

BACHELOR PAPER

Term paper submitted in partial fulfillment of the requirements for the degree of Bachelor of Science in Engineering at the University of Applied Sciences Technikum Wien - Degree Program Computer Science

Anomaly detection in data for data cleansing

By: David Zelenay

Student Number: 000000000000

Supervisor: Degree First Name Surname

Vienna, May 23, 2022

Declaration

“As author and creator of this work to hand, I confirm with my signature knowledge of the relevant copyright regulations governed by higher education acts (see Urheberrechtsgesetz /Austrian copyright law as amended as well as the Statute on Studies Act Provisions / Examination Regulations of the UAS Technikum Wien as amended).

I hereby declare that I completed the present work independently and that any ideas, whether written by others or by myself, have been fully sourced and referenced. I am aware of any consequences I may face on the part of the degree program director if there should be evidence of missing autonomy and independence or evidence of any intent to fraudulently achieve a pass mark for this work (see Statute on Studies Act Provisions / Examination Regulations of the UAS Technikum Wien as amended).

I further declare that up to this date I have not published the work to hand nor have I presented it to another examination board in the same or similar form. I affirm that the version submitted matches the version in the upload tool.“

Vienna, May 23, 2022

Signature

Kurzfassung

Diese Bachelorarbeit befasst sich mit dem Thema Anomalie-Erkennung mit einem Fokus auf Ausreißererkennung. Der Einstieg bildet ein Überblick, welche Merkmale Datenqualität ausmachen. Danach werden einige Methoden zur Erkennung von Ausreißern vorgestellt. Diese Methoden werden dann mit Daten aus der echten Welt getestet. Die Daten sind Wasserstandsdaten von Flüssen, welche verwendet werden, um Fluten rechtzeitig zu erkennen und betroffene Personen zu warnen. Bevor die Modelle mit den Daten getestet werden, wurde noch eine explorative Datenanalyse durchgeführt, um sich einen Überblick über die Daten zu verschaffen. Am Schluss werden die Ergebnisse der verschiedenen Ansätze Ausreißer zu erkennen, verglichen und ausgewertet. Um die Ergebnisse vergleichen zu können, wurde eine passende Kennzahl definiert, welche Auskunft über die Leistung des Modells gibt.

Abstract

This bachelor thesis analyses anomaly detection methods, focusing on outlier detection for time series data. In the beginning, an overview of the characteristics that make up data quality is outlined. Afterwards, some methods for outlier detection are presented. These Methods are then tested with real-world data. The data used for testing consists of the water levels of different rivers. This data is originally used to warn people about upcoming floods in their area of interest. However, before testing the models with the data, explorative data analysis was conducted to provide an overview of the data. In the end, the results of the different approaches to detect outliers are compared and evaluated. To be able to compare the models, a suitable performance metric was defined.

Keywords: anomaly detection, data cleansing, data cleaning, data quality, outlier detection

Contents

1	Overview	1
1.1	Introduction	1
1.2	Research Question	1
1.3	Research Method	1
2	Data Quality	2
2.1	Features of Data Quality	2
2.2	Improving Data Quality	3
2.3	Data Cleaning & Data Cleansing Approaches	4
2.4	Data Cleaning	4
2.5	Data Cleansing	5
3	Outlier Detection	6
3.1	Outlier Types	6
3.2	Outlier Detection Approaches	9
3.3	Threshold-based Outlier Detection	10
3.4	Outlier Detection using z-score	11
3.5	Outlier Detection using modified z-score	12
4	Outlier Detection based on Water Level Data	13
4.1	What is the Goal?	13
4.2	How to retrieve the Data (Description of the API)	13
4.3	Overview of the Data	17
4.4	Manual Classification of Outliers	17
4.5	Explorative Data Analysis	18
4.5.1	Aghacashlaun Station	18
4.6	Outlier Detection performance Metrics	26
4.6.1	Confusion Matrix	26
4.6.2	Accuracy	26
4.6.3	Precision and Recall	27
4.6.4	F-score	28
4.7	Implementation of different Outlier Detection Approaches	28
4.7.1	Mean Threshold	29
4.7.2	Median Threshold	30
4.7.3	MAD Threshold	31

4.7.4	Z-Score	32
4.7.5	Modified z-score	33
4.7.6	Preprocessing the Data	34
4.7.7	Finding Parameters	34
4.8	Compare different Outlier Detection Approaches	35
4.8.1	Station Aghacashlaun, Aghacashlaun (36022-ie)	43
4.8.2	Station Murg, Frauenfeld (2386-ch)	44
4.8.3	Station Sieg, Betzdorf (2720050000-de)	46
4.8.4	Station Losse, Helsa (42960105-de)	48
4.8.5	Station Crana, Tullyarvan (39003-ie)	49
5	Conclusion	50
6	Future Work	51
7	Appendix	52
7.1	Python Code	52
7.1.1	Manual outlier detection	52
7.1.2	Explorative Data Analysis	54
7.1.3	Preprocessing	61
7.1.4	Outlier Detection Methods Implementation	62
7.1.5	Parameter Grid Search	65
7.1.6	Calculate Outlier Detection Results	67
	Bibliography	70
	List of Figures	73
	List of Tables	74
	List of source codes	75
	List of Abbreviations	76

1 Overview

1.1 Introduction

With the growing popularity of Internet of Things (IoT) and digitizing business processes, there is a growing amount of data available for analysis. In order to utilize the data from the IoT sensors, it needs to be preprocessed. One step of preprocessing is data cleaning (also referred to as data cleansing). The main goal of data cleansing is to increase the data quality and furthermore to detect and remove anomalies in the data. The quality requirements for the data can differ depending on the use case. Anomalies in sensor data are data points which do not picture the reality. For example, an anomaly of a temperature sensor would be if the sensor reads 0 °C and the real temperature is 23 °C. This paper specifically focuses on the removal of outliers for water level sensors. The data is provided by FloodAlert [1], which provides a service to warn people about floods, in their area of interest.

1.2 Research Question

The research questions are:

- What are common methods to detect outliers for time series data?
- How can outlier detection methods be compared, with a focus on water level data?
- Based on water levels from different rivers, which method is able to classify outliers most reliably?

1.3 Research Method

This thesis will provide an overview and comparison of different approaches to detect outliers. It will focus on time series data, especially water level measurements of rivers. To introduce the topic a general overview of data quality, data cleansing, data cleaning and outlier types is provided. For the theoretical parts of the chapters, literature research was conducted. After gathering knowledge on different outlier detection approaches they were implemented in Python. To test the performance a suitable performance metric needs to be chosen. To use real-world data to classify outliers, the water levels from different measurement stations were taken. In the end, the performance of the different approaches to detect outliers is compared.

2 Data Quality

2.1 Features of Data Quality

This section will provide a few examples of key features of data quality.

Completeness

Data completeness describes the wholeness of data. If there are certain aspects of data missing the data is not complete. For example, if each data point of a sensor includes the date, time and production speed, the data is not complete, if one of those features is missing or not entire, this data point is not complete. [2,3]

Accuracy

The accuracy of data describes the exactness. Examples of possible data which decrease the accuracy are outliers or time shifts. Usually, the accuracy of data is harder to measure than the completeness, consistency, structure or documentation. Due to the heterogeneity of sensor data (regarding numerical values like production speed or temperature, not categorical values like on/off) for each data point it is difficult to detect which values are genuine and which are sensor errors and therefore outliers. [2] Furthermore, an obstacle to accurate data are calculation errors made by computers, because the data type (e.g. float) is not always 100% exact. However, there are methods to reduce these errors. Nevertheless, this also needs to be considered when talking about accurate data. [4,5]

Consistency

One example of consistency would be if the data interval is equal. For example, there should be a data point every ten seconds. As soon as two data points are more than ten seconds apart from each other the data is not consistent anymore. [2]

Structure & Documentation

If the structure of the data is not homogeneous, it is very difficult to analyze in an automated way. As a result the data either needs to be structured from the beginning or a process needs to be fabricated, to structure the data automatically. Furthermore, documentation is required in order to structure and preprocess data.

Documentation of data might include data format (Comma Separated Values (CSV), parquet [6], Java Script Object Notation (JSON)), date format (e.g. ISO 8601 with UTC offset), valid value spans (e.g. temperature is only valid if it is between 100 and 400 °C) [2]

2.2 Improving Data Quality

This section will describe methods to improve data quality, based on the features elaborated in section 2.1.

Completeness

The most common methods to increase data completeness are statistical and deep learning-based approaches. The goal of these methods is to fill in the missing values of a dataset. An example of a statistical method is DynaMMo [7]. For Artificial Neuronal Networks (ANNs), Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) can be used to predict missing data. [3]

Accuracy

One approach to increase the accuracy of data is to define constraints for each value. E.g. When a machine cannot produce more than ten pieces per second, because it is physically not possible, the value could be limited to less than or equal to ten. However, limiting the values to a specific range might hide the fact that the machine has an error and is producing faulty products at a rate of 15 pieces per second. This is one of the reasons why more sophisticated outlier detection methods are used. [3]

Consistency

To facilitate consistent data, statistical smoothing or forecasting methods can be used. Examples of methods are AutoRegressive Integrated Moving Average (ARIMA) or Gaussian Process (GP). ANNs can also be used to unify the time series interval between data points. [3]

Structure & Documentation

The process of structuring heterogeneous and messy data is called data wrangling. In order to unify the structure of the data at least some documentation is required. Therefore the documentation of the data is fundamental to analyse or further process it.

2.3 Data Cleaning & Data Cleansing Approaches

There are two main methods when it comes to data cleaning or data cleansing. Ignoring faulty data or replacing it with a representative value. This paper will use the term data cleaning to describe the process of ignoring or deleting incorrect data and the term data cleansing to portray the process of replacing invalid data with representative values. Faulty, incorrect, invalid or wrong data is data which is inaccurate, incomplete or inconsistent.

Example sensor data: (Valid values for `production_speed` range from 0.00 to 2.00 meter(s) per minute)

ID	timestamp	production_speed (meter/minute)	machine_running
0	2021-12-01T12:00:00.000	1.56	True
1	2021-12-01T12:01:00.000	1.58	True
2	2021-12-01T12:02:00.000	3.50	True
3	2021-12-01T12:03:00.000	1.50	False
4	2021-12-01T12:04:00.000	1.50	True
5	2021-12-01T12:05:00.000	1.49	True

Table 1: Example of IoT sensor data

2.4 Data Cleaning

As already mentioned the approach for data cleaning is to ignore or delete faulty data. Depending on the use case either the entire data point needs to be ignored or just one value. The process of data cleaning will be shown with the example data pictured in Table 1. The first incorrect data point has the ID 2. This row is incorrect because the `production_speed` exceeds the maximum value of 2.00. Depending on the use case (e.g. summary of how long the machine has been running) it can make sense to just ignore the row `production_speed` and keep the value for `machine_running`. The second appearance of a faulty data point has the ID 3. This data point is incorrect since `machine_running` is False but the value of `production_speed` is not 0.00. In this case, it does not make sense to keep either of those values for further analysis, because it is impossible to determine which of the two columns is incorrect. A possible result after the data cleaning is shown in Table 2

ID	timestamp	production_speed (meter/minute)	machine_running
0	2021-12-01T12:00:00.000	1.56	True
1	2021-12-01T12:01:00.000	1.58	True
2	2021-12-01T12:02:00.000		True
3	2021-12-01T12:03:00.000		
4	2021-12-01T12:04:00.000	1.50	True
5	2021-12-01T12:05:00.000	1.49	True

Table 2: Example of IoT sensor data after cleaning

2.5 Data Cleansing

Data cleansing pursues a different approach. Incorrect data is not ignored but substituted by a representative value. This can only be done, when it is sure, that no future analysis steps depend on unaltered data. Since altering some data points is likely to change Key Performance Indicators (KPIs), which are calculated from the raw data. For example for the data point with the ID 2, several strategies could be followed. For example, the outlier value 3.50 could be replaced with the upper limit of the valid range, in this example 2.00, the value could also be replaced with the last valid value, in this example 1.58, or the value could be replaced with the average of the last n Values, for example with $\frac{1.56+1.58}{2} = 1.57$. For the data point with the ID 3 there are also different approaches. If the machine was indeed not running then it would make sense, to set the production_speed to 0.0, if short downtimes for this machine are very unlikely then the machine_running value could be set to True. A possible result after the data cleansing is shown in Table 3 [8]

ID	timestamp	production_speed (meter/minute)	machine_running
0	2021-12-01T12:00:00.000	1.56	True
1	2021-12-01T12:01:00.000	1.58	True
2	2021-12-01T12:02:00.000	2.00	True
3	2021-12-01T12:03:00.000	1.50	True
4	2021-12-01T12:04:00.000	1.50	True
5	2021-12-01T12:05:00.000	1.49	True

Table 3: Example of IoT sensor data after cleansing

3 Outlier Detection

3.1 Outlier Types

Outliers can be categorized as point outliers or subsequence outliers. Methods to detect outliers are described in section 3.2

Point Outliers

A point outlier is a single data point that strongly varies from the usual trend of the data points. [9] Examples of three point outliers are shown in Figure 1.

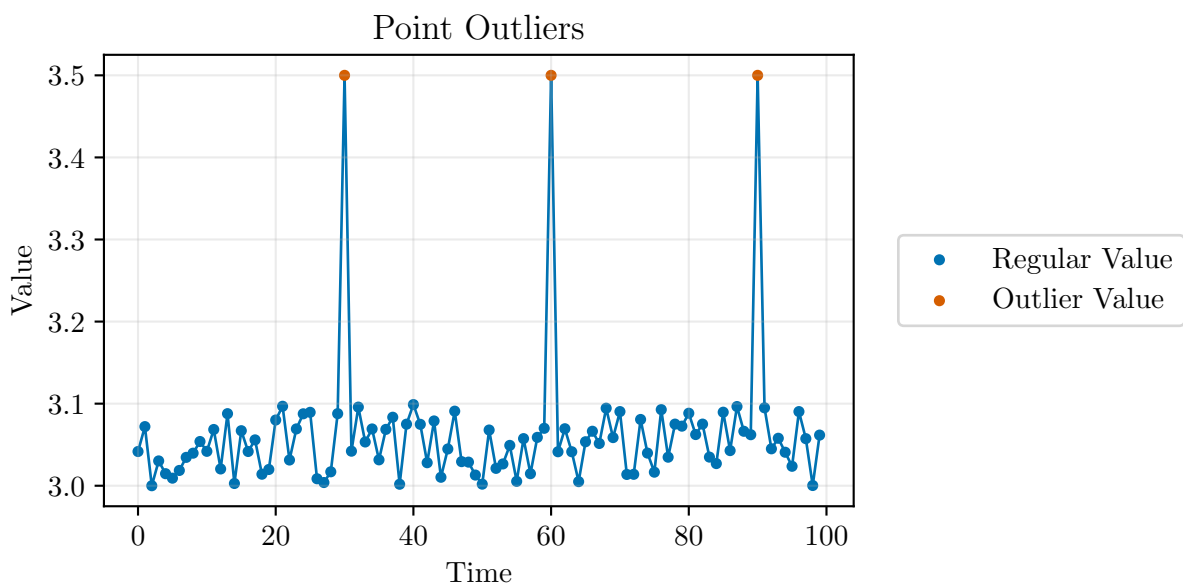


Figure 1: Examples of three point outliers

Subsequence outliers

Subsequence outliers are multiple consecutive data points that strongly vary from the usual trend of the data points. [9] In Figure 2, examples of subsequence outliers are shown.

Furthermore, outliers can be divided into local and global outliers.



Figure 2: Examples of three subsequence outliers

Local Outliers

A local outlier has a greater variance to its direct neighbouring data points (previous and next one) [9] In Figure 3 and Figure 4, examples of local outliers are shown. The first two outliers in Figure 3, have a great variance towards their direct neighbours, however, not to the values from Time 60 onwards.



Figure 3: Examples of four local point outliers



Figure 4: Examples of six local subsequence outliers

Global Outliers

Whereas a global outlier varies more in regard to all data points. Figure 1, 2 and 5 show picture examples of global outliers. [9]



Figure 5: Examples of four global point outliers

3.2 Outlier Detection Approaches

Outlier detection methods can be divided into the following groups

Statistical

For statistical outlier detection, historical data is taken to develop a model that pictures the expected behaviour of the data. An example of statistical outlier detection is the threshold-based method described in section 3.3 [10, 11]

Distance-based

For this approach, a distance metric needs to be defined, (e.g. Euclidean distance). Then each data point is compared to the data preceding it. The greater the distance between the current and previous data points, the greater the probability of an anomaly. [10–12]

Clustering

Clustering also requires a set of historical data in order to train the clustering model. Usually, the data is clustered into two clusters: normal data and anomalous data. Depending on the distance of a new data point to the "normal" and the "anomalous" cluster it is classified. [10–12]

Predictive

In this approach, a prediction model needs to be developed, which is based on previous data. The prediction of this model is then compared with the actual data point (new data, which was not used in training the model). If the actual data point differs too much from the prediction it is labelled as an anomaly. [10, 11]

Ensemble

As the word ensemble suggests, this is a collection of outlier detection methods that use a specific vote mechanism to determine whether a data point is faulty or normal. For example using the majority vote system and a statistical, distance-based and predictive method to detect outliers. If at least two methods flag a data point as an outlier the ensemble reports it as an outlier as well. If only one method reports it as an outlier the ensemble does not flag it as an anomaly. [10]

3.3 Threshold-based Outlier Detection

Threshold-based detection methods are able to identify outliers based on a given threshold τ . These methods can be described with the following formula

$$|x_t - \hat{x}_t| > \tau \quad [9] \quad (1)$$

Where x_t is the actual value and \hat{x}_t is the expected value and τ is a given threshold. Methods to calculate \hat{x}_t will be described in the following sections. Furthermore, \hat{x}_t can be calculated using the entire data series or with subsets (of equal length) of the entire data series. This means \hat{x}_t can be either calculated for the whole data series or just a segment. Depending on the sensitivity wanted for outlier detection an appropriate τ needs to be chosen. The greater τ is the fewer outliers will be detected. The smaller τ is the more outliers will be identified. [9]

Mean

$$\text{mean} = \bar{x} = \frac{1}{n} \sum_{t=0}^n x_t \quad (2)$$

Where n is the total number of samples. Using the mean as an expected value is not robust to outliers, because the mean is not as robust as the median in hindsight to outliers. To calculate the mean all data points of a series must be summed up and then divided by the number of data points.

Median

If n is odd:

$$\text{median}(x) = x_{(n+1)/2} \quad (3)$$

If n is even:

$$\text{median}(x) = \frac{x_{n/2} + x_{(n+1)/2}}{2} \quad (4)$$

Where x is a dataset of n elements ordered from smallest to largest

$(x_1 \leq x_2 \leq x_3 \leq \dots \leq x_{n-2} \leq x_{n-1} \leq x_n)$ [9] To calculate the median all values must be sorted from smallest to largest. If the number of data points is odd then the most centre data point is the Median (e.g. if the series consist of 7 values the third value is the median). If the number of data points is even, then the median is the mean of the two data points in the centre.

