
DS203 : EXERCISE 6

Nirav Bhattad

Last updated October 9, 2024

§1. Problem 1

1.1. Introduction

This report presents an analysis of the HT R Phase Current dataset. The primary goal is to identify an unstable period within the data and apply various outlier handling techniques, including imputation, trimming, capping, and robust estimation. The findings are illustrated with visualizations, and a statistical comparison of the methods is provided.

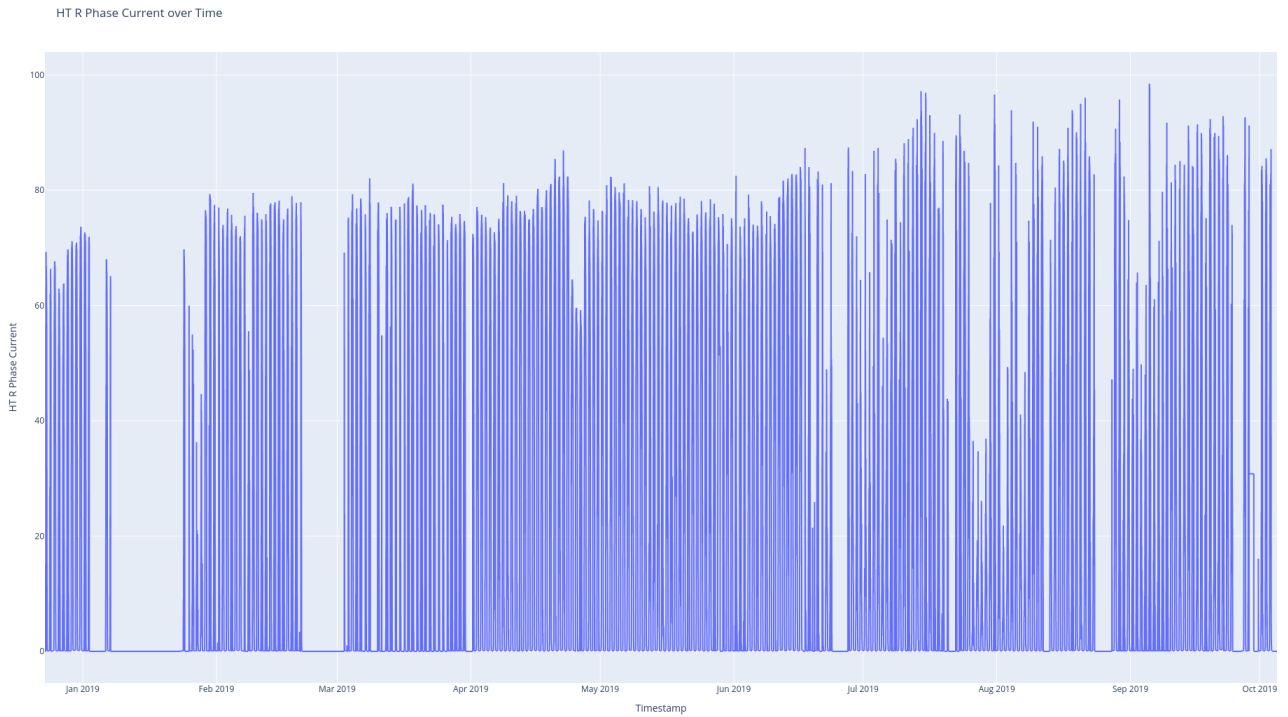


Figure 1: HT R Phase Current Dataset

1.2. Data Exploration

The dataset contains 1,000 records with a mean current of 16.912767, a standard deviation of 27.174448, and a range from 0 to 98.500000. The data is collected over a period of 9 Months, with a frequency of 1 record every 5 minutes.

```

1 import pandas as pd
2 import numpy as np
3 import plotly.express as px
4 import matplotlib.pyplot as plt
5 from scipy import stats
6 import os
7
8 # Create directory for saving images
9 if not os.path.exists('Images'):
10     os.makedirs('Images')
11
12 # Load data from CSV
13 data = pd.read_csv("e6-htr-current.csv", parse_dates=['Timestamp'], dayfirst=True)
14
15 # Part A: Perform EDA
16 print(data.describe())
17
18 # Plot the current over time using Plotly
19 fig = px.line(data, x='Timestamp', y='HT R Phase Current', title='HT R Phase Current over Time')
20 fig.write_image("Images/ht_r_phase_current.png", scale=1, width=1920, height=1080, format='png')
21 fig.show()
22
23 # Part B: Identify a 2-week unstable period
24 unstable_period = data[(data['Timestamp'] >= '2019-07-30') & (data['Timestamp'] <= '2019-08-14')].
25     copy()
26
27 # Plot the unstable period
28 plt.figure(figsize=(16, 10))
29 plt.plot(unstable_period['Timestamp'], unstable_period['HT R Phase Current'], color='blue')
30 plt.title('HT R Phase Current - Unstable Period')
31 plt.xlabel('Timestamp')
32 plt.ylabel('Current')
33 plt.savefig('Images/unstable_period.png', dpi=300)
34 plt.show()

```

Using the plot above, we can identify an unstable period between July 30, 2019, and August 14, 2019. This period will be used to evaluate the outlier handling techniques.

