

$S(i)$  close to 1  $\Leftrightarrow$  potential outlier



## Example: MNIST

$$X \in \{0, \dots, 255\}^{n \times d}$$

$n = 10,000$  samples

$$d = 28 \times 28 = 784$$

Outliers



Simple LOF



Isolation forest

# Outline

## ► **Algorithms**

Isolation metrics

Isolation tables

Isolation forests

## ► **Metrics**

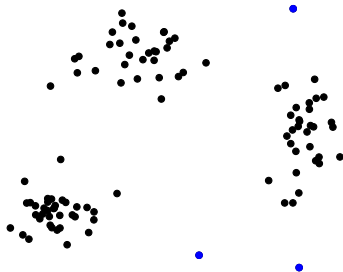
Nearest neighbors

Locally sensitive hashing

Ensemble methods

# Metric

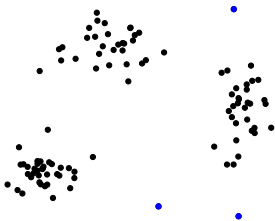
**Idea:** Use **annotated data** to assess the quality of anomaly detection



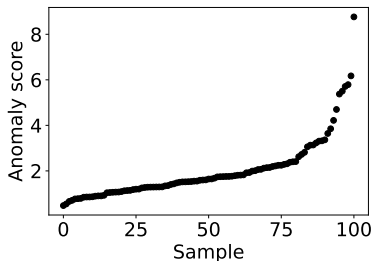
# Metric

**Issue:** The 3 algorithms (isolation metrics / tables / forests) provide an **anomaly score**  $S$ , not a binary decision

Samples



Scores



## Binary classification

Let  $y_1, \dots, y_n \in \{0, 1\}$  be the **true labels** (1 = anomaly)

Let  $\hat{y}_1 = 1_{\{s_1 \geq t\}}, \dots, \hat{y}_n = 1_{\{s_n \geq t\}} \in \{0, 1\}$  be the **predicted labels** at threshold  $t$

## Binary classification

Let  $y_1, \dots, y_n \in \{0, 1\}$  be the **true labels** (1 = anomaly)

Let  $\hat{y}_1 = 1_{\{s_1 \geq t\}}, \dots, \hat{y}_n = 1_{\{s_n \geq t\}} \in \{0, 1\}$  be the **predicted labels** at threshold  $t$

True positive rate (recall)

$$\frac{\# \text{ True positive}}{\# \text{ Positive}} = \frac{\sum_i 1_{\{\hat{y}_i=1, y_i=1\}}}{\sum_i 1_{\{y_i=1\}}}$$

## Binary classification

Let  $y_1, \dots, y_n \in \{0, 1\}$  be the **true labels** (1 = anomaly)

Let  $\hat{y}_1 = 1_{\{s_1 \geq t\}}, \dots, \hat{y}_n = 1_{\{s_n \geq t\}} \in \{0, 1\}$  be the **predicted labels** at threshold  $t$

True positive rate (recall)

$$\frac{\# \text{ True positive}}{\# \text{ Positive}} = \frac{\sum_i 1_{\{\hat{y}_i=1, y_i=1\}}}{\sum_i 1_{\{y_i=1\}}}$$

False positive rate (fall out)

$$\frac{\# \text{ False positive}}{\# \text{ Negative}} = \frac{\sum_i 1_{\{\hat{y}_i=1, y_i=0\}}}{\sum_i 1_{\{y_i=0\}}}$$

## A probabilistic view

Let  $y \in \{0, 1\}$  be the **true label** of a sample ( $1 = \text{anomaly}$ )

Let  $S$  be its score, and  $\hat{y} = 1_{\{S \geq t\}}$  the predicted label



## A probabilistic view

Let  $y \in \{0, 1\}$  be the **true label** of a sample (1 = anomaly)

Let  $S$  be its score, and  $\hat{y} = 1_{\{S \geq t\}}$  the predicted label

True positive rate (recall)

$$R(t) = P(S \geq t \mid y = 1)$$

## A probabilistic view

Let  $y \in \{0, 1\}$  be the **true label** of a sample (1 = anomaly)