Median Absolute Deviation (MAD)

The Median Absolute Deviation is defined as followed:

$$MAD = \text{median}(|x_t - \text{median}(x)|) \quad (5)$$

The MAD is a more robust (regarding outliers) way to calculate the deviation of a dataset. To calculate the MAD firstly the median of the dataset must be calculated. Then the absolute difference between x_t and the median of the dataset is calculated. The Median of all differences results in the MAD [13, 14]

3.4 Outlier Detection using z-score

The z-score, also known as the standard score, is the factor of how many standard deviations a data point differs from the mean. Using the standard score to detect outliers works best when the data is distributed normally. Because then it can be assumed, that, e.g. the top and bottom 0.5% are outliers and therefore every value with $|z| > \approx 2.576$ can be classified as an outlier. A visual representation of the z-score in a normal distribution is shown in Figure 6

$$z_t = \frac{x_t - \mu}{\sigma} \quad (6)$$

Where μ is the mean of the dataset and σ is the standard deviation. When the mean and the standard deviation of the entire dataset are not known, the mean and the standard deviation of a known sample can be used. [15–17]

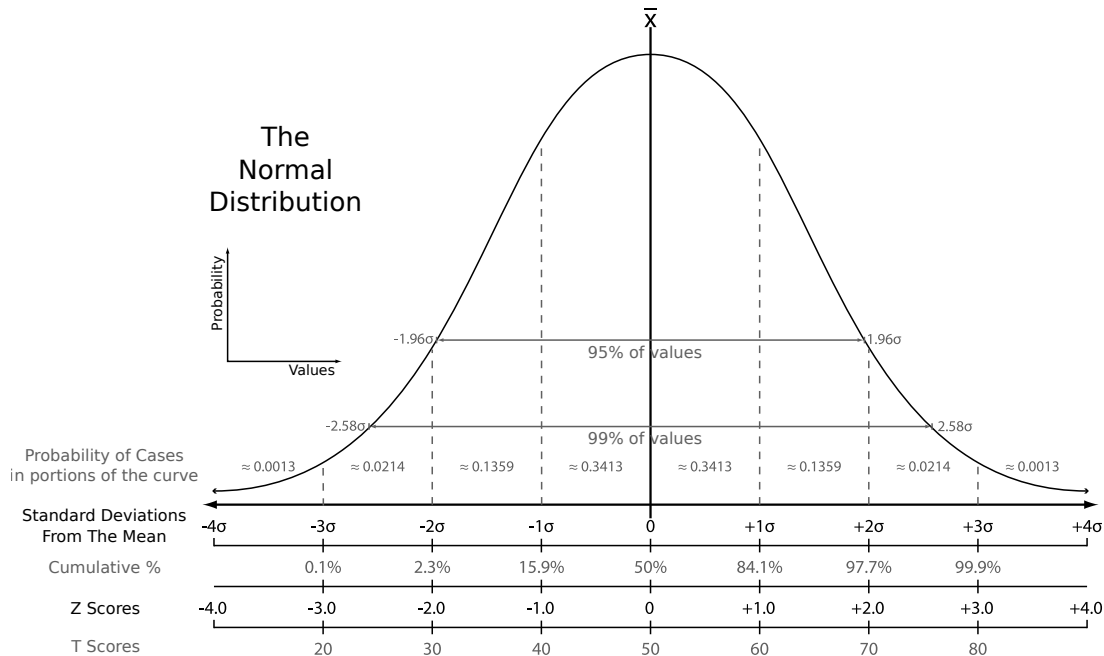


Figure 6: Standard Score in a normal distribution [18]

$$\sigma = s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (7)$$

[16,17]

To check whether the value x_t is an outlier the absolute of its z-score is compared against a threshold (τ) and if the absolute value of the z-score exceeds the threshold x_t is classified as an outlier.

$$|z_t| > \tau \quad (8)$$

3.5 Outlier Detection using modified z-score

Because the mean and standard deviation are not robust towards outliers, the z-score can be modified to use more robust metrics for the expected value and the variation of the values. To make the outlier detection with the z-score more robust, the mean can be replaced with the median and the standard deviation with the MAD or the Normalized Median Absolute Deviation (MADN). A possible formula for a modified z-score, as described in [19], could be:

$$m_t = \frac{|x_t - \text{median}(X)|}{MADN(X)} \quad (9)$$

Where the formula for the MADN is:

$$MADN(X) = \frac{MAD(X)}{0.6745} \quad (10)$$

The constant value 0.6745 is the 75th percentile of a standard normal distribution, which is equal to the MAD of a standard normal distribution with $\sigma = 1$. [19]

The classification of outliers using the modified z-score (m_t) is the same for the regular standard score described in section 3.4:

$$|m_t| > \tau \quad (11)$$

4 Outlier Detection based on Water Level Data

FloodAlert (or “Pegelalarm” in German) is a service, that provides water levels from about 30,000 measurement stations. It provides notifications to registered individuals when a certain water level threshold is reached. Thus warning people about possible floods in their area. [1]

Unfortunately, sometimes an incorrect water level is reported by the sensors, which are not maintained by the team of FloodAlert. FloodAlert just fetches the data from different maintainers of sensors. To reduce the frequency of false alarms, outliers should be identified and flagged. The theoretical parts of the previous sections are applied and tested on real-world data.

4.1 What is the Goal?

The ideal goal of FloodAlert would be to have a general model, that works on a variety of rivers. The model should be able to predict the probability of being an outlier for each water level value, without knowing future values. Thereby making it possible to classify incoming data points in real-time, without needing future data points to make a prediction. Furthermore, it would be ideal if the outlier detection method is able to adjust its parameters automatically for each river/water level measurement location. Because different rivers have different fluctuations in water level. Therefore each water level measurement station needs individual parameters. For example, one river regularly has an increase and decrease of 10 centimetres whereas for another river 10 cm of water level variance from one data point to the next is definitely an outlier.

Another approach would be to also take future values into consideration when predicting the probability of a value being an outlier. This method is probably easier and likely leads to better results. However, a method which does not take future values into consideration would be more useful for FloodAlert, since the values could be classified immediately.

4.2 How to retrieve the Data (Description of the API)

To make the access to the API easier SOBOS GmbH developed a Python wrapper which returns the requested data as a Pandas dataframe. [20] The code for the Python wrapper is available on GitHub [21]: https://github.com/SOBOS-GmbH/pegelalarm_public_pas_doc [22]

Listing 1 Request body to get API key

```
1 {  
2     "username": "myUsername",  
3     "password": "myPassword"  
4 }
```

In order to request data, an API key needs to be requested using credentials. To request the API key a POST request needs to be sent to this endpoint <https://api.pegelalarm.at/api/login>, where the request body contains the users' credentials as shown in Listing 1. To generate the actual API key, the key from the response needs to be hashed using the Keyed-Hashing for Message Authentication (HMAC) Algorithm. This is done using Python's built-in "hmac" module [23]. Afterwards, the HMAC is byte64 encoded, so it can be sent in the "X-AUTH-TOKEN" header field.

To get the unique identifier for a specific measurement station the list endpoint can be used. Which accepts three optional query parameters (qStationName - station name, qWater - water name and commonid - the unique identifier of a station) and returns a list of matching stations with metadata like coordinates, country or last water level. E.g. to retrieve the identifier of the Danube station located in Linz the URL is the following: <https://api.pegelalarm.at/api/station/1.1/list?qStationName=Linz&qWater=Donau>. Keep in mind in order to retrieve any data the header value "X-AUTH-TOKEN" must be set. An example response to this request is shown in Listing 2. Where the unique identifier is the "commonid" in line 9.

To retrieve historical data the history endpoint can be used. The request URL has the following structure: <https://api.pegelalarm.at/api/station/1.1/<unit>/<commonid>/history?<parameters>>. The unit can either be "height" or "flow". For this thesis, only height data was used. The following parameters can be set:

- **loadStartDate**: The start timestamp of the queried data.
- **loadEndDate**: The end timestamp of the queried data.
- **granularity**: The granularity of the response.

Possible values for the granularity are: "raw" (only available for the last 3 months of data), "hour", "day", "month", "year" or "era" (era returns one value for a given time) When requesting aggregated data (anything other than "raw"), the maximum value of the timespan is used. The timestamps are in the following format: "<DD>.<MM>.<YYYY>T<HH>:<MM>:<SS><+-UTCOffset>", where the '+' is URL encoded to "%2B". E.g.: "31.03.2022T13:35:40%2B0200". If no parameters are provided the API returns the last few data points.

An example request would be: <https://api.pegelalarm.at/api/station/1.1/height/207068-at/history?loadStartDate=01.03.2022T13:35:40%2B0200&loadEndDate=01.03.2022T18:00:00%2B0200&granularity=hour> The result of this request is shown in Listing 3.

Listing 2 Example response of the list endpoint

```
1 {
2     "status": {
3         "code": 200
4     },
5     "payload": {
6         "stations": [
7             {
8                 "name": "Donau / Linz / at",
9                 "commonid": "207068-at",
10                "country": "Österreich",
11                "stationName": "Linz",
12                "water": "Donau",
13                "region": "Oberösterreich",
14                "latitude": 48.306915712282,
15                "longitude": 14.284689597541,
16                "positionKm": 2135.17,
17                "altitudeM": 247.74,
18                "defaultWarnValueCm": 550.0,
19                "defaultAlarmValueCm": 630.0,
20                "data": [
21                    {
22                        "type": "height in cm",
23                        "value": 358.0,
24                        "requestDate": "19.04.2022T14:59:51+0200",
25                        "sourceDate": "19.04.2022T14:45:00+0200"
26                    }
27                ],
28                "trend": 10,
29                "situation": 10,
30                "visibility": "PUBLIC",
31                "stationType": "surfacewater"
32            }
33        ]
34    }
35 }
```

Listing 3 Example response of historical water level data for one station

```
1 {
2   "status": {
3     "code": 200
4   },
5   "payload": {
6     "history": [
7       {
8         "value": 360.0,
9         "sourceDate": "01.03.2022T13:00:00+0100"
10      },
11      {
12        "value": 360.0,
13        "sourceDate": "01.03.2022T14:00:00+0100"
14      },
15      {
16        "value": 359.0,
17        "sourceDate": "01.03.2022T15:00:00+0100"
18      },
19      {
20        "value": 361.0,
21        "sourceDate": "01.03.2022T16:00:00+0100"
22      },
23      {
24        "value": 362.0,
25        "sourceDate": "01.03.2022T17:00:00+0100"
26      }
27    ]
28  }
29 }
```

4.3 Overview of the Data

To test and tune the outlier detection methods the entire history of data was used. The granularity of the data usually is hourly up to the last three months of data. There are a few exceptions, where the granularity is not hourly, because either one or a few data points are missing. For the last three months, the granularity is “raw”, which is different for each station. Additionally, the base (where the sensor reports water level equals zero) also differs per station. Some stations have an arbitrary zero level whereas some stations use the sea level as a base. Thus the water level does not actually reflect the actual height, or depth, of the water. Instead, it represents the relative water level to a specific height.

The following water level measurement stations were used to test the outlier detection methods:

- **Station Aghacashlaun, Aghacashlaun (36022-ie)**

This station, on the river Aghacashlaun and in the area of Aghacashlaun, is located in the north of Ireland. The coordinates of this station are 54.03647, -7.94672 (Latitude (Lat), Longitude (Long)). [24]

- **Station Murg, Frauenfeld (2386-ch)**

The station, on the river Murg and in the area of Frauenfeld, is located in northern Switzerland. The coordinates of this station are 47.56852, 8.89432 (Lat, Long). [25]

- **Station Sieg, Betzdorf (2720050000-de)**

The station, on the river Sieg and in the area of Betzdorf, is located in western Germany. The coordinates of this station are 50.79332, 7.86390 (Lat, Long). [26]

- **Station Losse, Helsa (42960105-de)**

The station, on the river Losse and in the area of Helsa, is located in central Germany. The coordinates of this station are 51.25574, 9.68537 (Lat, Long). [27]

- **Station Crana, Tullyarvan (39003-ie)**

The station, on the river Crana and in the area of Tullyarvan, is located in the north of Ireland. The coordinates of this station are 55.14356, -7.45239 (Lat, Long). [28]

4.4 Manual Classification of Outliers

In order to speed up the manual labelling of outliers, a program was written. The program is a Plotly Dash [29] web application which displays the water level data as a scatter chart. By clicking the data points in the chart the user is able to toggle the data point as an outlier or back to a regular value. In Listing 9 the source code of the Dash application is shown. In Figure 7 the Website of the Python app is shown. The red dots between October 25 and October 27 are already classified as outliers. Below the chart a range slider is located, to move the zoomed-in view horizontally. The refresh button on the left side refreshes the graph, thus updating the

colour and label of previously selected outliers. Furthermore, it saves the data to a parquet file. The upper and lower limit of the y-axis also gets updated, when pressing the refresh button. The limits are automatically set to the lowest and highest regular value. This was implemented because some datasets had outliers with a huge difference from the regular data. Without the automatic scaling of the y-axis, detecting other outliers with a smaller difference was not possible.

Manual Outlier Selection

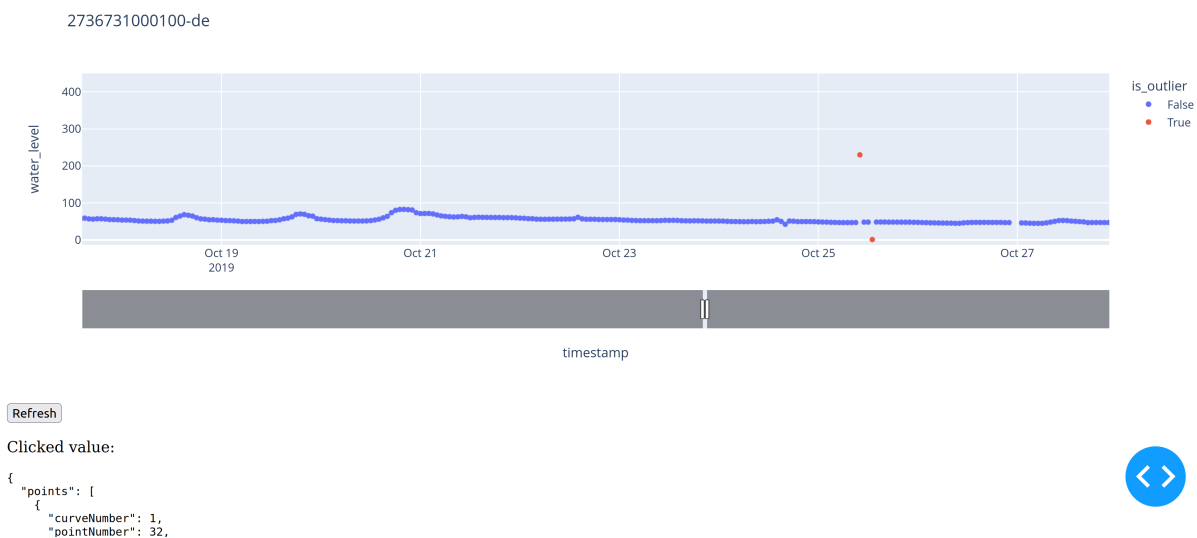


Figure 7: Dash Webapp to classify outliers

4.5 Explorative Data Analysis

4.5.1 Aghacashlaun Station

Table 4 shows seven number summary for all data points (outliers and regular values). The column “water_level_diff” is the difference in water level between two consecutive data points. The column “timedelta” is the time difference between two consecutive data points in hours. Table 5 and Table 6 picture seven number summaries for regular and only outlier values. The standard deviation is clearly higher for outlier values when compared to the standard deviation of regular values.

In Figure 8 the class distribution is shown. The classes are not balanced, there are far more regular values than outliers. This needs to be kept in mind when choosing a performance metric.

	water_level	water_level_diff	timedelta
count	27189.000000	27188.000000	27188.000000
mean	36.009754	-0.000077	0.809990
std	14.716286	7.220175	1.256409
min	0.000000	-107.500000	0.000000
25%	26.400000	-0.300000	1.000000
50%	31.700000	-0.100000	1.000000
75%	40.300000	0.000000	1.000000
max	190.000000	94.800000	178.000000

Table 4: Seven number summary of Aghacashlaun - Aghacashlaun (all values)

	water_level	water_level_diff	timedelta
count	26544.000000	26543.000000	26543.000000
mean	35.381205	-0.411276	0.808311
std	13.763332	5.165770	1.268532
min	20.000000	-107.500000	0.000000
25%	26.400000	-0.300000	1.000000
50%	31.400000	-0.100000	1.000000
75%	39.800000	0.000000	1.000000
max	151.600000	51.900000	178.000000

Table 5: Seven number summary of Aghacashlaun - Aghacashlaun (regular values)

	water_level	water_level_diff	timedelta
count	645.000000	645.000000	645.000000
mean	61.876744	16.921550	0.879070
std	25.476859	28.411003	0.560844
min	0.000000	-90.000000	0.000000
25%	40.000000	-10.000000	1.000000
50%	60.000000	12.500000	1.000000
75%	80.000000	39.400000	1.000000
max	190.000000	94.800000	9.000000

Table 6: Seven number summary of Aghacashlaun - Aghacashlaun (outlier values)

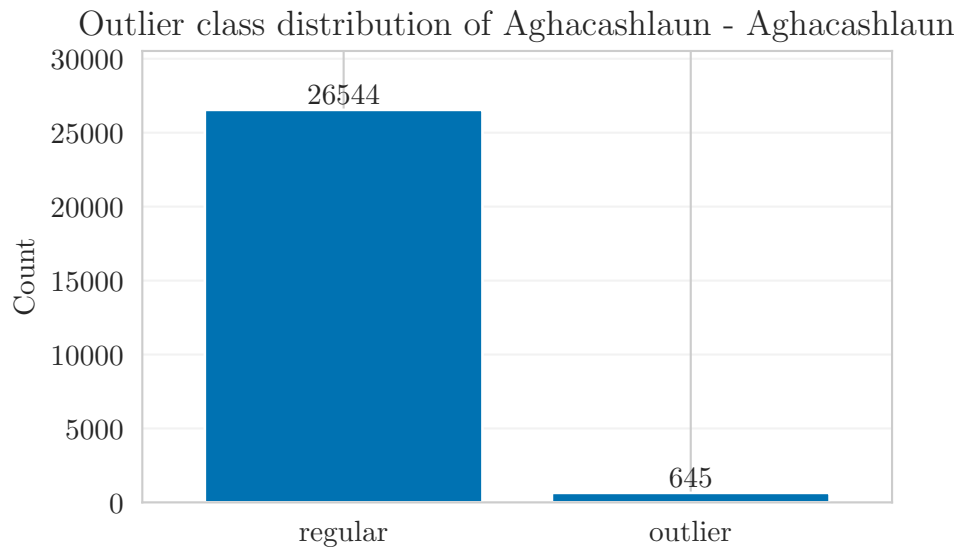


Figure 8: Class distribution of Aghacashlaun - Aghacashlaun

Figure 9 shows boxplots for all values, regular values and outlier values. The boxplots show, that the upper and lower whisker boundaries cannot be used to classify outliers reliably. Since there are many regular values that exceed the upper limit of the whisker, which is 1.5 times the InterQuartile Range (IQR).

In Figure 10 histograms of the different water levels are shown. The y-axis is scaled logarithmically, so less common values are also visible. The histogram shows, that values between 20 and 40 are most common and the higher the water level is the less common it becomes. Figure 11 is quite similar, but instead of the water level, the difference in water level between two consecutive data points is used.

Figure 12 shows the histogram of the time difference, in hours, between each data point. The largest gap in the data is 178 hours, which is about one week.

In Figure 13 the water level over time is shown. This chart shows, that the water level usually oscillates between 20 and 100. Figure 14 displays an example of subsequence outliers. Within one hour the water level rises by double the previous value. This is very unlikely, thus these values were classified as outliers. The reason, for multiple consecutive outliers, with decreasing water levels over time, could be, that the station does not report the actual water level but a moving average over the last few recorded values. This would at least explain the pattern of outliers shown in Figure 14.