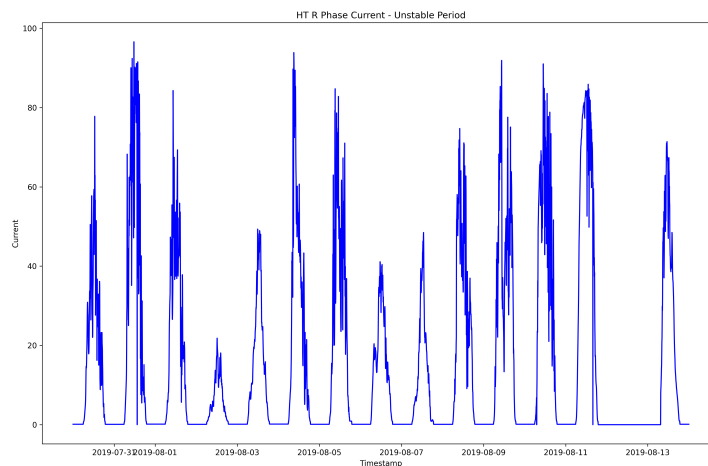


Figure 2: HT R Phase Current - Unstable Period

1.3. Removing outliers, smoothening, and imputing missing data

Here we apply four outlier handling techniques to the unstable period: imputation, trimming, capping, and robust estimation. The results are visualized to compare the effectiveness of each method.

```

1 # Method 1: Imputation (Replacing outlier values with the mean)
2 mean_value = unstable_period['HT R Phase Current'].mean()
3 median_value = unstable_period['HT R Phase Current'].median()
4 unstable_period['Current_Imputed_Mean'] = unstable_period['HT R Phase Current'].copy()
5 unstable_period.loc[(unstable_period['HT R Phase Current'] > unstable_period['HT R Phase Current'].
    quantile(0.95)) | (unstable_period['HT R Phase Current'] < unstable_period['HT R Phase Current'
    ].quantile(0.05)), 'Current_Imputed_Mean'] = mean_value
6
7 unstable_period['Current_Imputed_Median'] = unstable_period['HT R Phase Current'].copy()
8 unstable_period.loc[(unstable_period['HT R Phase Current'] > unstable_period['HT R Phase Current'].
    quantile(0.95)) | (unstable_period['HT R Phase Current'] < unstable_period['HT R Phase Current'
    ].quantile(0.05)), 'Current_Imputed_Median'] = median_value
9
10 # Method 2: Trimming (Removing outliers)
11 q_low = unstable_period['HT R Phase Current'].quantile(0.1)
12 q_high = unstable_period['HT R Phase Current'].quantile(0.9)
13 unstable_period_trimmed = unstable_period[(unstable_period['HT R Phase Current'] >= q_low) & (
    unstable_period['HT R Phase Current'] <= q_high)]
14
15 # Method 3: Capping (Setting a cap on the maximum and minimum values)
16 max_value = unstable_period['HT R Phase Current'].quantile(0.95)
17 min_value = unstable_period['HT R Phase Current'].quantile(0.05)
18 unstable_period['Current_Capped'] = unstable_period['HT R Phase Current'].copy()
19 unstable_period['Current_Capped'] = np.where(unstable_period['Current_Capped'] > max_value,
    max_value, unstable_period['Current_Capped'])
20 unstable_period['Current_Capped'] = np.where(unstable_period['Current_Capped'] < min_value,
    min_value, unstable_period['Current_Capped'])
21
22 # Method 4: Robust Estimation (Using RANSAC regression)
23 from sklearn.linear_model import RANSACRegressor, LinearRegression
24
25 # Prepare the data for RANSAC
26 X = unstable_period['Timestamp'].map(pd.Timestamp.timestamp).values.reshape(-1, 1) # Convert
    datetime to timestamp
27 y = unstable_period['HT R Phase Current'].values
28
29 # Fit RANSAC
30 model = RANSACRegressor(LinearRegression()).fit(X, y)
31 unstable_period['Current_Robust'] = model.predict(X)

```

I have used the following techniques to handle the outliers:

- **Imputation:** Replacing outlier values with the mean or median of the data.
- **Trimming:** Removing outliers from the dataset.
- **Capping:** Setting a cap on the maximum and minimum values.
- **Robust Estimation:** Using RANSAC regression to estimate the data without outliers.

