

COMPSCI 762 Tutorial 11

Tutorial 11 – Anomaly Detection and Data Stream Mining

Luke Chang

June 2021

The University of Auckland

Anomaly Detection

Data Stream Mining

Anomaly Detection

Types of Anomaly

- **Global outlier (Point anomaly):** deviates significantly from the rest of the data set. The simplest type of outliers.
- **Contextual outlier (Conditional outlier):** deviates significantly with respect to a specific context of the object.
 - **Contextual attributes:** define the object's context.
 - **Behavioral attributes:** define the object's characteristics, and are used to evaluate whether the object is an outlier in the context.

Example

A temperature sensor measures 4°C in May. It is a perfectly normal reading in Wellington, but it might be an outlier if the location is New York. The location and the date are **contextual attributes**, and the temperature is a **behavioral attribute**.

- **Collective outliers:** the objects as a whole deviate significantly from the entire data set.

Types of Anomaly

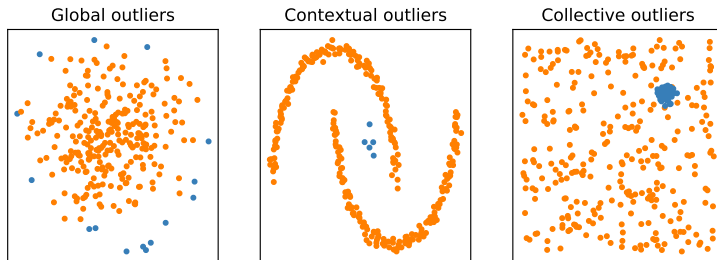
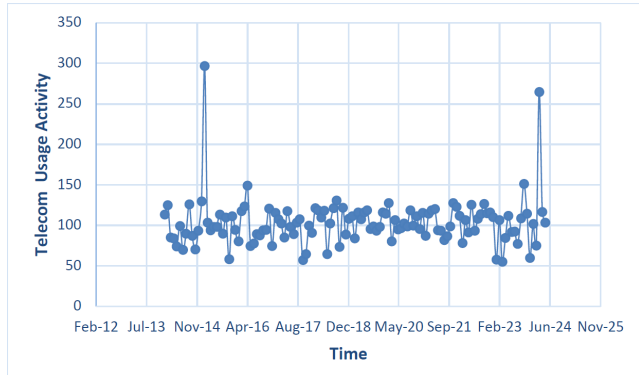


Figure 1: Types of outliers

- **Fig 1a:** Global outliers
- **Fig 1b:** Contextual outliers – Given the dataset has two clusters, each one has a moon shape.
- **Fig 1c:** Blue points are collective outliers because the density of those points is much higher than the rest.

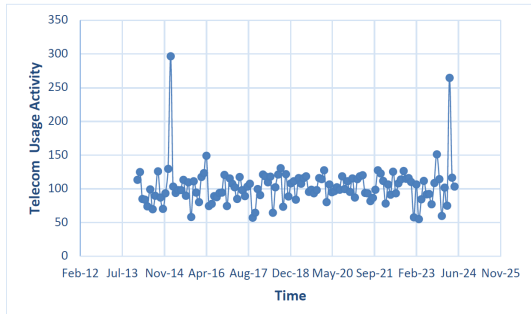
2020 361 Exam Question 4 – Outlier / Anomaly Detection

The following figure shows the monthly telecom usage activity in a specific area across several years. You were tasked to identify whether outliers exist in the data.



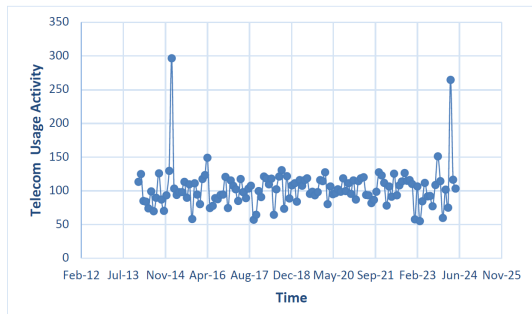
What outlier detection technique would you use to identify whether outliers exist for this case?

2020 361 Exam Question 4 – Outlier / Anomaly Detection



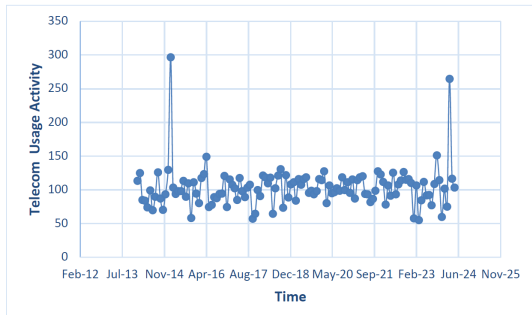
- No obvious periodic pattern;
- No dramatic density changes;
- Unlikely to have contextual outlier and collective outlier;
- We apply the **statistical approach – parametric method** to find global outliers.
- There are multiple ways to solve this problem.