The code for the explorative data analysis can be found in the appendix (Listing 11).

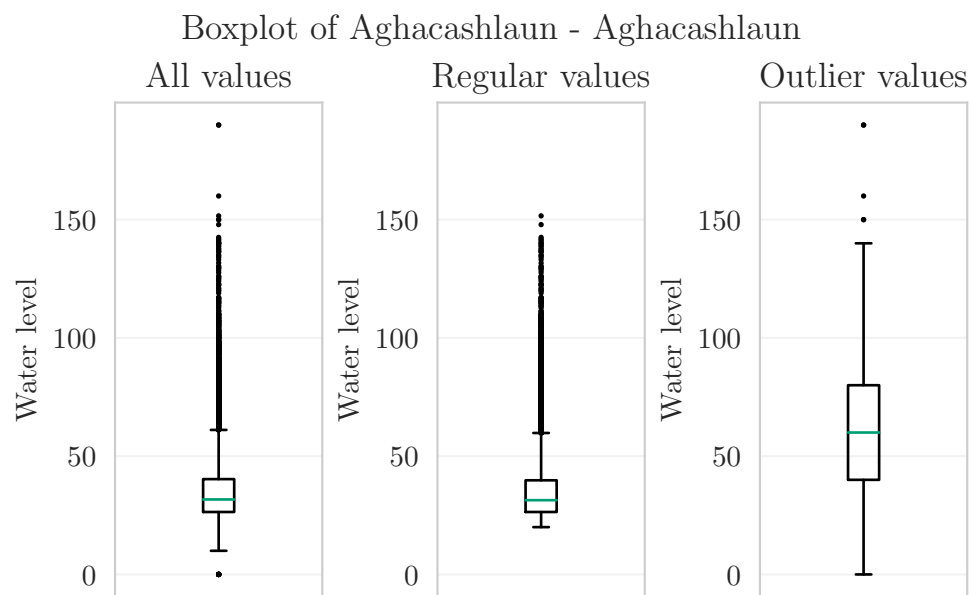
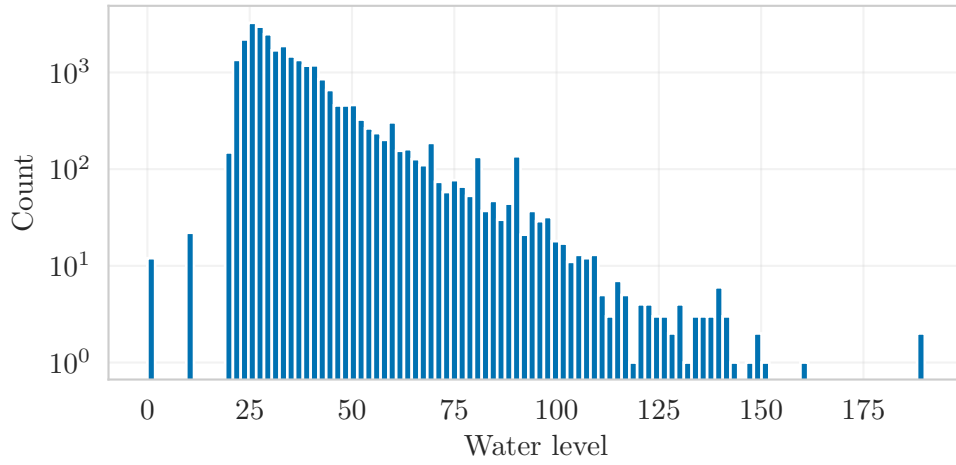
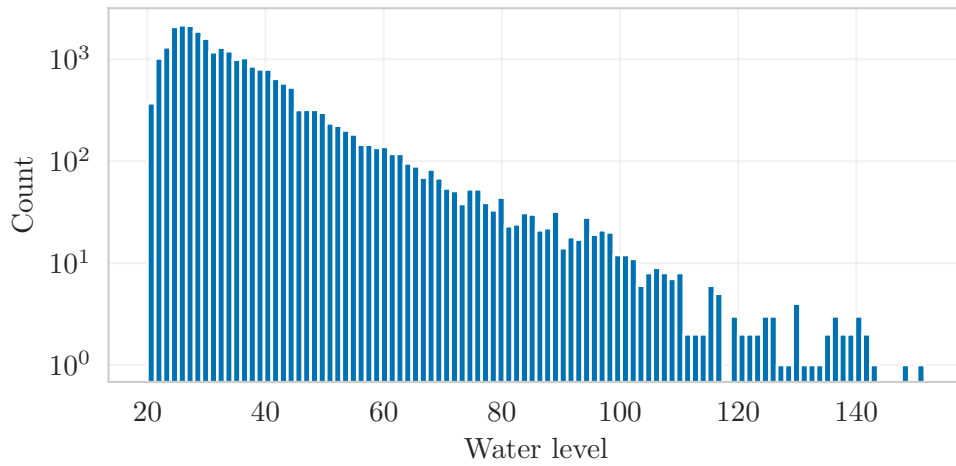


Figure 9: Boxplot of Aghacashlaun - Aghacashlaun

Histogram of the water Level of Aghacashlaun - Aghacashlaun
All Values



Regular values



Outlier values

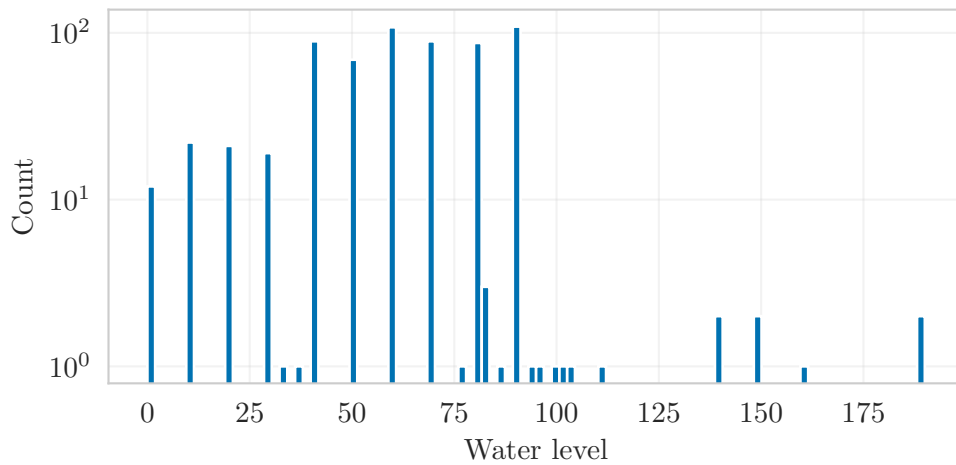


Figure 10: Histogram of the water level of Aghacashlaun - Aghacashlaun

Histogram of the water level delta of Aghacashlaun - Aghacashlaun

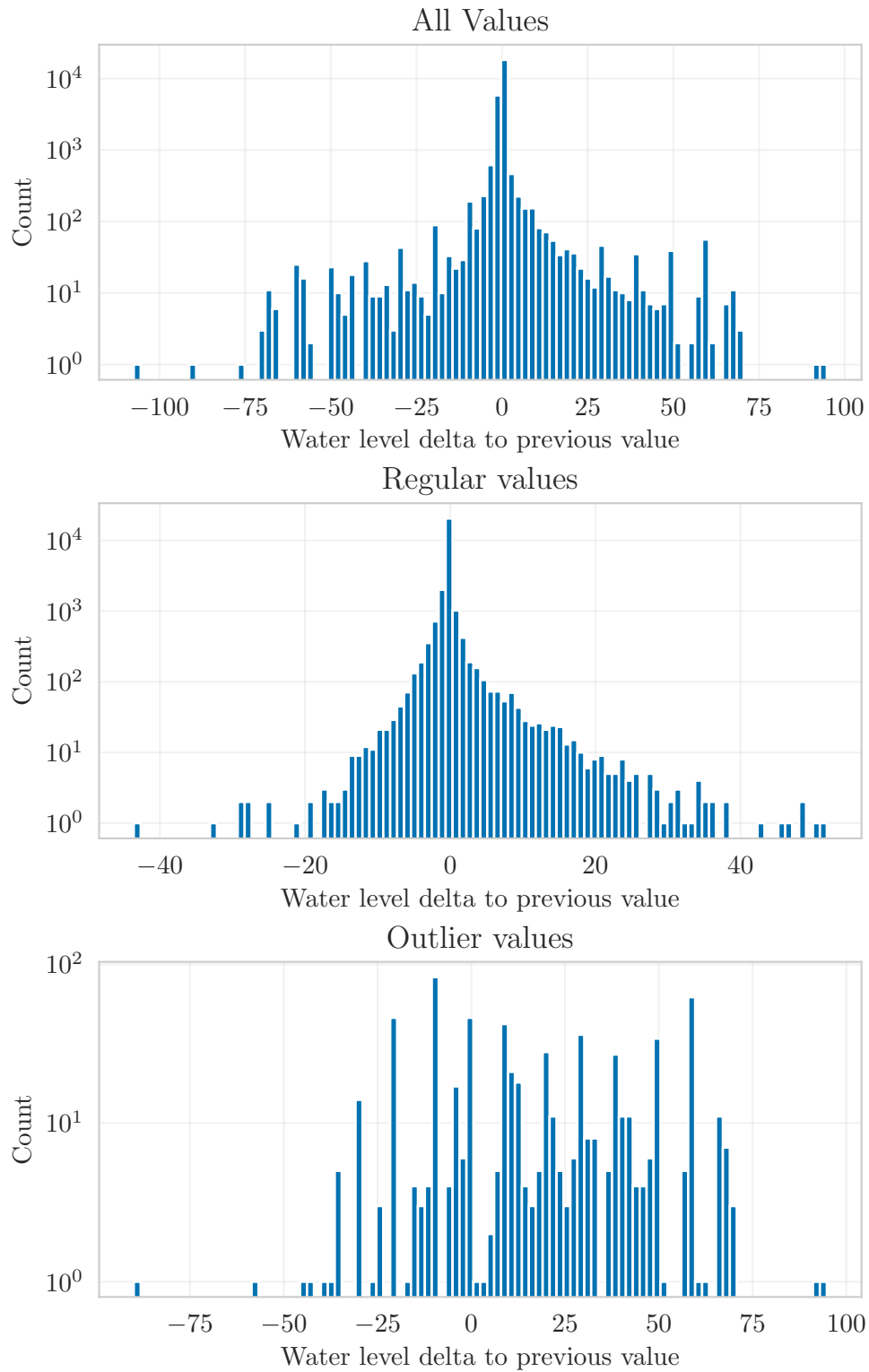
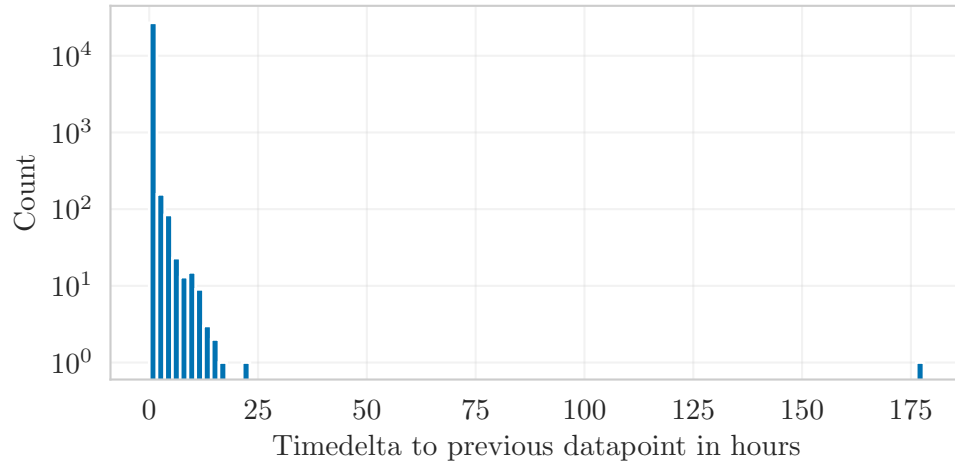
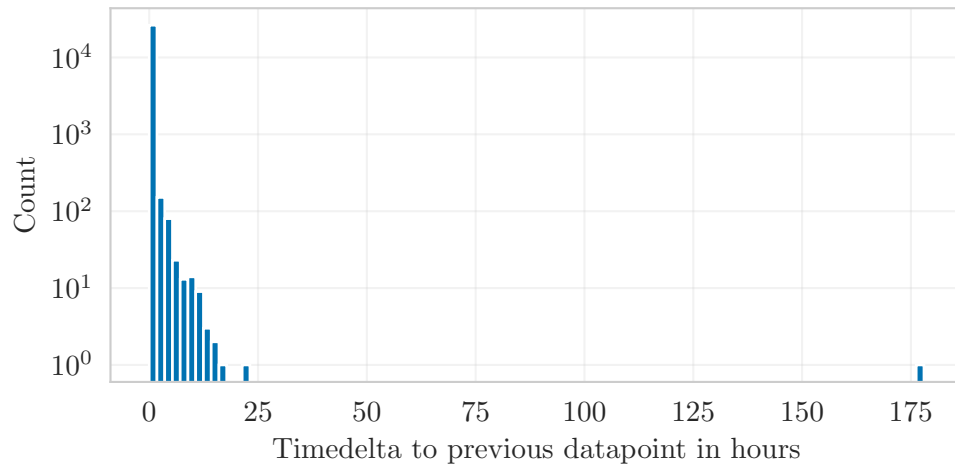


Figure 11: Histogram of the water level delta of Aghacashlaun - Aghacashlaun

Histogram of the timedelta between values of Aghacashlaun - Aghacashlaun
All Values



Regular values



Outlier values

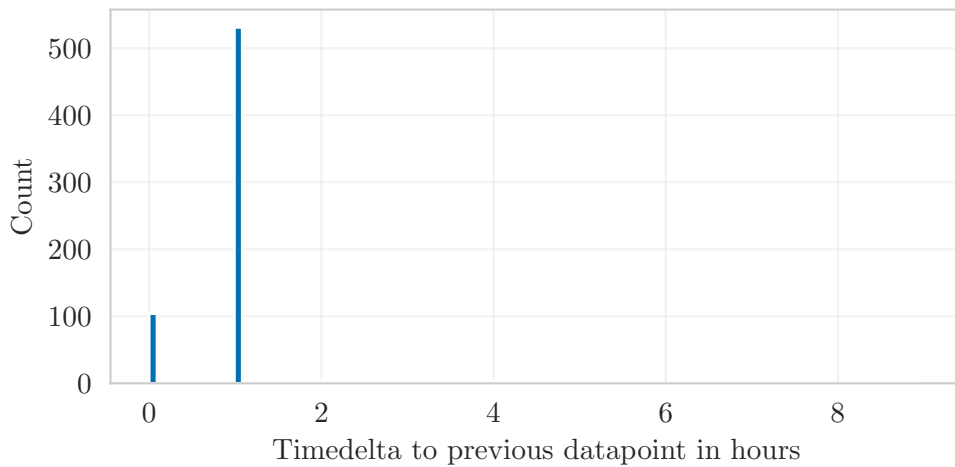


Figure 12: Histogram of the time delta (in hours) of Aghacashlaun - Aghacashlaun

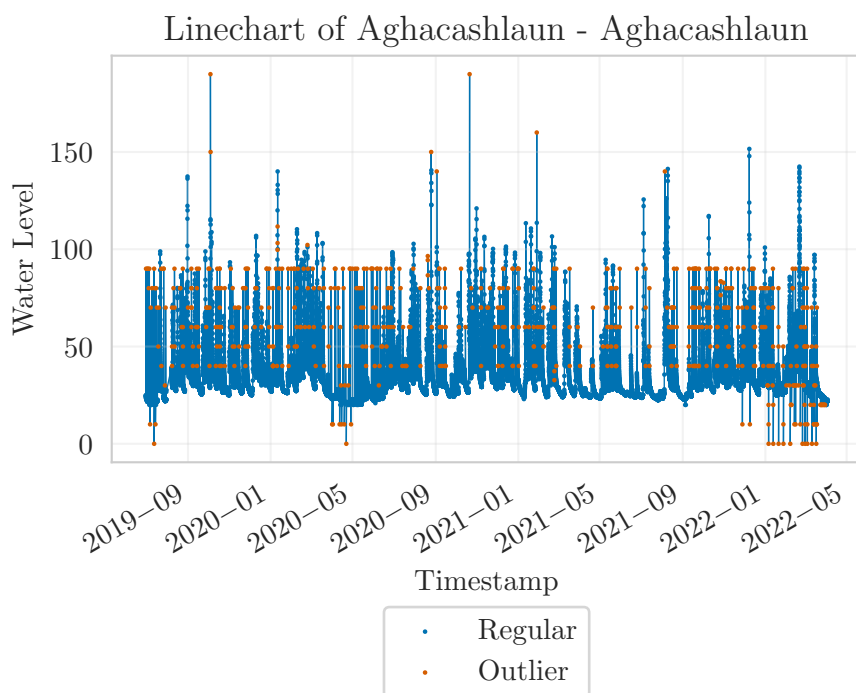


Figure 13: Line chart of Aghacashlaun - Aghacashlaun (whole data)

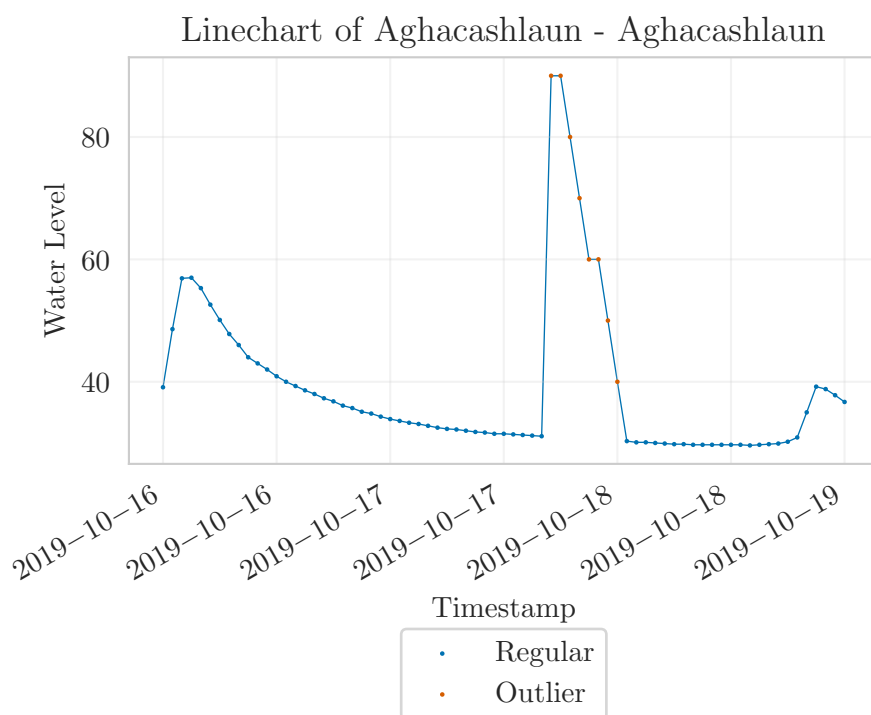


Figure 14: Line chart of Aghacashlaun - Aghacashlaun (2019-10-16 - 2019-10-19)

4.6 Outlier Detection performance Metrics

4.6.1 Confusion Matrix

To compare the performance of two classification methods a confusion matrix, shown in Table 7, can be used to provide an overview. A confusion matrix consists of the following elements (explained on the basis of outlier detection):

- True Positive (TP): actual class: outlier, predicted class: outlier
- False Negative (FN): actual class: outlier, predicted class: regular
- True Negative (TN): actual class: regular, predicted class: regular
- False Positive (FP): actual class: regular, predicted class: outlier

		Predicted label	
		Positive (P)	Negative (N)
True/ actual label	Positive	True positive	False negative
	Negative	False positive	True negative

Table 7: Example of a binary confusion matrix

The green cells represent the correctly classified values (TP & TN) and the red cells represent the incorrectly classified values (FN & FP). The confusion matrix is the basis of many more performance metrics.