Let  $S$  be its score, and  $\hat{y} = 1_{\{S \geq t\}}$  the predicted label

True positive rate (recall)

$$R(t) = P(S \geq t \mid y = 1)$$

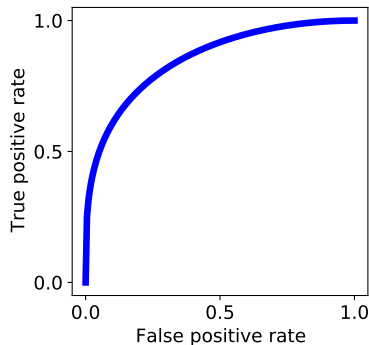
False positive rate (fall out)

$$F(t) = P(S \geq t \mid y = 0)$$

# The ROC<sup>1</sup> curve

## Definition

Plot of **Recall** against **Fall out** when the threshold  $t$  goes from  $-\infty$  to  $+\infty$



---

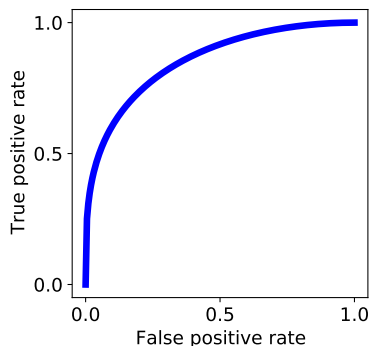
<sup>1</sup>Receiver Operating Characteristic

# The ROC AUC score

## Definition

The ROC AUC (Area Under Curve) score is:

$$\text{AUC} = \int_0^1 \text{ROC}(u) du \in [0, 1]$$



# Probabilistic interpretation of the ROC AUC score

## Proposition

The ROC AUC score is the probability that the score is consistent for a **random pair** of samples  $i, j$  with distinct labels:

$$\text{AUC} = P(S_i > S_j | y_i = 1, y_j = 0)$$

# Probabilistic interpretation of the ROC AUC score

## Proposition

The ROC AUC score is the probability that the score is consistent for a **random pair** of samples  $i, j$  with distinct labels:

$$\text{AUC} = P(S_i > S_j | y_i = 1, y_j = 0)$$

## Notes:

- ▶ We have  $\text{AUC} = \frac{1}{2}$  for a **random** prediction
- ▶ Alternative metric = **Mean Average Precision**  
→ AUC of Precision against Recall

# Summary

## Anomaly detection

### Algorithms

- ▶ Isolation metrics      Nearest neighbors
- ▶ Isolation tables      Locally sensitive hashing
- ▶ Isolation forests      Ensemble method

The ROC AUC score, a **quality metric** for annotated data

# 8 8 2 + 2 4 4 1 1  
# 1 8 1 2 1 2 1 2  
/ 4 3 6 1 8 ( 2 1 8  
2 3 4 1 2 5 3 7 /  
4 3 / 3 7 7 9 1 6 7  
1 1 9 / 7 7 8 8 / 1  
8 1 1 7 1 7 5 0 2 3  
1 8 \ 7 2 3 1 1 8 4  
4 1 0 8 3 1 8 1 1 7  
( 1 8 4 1 4 5 / 4: 1

2 2 2 6 7 6 5 0 7 5  
8 6 5 7 7 9 7 6 7 6  
6 7 0 2 5 0 2 2 7 2  
6 5 4 7 6 2 0 3 7 5  
7 5 4 9 4 7 5 0 7 3  
2 2 2 8 2 2 7 7 7  
2 4 3 5 7 4 9 0 6 3  
4 2 4 0 4 8 0 7 2 7  
5 6 7 7 9 2 4 2 3 3  
7 5 3 2 7 0 6 3

0 0 8 7 0 0 5 0 4 0  
4 0 0 7 6 0 5 7 6 0  
6 5 4 5 6 0 0 3 9 5  
0 0 0 8 0 6 0 0 2 2  
3 4 6 7 5 0 2 6 0 0  
0 5 3 5 6 5 0 0 2 4  
6 0 0 5 7 5 9 6 0 8  
4 6 0 0 4 6 6 9 4 2  
0 5 0 4 0 5 7 6 6 8  
0 0 0 7 7 7 0 6 0 0