

Analysis and Denoising of Time Series Data Using Unsupervised Methods

Shafkat Waheed¹⁺, B. M. Raihanul Haque¹⁺, Abir Roy¹, Mst. Shapna Akter¹, Sumit Ranjan Chakraborty², Shakran Hayat², and M.R.C. Mahdy^{1,2*}

¹ *Department of Electrical & Computer Engineering, North South University,
Bashundhara, Dhaka, 1229, Bangladesh.*

² *Pi Labs Bangladesh Ltd., Eden Center, 2/1/E, Toyenbee Rd, Dhaka 1000, Bangladesh.*

+ Equal contribution

*Corresponding Author: mahdy.chowdhury@northsouth.edu

ABSTRACT:

A dataset, collected under an industrial setting, often contains a significant portion of noises. In many cases, using trivial filters is not enough to retrieve useful information i.e., accurate value without the noise. One such data is time-series sensor readings collected from moving vehicles containing fuel information. Due to the noisy dynamics and mobile environment, the sensor readings can be very noisy. Denoising such a dataset is a prerequisite for any useful application. It has led us to develop a system that can remove noise and keep the original value and help vehicle industry, fuel station, and power-plant station that require fuel. In this work, we have only considered the value of fuel level, and we have come up with a unique solution to filter out the noise of high magnitudes using several algorithms such as interpolation, extrapolation, spectral clustering, agglomerative clustering, wavelet analysis, and median filtering. We have also employed peak detection and peak validation algorithms to detect fuel refill and consumption in charge-discharge cycles. We have used the R-squared metric to evaluate our model, and it is 98%. In most cases, the difference between detected value and real value remains within the range of $\pm 1L$.

Keywords: Interpolation, Extrapolation, Spectral clustering, Agglomerative clustering, Wavelet analysis, Median filtering.

1. Introduction

Data denoising or removing noise from a dataset is a challenging and fascinating topic for the researchers. Each dataset is different in terms of data distribution and noise. Dataset might be a set of images or audio signals or any other generic data. These noises are generated due to device malfunctioning, human limitation, machine limitation, improper approaches of collecting and preserving data, etc.

However, in this work, the focus revolves around the industrial dataset that has issues like the severely noisy environment, and this dataset contains fuel level information found in a fuel tank. Whether this is a vehicle industry or fuel station or a power plant, which is being operated using fuel, they have almost similar problems. Since this is a time-series data, any given point within a certain time frame either indicates charging or discharging, and noise present in the dataset gives a false representation of peak or consumption. Here, false representation stands for noisy data that plagues the actual value. For instance, at any given time, the value obtained from sensors could be 40 liters, whereas the actual value could be of 60 or 30 liters. Using a typical filter while cleaning the dataset, causes data loss, and fails to remove noise properly. In an attempt to reduce noise from time-series data representing the fuel volume of a vehicle; some related works have been discussed below.

So far, most of the researches has been conducted on medical images [1-5]. In all these papers, the dataset comprises MRI images, X-ray images, CT scan images, PET scan images, and so on. Some of the proposed methods incorporate NN (neural network), SVM (support vector machine), denoising autoencoder, etc. The problem with these methods is that: in the case of a complete unsupervised scenario, all the proposed models, especially NN, perform poorly. In our case, implementing a device to a car to get accurate real-time values has not been feasible. Therefore, NN has not been an option to deduce real value. For similar reasons, the use of LSTM (long short term memory), GAN (generative adversarial network), RNN (recurrent neural network), and other classifiers, which have been proposed by some authors [6-8], have not been possible either.

Another notable area where denoising techniques have been implemented is audio signals [9-12]. Audio signals usually come from speech, background noises, or any other automated source (electronic, electromagnetic, acoustic). The suggested methodology is heavily dependent on filters such as Kalman filter, Butter Worth, Chebyshev, Elliptical Filters, and so on. In the medical sector, ECG (electrocardiogram) signals are also cleaned using filters such as low pass filter [13]. Apart from these, deep learning methodologies [9] and non-linear diffusion [12] have also been considered. We have already ruled out a deep learning technique due to the unsupervised nature of the data, which we have dealt in this work. The other problem is that: relying just on filters produces poor results due to high regularization. In

contrast, our model tends to be generic, which means it will work efficiently on any vehicle with ambiguous noise patterns.

It is evident that: the contemporary techniques are focused on denoising images or audio signals. Although denoising techniques are applied in the realm of IoT (internet of things) by using ANN (artificial neural network) [14], ToF (Time-of-Flight) data by using GAN [15], and vibration sensor data by using TFM (time-frequency manifold) [16], none of the solutions seems suited for our dataset due to unsupervised nature of the data. In our case, we had to tackle two parallel problems; one is to minimize data loss, and another is to get the actual value in a completely unsupervised manner.

Among the few works in literature, Manmohan et al. [17] have thoroughly studied and implemented the process of wavelet analysis. They have also proposed a method which includes the median filter [17]. They have suggested that: ordinary median filter indicates an improvement against other filtering techniques i.e., the mean filter. Also, two thresholding methods and wavelet allow them to retrieve original features of the image effectively. Additionally, implementing median filtering is suggested in reference [18]. This inspires us to implement wavelet analysis and the median filter in our dataset.

Among many other techniques, clustering has been implemented in amplicon sequencing by Kaisa Koskinen et al. [19]. The impact of clustering is huge on the number of operational taxonomic units (OTUs). Priyam Chatterjee et al. have suggested the positive impact of clustering in data denoising as well [20]. Therefore, clustering has been considered as an imperative option for our dataset. For time series and other types of data, such as for predicting time-averaged force [21], forecasting price of the market [22-23], or classify visual pollutants from an image [24], machine learning and deep learning can also serve as highly effective tools. But for our case, due to the aforementioned reasons, a different approach has been implemented.

To test the efficiency of our model, we have been provided with a specific industrial dataset. This dataset contains the charging-discharging cycle value of the moving vehicle. All the vehicles have been operative on the street in Bangladesh. Due to the poor condition of a few streets in Bangladesh, noise may have been prominent in the dataset. Since, barely any work has been done in denoising time-series data containing fuel information retrieved from the industrial domain; therefore, to accomplish this, we have tried many methodologies to find out the optimal solution. Our system has started with simple interpolation, extrapolation, and white noise removal techniques. Here, by white noise, we are referring to random noise present in a dataset whose mean value is zero and has a finite variance [25]. Then two clustering methods [26-27] have been brought to isolate noise from real value. Wavelet analysis [28-30] has cleaned it further through wavelet transformation and inverse transformation. Finally, median filtering

[31] has been used to give the final push. Furthermore, many custom algorithms have been used to arrange previously mentioned algorithms and compare different datasets obtained from those algorithms to narrow down the refill value.