4.6.2 Accuracy

To make the comparison easier, the confusion matrix can be aggregated into one single value as a metric. There are numerous ways how this single metric can be calculated. One of the most common and basic methods is the Accuracy:

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{TP + TN}{TP + FN + TP + FP} = \frac{TP + TN}{P + N} \quad (12)$$

The accuracy provides information on which percentage of values are correctly classified. However, if the classes are not equally distributed, the accuracy is a bad metric. This is especially true for heavily imbalanced classes. For example, if the dataset has 1 000 000 positive values and 10 negative, when optimizing for a high accuracy it can happen that the classifier predicts every value as positive and still achieves a very high accuracy:

$$Accuracy = \frac{1\,000\,000 + 0}{1\,000\,000 + 10} = 0.99999 \quad (13)$$

Thus for measuring the performance of outliers, the accuracy is not a suitable metric.

4.6.3 Precision and Recall

Figure 15 shows a visual comparison between precision and recall. The precision provides information about what percentage of positive predictions were actually correct. Whereas the recall provides information about what percentage of all positive values were actually predicted as such. Precision and Recall are used to calculate the F-Score described in subsection 4.6.4.

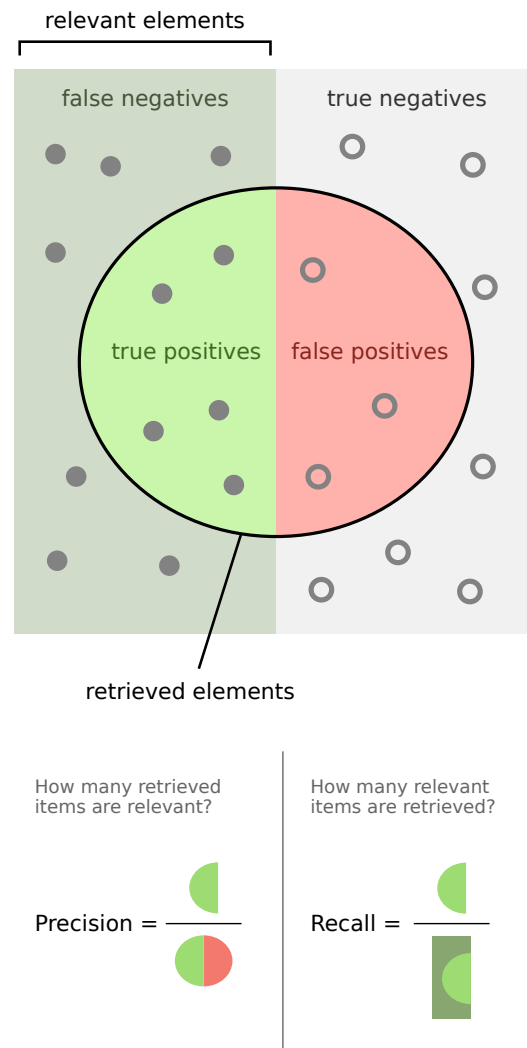


Figure 15: Precision vs Recall [30]

Precision

$$Precision = \frac{TP}{TP + FP} \quad (14)$$

Recall

$$Recall = \frac{TP}{TP + FN} \quad (15)$$

4.6.4 F-score

The $F_1 - score$ is the harmonic mean of the precision and recall.

$$F_1 - score = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}} = \frac{2}{\frac{Precision+Recall}{Precision \cdot Recall}} = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall} \quad (16)$$

[31, 32]

The $F_\beta - score$ is more sophisticated compared to the $F_1 - score$ where an additional parameter (β) needs to be chosen, which is used to weight the precision.

$$F_\beta - score = \frac{(\beta^2 + 1) \cdot Precision \cdot Recall}{\beta^2 \cdot Precision + Recall} \quad (17)$$

[32, 33]

The $F_1 - score$ uses 1 for β so the precision and the recall are equally weighted.

4.7 Implementation of different Outlier Detection Approaches

The input data structure for the different approaches is a pandas DataFrame. [20] The structure of the DataFrame is equal for all measurement stations. It contains the following columns:

- `water_level`: the water level of the river in cm (datatype: float64)
- `timestamp`: the date and time of the data point in UTC (datatype: datetime64[ns, UTC])
- `is_outlier`: true or false depending on if the value is an outlier (datatype: bool)

To test the performance of the models the column `is_outlier` contains the ground truth, which was manually labelled. An example of the input data can be seen in Table 8.

water_level	timestamp	is_outlier
24.9	2019-06-30 15:00:00+00:00	False
24.9	2019-06-30 16:00:00+00:00	False
24.8	2019-06-30 17:00:00+00:00	False
24.4	2019-06-30 18:00:00+00:00	False
24.4	2019-06-30 19:00:00+00:00	False
24.3	2019-06-30 20:00:00+00:00	False
24.1	2019-06-30 21:00:00+00:00	False
24.0	2019-06-30 22:00:00+00:00	False
90.0	2019-06-30 23:00:00+00:00	True

Table 8: First 9 values of Aghacashlaun - Aghacashlaun

To speed up the process of calculating the result of the outlier detection method ($|x_t - \hat{x}_t|$, the z-score or the modified z-score) is all done in one step and in the second step the result of this calculation is compared against a range of thresholds.

4.7.1 Mean Threshold

Calculating \hat{x}_t using the mean for a pandas DataFrame [20] is quite straightforward and is shown in Listing 4. To reduce the amount of duplicate code, the imports are only included in Listing 4.

Listing 4 The first step of classifying outliers using the mean

```
1 from typing import Union
2
3 import numpy as np
4 import pandas as pd
5
6
7 def mean_outlier_detection(input_df: pd.DataFrame,
8                             window: Union[int, None],
9                             center_window: bool):
10     """
11     Detects outliers in a dataframe using a (moving) average.
12     :param input_df: the input dataframe where the values are stored
13                       in the column water_level
14     :param window: the size of the window, None if no window should
15                     be used
16     :param center_window: whether the window should be centred or not
17     :return: a copy of the input dataframe where the column result
18              should be compared to a threshold to detect outliers
19     """
20     od_df = input_df.copy()
21     if window is None:
22         od_df['x_hat'] = od_df['water_level'].mean()
23     else:
24         od_df['x_hat'] = \
25             od_df['water_level'].rolling(window=window,
26                                         center=center_window,
27                                         min_periods=1).mean()
28     od_df['result'] = np.abs(od_df['water_level'] - od_df['x_hat'])
29     return od_df
```

4.7.2 Median Threshold

The calculation of the median is quite similar to the mean. It is shown in Listing 5.

Listing 5 The first step of classifying outliers using the median

```
1 def median_outlier_detection(input_df: pd.DataFrame,
2                             window: Union[int, None],
3                             center_window: bool):
4     """
5     Detects outliers in a dataframe using a (moving) average.
6     :param input_df: the input dataframe where the values are stored
7                     in the column water_level
8     :param window: the size of the window, None if no window should
9                     be used
10    :param center_window: whether the window should be centred or not
11    :return: a copy of the input dataframe where the column result
12            should be compared to a threshold to detect outliers
13    """
14    od_df = input_df.copy()
15    if window is None:
16        od_df['x_hat'] = od_df['water_level'].median()
17    else:
18        od_df['x_hat'] = \
19            od_df['water_level'].rolling(window=window,
20                                       center=center_window,
21                                       min_periods=1).median()
22    od_df['result'] = np.abs(od_df['water_level'] - od_df['x_hat'])
23    return od_df
```

4.7.3 MAD Threshold

In Listing 6 an example of an implementation for calculating the MAD is provided.

Listing 6 The first step of classifying outliers using the MAD

```
1 def mad_outlier_detection(input_df: pd.DataFrame,
2                           window: Union[int, None],
3                           center_window: bool):
4     """
5     Detects outliers in a dataframe using a (moving) average.
6     :param input_df: the input dataframe where the values are stored
7                     in the column water_level
8     :param window: the size of the window, None if no window should
9                     be used
10    :param center_window: whether the window should be centred or not
11    :return: a copy of the input dataframe where the column result
12            should be compared to a threshold to detect outliers
13    """
14    od_df = input_df.copy()
15    if window is None:
16        od_df['x_hat'] = np.median(
17            np.abs(od_df['water_level'] - np.median(
18                od_df['water_level'])))
19    else:
20        od_df['x_hat'] = \
21            od_df['water_level'].rolling(window=window,
22                                         center=center_window,
23                                         min_periods=1).apply(
24                lambda x: np.median(np.abs(x - np.median(x))))
25    od_df['result'] = np.abs(od_df['water_level'] - od_df['x_hat'])
26    return od_df
```

4.7.4 Z-Score

Listing 7 shows how to calculate the z-score in Python.

Listing 7 The first step of classifying outliers using the z-score

```
1 def z_score_outlier_detection(input_df: pd.DataFrame,
2                               window: Union[int, None],
3                               center_window: bool):
4     """
5     Detects outliers in a dataframe using a (moving) average.
6     :param input_df: the input dataframe where the values are stored
7                       in the column water_level
8     :param window: the size of the window, None if no window should
9                     be used
10    :param center_window: whether the window should be centred or not
11    :return: a copy of the input dataframe where the column result
12             should be compared to a threshold to detect outliers
13    """
14    od_df = input_df.copy()
15    if window is None:
16        od_df['mean'] = od_df['water_level'].mean()
17        od_df['std'] = od_df['water_level'].std()
18    else:
19        od_df['mean'] = \
20            od_df['water_level'].rolling(window=window,
21                                         center=center_window,
22                                         min_periods=1).mean()
23        od_df['std'] = \
24            od_df['water_level'].rolling(window=window,
25                                         center=center_window,
26                                         min_periods=1).std()
27    od_df['result'] = \
28        (od_df['water_level'] - od_df['mean']).divide(od_df['std'])
29    return od_df
```

4.7.5 Modified z-score

The code for calculating the result of the modified z-score (using the MADN) is shown in Listing 8.

Listing 8 The first step of classifying outliers using the modified z-score (MADN-z-score)

```
1 def madn_z_score_outlier_detection(input_df: pd.DataFrame,
2                                   window: Union[int, None],
3                                   center_window: bool):
4     """
5     Detects outliers in a dataframe using a (moving) average.
6     :param input_df: the input dataframe where the values are stored
7                       in the column water_level
8     :param window: the size of the window, None if no window should
9                     be used
10    :param center_window: whether the window should be centred or not
11    :return: a copy of the input dataframe where the column result
12             should be compared to a threshold to detect outliers
13    """
14    od_df = input_df.copy()
15    if window is None:
16        od_df['median'] = od_df['water_level'].median()
17        od_df['mad'] = np.median(
18            np.abs(od_df['water_level'] - od_df['median']))
19        od_df['madn'] = od_df['mad'] / 0.6745
20    else:
21        od_df['median'] = \
22            od_df['water_level'].rolling(window=window,
23                                         min_periods=1,
24                                         center=center_window).median()
25        od_df['mad'] = \
26            od_df['water_level'].rolling(window=window,
27                                         min_periods=1,
28                                         center=center_window).apply(
29                lambda x: np.median(np.abs(x - np.median(x))))
30        od_df['madn'] = od_df['mad'] / 0.6745
31    od_df['result'] = \
32        (od_df['water_level'] - od_df['median']).abs() \
33        .divide(od_df['madn'])
34    return od_df
```

4.7.6 Preprocessing the Data

In the preprocessing step outliers which extremely vary from the usual trend of the other values were removed. This was done by defining upper and lower limits for the data. If a value is not inside this limit it is removed from the dataset. The thought behind preprocessing the data was to increase the model performance by removing extreme outlier values from the beginning. The disadvantage of the preprocessing step is, that two limits need to be defined for each measurement station. The limits need to be chosen carefully. On the one hand, if the valid value range is too large, no extreme outliers are removed and the preprocessing is useless, on the other hand, if the valid value range is too small the preprocessing might remove valid values, which would indicate a possible flood. Thus it is also a good idea to regularly check and maybe update those limits.

The Python code for preprocessing the data can be found in the appendix in Listing 12.

4.7.7 Finding Parameters

To find the best parameters a grid search was used. The first step was to calculate the results for different methods, window sizes and types (centred and non-centred). For each unique parameter combination, the result dataframe was stored as a file, where the filename provided information about the parameters. An example for a filename would be “12_cw_median.parquet”, where “12” is the window size, “cw” stands for center window (“nocw” is for no center window) and “median” is the method used. This was done for every unique combination of parameters. The possible values for each parameter, that was used, are listed below.

- **normalized**: yes, no
- **preprocessed**: yes, no
- **window size**: None, 2-51 (in steps of one)
- **centered window**: yes, no
- **method**: mean, median, mad, z-score, modified z-score (madn-z-score)
- **common-id**: “36022-ie”, “39003-ie”, “2386-ch”, “42960105-de”, “2720050000-de”

Due to the large number of unique combinations of parameters the size of the resulting files was quite large (multiple Gigabytes). The grid search was executed in parallel using multiprocessing, in order to speed up the process.

After the results of each parameter were calculated, for each file a range of thresholds was tested and the performance (confusion matrix and $F_1 - score$) of this threshold and the parameters used were saved to a DataFrame. This was also conducted in parallel to decrease the runtime.

The code of the grid search for different parameters can be found in the appendix in Listing 14 and the code to get the prediction for different thresholds in Listing 15.

4.8 Compare different Outlier Detection Approaches

As expected, using a centred window yields the best $F_1 - score$. Furthermore, the methods with the best performance are the median and the MADN z-score. While the median performed best with smaller window sizes (5-7) the modified z-score performed best with larger window sizes (22-28). Using the MAD to calculate the \hat{x}_t yields the worst performance.

Table 9 shows the top 5 average $F_1 - scores$ of all stations tested. When using one set of parameters for different measurement stations the median performed best. Using the mean for \hat{x}_t delivered the second-best performance when comparing shared parameters between all tested stations. However the average $F_1 - score$ of the mean is only 0.47794, it was reached by using a centred window of size 5, a threshold of 18.87960 and not normalizing the data. Due to the heterogeneity of the fluctuations of the water levels for the different stations, it is not recommended to use the same parameters for different stations. Similar or equal parameters should only be used when the water levels of the two measurement stations behave similarly. Additionally, when using the same parameters for different stations the performance of those parameters should be looked at per station and not as an average of the $F_1 - scores$. This hinders the fact that one model performs perfectly ($F_1 - score = 1$) and the other very poorly ($F_1 - score = 0.5$).

window_size	center_window	normalized	threshold	model_type	average_f1_score
3.0	True	False	6.628763	median	0.725722
3.0	True	False	6.959866	median	0.725293
3.0	True	False	6.297659	median	0.723374
3.0	True	False	8.284281	median	0.719480
3.0	True	False	8.615385	median	0.719397

Table 9: Best parameters of the average F1-score of all stations tested

Using the preprocessing (described in subsection 4.7.6) usually resulted in lower performance, when comparing the highest $F_1 - scores$ for both datasets (preprocessed and not preprocessed). The reason for this is, that the outliers removed by the two-sided filter were detected as outliers anyway. So just the number of TPs was reduced. This resulted in a slightly lower overall performance since the outliers were completely removed from the dataset in the preprocessing step.

Mean

Figure 16 shows the best performing outlier detection using mean. The F_1 -score is about 0.68, the *Precision* is about 0.73 and the *Recall* is about 0.65.

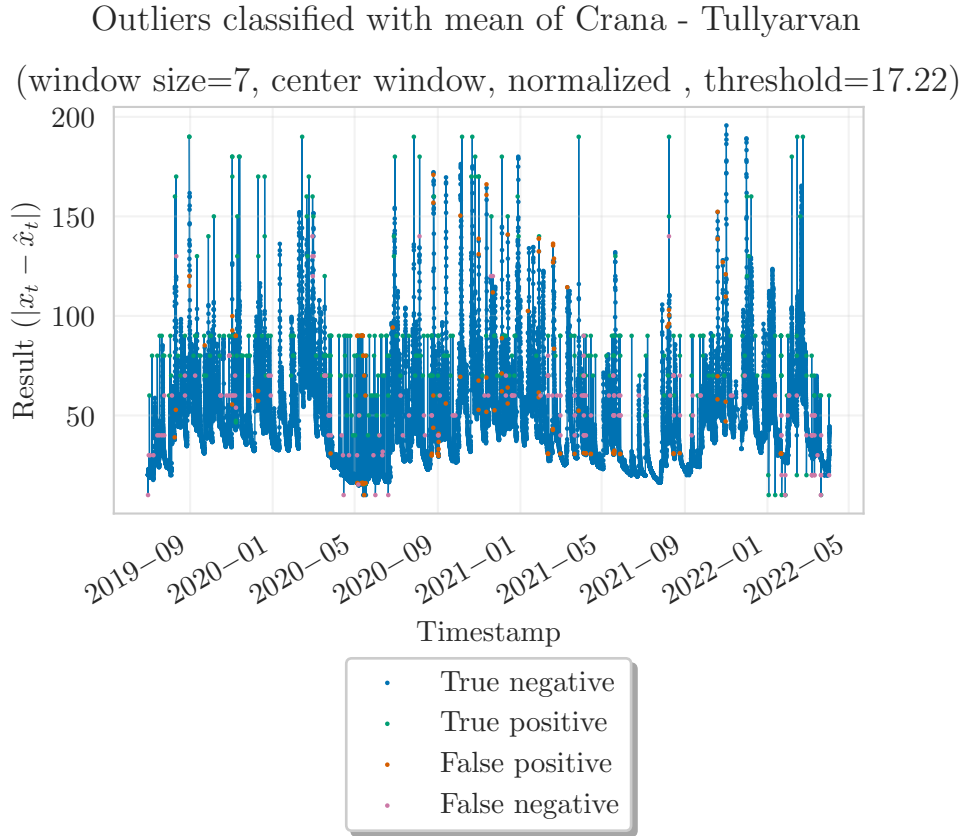


Figure 16: Best performance of outlier detection using mean (Crana - Tullyarvan)

Median

Figure 17 shows the best performing outlier detection using the median. The F_1 - score is about 0.91, the *Precision* is about 0.93 and the *Recall* is about 0.89.