2020 361 Exam Question 4 – Outlier / Anomaly Detection



- **Parametric method:** Assume the samples are drawn from a normal distribution, estimate the maximum likelihood of mean μ and standard deviation σ .
- We know that the $\mu \pm 3\sigma$ region contains 99.7%. According to z-score, if $z = \frac{|x_i - \bar{x}|}{s} > 3$, the sample x_i is generated by the normal distribution is less than $\frac{0.3}{2} = 0.15\%$.
- Note that we do not know the population mean and standard deviation, we only have sample mean and sample standard deviation.

2020 361 Exam Question 4 – Outlier / Anomaly Detection



- Based on observation, each grid contains around 14 points, we have 110 data points in total.
- $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \approx 100$
- $s = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2} \approx 20$
- $\frac{|300-100|}{20} = 10 > 3$, an outlier
- $\frac{|270-100|}{20} = 8.5 > 3$, an outlier
- $\frac{|150-100|}{20} = 2.5 < 3$, not an outlier
- $\frac{|55-100|}{20} = 2.25 < 3$, not an outlier

Tutorial Question 1 – Anomaly Detection

Question 1

You are given the following list of 2D data points. $[[1; 1]; [1; 2]; [2; 2]; [2; 1]; [3; 3]; [2; 5]; [2; 3]]$ If you had to select one point to be anomalous, which would you pick? Explain your answer. Please link your explanation to an anomaly detection technique.

There are multiple ways to solve this problem.

Let's use **distance-based outlier detection**.

Tutorial Question 1 – Anomaly Detection

Distance-based outlier detection

	1,1	1,2	2,2	2,1	3,3	2,5	2,3
1,1	0	1	2	1	4	5	3
1,2	1	0	1	2	3	4	2
2,2	2	1	0	1	2	3	1
2,1	1	2	1	0	3	4	2
3,3	4	3	2	3	0	3	1
2,5	5	4	3	4	3	0	2
2,3	3	2	1	2	1	2	0

Table 1: Manhattan Distance Matrix

$$\frac{|\{o' | \text{dist}(o, o') \leq r\}|}{|D|} \leq \pi$$

where $|D|$ is the *cardinality* of the set D .

We need to define the hyperparameters: r and π .

The number of data points, $|D|$, is 7.

Let $r = 2$,

	# of objects within r	Divide by $ D $
1,1	3	0.43
1,2	4	0.57
2,2	5	0.71
2,1	4	0.57
3,3	2	0.29
2,5	1	0.14
2,3	5	0.71

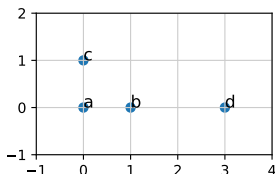
Table 2: The # of neighbours within r

If we select only one point as an outlier, we can set π to any value between 0.14 and 0.29, e.g. 0.15. Therefore, we identify [2;5] is an outlier.

Tutorial Question 2 – Density-based Approach: Local Outlier Factor (LOF)

Question 2

Consider a set of points $(0,0)$, $(1,0)$, $(0,1)$, $(3,0)$. Calculate the LOF score for the points using Manhattan distance and k is 2.



	a (0,0)	b (1,0)	c (0,1)	d (3,0)
a (0,0)	0	1	1	3
b (1,0)	1	0	2	2
c (0,1)	1	2	0	4
d (3,0)	3	2	4	0

Table 3: Manhattan distance matrix

	$\text{dist}_2(o)$	$N_2(o)$
a (0,0)	1	2
b (1,0)	2	3
c (0,1)	2	2
d (3,0)	3	2

Table 4: Distance between data point o and its k -th nearest neighbour ($k = 2$)

We denote the set of k -nearest neighbours as $N_k(o)$:

$$N_k(o) = \{o' | o' \in D, \text{dist}(o, o') \leq \text{dist}_k(o)\}$$

Note: $|N_k(o)|$ may contain more than k objects, because objects may have same distance.

Tutorial Question 2 – LOF

Let o' be a neighbour of o , to avoid $\text{dist}(o, o')$ is too small, we compute the reachability distance from o to o' :

$$\text{reachdist}_k(o' \leftarrow o) = \max\{\text{dist}_k(o), \text{dist}(o, o')\}$$

Note that reachability distance is not symmetric, thus

$$\text{reachdist}_k(o \leftarrow o') \neq \text{reachdist}_k(o' \leftarrow o)$$

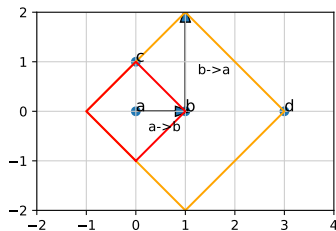


Figure 2: $\text{reachdist}_2(b \leftarrow a) \neq \text{reachdist}_2(a \leftarrow b)$