Apart from the unorthodox noises caused by the reason mentioned above, there is also an unusual noise present in the dataset. Usually, when the engine or machine is switched off, there is no charging or discharging occurring within a certain time period. This yields a stationary or constant set of values within that time frame. Our system has required to deal with this issue, and for most of the parts, it has accomplished successfully.

At first, we have developed our model using a concatenated dataset containing information from multiple vehicles. Then, to test it further we have been provided with an additional dataset of various vehicles. The initial dataset contains over 100,000 data points, which have been narrowed down to 37 peak values. To evaluate our result properly, we have used the R-squared metric and compared our result with previous work that has been done on similar datasets [32-38]. The R-squared value of the final result is 0.98, and RMSE is 1.49. Our system detects not only the refill of fuel but also the consumption rate while minimizing data loss in a very noisy industrial environment.

2. Methodology

All the data have been collected from Pi Labs Bangladesh Ltd. Initially, from 9 vehicles, roughly 1,00,000 data points have been collected and used for our system. Data flow is shown in Fig.1, which demonstrates the cycle starting from raw data to processed data. It begins with employing industrial settings in a moving vehicle. Then some generic transmission occurs. Our algorithm can reduce noise from the time-series data while it is streamed from source to destination through a server.

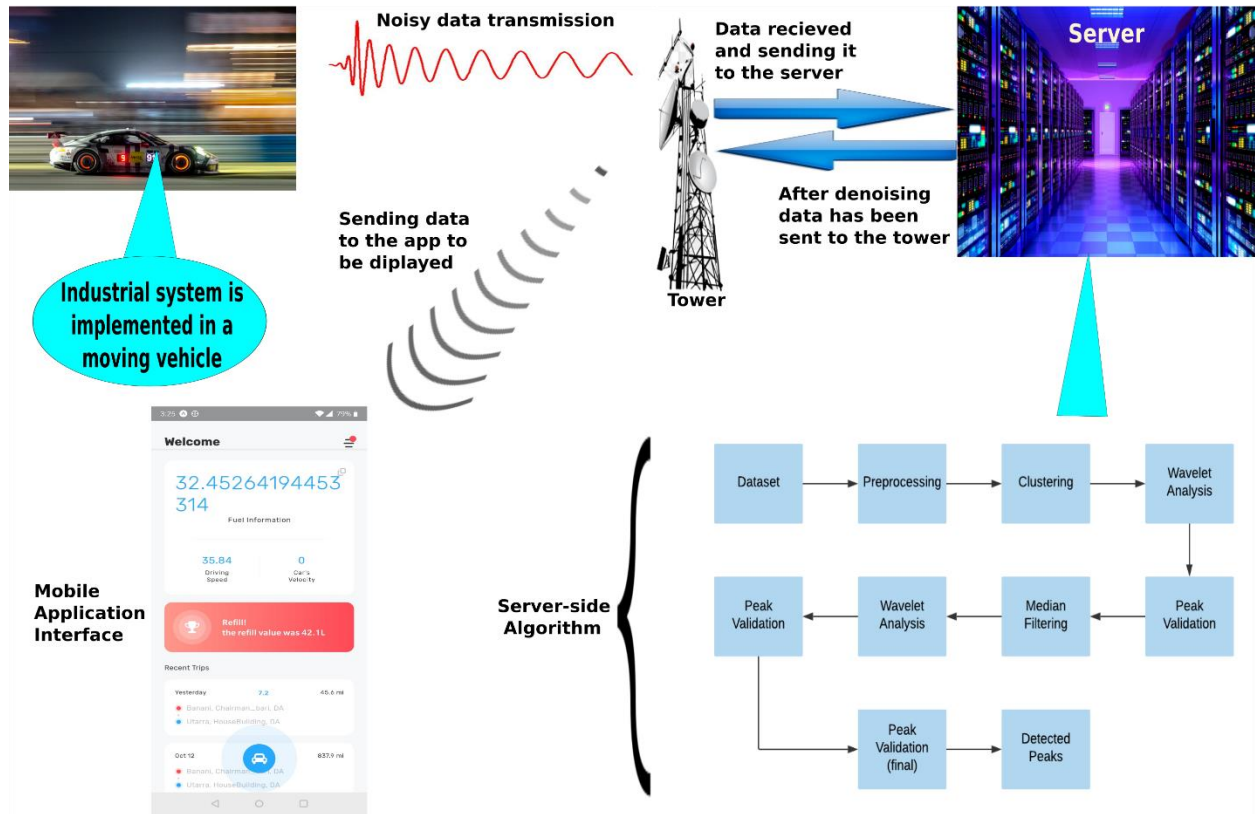


Fig.1 A general overview of the data cycle. The industrial setting is implanted in a vehicle, and it collects and sends data to the server. Inside the server, our algorithms are implemented where the data traverse through each phase. After successful peak detection, the data is sent to the mobile application (beta version) from where the result can be displayed.

The scope of this project revolves around the algorithms. We wanted to test whether real-time data can be displayed, and this is why a mobile application prototype has been built. An enhanced overview of our algorithm is represented in Diagram.1. Our system begins cleaning the dataset with simple extrapolation and interpolation. This removes some trivial types of noises i.e., white noise. Then two sets of clustering methodologies (spectral and agglomerative) isolate noise from the actual value. After that, the dataset is stored and sent to the next phase to perform wavelet analysis. This generates two datasets to be compared and keep the common values. A similar strategy is undertaken for median filtering and wavelet analysis. Finally, when all prior operation is done, we have performed a final validation to get the peak values and consumption rate.