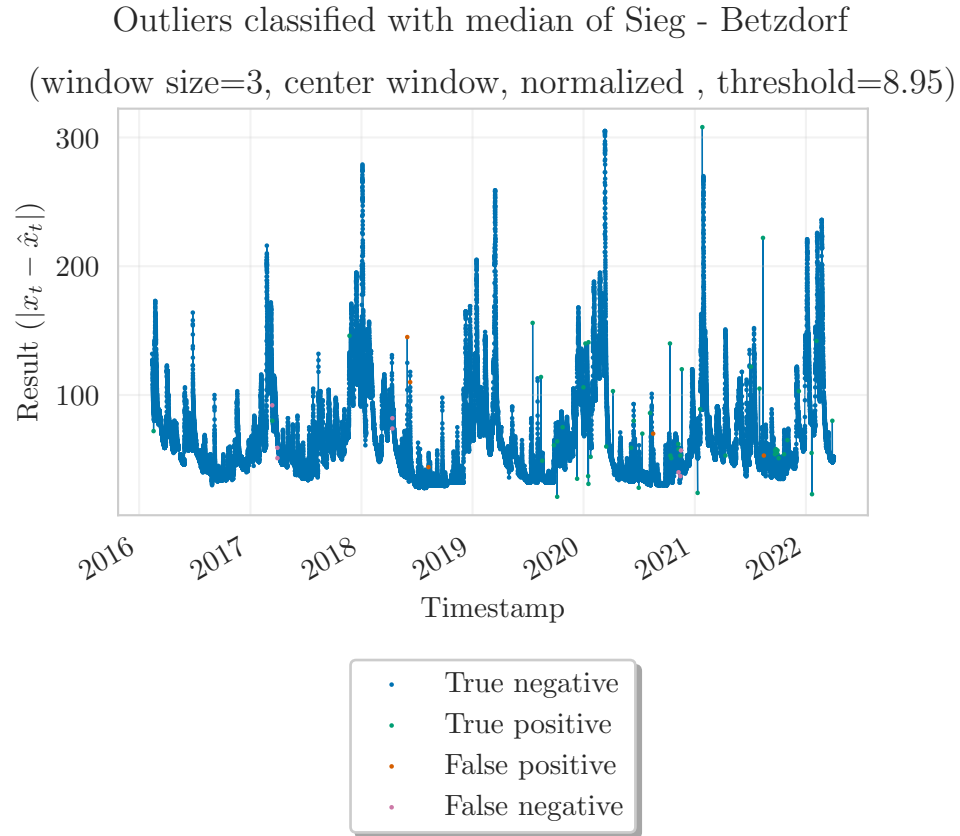


Figure 17: Best performance of outlier detection using median (Sieg - Betzdorf)

MAD

Figure 18 shows the best performing outlier detection using the MAD. The F_1 - score is about 0.45, the *Precision* is about 0.53 and the *Recall* is 0.4.

Outliers classified with mad of Losse - Helsa
(window size=2, no center window, not normalized , threshold=9.28)

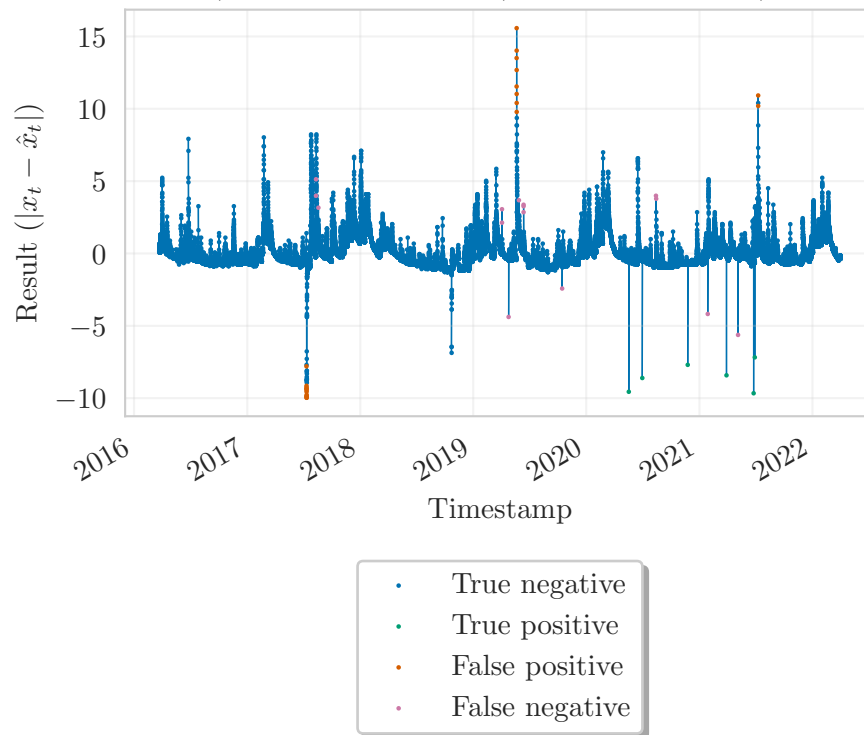


Figure 18: Best performance of outlier detection using MAD (Losse - Helsa)

Z-score

Figure 19 shows the best performing outlier detection using z-score. The F_1 - score is about 0.68, the *Precision* is about 0.80 and the *Recall* is about 0.59. In Figure 22 examples for TP, TN, FN and FP are shown.

Outliers classified with z-score of Crana - Tullyarvan
(window size=26, center window, not normalized , threshold=1.99)

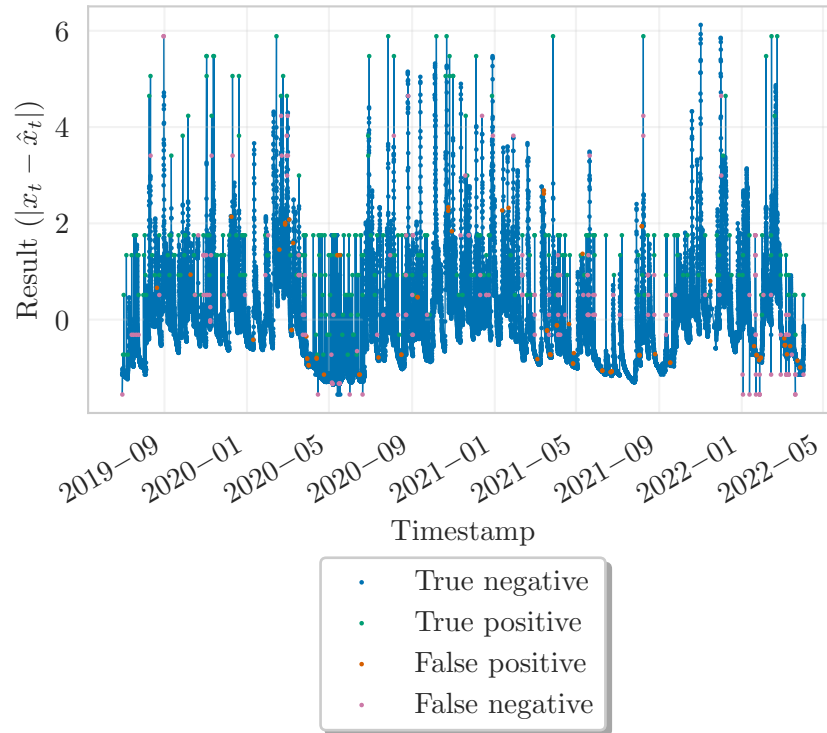


Figure 19: Best performance of outlier detection using z-score (Crana - Tullyarvan)

Outliers classified with z-score of Crana - Tullyarvan zoomed in
(window size=26, center window, not normalized , threshold=1.99)

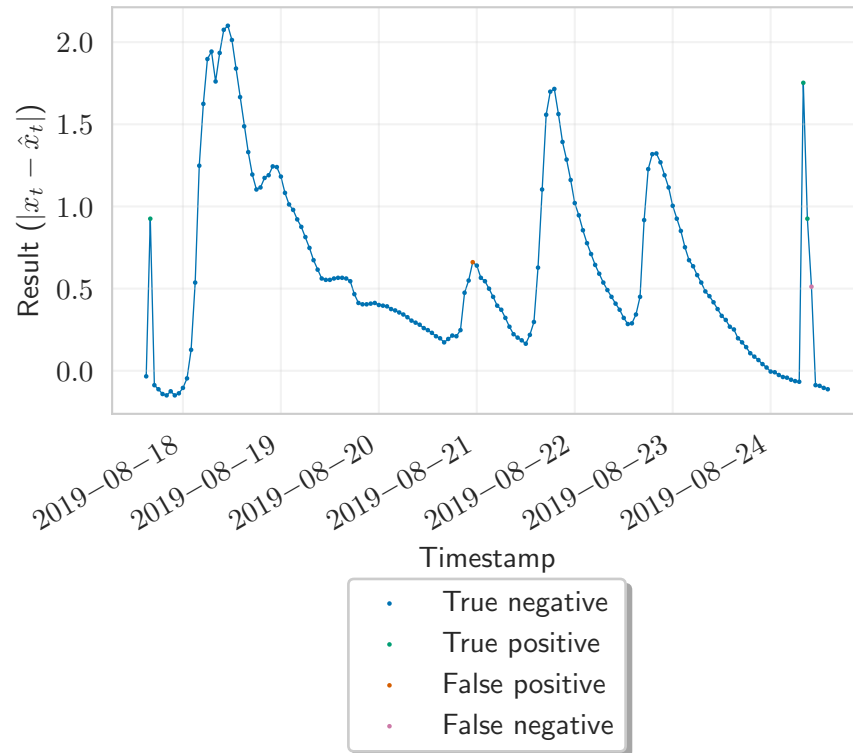


Figure 20: Best performance of outlier detection using z-score (Crana - Tullyarvan) zoomed in

MADN-z-score

Figure 21 shows the best performing outlier detection using MADN-z-score. The F_1 - score is about 0.79, the *Precision* is about 0.75 and the *Recall* is about 0.84. In Figure 22 examples for TP, TN, FN and FP are shown.

Outliers classified with madn-z-score of Aghacashlaun - Aghacashlaun
(window size=24, center window, not normalized , threshold=3.65)

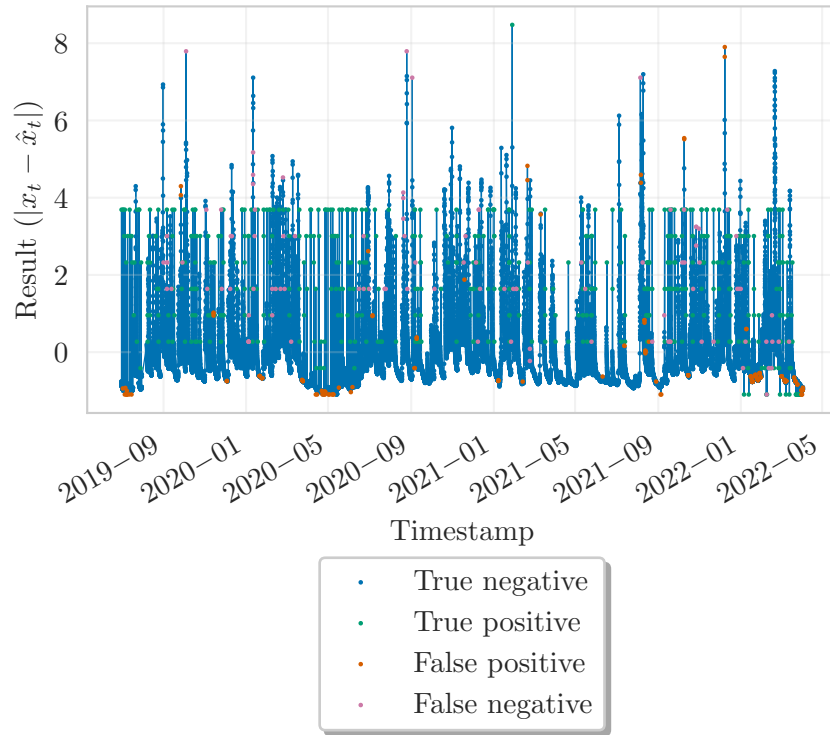


Figure 21: Best performance of outlier detection using MADN-z-score (Aghacashlaun - Aghacashlaun)

Outliers classified with madn-z-score of Aghacashlaun - Aghacashlaun zoomed in
(window size=24, center window, not normalized , threshold=3.65)

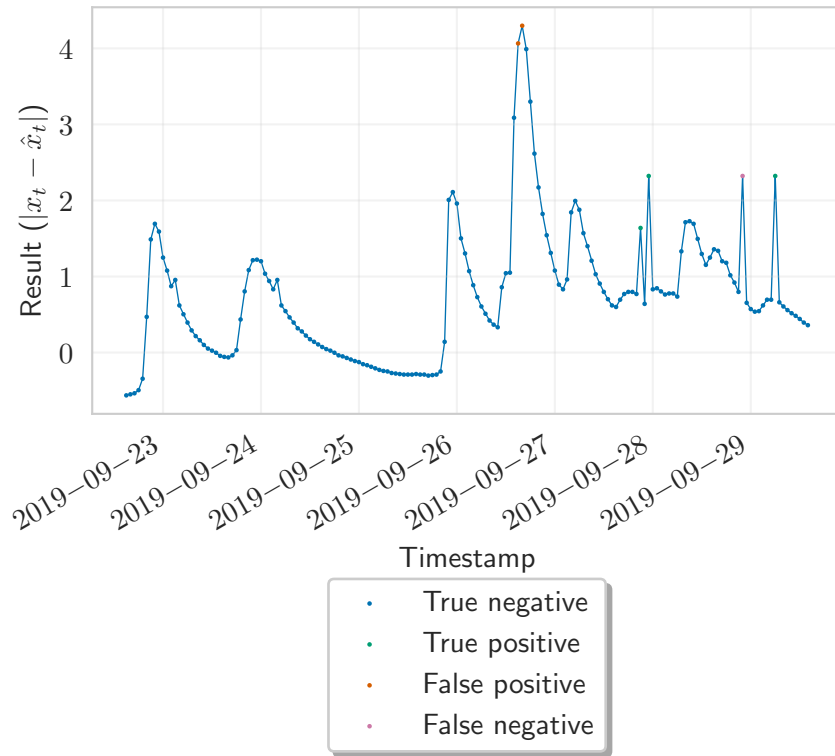


Figure 22: Best performance of outlier detection using MADN-z-score (Aghacashlaun - Aghacashlaun) zoomed in

4.8.1 Station Aghacashlaun, Aghacashlaun (36022-ie)

Table 10 shows the top three predictions per model for normalized and not normalized data. It clearly pictures, that methods using a centred moving window perform better. It is also interesting to see that the MADN z-score and the regular z-score performed about equally well, regardless of the fact, that the data is normalized or not.

window_size	center_window	normalized	threshold	model_type	f1_score
5.0	True	False	7.622074	median	0.812550
5.0	True	False	8.284281	median	0.812348
5.0	True	False	7.953177	median	0.811688
24.0	True	True	3.648829	madn-z-score	0.789205
24.0	True	False	3.648829	madn-z-score	0.789205
22.0	True	True	3.648829	madn-z-score	0.787572
22.0	True	False	3.648829	madn-z-score	0.787572
26.0	True	True	4.311037	madn-z-score	0.786415
26.0	True	False	4.311037	madn-z-score	0.786415
9.0	True	True	1.000000	median	0.759430
7.0	True	True	1.000000	median	0.750670
8.0	True	True	1.000000	median	0.749358
24.0	True	False	1.993311	z-score	0.626321
32.0	True	True	2.324415	z-score	0.622718
32.0	True	False	2.324415	z-score	0.621457
23.0	True	True	1.993311	z-score	0.621359
27.0	True	True	2.324415	z-score	0.620833
27.0	True	False	2.324415	z-score	0.620833
9.0	True	False	16.892977	mean	0.619647
10.0	True	False	18.879599	mean	0.619173
10.0	True	False	19.210702	mean	0.618538
11.0	True	True	1.331104	mean	0.611993
9.0	True	True	1.331104	mean	0.611418
10.0	True	True	1.331104	mean	0.610659
24.0	True	True	1.993311	mad	0.404255
26.0	True	True	1.993311	mad	0.404000
28.0	True	True	1.993311	mad	0.403183
22.0	True	False	66.889632	mad	0.401679
22.0	True	False	67.220736	mad	0.401126
22.0	True	False	67.551839	mad	0.400856

Table 10: Top predictions summary of Aghacashlaun - Aghacashlaun

4.8.2 Station Murg, Frauenfeld (2386-ch)

In Table 11 the top three predictions per model for normalized and not normalized data are shown. The outlier detection for this station has the worst performance among those tested. The result of this model is shown in Figure 23. Every value above the threshold is classified as an outlier. Every value below the threshold is classified as a regular value.

window_size	center_window	normalized	threshold	model_type	f1_score
3.0	True	False	5.966555	median	0.657895
3.0	True	False	5.635452	median	0.649351
3.0	True	False	6.628763	median	0.639175
3.0	True	False	4.642140	mean	0.447917
3.0	True	False	4.973244	mean	0.443804
3.0	True	False	5.304348	mean	0.429907
50.0	True	True	3.979933	z-score	0.327273
50.0	True	False	3.979933	z-score	0.327273
49.0	True	True	3.648829	z-score	0.306931
49.0	True	False	3.648829	z-score	0.306931
48.0	True	False	3.648829	z-score	0.305419
50.0	True	True	3.648829	z-score	0.303922
47.0	True	True	1.000000	median	0.120690
42.0	False	True	11.926421	median	0.120690
42.0	False	True	8.946488	median	0.120690
49.0	False	True	24.177258	mean	0.120690
34.0	True	True	9.277592	mean	0.120690
34.0	True	True	6.959866	mean	0.120690
NaN	True	True	97.351171	madn-z-score	0.120690
NaN	True	True	76.491639	madn-z-score	0.120690
NaN	True	True	74.505017	madn-z-score	0.120690
NaN	True	False	97.351171	madn-z-score	0.120690
NaN	True	False	76.491639	madn-z-score	0.120690
NaN	True	False	74.505017	madn-z-score	0.120690
7.0	True	True	1.000000	mad	0.120690
11.0	True	True	1.993311	mad	0.120690
34.0	False	True	24.839465	mad	0.120690
7.0	True	False	1.000000	mad	0.004301
27.0	True	False	4.973244	mad	0.004301
27.0	True	False	5.635452	mad	0.004301

Table 11: Top predictions summary of Murg - Frauenfeld

Result of median based outlier detection of Murg - Frauenfeld
 (window size=3, center window, normalized , threshold=5.97)

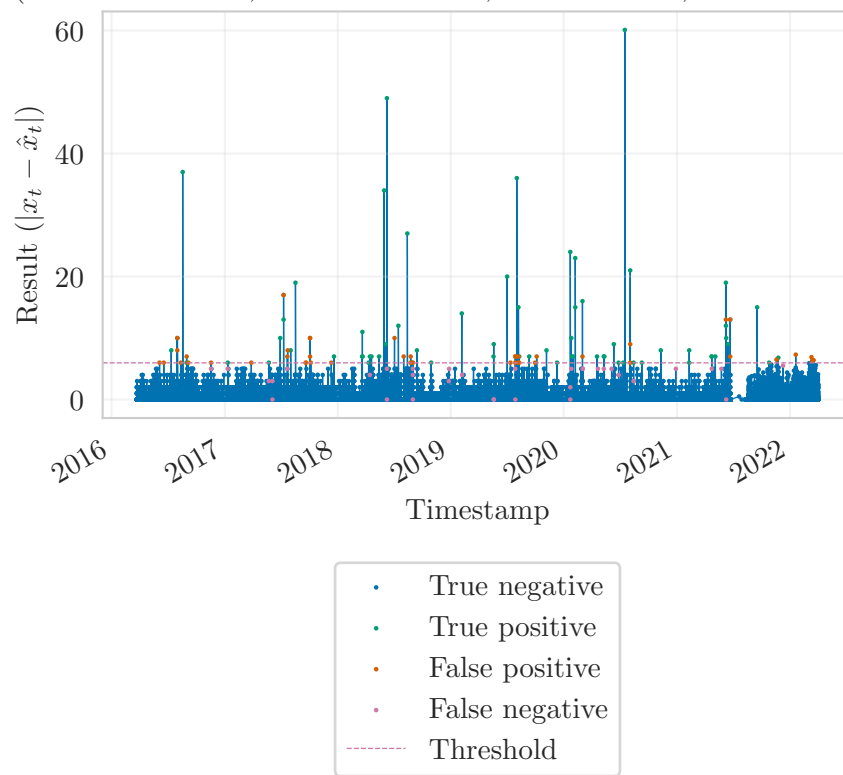


Figure 23: Result of median-based outlier detection Murg - Frauenfeld

4.8.3 Station Sieg, Betzdorf (2720050000-de)

In Table 12 the top three predictions per model for normalized and not normalized data are shown. With an F_1 - score of 0.905109 this station had the best performance. The result of this model is shown in Figure 24. Every value above the threshold is classified as an outlier. Between 2018 and 2019 an example of an FP can be seen (the orange values above the threshold).