Tutorial Question 2 – LOF

$\text{reachdist}_2(b \leftarrow a)$	$= \max\{\text{dist}_2(a), \text{dist}(a, b)\}$	$= \max\{1, 1\}$	$= 1$
$\text{reachdist}_2(c \leftarrow a)$	$= \max\{\text{dist}_2(a), \text{dist}(a, c)\}$	$= \max\{1, 1\}$	$= 1$
$\text{reachdist}_2(d \leftarrow a)$	$= \max\{\text{dist}_2(a), \text{dist}(a, d)\}$	$= \max\{1, 3\}$	$= 3$
$\text{reachdist}_2(a \leftarrow b)$	$= \max\{\text{dist}_2(b), \text{dist}(b, a)\}$	$= \max\{2, 1\}$	$= 2$
$\text{reachdist}_2(c \leftarrow b)$	$= \max\{\text{dist}_2(b), \text{dist}(b, c)\}$	$= \max\{2, 2\}$	$= 2$
$\text{reachdist}_2(d \leftarrow b)$	$= \max\{\text{dist}_2(b), \text{dist}(b, d)\}$	$= \max\{2, 2\}$	$= 2$
$\text{reachdist}_2(a \leftarrow c)$	$= \max\{\text{dist}_2(c), \text{dist}(c, a)\}$	$= \max\{2, 1\}$	$= 2$
$\text{reachdist}_2(b \leftarrow c)$	$= \max\{\text{dist}_2(c), \text{dist}(c, b)\}$	$= \max\{2, 2\}$	$= 2$
$\text{reachdist}_2(d \leftarrow c)$	$= \max\{\text{dist}_2(c), \text{dist}(c, d)\}$	$= \max\{2, 4\}$	$= 4$
$\text{reachdist}_2(a \leftarrow d)$	$= \max\{\text{dist}_2(d), \text{dist}(d, a)\}$	$= \max\{3, 3\}$	$= 3$
$\text{reachdist}_2(b \leftarrow d)$	$= \max\{\text{dist}_2(d), \text{dist}(d, b)\}$	$= \max\{3, 2\}$	$= 3$
$\text{reachdist}_2(c \leftarrow d)$	$= \max\{\text{dist}_2(d), \text{dist}(d, c)\}$	$= \max\{3, 4\}$	$= 4$

$$\text{reachdist}_k(o \leftarrow o') \neq \text{reachdist}_k(o' \leftarrow o)$$

Tutorial Question 2 – LOF

The *Local Reachability Density* (LRD) of an object o is defined as

$$\text{lrd}_k(o) = \frac{|N_k(o)|}{\sum_{o' \in N_k(o)} \text{reachdist}_k(o \leftarrow o')}$$

$\text{lrd}_2(a)$	$= N_2(a) / (\text{reachdist}_2(a \leftarrow b) + \text{reachdist}_2(a \leftarrow c))$	$= 2 / (2 + 2)$	$= 0.5$
$\text{lrd}_2(b)$	$= N_2(b) / (\text{reachdist}_2(b \leftarrow a) + \text{reachdist}_2(b \leftarrow c) + \text{reachdist}_2(b \leftarrow d))$	$= 3 / (1 + 2 + 3)$	$= 0.5$
$\text{lrd}_2(c)$	$= N_2(c) / (\text{reachdist}_2(c \leftarrow a) + \text{reachdist}_2(c \leftarrow b))$	$= 2 / (1 + 2)$	$= 0.667$
$\text{lrd}_2(d)$	$= N_2(d) / (\text{reachdist}_2(d \leftarrow a) + \text{reachdist}_2(d \leftarrow b))$	$= 2 / (3 + 2)$	$= 0.4$

Tutorial Question 2 – LOF

The *Local Outlier Factor* (LOF) of an object o is

$$\text{LOF}_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{\text{lrd}_k(o')}{\text{lrd}_k(o)}}{|N_k(o)|} = \frac{\sum_{o' \in N_k(o)} \text{lrd}_k(o')}{|N_k(o)| \cdot \text{lrd}_k(o)}$$

$\text{LOF}_2(a)$	$= (\text{lrd}_2(b) + \text{lrd}_2(c)) / (N_2(a) \cdot \text{lrd}_2(a))$	$= (0.5 + 0.667) / (2 \times 0.5)$	$= 1.167$
$\text{LOF}_2(b)$	$= (\text{lrd}_2(a) + \text{lrd}_2(c) + \text{lrd}_2(d)) / (N_2(b) \cdot \text{lrd}_2(b))$	$= (0.5 + 0.667 + 0.4) / (3 \times 0.5)$	≈ 1.045
$\text{LOF}_2(c)$	$= (\text{lrd}_2(a) + \text{lrd}_2(b)) / (N_2(c) \cdot \text{lrd}_2(c))$	$= (0.5 + 0.5) / (2 \times 0.667)$	≈ 0.750
$\text{LOF}_2(d)$	$= (\text{lrd}_2(a) + \text{lrd}_2(b)) / (N_2(d) \cdot \text{lrd}_2(d))$	$= (0.5 + 0.5) / (2 \times 0.4)$	$= 1.25$