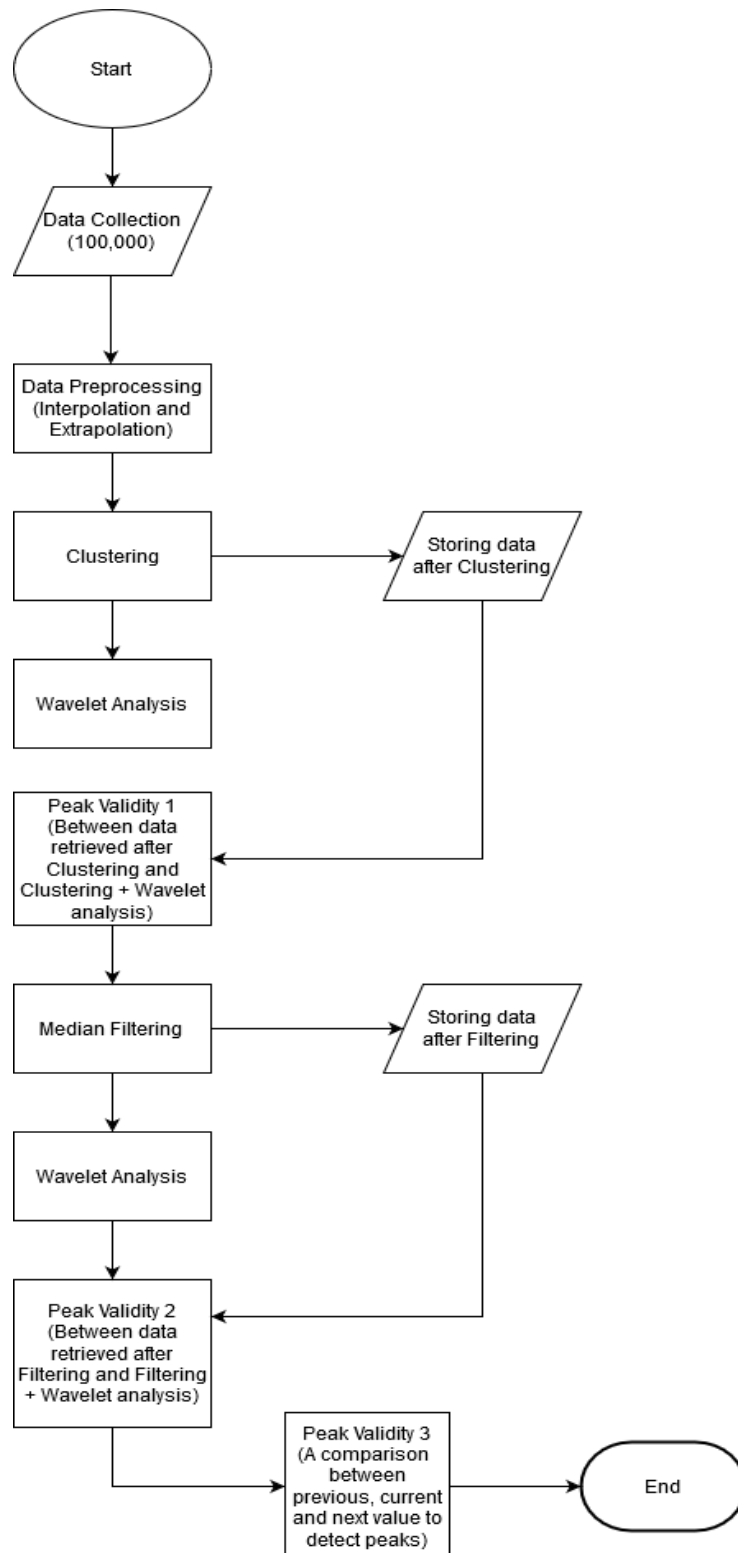


Diagram.1 After receiving the data, it has been sent for preprocessing, where we have used a few data wrangling techniques to deal with missing values and white noises. This data is then stored, and data is sent to the next phase for wavelet analysis. These two sets of data are being compared to eliminate redundant data points. A similar process can be observed for median filtering and wavelet transformation. Lastly, a final comparison is made to detect peaks.

Initially, the dataset is comprised of a massive amount of noise. As a result, it has been hard to isolate them from the actual value just through observation. Fig.2 (a) gives a generic outlook of the data set.