window_size	center_window	normalized	threshold	model_type	f1_score
3.0	True	False	8.284281	median	0.905109
3.0	True	False	8.615385	median	0.905109
3.0	True	False	8.946488	median	0.905109
29.0	True	False	3.648829	z-score	0.545455
23.0	True	False	3.317726	z-score	0.545455
27.0	True	False	3.648829	z-score	0.542373
22.0	True	True	3.317726	z-score	0.533333
28.0	True	True	3.648829	z-score	0.533333
20.0	True	True	3.317726	z-score	0.516667
3.0	True	False	7.290970	mean	0.510288
3.0	True	False	6.959866	mean	0.510121
3.0	True	False	6.628763	mean	0.503937
NaN	False	True	10.602007	madn-z-score	0.133333
NaN	True	True	10.602007	madn-z-score	0.133333
NaN	False	False	10.602007	madn-z-score	0.133333
NaN	True	False	10.602007	madn-z-score	0.133333
NaN	False	True	10.270903	madn-z-score	0.119048
NaN	False	False	10.270903	madn-z-score	0.119048
47.0	True	True	1.000000	median	0.108108
29.0	False	True	1.000000	median	0.108108
50.0	False	True	2.324415	median	0.108108
NaN	False	True	1.662207	mean	0.108108
NaN	True	True	1.662207	mean	0.108108
NaN	True	True	1.000000	mean	0.108108
7.0	True	True	1.000000	mad	0.108108
31.0	True	True	1.000000	mad	0.108108
25.0	False	True	2.324415	mad	0.108108
49.0	True	False	100.000000	mad	0.006390
50.0	True	False	100.000000	mad	0.006379
47.0	True	False	100.000000	mad	0.006360

Table 12: Top predictions summary of Sieg - Betzdorf

Result of median based outlier detection of Sieg - Betzdorf
 (window size=3, center window, normalized , threshold=8.95)

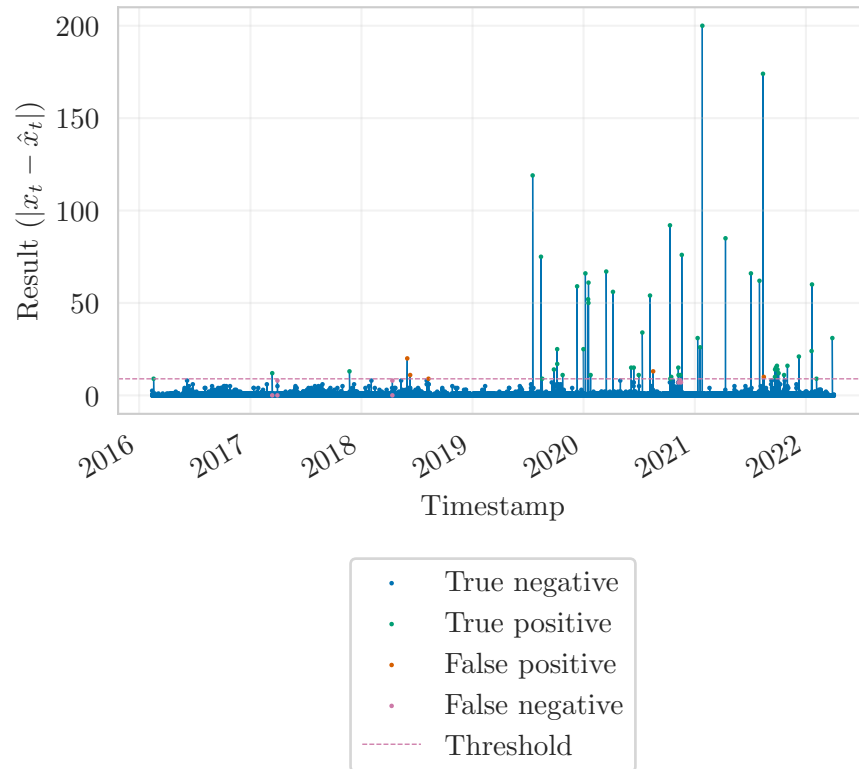


Figure 24: Result of median based outlier detection Sieg - Betzdorf

4.8.4 Station Losse, Helsa (42960105-de)

In Table 13 the top three predictions per model for normalized and not normalized data are shown.

window_size	center_window	normalized	threshold	model_type	f1_score
6.0	True	True	2.655518	median	0.731707
4.0	True	False	25.170569	median	0.731707
4.0	True	False	26.494983	median	0.731707
4.0	True	False	26.163880	median	0.731707
3.0	True	True	1.000000	median	0.720000
4.0	True	True	2.324415	median	0.714286
19.0	True	True	5.966555	mean	0.571429
17.0	True	True	5.635452	mean	0.571429
15.0	True	True	5.304348	mean	0.571429
15.0	True	False	62.585284	mean	0.571429
16.0	True	False	62.585284	mean	0.571429
16.0	True	False	59.605351	mean	0.571429
2.0	False	True	9.277592	mad	0.454545
2.0	True	True	9.277592	mad	0.454545
2.0	False	True	9.939799	mad	0.439024
NaN	True	True	18.879599	z-score	0.275862
NaN	True	True	17.555184	z-score	0.275862
NaN	False	True	16.892977	z-score	0.275862
NaN	False	False	19.541806	z-score	0.275862
NaN	False	False	16.561873	z-score	0.275862
NaN	False	False	15.899666	z-score	0.275862
NaN	False	True	30.137124	madn-z-score	0.275862
NaN	False	True	32.785953	madn-z-score	0.275862
NaN	False	True	35.103679	madn-z-score	0.275862
NaN	False	False	30.137124	madn-z-score	0.275862
NaN	False	False	32.785953	madn-z-score	0.275862
NaN	False	False	35.103679	madn-z-score	0.275862
NaN	True	False	100.000000	mad	0.005068
NaN	False	False	100.000000	mad	0.005068
NaN	True	False	99.668896	mad	0.004548

Table 13: Top predictions summary of Losse - Helsa

4.8.5 Station Crana, Tullyarvan (39003-ie)

In Table 14 the top three predictions per model for normalized and not normalized data are shown.

window_size	center_window	normalized	threshold	model_type	f1_score
5.0	True	False	8.615385	median	0.859035
5.0	True	False	9.277592	median	0.858500
5.0	True	False	8.946488	median	0.857685
34.0	True	True	3.648829	madn-z-score	0.736842
34.0	True	False	3.648829	madn-z-score	0.736842
34.0	True	True	4.311037	madn-z-score	0.733962
34.0	True	False	4.311037	madn-z-score	0.733962
32.0	True	True	4.642140	madn-z-score	0.731707
32.0	True	False	4.642140	madn-z-score	0.731707
11.0	True	True	1.000000	median	0.723699
10.0	True	True	1.000000	median	0.723059
9.0	True	True	1.000000	median	0.720189
7.0	True	False	17.224080	mean	0.684211
7.0	True	False	16.230769	mean	0.683398
7.0	True	False	16.892977	mean	0.683267
26.0	True	True	1.993311	z-score	0.677668
26.0	True	False	1.993311	z-score	0.677668
27.0	True	True	1.993311	z-score	0.672489
27.0	True	False	1.993311	z-score	0.672489
25.0	True	True	1.993311	z-score	0.669633
25.0	True	False	1.993311	z-score	0.669633
8.0	True	True	1.000000	mean	0.643868
9.0	True	True	1.000000	mean	0.636678
10.0	True	True	1.000000	mean	0.633596
36.0	True	False	85.431438	mad	0.180556
34.0	True	False	85.431438	mad	0.180149
36.0	True	False	85.100334	mad	0.179840
30.0	True	True	1.662207	mad	0.173418
14.0	False	True	1.662207	mad	0.173355
29.0	True	True	1.662207	mad	0.173149

Table 14: Top predictions summary of Crana - Tullyarvan

5 Conclusion

This bachelor thesis provides an overview of the topic of anomaly detection, especially outlier detection for time series data. It describes key features of data quality and introduces the topic of data cleaning and data cleansing. Furthermore, this paper provides a general overview of outlier detection approaches. After a theoretical overview of different outlier detection approaches, they are tested on water level data from different rivers.

The overall best performance, across different water level measurement stations, was achieved by using the median threshold-based outlier detection method with a centred window, which has a size of three and a threshold of about 6.6. The median threshold-based outlier detection also delivered the highest $F_1 - score$ (0.905). Using the mean to calculate \hat{x}_t is not recommended since the mean is not robust against outliers. Using the MAD with the threshold-based outlier detection resulted in the lowest $F_1 - score$, with the best score only being about 0.45. The second-best result was achieved by using the modified z-score described in section 3.5. For the stations tested the approach using the median delivered the best performances. However, this does not mean, that this will be true for all stations. It has to be assessed for each station individually which model is able to detect outliers the best. Furthermore, it depends on the use case if higher precision or recall is required. Depending on that, β for the $F_\beta - score$ needs to be chosen accordingly. For the tests, the $F_1 - score$ was used since precision and recall are equally important. In addition preprocessing the data by setting upper and lower boundaries and removing data points which exceed those limits, did not improve the performance of the models, on the contrary, the performance was worse. Because the extreme outliers were mostly detected anyways, thus fewer outliers were detected when setting upper and lower limits, which resulted in lower performance.

6 Future Work

There are countless outlier detection approaches, which can be applied for time series data, that have not been covered in detail in this paper. For example, one approach could be setting a maximum gradient for both directions (one for falling and one for rising values) for each measurement station. If the water level measurements exceed the maximum gradient then the value is classified as an outlier. Another approach would be to use ANNs with either LSTM, GRU or an autoencoder architecture. Additionally prediction and classification ANNs could be compared to see, which approach performs better.

7 Appendix

7.1 Python Code

7.1.1 Manual outlier detection

```
1 import json
2
3 import numpy as np
4 import pandas as pd
5 import plotly.express as px
6 from dash import Dash, dcc, html
7 from dash.dependencies import Input, Output
8
9 app = Dash(__name__)
10
11 stations_dict = pd.read_csv('./data/stations.csv').groupby(
12     ['common_id']).first().to_dict('index')
13
14 common_id = '2736731000100-de'
15 df = pd.read_parquet(f'./data/{common_id}_outliers_classified.parquet')
16 df.info()
17 fig = px.scatter(df, x='timestamp', y='water_level',
18                 title=f'{common_id}', color='is_outlier')
19 fig.update_layout(
20     xaxis=dict(
21         rangeselector=dict(
22             buttons=list([
23                 dict(count=1,
24                     step="day",
25                     stepmode="backward"),
26             ])
27         ),
28         rangeslider=dict(
29             visible=True
30         ),
31         type="date"
32     )
33 )
34 app.layout = html.Div([
35     html.H1('Manual Outlier Selection'),
36     dcc.Graph(
37         id='water-level-graph',
38         figure=fig
39     ),
40     html.Button('Refresh', id='refresh-btn', n_clicks=0),
41     html.Div([
42         dcc.Markdown('Clicked value:'),
43         html.Pre(id='click-data'),
```

```

44     ]),
45 ])
46
47
48 def toggle_outlier(timestamp: str):
49     df.loc[(df['timestamp'] == timestamp), 'is_outlier'] = np.invert(
50         df.loc[(df['timestamp'] == timestamp), 'is_outlier'])
51
52
53 @app.callback(
54     Output('click-data', 'children'),
55     Input('water-level-graph', 'clickData'))
56 def display_click_data(clickData):
57     if clickData is not None:
58         toggle_outlier(clickData['points'][0]['x'])
59     return json.dumps(clickData, indent=2)
60
61
62 @app.callback(
63     Output('water-level-graph', 'figure'),
64     Input('refresh-btn', 'n_clicks'))
65 def update_output(n_clicks):
66     fig = px.scatter(df, x='timestamp', y='water_level',
67                     title=f'{common_id}', color='is_outlier')
68     df.to_parquet(f'./data/{common_id}_outliers_classified.parquet')
69     fig.update_layout(
70         xaxis=dict(
71             rangeselector=dict(
72             ),
73             rangeslider=dict(
74                 visible=True
75             ),
76             type="date"
77         )
78     )
79     fig.update_layout(
80         yaxis_range=[df.loc[df['is_outlier'] == False,
81                          'water_level'].min() - 5,
82                     df.loc[df['is_outlier'] == False,
83                          'water_level'].max() + 5])
84     return fig
85
86
87 if __name__ == '__main__':
88     app.run_server(debug=True)

```

Listing 9: Manual Outlier Detection Webapp using Dash [34]

7.1.2 Explorative Data Analysis

```
1  # from https://jwalton.info/Embed-Publication-Matplotlib-Latex/
2  def set_size(width, fraction=1, subplots=(1, 1)):
3      """Set figure dimensions to avoid scaling in LaTeX.
4
5      Parameters
6      -----
7      width: float or string
8          Document width in points, or string of predined document type
9      fraction: float, optional
10         Fraction of the width which you wish the figure to occupy
11      subplots: array-like, optional
12         The number of rows and columns of subplots.
13
14      Returns
15      -----
16      fig_dim: tuple
17          Dimensions of figure in inches
18      """
19      if width == 'thesis':
20          width_pt = 364
21          height_pt = 595
22      else:
23          width_pt = width
24          height_pt = 595
25
26      # Width of figure (in pts)
27      fig_width_pt = width_pt * fraction
28      # Convert from pt to inches
29      inches_per_pt = 1 / 72.27
30
31      # Golden ratio to set aesthetic figure height
32      # https://disq.us/p/2940ij3
33      golden_ratio = (5 ** .5 - 1) / 2
34
35      # Figure width in inches
36      fig_width_in = fig_width_pt * inches_per_pt
37      # Figure height in inches
38      fig_height_in = min(
39          fig_width_in * golden_ratio * (subplots[0] / subplots[1]),
40          height_pt * inches_per_pt)
41
42      return (fig_width_in, fig_height_in)
```

Listing 10: Helper method to size plots [34]

```
1  # %%
2  import os
3
4  import matplotlib.pyplot as plt
5  import numpy as np
6  import pandas as pd
7
8  from shared_logic import *
9
10 random_seed = 1
11 np.random.seed(random_seed)
```

```

12
13 # %%
14 plt.style.use('seaborn-colorblind')
15 # https://jwalton.info/Embed-Publication-Matplotlib-Latex/
16 tex_fonts = {
17     # Use LaTeX to write all text
18     "text.usetex": True,
19     "font.family": "serif",
20     # Use 11pt font in plots, to match 11pt font in document
21     "axes.labelsize": 11,
22     "font.size": 11
23 }
24 plt.rcParams.update(tex_fonts)
25
26 # %%
27 stations_df = pd.read_csv('./data/stations.csv')
28 stations_dict = stations_df.groupby(['common_id']).first().to_dict('index')
29
30 # %%
31 common_id = '36022-ie'
32 # common_id = 'auto-1003803'
33 tex_plots_path = f'../bachelor-thesis/plots/pdfs/{common_id}/'
34 tex_table_path = f'../bachelor-thesis/tables/{common_id}/'
35 if not os.path.exists(tex_plots_path):
36     os.makedirs(tex_plots_path)
37 if not os.path.exists(tex_table_path):
38     os.makedirs(tex_table_path)
39
40 df = pd.read_parquet(
41     f'./data/classified_raw/{common_id}_outliers_classified.parquet')
42
43 # %%
44 df.info()
45
46 # %%
47 print(df.describe().to_latex())
48 # pd.Styler.to_latex(df.describe())
49
50 # %%
51 fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
52 bars = plt.bar(['regular', 'outlier'],
53               df['is_outlier'].value_counts().to_numpy())
54 plt.margins(0.1, 0.15)
55 ax.bar_label(bars)
56 # plt.xlabel('Is an outlier?')
57 plt.ylabel('Count')
58 plt.title(f'Outlier class distribution of '
59           f'{stations_dict[common_id]["water_name"]} - '
60           f'{stations_dict[common_id]["station_name"]}')
61 plt.grid(alpha=0.25, axis='y')
62 plt.savefig(f'{tex_plots_path}outlier_class_distribution_{common_id}.pdf',
63           format='pdf', bbox_inches='tight')
64
65 # %%
66 fig, (ax1, ax2, ax3) = plt.subplots(1, 3, constrained_layout=True,
67                                     figsize=set_size('thesis'))
68 # fig.subplots_adjust(hspace=0.25)

```

```

69
70 ax1.set_title('All values')
71 ax1.boxplot(df['water_level'].to_numpy(), vert=True,
72             flierprops={'marker': '.', 'markersize': 2})
73 ax1.set_ylabel('Water level')
74 # Hide x axis
75 ax1.get_xaxis().set_visible(False)
76 ax1.grid(alpha=0.25)
77
78 ax2.set_title('Regular values')
79 ax2.boxplot(df.loc[~df['is_outlier'], 'water_level'].to_numpy(), vert=True,
80             flierprops={'marker': '.', 'markersize': 2})
81 ax2.set_ylabel('Water level')
82 # Hide x axis
83 ax2.get_xaxis().set_visible(False)
84 ax2.grid(alpha=0.25)
85
86 ax3.set_title('Outlier values')
87 ax3.boxplot(df.loc[df['is_outlier'], 'water_level'].to_numpy(), vert=True,
88             flierprops={'marker': '.', 'markersize': 2})
89 ax3.set_ylabel('Water level')
90 # Hide x axis
91 ax3.get_xaxis().set_visible(False)
92 ax3.grid(alpha=0.25)
93
94 plt.suptitle(f'Boxplot of {stations_dict[common_id]["water_name"]} - '
95             f'{stations_dict[common_id]["station_name"]}')
96 plt.setp((ax1, ax2, ax3), ylim=ax1.get_ylim())
97
98 plt.savefig(f'{tex_plots_path}boxplot_{common_id}.pdf', format='pdf',
99             bbox_inches='tight')
100
101 # %%
102 fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
103 plt.boxplot(df['water_level'].to_numpy(), vert=False,
104             flierprops={'marker': '.', 'markersize': 2})
105 plt.xlabel('Water level')
106 # Hide y axis
107 ax.get_yaxis().set_visible(False)
108 plt.title(f'Boxplot of {stations_dict[common_id]["water_name"]} - '
109          f'{stations_dict[common_id]["station_name"]} (all values)')
110 plt.grid(alpha=0.25)
111 plt.savefig(f'{tex_plots_path}boxplot_{common_id}_all.pdf', format='pdf',
112             bbox_inches='tight')
113
114 # %%
115 fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
116 plt.boxplot(df.loc[~df['is_outlier'], 'water_level'].to_numpy(), vert=False,
117             flierprops={'marker': '.', 'markersize': 2})
118 plt.xlabel('Water level')
119 # Hide y axis
120 ax.get_yaxis().set_visible(False)
121 plt.title(f'Boxplot of {stations_dict[common_id]["water_name"]} - '
122          f'{stations_dict[common_id]["station_name"]} (regular values)')
123 plt.grid(alpha=0.25)
124 plt.savefig(f'{tex_plots_path}boxplot_{common_id}_regular.pdf', format='pdf',
125             bbox_inches='tight')

