# READ ME Data Version 0.

## Content

## Introduction

This is the third version of the dataset. This is a bugfix. We had a problem with the normalization of the traffic value.
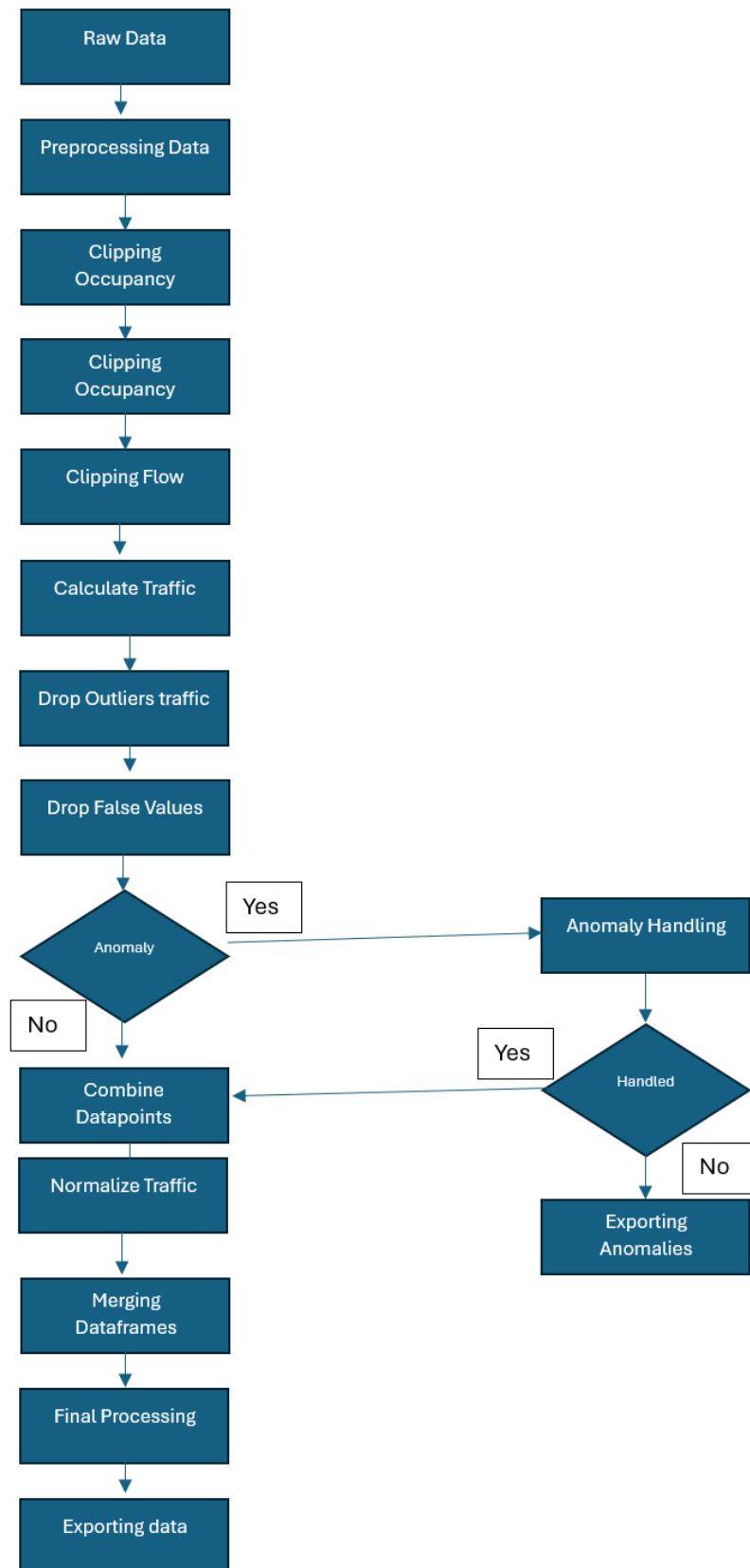
## Release

Samuel Paul 10.11.2024 21:00

# Changes

There are changes to the data data preparation.