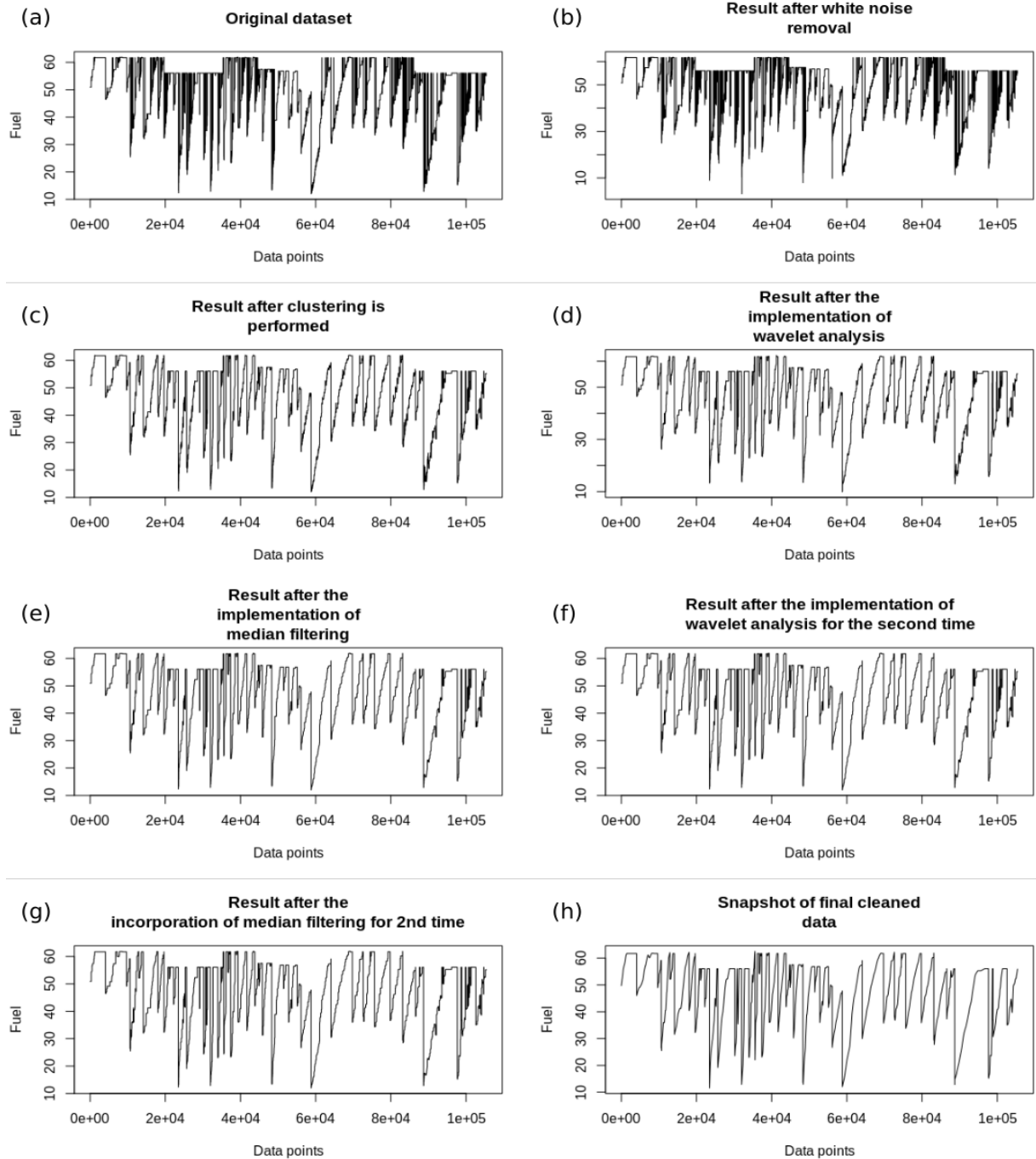


Fig.2 The outlook of our dataset after each step from initial to processed. The congested regions are the indication of data noise. (a)-(b) Raw data and data after white noise removal, respectively. In both cases, heavy noise is present in the dataset. (c)-(g) Graphical representation of data after clustering, wavelet analysis (1st), median filtering (1st), wavelet analysis (2nd), and median filtering (2nd). A significant portion of the noise is removed in each step. (h) Processed data, containing 37 peaks.

In this figure, from right to left, we see downward spikes. This slow decrease indicates the consumption rate over a certain time period. The sharp jump from the spike represents refuel. Each segment in this graph illustrates the individual vehicle's fuel level.

2.1 Data characteristics

Before designing a model, it is important to know some characteristics of the data so that similar work can be replicated, or the designed module can be implemented on a similar data set. The goal was to plot a histogram within a certain frequency. As seen in Fig.3, there are multiple intervals and corresponding frequency where data points are grouped within a certain range.

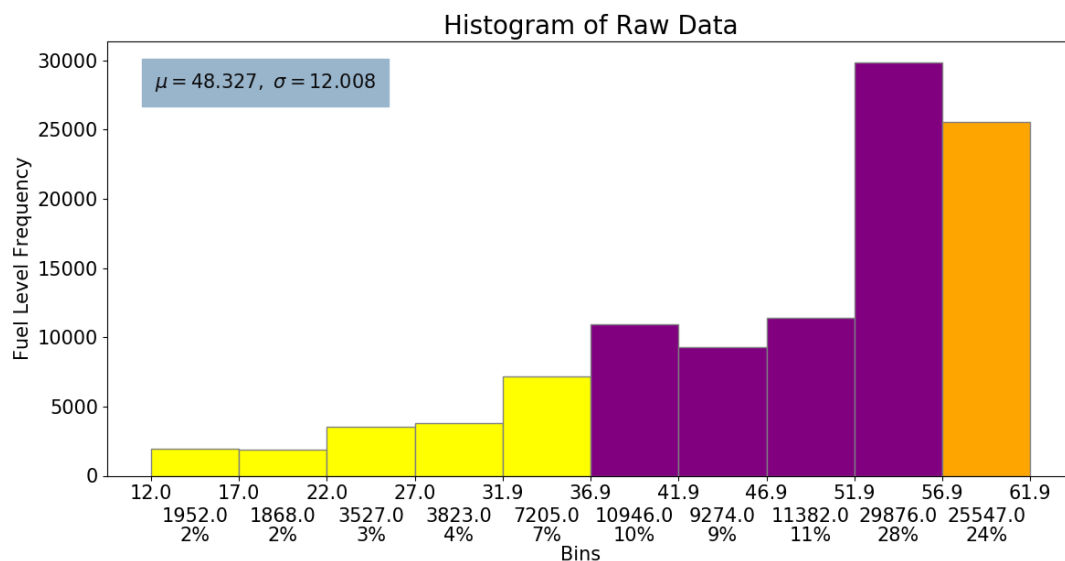


Fig.3 Histogram shows the distribution of the dataset

To understand the characteristics better, we have calculated the mean and median of the data along with the 1st and 3rd quartile range. Our objective was to find out whether our median is closed to the 1st or 3rd quartile range. Our mean = 48.327, standard deviation = 12.008, median = 52.316, 1st quartile range = 40.790 and 3rd quartile range = 56.837. Obviously, the median is close to the 3rd quartile range, and the mean is less than the median. We have calculated the moments of the dataset to identify the distribution. But we could not find any specific distribution related to the dataset, which is why we have undertaken an unsupervised approach to denoise the dataset.

2.2 Data Preprocessing

Initially, some rows contained zero values. This could contaminate the results of further processes, and therefore data needed to be tackled in an efficient manner where data replacement conforms to the pattern

of non-zero values. Preprocessing involves white noise removal using various data wrangling techniques and applying interpolation-extrapolation to recover the lost points. Also, we have differentiated the dataset to create a new feature.

Interpolation is a method that can be achieved using discrete data set and by forming general formula within a certain range. In our case, linear interpolation was used, which can be expressed as $y = y_1 + \left(\frac{y_2 - y_1}{x_2 - x_1}\right)(x - x_1)$. Here, unknown points are denoted as y for a certain value of x .

Extrapolation operates outside of the range of the observation. It is subjected to higher uncertainty, as opposed to interpolation, since it fills out data based on the relation of other data points. In this research, midpoint extrapolation was incorporated, denoted as $(x, y) = \left(\frac{x_1 + x_2}{2}, \frac{y_1 + y_2}{2}\right)$, and which is operative for one unknown point.

Interpolating and extrapolating data points were achieved through the following process. At the beginning, it is checked whether the data points contain 0 values or not. If it does so, all those values are set to NULL. Finally, through interpolation and extrapolation, new data has been generated and is being stored.

Apart from that, a general form of noise also known as random noise, has been removed from the data set. The main characteristic of the noise includes the equal intensity of the signal at different frequencies. If the mean equals zero, then data is considered white noise, as data is spread across the negative and positive portion of the graph in an equal manner. The white noise is excluded from the dataset by using the following formula iteratively: $\alpha = (x_i - x_{i+1}) + (x_i - x_{i-1})$. However, there is a strict condition for white noise removal, which is, $x_i \neq x_{i+1}$ and $x_i \neq x_{i-1}$. The iterative process can be represented as $y_i = x_i \alpha$. After successful completion of preprocessing, we have gotten a visual representation of the data set, as presented in Fig.2 (b).