```

```

126
127 # %%
128 fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
129 plt.boxplot(df.loc[df['is_outlier'], 'water_level'].to_numpy(), vert=False,
130             flierprops={'marker': '.', 'markersize': 2})
131 plt.xlabel('Water level')
132 # Hide y axis
133 ax.get_yaxis().set_visible(False)
134 plt.title(f'Boxplot of {stations_dict[common_id]["water_name"]} - '
135          f'{stations_dict[common_id]["station_name"]} (outlier values)')
136 plt.grid(alpha=0.25)
137 plt.savefig(f'{tex_plots_path}boxplot_{common_id}_outlier.pdf', format='pdf',
138            bbox_inches='tight')
139
140 # %%
141 fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
142 plt.violinplot(df['water_level'].to_numpy(), vert=False, showmeans=True,
143               showmedians=True, showextrema=True)
144 plt.xlabel('Water level')
145 # Hide y axis
146 ax.get_yaxis().set_visible(False)
147 plt.title(f'Violinplot of {stations_dict[common_id]["water_name"]} - '
148          f'{stations_dict[common_id]["station_name"]}')
149 plt.grid(alpha=0.25)
150
151 # %%
152 fig, (ax1, ax2, ax3) = plt.subplots(3, 1, constrained_layout=True,
153                                     figsize=set_size('thesis', subplots=(3, 1)))
154 ax1.set_title('All Values')
155 ax1.hist(df['water_level'].to_numpy(), bins=100)
156 ax1.set_xlabel('Water level')
157 ax1.set_yscale('log')
158 ax1.set_ylabel('Count')
159 ax1.grid(alpha=0.25)
160
161 ax2.set_title('Regular values')
162 ax2.hist(df.loc[~df['is_outlier'], 'water_level'].to_numpy(), bins=100)
163 ax2.set_xlabel('Water level')
164 ax2.set_ylabel('Count')
165 ax2.set_yscale('log')
166 ax2.grid(alpha=0.25)
167
168 ax3.set_title('Outlier values')
169 ax3.hist(df.loc[df['is_outlier'], 'water_level'].to_numpy(), bins=100)
170 ax3.set_xlabel('Water level')
171 ax3.set_ylabel('Count')
172 ax3.set_yscale('log')
173 ax3.grid(alpha=0.25)
174
175 plt.suptitle(f'Histogram of the water Level of '
176             f'{stations_dict[common_id]["water_name"]} - '
177             f'{stations_dict[common_id]["station_name"]}')
178 plt.savefig(f'{tex_plots_path}water_level_histogram_{common_id}.pdf',
179            format='pdf', bbox_inches='tight')
180
181 # %%
182 fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))

```

```

183 plt.hist(df['water_level'].to_numpy(), bins=100)
184 plt.xlabel('Water level')
185 plt.ylabel('Count')
186 plt.title(f'Histogram of {stations_dict[common_id]["water_name"]} - '
187           f'{stations_dict[common_id]["station_name"]}')
188 plt.grid(alpha=0.25)
189 # plt.show()
190 plt.savefig(f'{tex_plots_path}histogram_{common_id}.pdf', format='pdf',
191            bbox_inches='tight')
192
193 # %%
194 set_size('thesis', subplots=(3, 1))
195
196 # %%
197 df['water_level_diff'] = df['water_level'].diff()
198 fig, (ax1, ax2, ax3) = plt.subplots(3, 1, constrained_layout=True,
199                                   figsize=set_size('thesis', subplots=(3, 1)))
200 ax1.set_title('All Values')
201 ax1.hist(df['water_level_diff'].to_numpy(), bins=100)
202 ax1.set_yscale('log')
203 ax1.set_xlabel('Water level delta to previous value')
204 ax1.set_ylabel('Count')
205 ax1.grid(alpha=0.25)
206
207 ax2.set_title('Regular values')
208 ax2.hist(df.loc[~df['is_outlier'], 'water_level'].diff().to_numpy(), bins=100)
209 ax2.set_yscale('log')
210 ax2.set_xlabel('Water level delta to previous value')
211 ax2.set_ylabel('Count')
212 ax2.grid(alpha=0.25)
213
214 ax3.set_title('Outlier values')
215 ax3.hist(df.loc[df['is_outlier'], 'water_level_diff'].to_numpy(), bins=100)
216 ax3.set_xlabel('Water level delta to previous value')
217 ax3.set_ylabel('Count')
218 ax3.set_yscale('log')
219 ax3.grid(alpha=0.25)
220
221 plt.suptitle(f'Histogram of the water level delta of '
222             f'{stations_dict[common_id]["water_name"]} - '
223             f'{stations_dict[common_id]["station_name"]}')
224 plt.savefig(f'{tex_plots_path}water_level_delta_histogram_{common_id}.pdf',
225            format='pdf', bbox_inches='tight')
226
227 # %%
228 # get time difference between measurements in hours
229 df['timedelta'] = df['timestamp'].diff().astype('timedelta64[h]').astype(float)
230 fig, (ax1, ax2, ax3) = plt.subplots(3, 1, constrained_layout=True,
231                                   figsize=set_size('thesis', subplots=(3, 1)))
232 ax1.set_title('All Values')
233 ax1.hist(df['timedelta'].to_numpy(), bins=100)
234 ax1.set_yscale('log')
235 ax1.set_xlabel('Timedelta to previous datapoint in hours')
236 ax1.set_ylabel('Count')
237 ax1.grid(alpha=0.25)
238
239 ax2.set_title('Regular values')

```



```

240 ax2.hist(df.loc[~df['is_outlier'], 'timedelta'].to_numpy(), bins=100)
241 ax2.set_yscale('log')
242 ax2.set_xlabel('Timedelta to previous datapoint in hours')
243 ax2.set_ylabel('Count')
244 ax2.grid(alpha=0.25)
245
246 ax3.set_title('Outlier values')
247 ax3.hist(df.loc[df['is_outlier'], 'timedelta'].to_numpy(), bins=100)
248 ax3.set_xlabel('Timedelta to previous datapoint in hours')
249 ax3.set_ylabel('Count')
250 ax3.grid(alpha=0.25)
251
252 plt.suptitle(f'Histogram of the timedelta between values of '
253             f'{stations_dict[common_id]["water_name"]} - '
254             f'{stations_dict[common_id]["station_name"]}')
255 plt.savefig(f'{tex_plots_path}time_delta_histogram_{common_id}.pdf',
256           format='pdf', bbox_inches='tight')
257
258 # %%
259 from scipy.stats import kde
260
261 density = kde.gaussian_kde(df['water_level'].to_numpy())
262 x = np.linspace(np.floor(df['water_level'].min()),
263               np.ceil(df['water_level'].max()), 300)
264 y = density(x)
265 plt.hist(df['water_level'], bins=100, density=True)
266 plt.plot(x, y)
267 plt.title("Density Plot of the data")
268 plt.show()
269
270 # %%
271 df.describe().to_latex(f'{tex_table_path}/{common_id}-7-number-summary-all.tex',
272                      position='htp',
273                      label=f'table:{common_id}-7-number-summary-all',
274                      caption=f'Seven number summary of '
275                             f'{stations_dict[common_id]["water_name"]} - '
276                             f'{stations_dict[common_id]["station_name"]} '
277                             f'(all values)')
278
279 # %%
280 df.loc[~df['is_outlier']].describe().to_latex(
281     f'{tex_table_path}/{common_id}-7-number-summary-regular.tex',
282     position='htp',
283     label=f'table:{common_id}-7-number-summary-regular',
284     caption=f'Seven number summary of {stations_dict[common_id]["water_name"]} '
285            f'- {stations_dict[common_id]["station_name"]} (regular values)')
286
287 # %%
288 df.loc[df['is_outlier']].describe().to_latex(
289     f'{tex_table_path}/{common_id}-7-number-summary-outlier.tex',
290     position='htp',
291     label=f'table:{common_id}-7-number-summary-outlier',
292     caption=f'Seven number summary of {stations_dict[common_id]["water_name"]} '
293            f'- {stations_dict[common_id]["station_name"]} (outlier values)')
294
295 # %%
296 df.loc[df['is_outlier']].describe().to_latex(

```

```

297     f'{tex_table_path}/{common_id}-7-number-summary-outlier.tex',
298     position='htp',
299     label=f'table:{common_id}-7-number-summary-outlier',
300     caption=f'Seven number summary of {stations_dict[common_id]["water_name"]} '
301             f'- {stations_dict[common_id]["station_name"]} (outlier values)'
302
303     # %%
304     fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
305     plt.plot(df['timestamp'], df['water_level'], linewidth=0.5, zorder=-1)
306     # plt.scatter(df['timestamp'], df['water_level'], s=0.5)
307     plt.scatter(df.loc[~df['is_outlier'], 'timestamp'],
308                df.loc[~df['is_outlier'], 'water_level'], s=0.5, c='C0',
309                label='Regular')
310
311     plt.scatter(df.loc[df['is_outlier'], 'timestamp'],
312                df.loc[df['is_outlier'], 'water_level'], s=0.5, c='C2',
313                label='Outlier')
314
315     plt.gcf().autofmt_xdate()
316     plt.xlabel('Timestamp')
317     plt.ylabel('Water Level')
318     plt.title(f'Linechart of {stations_dict[common_id]["water_name"]} - '
319             f'{stations_dict[common_id]["station_name"]}')
320     plt.grid(alpha=0.25)
321     ax.legend(loc='upper center', bbox_to_anchor=(0.5, -0.32))
322
323     # plt.show()
324     plt.savefig(f'{tex_plots_path}linechart_{common_id}.pdf', format='pdf',
325                bbox_inches='tight')
326
327     # %%
328     from matplotlib.dates import DateFormatter
329
330     fig, ax = plt.subplots(1, 1, figsize=set_size('thesis'))
331     start_date = '2019-10-16'
332     end_date = '2019-10-19'
333     df_slice = df.loc[
334         (df['timestamp'] >= start_date) & (df['timestamp'] <= end_date)]
335     plt.plot(df_slice['timestamp'], df_slice['water_level'], linewidth=0.5,
336             zorder=-1)
337     # plt.scatter(df['timestamp'], df['water_level'], s=0.5)
338     plt.scatter(df_slice.loc[~df_slice['is_outlier'], 'timestamp'],
339                df_slice.loc[~df_slice['is_outlier'], 'water_level'], s=0.5, c='C0',
340                label='Regular')
341
342     plt.scatter(df_slice.loc[df_slice['is_outlier'], 'timestamp'],
343                df_slice.loc[df_slice['is_outlier'], 'water_level'], s=0.5, c='C2',
344                label='Outlier')
345
346     ax.fmt_xdata = DateFormatter('%Y-%m-%d')
347     import matplotlib.dates as mdates
348
349     ax.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
350     plt.gcf().autofmt_xdate()
351     plt.xlabel('Timestamp')
352     plt.ylabel('Water Level')
353     plt.title(f'Linechart of {stations_dict[common_id]["water_name"]} - '

```

```

354         f'{stations_dict[common_id]["station_name"]}')
355 plt.grid(alpha=0.25)
356 ax.legend(loc='upper center', bbox_to_anchor=(0.5, -0.37))
357
358 # plt.show()
359 plt.savefig(f'{tex_plots_path}slice_linechart_{common_id}.pdf', format='pdf',
360             bbox_inches='tight')

```

Listing 11: Explorative Data Analysis [34]

7.1.3 Preprocessing

```

1  # %%
2  import os
3
4  import numpy as np
5  import pandas as pd
6
7  random_seed = 1
8  np.random.seed(random_seed)
9
10 # %%
11 stations_df = pd.read_csv('./data/stations.csv')
12 stations_dict = stations_df.groupby(['common_id']).first().to_dict('index')
13 stations_dict['2386-ch']['lower_limit'] = 38_000.0
14 stations_dict['2386-ch']['upper_limit'] = 40_000.0
15 stations_dict['2720050000-de']['lower_limit'] = 20.0
16 stations_dict['2720050000-de']['upper_limit'] = 500.0
17 stations_dict['36022-ie']['lower_limit'] = 20.0
18 stations_dict['36022-ie']['upper_limit'] = 175.0
19 stations_dict['39003-ie']['lower_limit'] = 10.0
20 stations_dict['39003-ie']['upper_limit'] = 225.0
21 stations_dict['42960105-de']['lower_limit'] = -10.0
22 stations_dict['42960105-de']['upper_limit'] = 275.0
23 stations_dict['auto-1003803']['lower_limit'] = -10.0
24 stations_dict['auto-1003803']['upper_limit'] = 275.0
25
26 # %%
27 for common_id, station_dict in stations_dict.items():
28     fp = f'./data/classified_raw/{common_id}_outliers_classified.parquet'
29     if not os.path.exists(fp):
30         continue
31     raw_classified_df = pd.read_parquet(fp)
32     print(f'Removing for {common_id}, lower limit '
33           f'{station_dict["lower_limit"]}, upper limit '
34           f'{station_dict["upper_limit"]}')
35     print('Removing following rows:')
36     mask = (raw_classified_df['water_level'] < station_dict['lower_limit']) | (
37             raw_classified_df['water_level'] > station_dict['upper_limit'])
38     print(raw_classified_df[mask])
39     raw_classified_df[~mask].to_parquet(
40         f'./data/classified/{common_id}_outliers_classified.parquet')
41     print()

```

Listing 12: Data preprocessing code in Python [34]

7.1.4 Outlier Detection Methods Implementation

```
1  from typing import Union
2
3  import numpy as np
4  import pandas as pd
5
6
7  def mean_outlier_detection(input_df: pd.DataFrame,
8                             window: Union[int, None],
9                             center_window: bool):
10
11     """
12     Detects outliers in a dataframe using a (moving) average.
13     :param input_df: the input dataframe where the values are stored
14                       in the column water_level
15     :param window: the size of the window, None if no window should
16                    be used
17     :param center_window: whether the window should be centered or not
18     :return: a copy of the input dataframe where the column result
19              should be compared to a threshold to detect outliers
20     """
21     od_df = input_df.copy()
22     if window is None:
23         od_df['x_hat'] = od_df['water_level'].mean()
24     else:
25         od_df['x_hat'] = \
26             od_df['water_level'].rolling(window=window,
27                                          center=center_window,
28                                          min_periods=1).mean()
29     od_df['result'] = np.abs(od_df['water_level'] - od_df['x_hat'])
30     return od_df
31
32 def median_outlier_detection(input_df: pd.DataFrame,
33                              window: Union[int, None],
34                              center_window: bool):
35
36     """
37     Detects outliers in a dataframe using a (moving) average.
38     :param input_df: the input dataframe where the values are stored
39                       in the column water_level
40     :param window: the size of the window, None if no window should
41                    be used
42     :param center_window: whether the window should be centered or not
43     :return: a copy of the input dataframe where the column result
44              should be compared to a threshold to detect outliers
45     """
46     od_df = input_df.copy()
47     if window is None:
48         od_df['x_hat'] = od_df['water_level'].median()
49     else:
50         od_df['x_hat'] = \
51             od_df['water_level'].rolling(window=window,
52                                          center=center_window,
53                                          min_periods=1).median()
54     od_df['result'] = np.abs(od_df['water_level'] - od_df['x_hat'])
55     return od_df
56
```

```

57 def mad_outlier_detection(input_df: pd.DataFrame,
58                           window: Union[int, None],
59                           center_window: bool):
60     """
61     Detects outliers in a dataframe using a (moving) average.
62     :param input_df: the input dataframe where the values are stored
63                     in the column water_level
64     :param window: the size of the window, None if no window should
65                   be used
66     :param center_window: whether the window should be centered or not
67     :return: a copy of the input dataframe where the column result
68             should be compared to a threshold to detect outliers
69     """
70     od_df = input_df.copy()
71     if window is None:
72         od_df['x_hat'] = np.median(
73             np.abs(od_df['water_level'] - np.median(
74                 od_df['water_level'])))
75     else:
76         od_df['x_hat'] = \
77             od_df['water_level'].rolling(window=window,
78                                         center=center_window,
79                                         min_periods=1).apply(
80                 lambda x: np.median(np.abs(x - np.median(x))))
81     od_df['result'] = np.abs(od_df['water_level'] - od_df['x_hat'])
82     return od_df
83
84
85 def z_score_outlier_detection(input_df: pd.DataFrame,
86                               window: Union[int, None],
87                               center_window: bool):
88     """
89     Detects outliers in a dataframe using a (moving) average.
90     :param input_df: the input dataframe where the values are stored
91                     in the column water_level
92     :param window: the size of the window, None if no window should
93                   be used
94     :param center_window: whether the window should be centered or not
95     :return: a copy of the input dataframe where the column result
96             should be compared to a threshold to detect outliers
97     """
98     od_df = input_df.copy()
99     if window is None:
100         od_df['mean'] = od_df['water_level'].mean()
101         od_df['std'] = od_df['water_level'].std()
102     else:
103         od_df['mean'] = \
104             od_df['water_level'].rolling(window=window,
105                                         center=center_window,
106                                         min_periods=1).mean()
107         od_df['std'] = \
108             od_df['water_level'].rolling(window=window,
109                                         center=center_window,
110                                         min_periods=1).std()
111     od_df['result'] = \
112         (od_df['water_level'] - od_df['mean']).divide(od_df['std'])
113     return od_df

```

```

114
115
116 def delta_z_score_outlier_detection(input_df: pd.DataFrame,
117                                     window: Union[int, None],
118                                     center_window: bool):
119     """
120     Detects outliers in a dataframe using a (moving) average.
121     :param input_df: the input dataframe where the values are stored
122                       in the column water_level
123     :param window: the size of the window, None if no window should
124                     be used
125     :param center_window: whether the window should be centered or not
126     :return: a copy of the input dataframe where the column result
127              should be compared to a threshold to detect outliers
128     """
129     od_df = input_df.copy()
130     od_df['water_level_delta'] = od_df['water_level'].diff().fillna(0)
131     if window is None:
132         od_df['mean'] = od_df['water_level_delta'].mean()
133         od_df['std'] = od_df['water_level_delta'].std()
134     else:
135         od_df['mean'] = \
136             od_df['water_level_delta'].rolling(window=window,
137                                                 center=center_window,
138                                                 min_periods=1).mean()
139         od_df['std'] = \
140             od_df['water_level_delta'].rolling(window=window,
141                                                 center=center_window,
142                                                 min_periods=1).std()
143     od_df['result'] = \
144         (od_df['water_level_delta'] - od_df['mean']).divide(
145             od_df['std'])
146     return od_df
147
148
149 def madn_z_score_outlier_detection(input_df: pd.DataFrame,
150                                    window: Union[int, None],
151                                    center_window: bool):
152     """
153     Detects outliers in a dataframe using a (moving) average.
154     :param input_df: the input dataframe where the values are stored
155                       in the column water_level
156     :param window: the size of the window, None if no window should
157                     be used
158     :param center_window: whether the window should be centered or not
159     :return: a copy of the input dataframe where the column result
160              should be compared to a threshold to detect outliers
161     """
162     od_df = input_df.copy()
163     if window is None:
164         od_df['median'] = od_df['water_level'].median()
165         od_df['mad'] = np.median(
166             np.abs(od_df['water_level'] - od_df['median']))
167         od_df['madn'] = od_df['mad'] / 0.6745
168     else:
169         od_df['median'] = \
170             od_df['water_level'].rolling(window=window,