Tutorial Question 2 – LOF

```
import numpy as np
from sklearn.neighbors import LocalOutlierFactor

data = np.array([[0,0],[1,0],[0,1],[3,0]])
# The neighbours include itself.
clf = LocalOutlierFactor(n_neighbors=3, metric='manhattan', novelty=False)
pred = clf.fit_predict(data)
print(pred)
# [1 1 1 1]
print(clf.negative_outlier_factor_)
# [-1.04377104 -1.18148148 -0.90606061 -0.90606061]
```

Tutorial Question 3 – Supervised

Question 3

What is a difference, between a supervised, semi-supervised, and unsupervised anomaly detection techniques?

Supervised

- With labels, i.e., normal objects vs. outliers
- Binary classification problem
- Extremely imbalanced, majority of the objects are normal.

Unsupervised

- Labels are not available
- Clustering based methods

Semi-supervised

- Some labeled examples is feasible, the number of such labeled examples is small.
- Combining classification-based and clustering based methods
- Use the one-class model to label normal objects, e.g. One-class SVM

Data Stream Mining

Tutorial Question 6 – Data Stream

Question 6

Give examples of three requirements for data stream mining that makes it different from regular data mining.

Examples:

- Web data, e.g. Amazon web orders, video steaming (real-time decoding)
- IoT sensors, e.g. Real-time temperature monitoring, blood pressure monitors, step monitor on a smartwatch
- Finance, e.g. Real-time stock price data
- Transportation and supply-chain, e.g. traffic information, pack-house logistic

Requirements:

- Process one instance at a time, and only inspect it once
- Time and space constraint
- Need to adapt changes (Concept drift)

Tutorial Question 3 – Sliding Window

Question 3

When would you use a sliding window technique as compared to a reservoir sampling technique.

Reservoir Sampling

- A randomized algorithm for choosing a *simple random sample* without replacement of k items from a population of unknown size in a single pass over the items.
- No matter how many stream elements t have been read so far, each of them is currently in the reservoir with the same probability, k/t .
- Maintains a fixed-size uniform random sample

Tutorial Question 3 – Sliding Window

Question 3

When would you use a sliding window technique as compared to a reservoir sampling technique.

Sliding Window

- Only consider the last n samples
- Clear way to bound memory
- Natural in applications: Emphasizes the most recent data
- Do not keep the old samples outside the current window

Concept Drift

Example

You use a real-time temperature sensor to predict the weather in Wellington. The model is trained from the data which are gathered from Wellington. The model shows high performance, but there is a dramatic performance drop when using it to predict the weather in New York.

- **Concept Drift** describes a change in the distribution of input data and labels over time.
- It may cause prediction accuracy to deteriorate over time.

Types of concept drift:

- **Real** concept drift: The relationship between the input data and label has changed. e.g. The concept of the keyword – "Doge" on Twitter.
- **Virtual** concept drift: The distribution of input data has changed, but not the decision boundary. e.g. School holidays, less sales from the food court, the distribution changes, but the purchase pattern still holds.

Drift Detection Method (DDM)

- DDM monitors the number of errors produced by a model learned on the previous stream items.
- Let P_t denote the error rate of the predictor at time t . Given the number of errors in a sample of t examples, its standard deviation at time t is given by:

$$S_t = \sqrt{P_t(1 - P_t)/t}$$

- DDM stores the smallest value of the error rates P_{\min} observed up to time t , and the standard deviation S_{\min} at that point.
- **Warning:** If $P_t + S_t \geq P_{\min} + 2S_{\min}$, **new samples** are stored in anticipation of a possible change.
- **Change:** If $P_t + S_t \geq P_{\min} + 3S_{\min}$, a **new model** is built using the examples stored since the warning. Reset P_{\min} and S_{\min} .

Tutorial Question 1 – Drift Detection Method (DDM)

Question 1

DDM has the drawback that it may take a long time to react to changes after a long period without change. Suggest a couple of ways to fix this, possibly at the cost of introducing some parameters.

- DDM has difficulties when the change is **slowly** gradual.
- The examples will be stored for long time, the drift level can take too much time to trigger and the examples memory can be exceeded.
- *Early Drift Detection Method* (EDDM) by Baena-Garcia, et al. (2006). The idea is to keep track of the **average distance between two errors** instead of only the error rate. It assumes that the distance between two occurrences of classification errors increases as a stable concept is being learnt.