2.3 Clustering

Clustering is an unsupervised learning method. It is used to group data according to their respective characteristics. Any data in a single cluster possesses a similar characteristic as every other data point in the same cluster. Right after the preprocessing is done, two clustering techniques have been incorporated. These clustering methods work together to separate noisy data from correct data.

Spectral clustering is best suited when the variation of data is not much. It clusters the data based on the density of data points. The state of being close proximity to each other is called affinity, and this phenomenon can be described by the affinity matrix. Different vectors from this matrix can be

extracted using principal component analysis (PCA), which later leads Eigenvectors to be formed. These vectors are referred to as feature vectors of each object of affinity or Laplacian matrix.

Hierarchical clustering, also known as Agglomerative clustering, generates a tree where roots are the lowest point of the data-set, and leaf nodes consist of values that are greater than the values of the root node. It uses the bottom-up approach to group the data points in a hierarchical manner. Every data point groups together in its own cluster. This goes on for all data points, and these clusters are then joined using a greedy approach. The greedy approach involves merging two most similar clusters together.

To perform the clustering on the entire dataset we have incorporated the following formula, $y_{s,h} = \alpha f_s(x_i) + (1 - \alpha)f_h(x_i)$. Here, f_s and f_h represent the mechanism of spectral and agglomerative clustering respectively as functions. At the start of clustering, we set a constant value denoted as threshold T whose value is 0.1. Then standard deviation is being calculated based on all the data points. As agglomerative clustering uses a dendrogram to cluster that data it works well on the dataset having higher deviational value. As opposed to that, spectral clustering checks on the affinity of the data points to group them; That is why it is better suited for dataset having lower deviational value. Due to variation of standard deviation after each iteration, we used two clustering. If T is greater than standard deviation, then $\alpha = 0$ and data points are grouped using Agglomerative clustering, and if it is less then, α is set to 1, and data are grouped together using spectral clustering. Each iterative step is incremented by 200. After clustering, many data points get eliminated as part of the noise removal technique. Interpolation and extrapolation are used to fill out those data points. Fig.2 (c) is a visual representation of the effect of clustering on the data set.

2.4 Wavelet Analysis

Unlike any other Fourier transformation methodology, wavelet transformation is used for transforming data represented in the time domain to the frequency domain. In previous work, we have found that it has the ability to compress an image efficiently. By managing factors like shifting and scaling, it can decompose an image to multiple lower resolution images. The waves have features like varying frequency, limited duration, and zero average value. This is also eligible to remove high-frequency noise from any time-series data with non-continuous peaks. The implementation of wavelets revolves around implementing two different transformations and incorporating one threshold function. These transformations are wavelet transformation and inverse wavelet transformation. The wavelet transformation is achieved by the following formula, $C(\tau, s) = \frac{1}{\sqrt{s}} \int_t f(t) \psi^*\left(\frac{t-\tau}{s}\right) dt$.

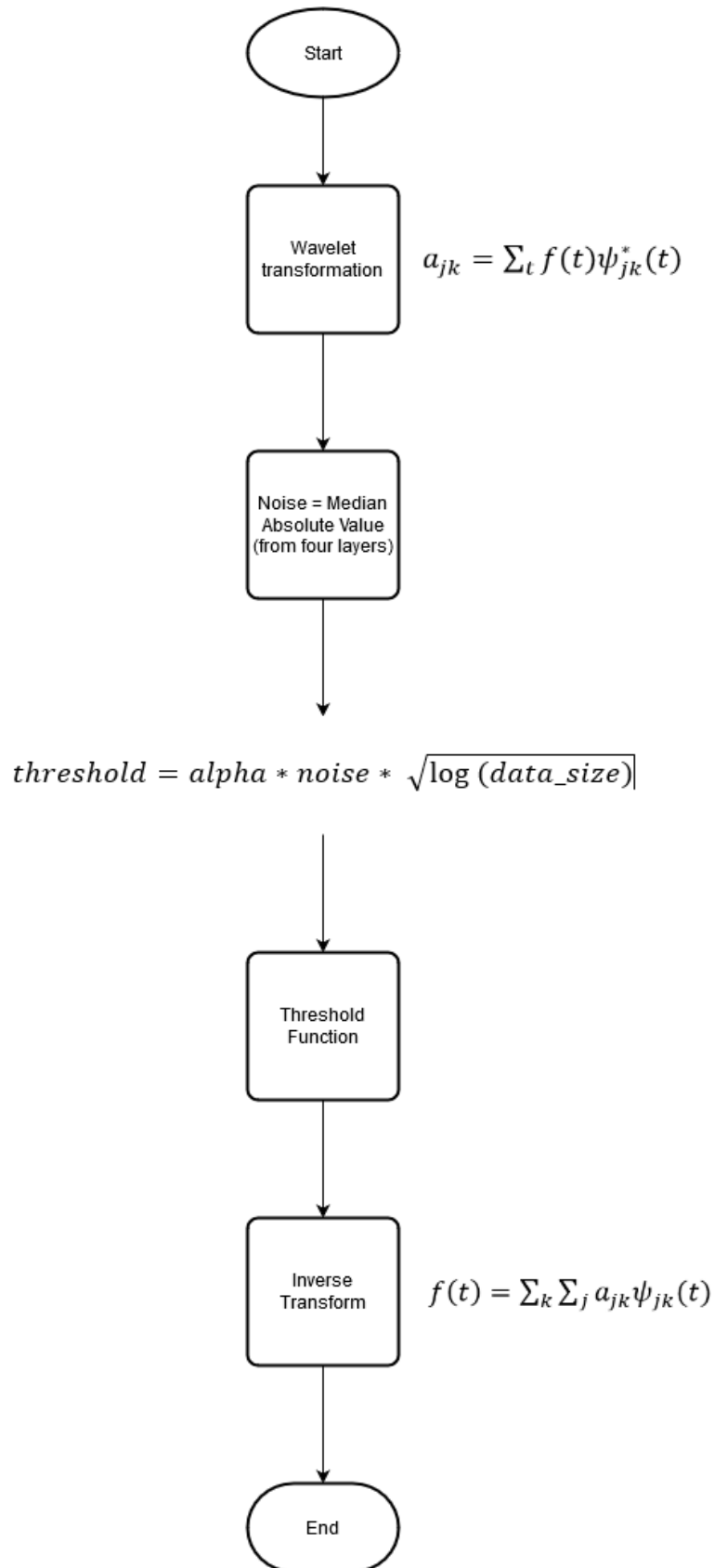


Diagram-2 An illustration of wavelet analysis containing two transformations.