```

```

171                                     min_periods=1,
172                                     center=center_window).median()
173     od_df['mad'] = \
174         od_df['water_level'].rolling(window=window,
175                                     min_periods=1,
176                                     center=center_window).apply(
177             lambda x: np.median(np.abs(x - np.median(x))))
178     od_df['madn'] = od_df['mad'] / 0.6745
179     od_df['result'] = \
180         (od_df['water_level'] - od_df['median']).abs() \
181         .divide(od_df['madn'])
182     return od_df

```

Listing 13: Implementation of Outlier Detection Methods [34]

7.1.5 Parameter Grid Search

```

1  import itertools
2  import multiprocessing as mp
3  import os
4  from datetime import datetime
5  from typing import List
6
7  import numpy as np
8  import pandas as pd
9  from sklearn.preprocessing import StandardScaler
10
11 import outlier_detection_methods as odm
12
13 random_seed = 1
14 np.random.seed(random_seed)
15
16
17 def threshold_outlier_prediction(input_df, window, center_window,
18                                method):
19     od_df = input_df.copy()
20     if method == 'mean':
21         od_df = odm.mean_outlier_detection(od_df, window,
22                                           center_window)
23     elif method == 'median':
24         od_df = odm.median_outlier_detection(od_df, window,
25                                             center_window)
26     elif method == 'mad':
27         od_df = odm.mad_outlier_detection(od_df, window,
28                                           center_window)
29     elif method == 'z-score':
30         od_df = odm.z_score_outlier_detection(od_df, window,
31                                              center_window)
32     elif method == 'delta-z-score':
33         od_df = odm.delta_z_score_outlier_detection(od_df, window,
34                                                    center_window)
35     elif method == 'madn-z-score':
36         od_df = odm.madn_z_score_outlier_detection(od_df, window,
37                                                    center_window)
38     else:

```

```

39         raise ValueError(f'Method ({method}) not supported')
40     # od_df['result'] = od_df['result'].replace(np.inf, np.nan)
41     # od_df['result'] = od_df['result'].replace(-np.inf, np.nan)
42     od_df['result'] = od_df['result'].fillna(0)
43     od_df['result'] = od_df['result'].replace(np.inf, np.finfo(
44         np.float64).max)
45     od_df['result'] = od_df['result'].replace(-np.inf, np.finfo(
46         np.float64).min)
47     return {'window_size': window, 'center_window': center_window,
48           'method': method, 'df': od_df}
49
50
51 def run_grid_search_parallely(df: pd.DataFrame, windows: list,
52                               center_windows: list,
53                               methods: List[str],
54                               target_dir: str):
55     with mp.Pool(processes=12) as executor:
56         results = executor.starmap(threshold_outlier_prediction,
57                                   itertools.product([df], windows,
58                                                       center_windows,
59                                                       methods))
60
61     if not os.path.exists(target_dir):
62         os.makedirs(target_dir)
63     for res in results:
64         res['df'].to_parquet(
65             f'{target_dir}{res["window_size"]}_'
66             f'{"cw" if res["center_window"] else "nocw"}_'
67             f'{res["method"]}.parquet')
68
69
70 stations_df = pd.read_csv('./data/stations.csv')
71 stations_dict = stations_df.groupby(['common_id']).first().to_dict(
72     'index')
73
74 common_ids = ['36022-ie', '39003-ie', '2386-ch', '42960105-de',
75              '2720050000-de']
76 methods = ['median', 'mean', 'mad', 'z-score', 'madn-z-score', 'delta-z-score']
77 windows = [None] + list(range(2, 51))
78 center_windows = [False, True]
79 for common_id in common_ids:
80     print(
81         f'{datetime.now().isoformat()} - Processing {common_id} (regular, not preprocessed)')
82     tex_plots_path = f'../bachelor-thesis/plots/pdfs/{common_id}/'
83     df = pd.read_parquet(
84         f'data/classified_raw/{common_id}_outliers_classified.parquet')
85     run_grid_search_parallely(df, windows, center_windows, methods,
86                               f'./data/predictions/raw/regular/{common_id}/')
87     print(
88         f'{datetime.now().isoformat()} - Processing {common_id} (regular, preprocessed)')
89     df = pd.read_parquet(
90         f'data/classified/{common_id}_outliers_classified.parquet')
91     run_grid_search_parallely(df, windows, center_windows, methods,
92                               f'./data/predictions/raw_preprocessed/regular/{common_id}/')
93
94     print(
95         f'{datetime.now().isoformat()} - Processing {common_id} (normalized, not preprocessed)')

```



```

96     df = pd.read_parquet(
97         f'data/classified_raw/{common_id}_outliers_classified.parquet')
98     scaler = StandardScaler()
99     df['water_level'] = scaler.fit_transform(df[['water_level']])
100    run_grid_search_parallelly(df, windows, center_windows, methods,
101                              f'./data/predictions/raw/normalized/{common_id}/')
102    print(
103        f'{datetime.now().isoformat()} - Processing {common_id} (normalized, preprocessed)')
104    df = pd.read_parquet(
105        f'data/classified/{common_id}_outliers_classified.parquet')
106    scaler = StandardScaler()
107    df['water_level'] = scaler.fit_transform(df[['water_level']])
108    run_grid_search_parallelly(df, windows, center_windows, methods,
109                              f'./data/predictions/raw_preprocessed/normalized/{common_id}/')

```

Listing 14: Threshold based outlier detection [34]

7.1.6 Calculate Outlier Detection Results

```

1  # %%
2  import glob
3  import multiprocessing as mp
4
5  import matplotlib.pyplot as plt
6  import numpy as np
7  import pandas as pd
8  # from IPython.core.display_functions import display
9  import sklearn.metrics as metrics
10
11  random_seed = 1
12  np.random.seed(random_seed)
13
14  # %%
15  plt.style.use('seaborn-colorblind')
16
17  # from https://jwalton.info/Embed-Publication-Matplotlib-Latex/
18  tex_fonts = {
19      # Use LaTeX to write all text
20      "text.usetex": True,
21      "font.family": "serif",
22      # Use 11pt font in plots, to match 11pt font in document
23      "axes.labelsize": 11,
24      "font.size": 11
25  }
26  plt.rcParams.update(tex_fonts)
27  # tex_plots_path = f'./bachelor-thesis/plots/pdfs/{common_id}/'
28
29
30  # %%
31
32  file_list = glob.glob('./data/predictions/raw/**/*.parquet', recursive=True)
33
34  predictions = []
35
36  for fp in file_list:

```

```

37     file_name = fp.split('/')[-1]
38     metadata = file_name.split('_')
39     # remove .parquet
40     metadata = [m.split('.')[0] for m in metadata]
41     # df = pd.read_parquet(fp)
42     metadata_dict = {
43         'file_path': fp,
44         'normalized': fp.split('/')[-3] == 'normalized',
45         'window_size': None if metadata[0] == 'None' else int(metadata[0]),
46         'center_window': metadata[1] == 'cw',
47         'model_type': metadata[2],
48         'common_id': fp.split('/')[-2]
49     }
50     predictions.append(metadata_dict)
51
52     # %%
53     len(predictions)
54
55
56     # %%
57     def get_prediction_summary(metadata_dict):
58         predictions_summary = []
59         pred_df = pd.read_parquet(metadata_dict['file_path'])
60         threshold_min = 1
61         threshold_max = 100
62         threshold_steps = 300
63         thresholds = np.linspace(threshold_min, threshold_max, threshold_steps)
64         y_true = pred_df['is_outlier'].astype(int).to_numpy()
65         m = pred_df['result'].to_numpy()
66         for threshold in thresholds:
67             y_pred = np.where(m > threshold, 1, 0)
68             tn, fp, fn, tp = metrics.confusion_matrix(y_true, y_pred).ravel()
69             f1_score = metrics.f1_score(y_true, y_pred, zero_division=0)
70
71             predictions_summary.append({
72                 'common_id': metadata_dict['common_id'],
73                 'window_size': metadata_dict['window_size'],
74                 'center_window': metadata_dict['center_window'],
75                 'model_type': metadata_dict['model_type'],
76                 'normalized': metadata_dict['normalized'],
77                 'threshold': threshold,
78                 'f1_score': f1_score,
79                 'tn': tn,
80                 'fp': fp,
81                 'fn': fn,
82                 'tp': tp,
83             })
84         return predictions_summary
85
86
87     # predictions_summary_df = pd.DataFrame(predictions_summary)
88     # predictions_summary_df.info()
89
90     # %%
91     with mp.Pool(processes=12) as executor:
92         results = executor.map(get_prediction_summary, predictions)
93         result_lst = [item for sublist in results for item in sublist]

```

```

94
95  # %%
96  predictions_summary_df = pd.DataFrame(result_lst)
97
98  # %%
99  predictions_summary_df
100
101  # %%
102  predictions_summary_df.info()
103
104  # %%
105  predictions_summary_df.to_parquet(
106      f'./data/predictions/predictions_summary.parquet')
107
108  # %%
109  for id in predictions_summary_df['common_id'].unique():
110      print(id)
111      df = predictions_summary_df[predictions_summary_df['common_id'] == id]
112      df.to_csv(f'./data/predictions/predictions_preprocessed_summary/{id}.csv',
113              index=False)

```

Listing 15: Calculate the outlier detection results for different thresholds [34]

Bibliography

- [1] J. Strassmayr, G. Öller, E. Bragante, and J. Öhlböck, "FloodAlert Water Levels and Warning." [Online]. Available: <https://pegelalarm.at/en/>
- [2] L. Cai and Y. Zhu, "The Challenges of Data Quality and Data Quality Assessment in the Big Data Era," *Data Science Journal*, vol. 14, no. 0, p. 2, May 2015.
- [3] S. Song and A. Zhang, "IoT Data Quality," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. New York, NY, USA: Association for Computing Machinery, Oct. 2020, pp. 3517–3518.
- [4] D. Goldberg, "What every computer scientist should know about floating-point arithmetic," *ACM Computing Surveys*, vol. 23, no. 1, pp. 5–48, Mar. 1991. [Online]. Available: <https://doi.org/10.1145/103162.103163>
- [5] "What Every Computer Scientist Should Know About Floating-Point Arithmetic." [Online]. Available: https://docs.oracle.com/cd/E19957-01/806-3568/ncg_goldberg.html
- [6] "Apache Parquet," Dec. 2021. [Online]. Available: <https://parquet.apache.org/>
- [7] L. Li, J. McCann, N. S. Pollard, and C. Faloutsos, "DynaMMo: Mining and summarization of coevolving sequences with missing values," in *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '09. New York, NY, USA: Association for Computing Machinery, Jun. 2009, pp. 507–516.
- [8] J. I. Maletic and A. Marcus, "Data Cleansing: Beyond Integrity Analysis," p. 10, 2000.
- [9] A. Blázquez-García, A. Conde, U. Mori, and J. A. Lozano, "A review on Outlier/Anomaly detection in time series data," *arXiv:2002.04236 [cs, stat]*, Feb. 2020.
- [10] A. A. Cook, G. Mısırlı, and Z. Fan, "Anomaly Detection for IoT Time-Series Data: A Survey," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6481–6494, Jul. 2020.
- [11] F. Giannoni, M. Mancini, and F. Marinelli, "Anomaly Detection Models for IoT Time Series Data," *arXiv:1812.00890 [cs, eess]*, Nov. 2018.
- [12] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Computing Surveys*, vol. 41, no. 3, pp. 15:1–15:58, Jul. 2009.
- [13] C. Leys, C. Ley, O. Klein, P. Bernard, and L. Licata, "Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median," *Journal of Experimental Social Psychology*, vol. 49, no. 4, pp. 764–766, Jul. 2013.

- [14] S. Mehrang, E. Helander, M. Pavel, A. Chieh, and I. Korhonen, "Outlier detection in weight time series of connected scales," in *2015 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, Nov. 2015, pp. 1489–1496.
- [15] "Detection of Spatial Outlier by Using Improved Z-Score Test." [Online]. Available: <https://ieeexplore-1ieeee-1org-100033cm6022a.han.technikum-wien.at/document/8862582/>
- [16] G. Teschl and S. Teschl, *Spezielle Stetige Verteilungen*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014, pp. 325–345. [Online]. Available: http://link.springer.com/10.1007/978-3-642-54274-9_12
- [17] P. J. Rousseeuw and M. Hubert, "Anomaly detection by robust statistics," *WIREs Data Mining and Knowledge Discovery*, vol. 8, no. 2, p. e1236, 2018. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/widm.1236>
- [18] "Standard score," *Wikipedia*, Apr. 2022. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Standard_score&oldid=1082326996
- [19] I. Bae and U. Ji, "Outlier Detection and Smoothing Process for Water Level Data Measured by Ultrasonic Sensor in Stream Flows," *Water*, vol. 11, no. 5, p. 951, May 2019. [Online]. Available: <https://www.mdpi.com/2073-4441/11/5/951>
- [20] "Pandas documentation — pandas 1.4.2 documentation." [Online]. Available: <https://pandas.pydata.org/pandas-docs/version/1.4/index.html>
- [21] "GitHub." [Online]. Available: <https://github.com/>
- [22] J. Strassmayr, "Pegelalarm API wrapper," SOBOS-GmbH, Apr. 2022. [Online]. Available: https://github.com/SOBOS-GmbH/pegelalarm_public_pas_doc
- [23] "Hmac — Keyed-Hashing for Message Authentication — Python 3.9.12 documentation." [Online]. Available: <https://docs.python.org/3.9/library/hmac.html>
- [24] "Station Aghacashlaun - Aghacashlaun (36022-ie)." [Online]. Available: <https://waterlevel.ie/0000036022/>
- [25] "Station Murg - Frauenfeld (2386-ch)." [Online]. Available: <https://www.hydrodaten.admin.ch/en/2386.html>
- [26] "Station Sieg - Betzdorf (2720050000-de)." [Online]. Available: <http://www.hochwasser-rlp.de/pegeluebersichten/einzelpegel/flussgebiet/sieg/darstellung/tabellarisch/pegel/BETZDORF>
- [27] "Station Losse - Helsa (42960105-de)." [Online]. Available: <http://www.hlnug.de/static/pegel/wiskiweb2/stations/42960105/station.html>

- [28] "Station Crana - Tullyarvan (39003-ie)." [Online]. Available: <https://waterlevel.ie/0000039003/>
- [29] "Dash Documentation & User Guide | Plotly." [Online]. Available: <https://dash.plotly.com/>
- [30] "Precision and recall," *Wikipedia*, Mar. 2022. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Precision_and_recall&oldid=1080323102
- [31] Y. Sasaki, "The truth of the F-measure," p. 5.
- [32] N. Chinchor, "MUC-4 evaluation metrics," in *Proceedings of the 4th Conference on Message Understanding*, ser. MUC4 '92. USA: Association for Computational Linguistics, Jun. 1992, pp. 22–29. [Online]. Available: <https://doi.org/10.3115/1072064.1072067>
- [33] A. A. Taha and A. Hanbury, "Metrics for evaluating 3D medical image segmentation: Analysis, selection, and tool," *BMC Medical Imaging*, vol. 15, no. 1, p. 29, Aug. 2015. [Online]. Available: <https://doi.org/10.1186/s12880-015-0068-x>
- [34] D. Zelenay, "Outlier detection for water level data," Apr. 2022. [Online]. Available: <https://github.com/cellularegg/bachelor-thesis-code>

List of Figures

Figure 1	Examples of three point outliers	6
Figure 2	Examples of three subsequence outliers	7
Figure 3	Examples of four local point outliers	7
Figure 4	Examples of six local subsequence outliers	8
Figure 5	Examples of four global point outliers	8
Figure 6	Standard Score in a normal distribution [18]	11
Figure 7	Dash Webapp to classify outliers	18
Figure 8	Class distribution of Aghacashlaun - Aghacashlaun	20
Figure 9	Boxplot of Aghacashlaun - Aghacashlaun	21
Figure 10	Histogram of the water level of Aghacashlaun - Aghacashlaun	22
Figure 11	Histogram of the water level delta of Aghacashlaun - Aghacashlaun	23
Figure 12	Histogram of the time delta (in hours) of Aghacashlaun - Aghacashlaun	24
Figure 13	Line chart of Aghacashlaun - Aghacashlaun (whole data)	25
Figure 14	Line chart of Aghacashlaun - Aghacashlaun (2019-10-16 - 2019-10-19)	25
Figure 15	Precision vs Recall [30]	27
Figure 16	Best performance of outlier detection using mean (Crana - Tullyarvan)	36
Figure 17	Best performance of outlier detection using median (Sieg - Betzdorf)	37
Figure 18	Best performance of outlier detection using MAD (Losse - Helsa)	38
Figure 19	Best performance of outlier detection using z-score (Crana - Tullyarvan)	39
Figure 20	Best performance of outlier detection using z-score (Crana - Tullyarvan) zoomed in	40
Figure 21	Best performance of outlier detection using MADN-z-score (Aghacashlaun - Aghacashlaun)	41
Figure 22	Best performance of outlier detection using MADN-z-score (Aghacashlaun - Aghacashlaun) zoomed in	42
Figure 23	Result of median-based outlier detection Murg - Frauenfeld	45
Figure 24	Result of median based outlier detection Sieg - Betzdorf	47

List of Tables

Table 1	Example of IoT sensor data	4
Table 2	Example of IoT sensor data after cleaning	5
Table 3	Example of IoT sensor data after cleansing	5
Table 4	Seven number summary of Aghacashlaun - Aghacashlaun (all values)	19
Table 5	Seven number summary of Aghacashlaun - Aghacashlaun (regular values)	19
Table 6	Seven number summary of Aghacashlaun - Aghacashlaun (outlier values)	19
Table 7	Example of a binary confusion matrix	26
Table 8	First 9 values of Aghacashlaun - Aghacashlaun	28
Table 9	Best parameters of the average F1-score of all stations tested	35
Table 10	Top predictions summary of Aghacashlaun - Aghacashlaun	43
Table 11	Top predictions summary of Murg - Frauenfeld	44
Table 12	Top predictions summary of Sieg - Betzdorf	46
Table 13	Top predictions summary of Losse - Helsa	48
Table 14	Top predictions summary of Crana - Tullyarvan	49

List of source codes

1	Request body to get API key	14
2	Example response of the list endpoint	15
3	Example response of historical water level data for one station	16
4	The first step of classifying outliers using the mean	29
5	The first step of classifying outliers using the median	30
6	The first step of classifying outliers using the MAD	31
7	The first step of classifying outliers using the z-score	32
8	The first step of classifying outliers using the modified z-score (MADN-z-score)	33
9	Manual Outlier Detection Webapp using Dash [34]	53
10	Helper method to size plots [34]	54
11	Explorative Data Analysis [34]	61
12	Data preprocessing code in Python [34]	61
13	Implementation of Outlier Detection Methods [34]	65
14	Threshold based outlier detection [34]	67
15	Calculate the outlier detection results for different thresholds [34]	69

List of Abbreviations

ANN Artificial Neuronal Network

ARIMA AutoRegressive Integrated Moving Average

CSV Comma Separated Values

FN False Negative

FP False Positive

GP Gaussian Process

GRU Gated Recurrent Unit

HMAC Keyed-Hashing for Message Authentication

IQR InterQuartile Range

IoT Internet of Things

JSON Java Script Object Notation

KPI Key Performance Indicator

LSTM Long Short-Term Memory

Lat Latitude

Long Longitude

MADN Normalized Median Absolute Deviation

MAD Median Absolute Deviation

N Negative

P Positive

TN True Negative

TP True Positive