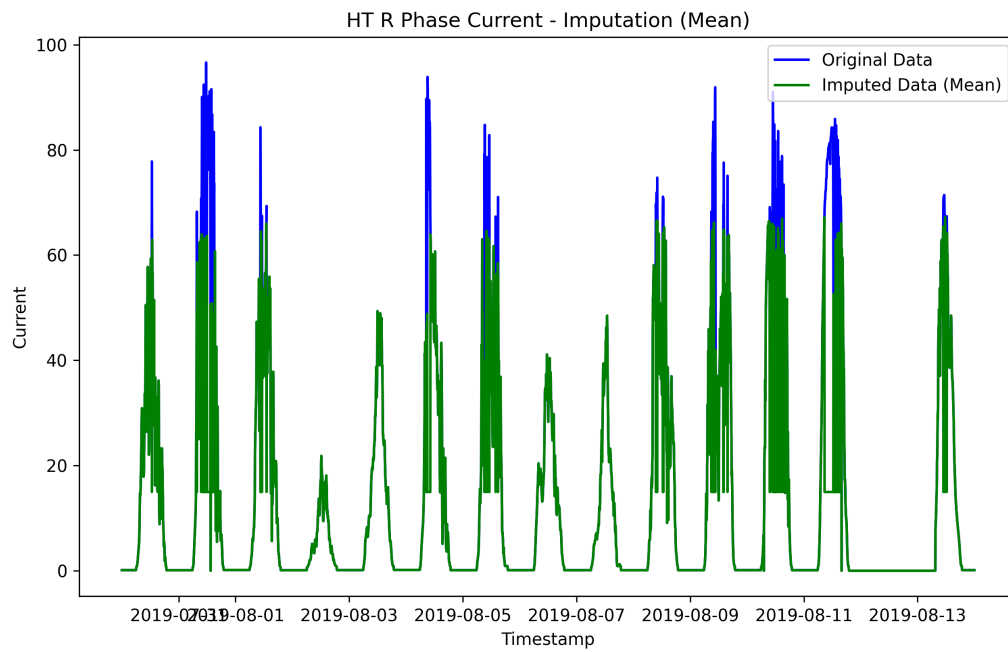


Figure 3: Imputed Mean Data

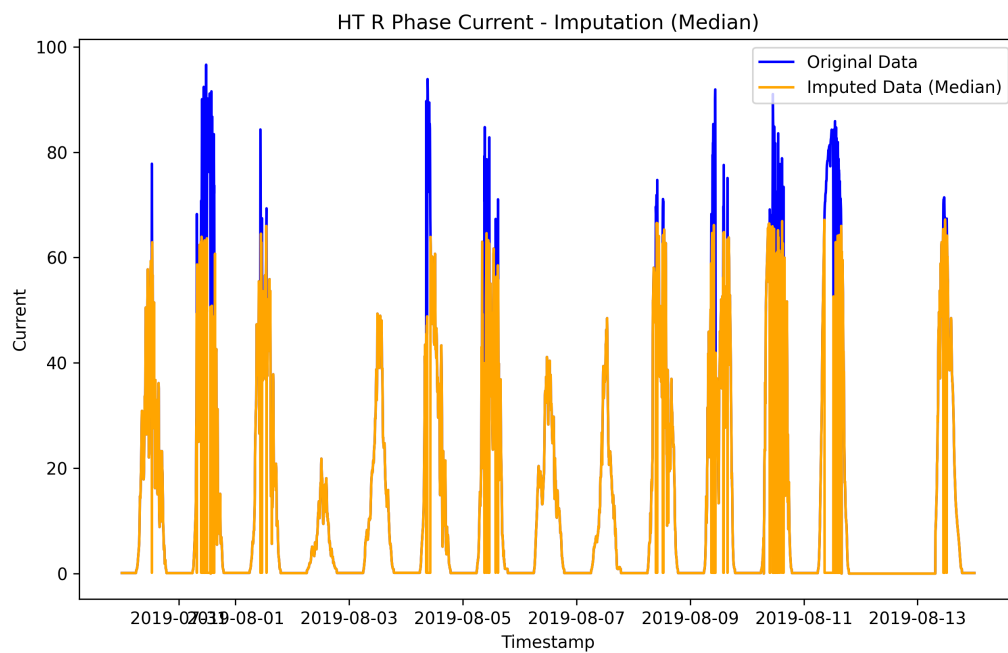


Figure 4: Imputed Median Data