Above is the formula for continuous wavelet transformation where τ and s are transition parameter and scale parameter respectively, $\frac{1}{\sqrt{s}}$ is normalization constant, and $\int_t f(t)\psi^*(\frac{t-\tau}{s})dt$ is the mother wavelet. The inverse operation is carried out by the following function: $f(t) = \frac{1}{\sqrt{s}} \int_{\tau} \int_s C(\tau, s)\psi(\frac{t-\tau}{s})d\tau ds$.

Discrete wavelet transformation is a bit straight-forward than this and, to state the obvious, free from the integral operation. The formula for discrete wavelet transformation is $a_{jk} = \sum_t f(t)\psi_{jk}^*(t)$ and inverse discrete wavelet transformation can be written as such, $f(t) = \sum_k \sum_j a_{jk}\psi_{jk}(t)$. These formulas provide simultaneous localization in time and scale, sparsity, adaptability and linear time complexity; which allow noise filtering, image compression, image fusion, recognition, image matching and retrieval efficiently. Finally, an additional threshold function can be represented as such to improve the proposed model, $threshold = \alpha * noise * \sqrt{\log(data_size)}$ and in this formula, α is a constant and $noise$ = absolute median value. We incorporated the discrete approach in our model.

This reserves the original signal but without the noise. After the implementation, the peak value was shifted by 5200 points. Diagram-2 provides the mechanism and Fig-2 (d) and (f) show the result after wavelet analysis.

2.5 Median Filtering

Median filtering is a non-linear filtering technique, and it removes noise from the data set while preserving the valuable data. In our case, the median filter is to traverse through the signal one by one. It generates a window to slide through the data entry, and this window is generated based on the pattern of the neighboring window. This replaces each entry with the median of neighboring entries and what remains is the peak. Fig.2 (e) and (f) show results after median filtering applied for the first time on the data set. Our data has one dimension. A window of the median filter contains the first few preceding and following entries. To carry out this task, we used this formula, $y(X_{top}, X_{bottom}) = M\{x_i, x_{i+1}\}$, where M =median, X_{top} =Initial point, X_{bottom} =Last point, and $\{x_i, x_{i+1}\} \in X$.

Once what algorithm to use is sorted out, we have been required to write our own custom algorithms where we can use clustering, filtering methodologies efficiently. It is necessary to choose data window and iterative steps carefully so that we can minimize the data loss. Our custom algorithms include one peak detection and three peak validation methods.

2.6 Peak detection

This is the part where peaks are detected. The algorithm uses the moving average to find any sudden changes in the data points. As we pass through data points, we calculate the simple moving average with $SMA = \frac{\sum x_i}{n}$. We used a window of 3 points to find any discrepancy in the data. So, for every three adjacent point $SMA_3 = \frac{\sum_{i=1}^{n=3} x_i}{3}$. Any deviation of more than 4 units is considered a peak, which means if $|x_i - SMA_3| > 4$, then x_i is a valid peak.

There are three peak validation methods in this system, among which two are almost identical. These validation methods validate the data points, whether they are valid peak points or not. The prior two peak validation methods include the following:

2.7 Peak validation (first and second)

At the start, peak validation data is sent to two different directions. One part is operative using clustering and peak detection algorithm. Another part is operative using a combination of clustering, wavelet analysis, and peak detection algorithm. This produces two different datasets. After that, each entry of one dataset is evaluated against the corresponding entry of other data set. If there is a similarity, we validate the peaks; otherwise, we discard it. This parallel processing gets the work done fast without interrupting each other. The formula for peak validation (first) is, $f_{rp} = |f_{pcp} - f_{pwp}| \leq 100$ where f_{rp} = valid peak points, f_{pcp} = data points obtained after performing peak detection on clustered data points = $SMA(y_{s,h}(x_i))$ and f_{pwp} = data points obtained after performing peak detection on data points retrieved from implementation of clustering and wavelet analysis (W_A) = $SMA(W_A(y_{s,h}(x_i)))$. In order to get valid peak f_{rp} must be less than or equal to 100. It is a comparison between values obtained by only after clustering and a combination of clustering and wavelet. Real peaks are validated based on similarity.

The second peak validation method, as mentioned earlier, works similarly. The only difference is that, in second peak validation, median filtering has been used instead of clustering. Likewise, it is a comparison between the two datasets. One is obtained after performing just filtering, and another is acquired through implementing a combined system of filtering and wavelet. So, the formula for second peak validation is, $f_{rp} = |f_{pcp} - f_{pmp}| \leq 100$ where f_{pmp} = data points obtained after performing peak detection on data points retrieved from the implementation of clustering and median filtering (M_F) = $SMA(M_F(y_{s,h}(x_i)))$. As part of the peak validation, the effect of wavelet transformation and median filtering, which is

implemented for the second time, can be found in Fig.2 (f) and Fig.2 (g), respectively.

2.8 Peak validation (final)

This peak validation method is operative on the final data set before producing the final result. Three consecutive data points, previous = $i-1$, current = i and next = $i+1$, are considered for data traversal.