## Data Preparation

- Datapoints are now combined before they are normalized.
- Datapoints are rounded not longer just converted to integers.

# UML

```
┌─────────────────┐
│    Raw Data     │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│Preprocessing Data│
└─────────────────┘
         │
         ▼
┌─────────────────┐
│    Clipping     │
│    Occupancy    │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│    Clipping     │
│    Occupancy    │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│  Clipping Flow  │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│ Calculate Traffic│
└─────────────────┘
         │
         ▼
┌─────────────────┐
│Drop Outliers traffic│
└─────────────────┘
         │
         ▼
┌─────────────────┐
│Drop False Values│
└─────────────────┘
         │
         ▼
      ◇ Anomaly ◇ ──── Yes ────▶ ┌─────────────────┐
         │                        │ Anomaly Handling │
        No                        └─────────────────┘
         │                                 │
         ▼                                 ▼
┌─────────────────┐        Yes        ◇ Handled ◇
│    Combine      │◀──────────────────────│
│   Datapoints    │                      No
└─────────────────┘                       │
┌─────────────────┐                       ▼
│ Normalize Traffic│             ┌─────────────────┐
└─────────────────┘             │    Exporting     │
         │                       │    Anomalies     │
         ▼                       └─────────────────┘
┌─────────────────┐
│    Merging      │
│   Dataframes    │
└─────────────────┘
         │
         ▼
┌─────────────────┐
│ Final Processing│
└─────────────────┘
         │
         ▼
┌─────────────────┐
│  Exporting data │
└─────────────────┘
```

# Parameters

In the create_dataset.py script there are a lot of parameters that can be changed. Does are the ones that got set for this version.

## Data Cleaning

Clipping occupancy outlier factor: 3

Clipping flow outlier factor: 3

Dropping traffic outlier factor: 2

## Anomaly Detection

Mean out of bound factor: 3

IQR to small, min IQR: 5

Not enough data, min datapoints: 4000

## Anomaly Handling

The anomaly handling has now only one function. There will be more in the future.

### Detectors with bad days

Minimum datapoints per day: 230

Minimum good days: 10

Minimum one good day for every weekday

# Script output

Starting script

Loading data from: C:\Users\samue\OneDrive\AIML\HS2024\Data Sicence Projekt\Data\London\London_UTD19.csv

Loading data from: C:\Users\samue\OneDrive\AIML\HS2024\Data Sicence Projekt\Data\London\London_detectors.csv

Data loaded

Preprocessing data

Errors found and dropped: 12234005

Preprocessing data took 12 seconds

Drop bad days

Total outliers detected and removed: 9504

Drop bad days took 177 seconds

Clipping outliers on occ

Total outliers clipped: 81

Clipping outliers on occ took 179 seconds

Clipping outliers on flow

Total outliers clipped: 57

Clipping outliers on flow took 160 seconds

Calculating traffic

Calculating traffic took 1 seconds

Droping outliers on traffic

Total outliers dropped: 661222

Droping outliers on traffic took 23 seconds

Drop false values

Total outliers detected and removed: 0

Drop false values took 145 seconds

Detecting anomalies

Anomalies detected based on IQR: 69

Anomalies detected based on IQR too small: 139

Anomalies detected based on not enough data: 1783

Detecting anomalies took 10 seconds

Handling anomalies

Anomalies with not enough data handled: 686

Total amount of dropeed anomalies: 1163

Handling anomalies took 366 seconds

Exporting anomalies to:  C:\Users\samue\OneDrive\AIML\HS2024\Data Sicence Projekt\Data

Exporting anomalies took 0 seconds

Combine datapoints

Combine datapoints took 6 seconds

Normalizing traffic

Normalizing traffic took 0 seconds

Merging dataframes

Merging dataframes took 1 seconds

Final processing

Final processing took 0 seconds

Exporting modified dataset to:  C:\Users\samue\OneDrive\AIML\HS2024\Data Sicence Projekt\Data

Exporting modified dataset took 14 seconds

Script finished

Total script execution time: 1130 seconds