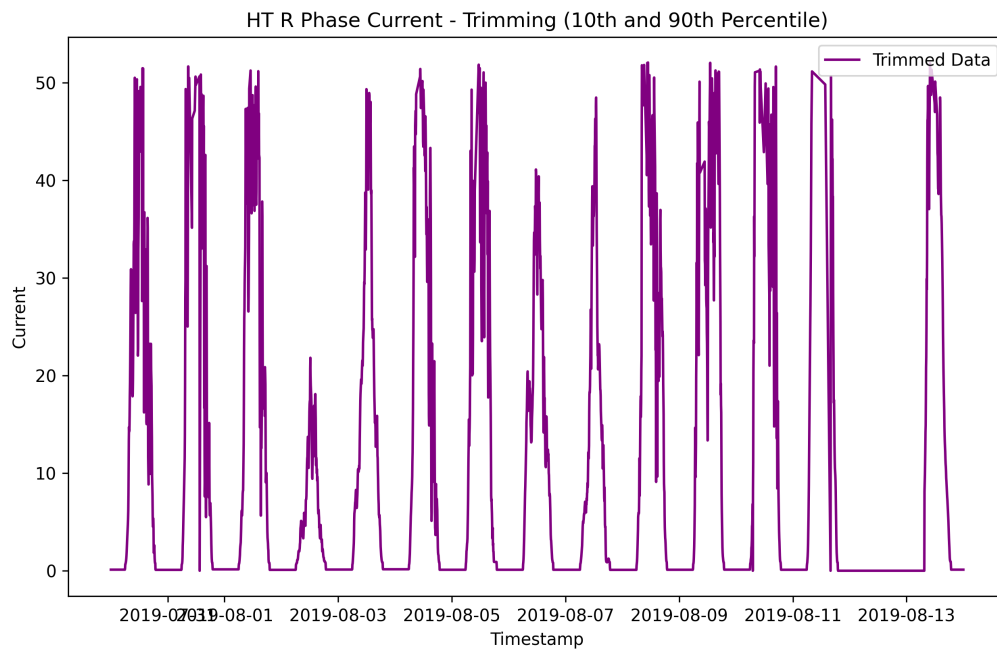


Figure 5: Trimmed Data

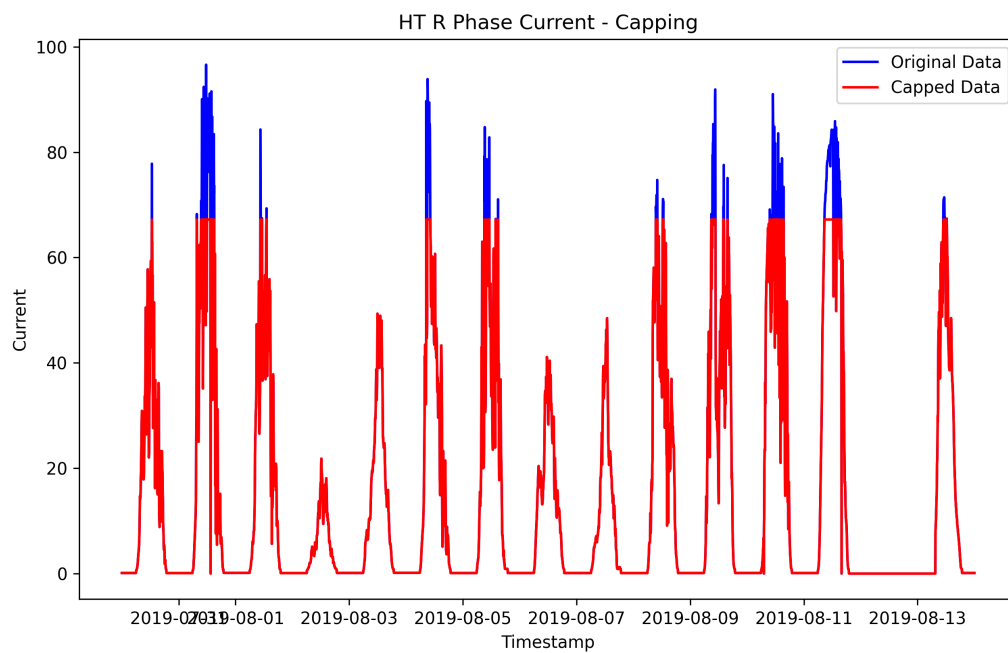


Figure 6: Capped Data

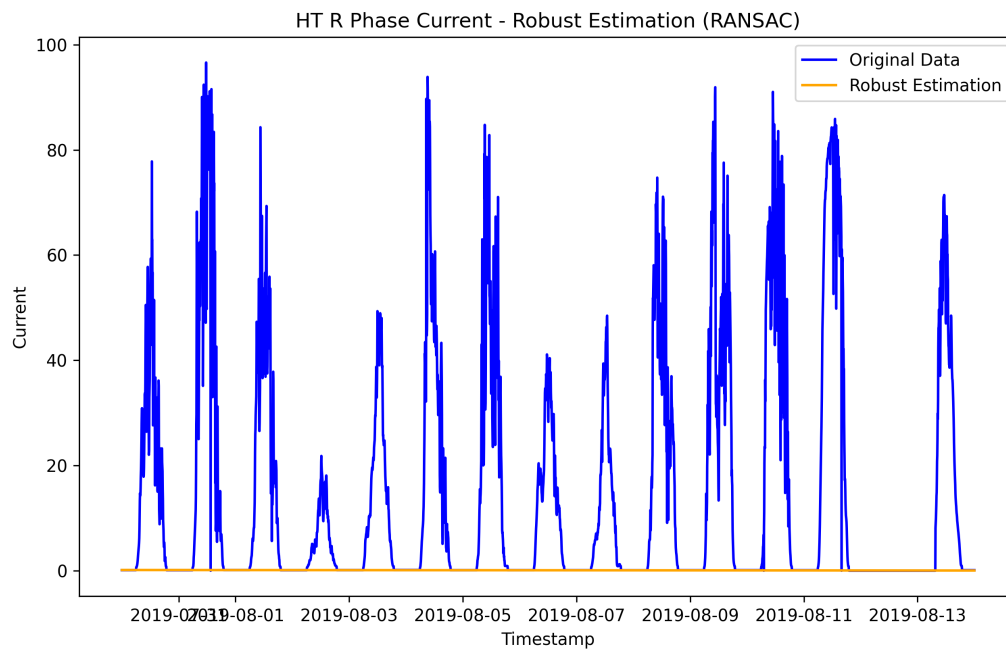