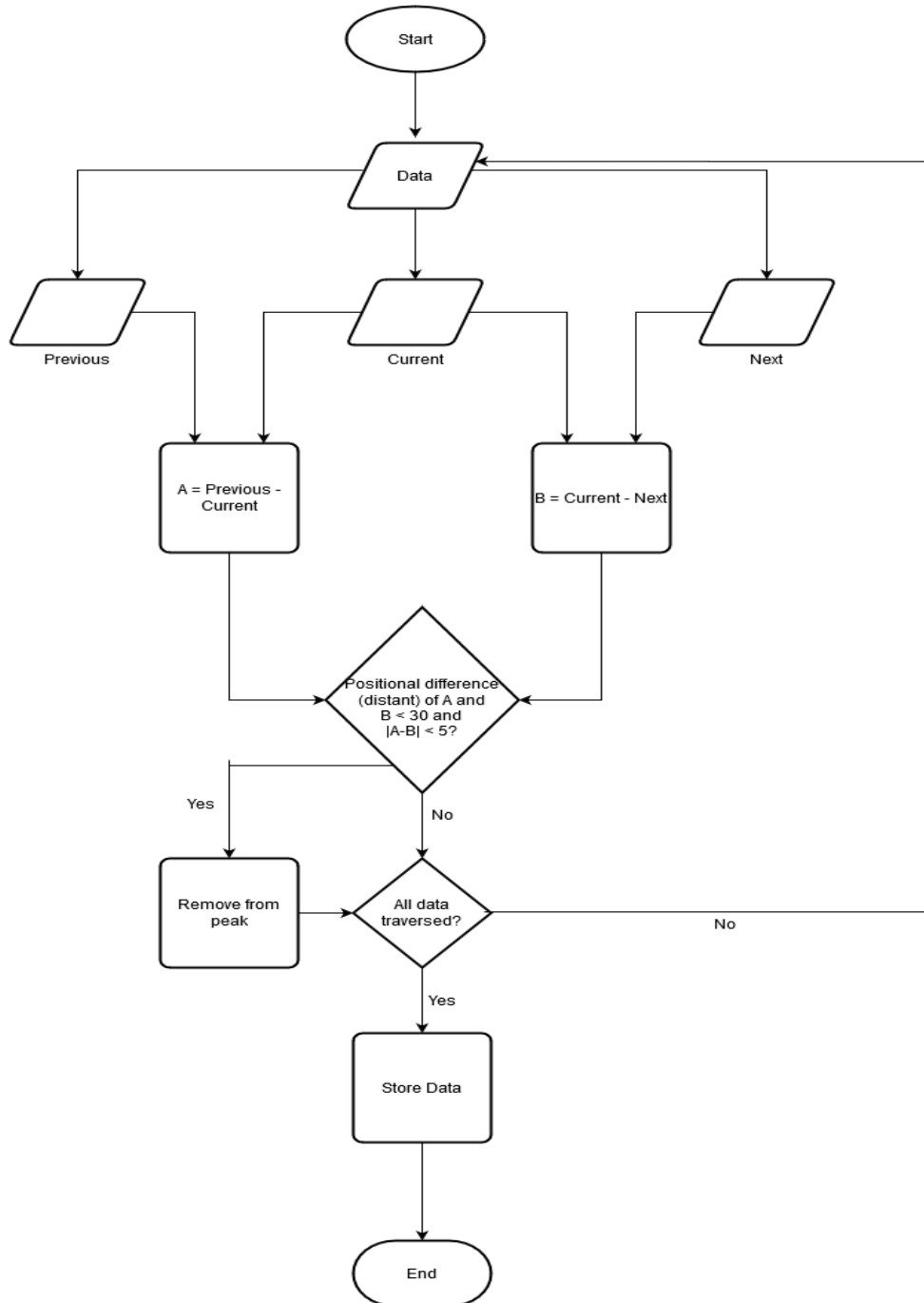


Diagram-3 Third and final peak validation method. It is a comparison between three adjacent data points. This exerts the final 37 data points, which correspond to the peak value.

We have defined two variables, such as $A = \text{previous} - \text{current}$ and $B = \text{current} - \text{next}$, to check the distance and difference of A and B. Distance refers to the position between A and B, and the difference is the value obtained by subtracting A and B. If the distance is less than 30 and the difference is less than 5, then the value is removed from the peak. This final procedure isolates the peak from the rest of the data to have a proper evaluation of peak detection. After this peak validation, we obtain 37 peak points. After getting the peaks, we can easily calculate the consumption rate by finding the difference between peaks. This final procedure isolates the peak from the rest of the data to have a proper evaluation of peak detection, as illustrated in Diagram-3.

3. Result and Analysis

Our paper has focused primarily on noise reduction and peak detection from the industrial dataset, which is retrieved from vehicles. To achieve this target, we have faced an intricate trade-off; detecting peaks while minimizing data loss. To get the actual value from the noisy time series dataset, we have implemented different algorithms mentioned in the previous section. Implementation of similar algorithms can be found in reference [32], where both wavelet transformation and adaptive denoising algorithms are implemented on non-linear (chaotic) time-series datasets. The performance of the wavelet analysis is poor as it has not been systematically studied for chaotic signals, and thus, no conclusive proof is given whether the wavelet-based noise reduction techniques truly are the best. Chaotic signals usually have a broad-band spectrum that overlaps with the spectrum of noise. In our case, after minimizing the noise using several algorithms, wavelet transformation is implemented as a catalyst to trim the noise further.

In another work [33] the authors have run wavelet transformation on the traffic volume time series dataset. However, they could not deal with past correlated samples due to the longer lags into the traffic volume pattern. Likewise, a threshold-based wavelet denoising algorithm is applied on hydrological time series data to acquire a detailed approximation of rainfall and runoff signals at various resolution levels, respectively [34]. It has been deduced that there are both positive and negative trends in each zone for the monsoon datasets, which shows the same periodicity, and thus, it affects the accuracy of acquiring approximate and detailed rainfall and runoff signals. An almost identical issue has emerged in other work [35], where the author used discrete wavelet transformation on hydrological time series data and found the limitation of positive and negative trends. In contrast, our approach has successfully removed all possible noises. In reference [36], different denoising and gap-filling methods such as interpolation, the Singular Spectrum Analysis, and the Lomb-Scargle algorithm. They have further

suggested to select diverse methods based on the problem type of each noise pattern of the dataset. Each of the algorithms has its strength and drawbacks i.e. the Lomb-Scargle method is sensitive in the presence of a strong trend in the raw data. This sensitivity may boost the power spectrum at low frequencies which may therefore mask other frequencies that would be significant in the absence of that limitation. Due to this issue, the method produces false, intermittent peaks in these systematic gaps, and the overall experiment shows the poor result. On the other hand, in our case, we achieved semi-accurate results by mashing up different algorithms according to the type of noises. Our two-layered noise reduction method has detected peaks while eliminating the noise based on a threshold value.

Other investigations have been done on the approach of peak detection in the time series data set. In one such work [37], the authors have proposed a continuous wavelet transform (CWT)-based pattern matching method for detecting strong and weak peaks without performing data preprocessing. Though they have efficiently detected the strong and weak peaks on the raw data but it can only be applicable for the dataset with default parameters, and the computational load is comparatively high. The experimental result they have achieved shows one isolated false positive identification from the spectrums. The work presented in reference [38] has shown several statistical ways to compute the time-series peak function. They have experimented on a dataset consisting of annual sunspot data for years 1700 to 2008, where the best result shows almost accurate results but failed to detect more than two peaks. Conversely, our algorithm is successfully operative in detecting multiple peaks.