Figure 7: Robust Data

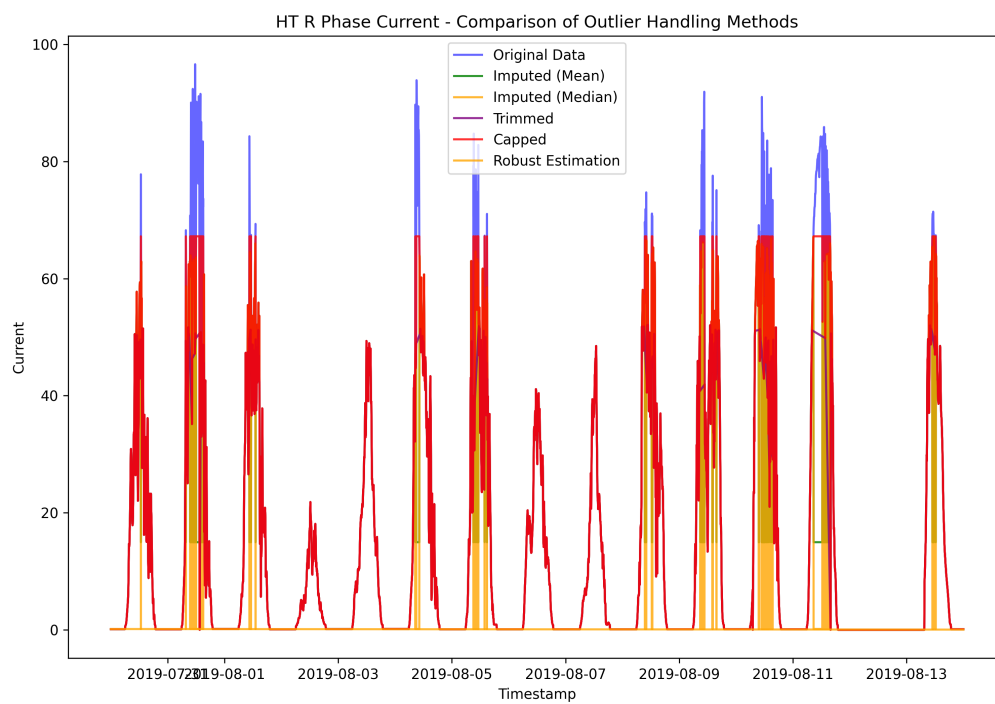


Figure 8: Outlier Handling Comparison

§2. Problem 2**§3. Problem 3**