Start_index_refill	Stop_index_refill	Deduced_value	Real_value	Error	Percentage_error
4041	4042	17.75168437	17.77	0.018315629965759	0.103070511906353
9709	9710	14.64695647	14.65	0.003043528614741	0.02077493934977
10663	10664	35.29107497	36.89	1.59892502522143	4.33430475798708
12938	12939	10.87529252	9.8	1.07529251780108	10.9723726306232
14083	14084	31.3662363	33.25	1.88376369997919	5.66545473677951
18052	18163	23.66448841	23.54	0.124488411240215	0.528837770774064
19687	19783	29.39111195	29.202177	0.188934954874462	0.646989280540495
22121	22122	14.69663734	14.66	0.036637341635856	0.24991365372344
23412	23413	46.69013676	47.14	0.449863239609741	0.954313193911202
25657	25658	35.89892004	36.8	0.901079957266589	2.44858684039834
30112	30113	32.85478607	32.79	0.064786069970658	0.197578743429882
31964	31990	52.91355821	45.71	7.20355820989442	15.7592610148642
33898	33899	37.63575852	37.915	0.279241484810619	0.736493432178871
35525	35551	39.18674741	38.46	0.726747413687697	1.88961886034243
36694	36695	8.34281661	8.342817	3.86526206597182E-07	4.63304189217122E-06
37283	37309	39.57391374	40.62	1.04608626490735	2.57529853497623
39543	39544	28.36952385	28.98	0.610476145318319	2.10654294450766
41653	41654	31.27637457	31.79	0.51362542510611	1.61568236900318
43770	43771	19.8699458	19.8	0.069945799087723	0.353261611554156
44859	44860	8.93436836	8.65	0.28436836419594	3.28749553983745
45817	45818	26.51803476	26.31	0.208034755485475	0.790706026170562
48269	48337	46.49246463	46.35	0.142464625747117	0.307367045840597
51489	51490	8.39503321	8.84	0.444966794162927	5.03356101994261
52715	52716	20.71302187	20.86	0.14697812776361	0.70459313405374
55116	55117	14.84071357	14.94	0.099286428987975	0.664567797777611
56075	56086	24.55695611	24.56	0.003043893387808	0.012393702719088
58861	58862	36.03948936	37.47	1.43051064281764	3.81774924691123
64186	64187	28.85588498	28	0.855884980066733	3.05673207166669
69723	69749	27.06560695	28	0.9343930504497	3.33711803732036
72642	72643	23.80248598	27.27	3.46751402431882	12.7154896381328
75681	75682	27.38444797	28.79	1.40555203108102	4.8820841649219
79819	79820	27.16472555	27	0.164725548989559	0.610094625887254
83096	83097	30.65766458	34.85	4.19233542465138	12.0296568856568
86449	86543	18.95143924	15.045849	3.90559023839059	25.9579252615827
88718	88744	44.74483446	44.7	0.044834460714711	0.100300806968033
97673	97674	41.87936796	41.97	0.090632040213478	0.215944818235592
102601	102602	21.01643509	24.07	3.05356491175409	12.6861857571836

Table-1 List of detected peaks as compared to the real value. For the most part, the error is minimum.

As seen in Table-1, for the most part, our module can detect peaks with an error of $\pm 1L$. Those values whose error rate is higher than 5 liters are due to the reason mentioned above. However, the R-squared score is 0.98, and RMSE is 1.49, which indicates that our module is highly efficient. The graphical representation of cleaned data after all noise removal can be seen in Fig.2 (h).

Due to hardware malfunction or some similar reason, at different time periods, the obtained data showed bizarre patterns, which has caused discrepancy in finding the correct result. In the highway, due to the high speed, the fuel consumption rate remains constant, as the device has its limitation. As a result, there is a huge data loss during this time period. If the dataset collected has no data loss, then the result can be significantly improved where the accuracy may exceed ours.

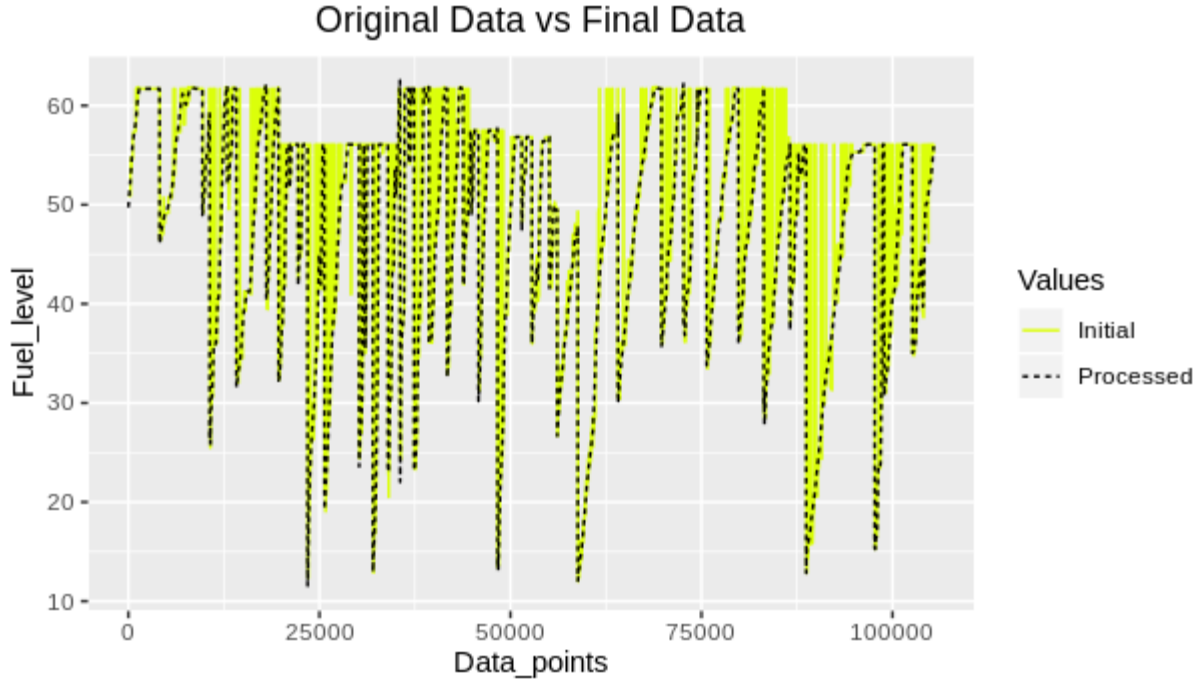


Fig.4 A comparison between our initial dataset and processed dataset. The data points are narrowed down, eliminating all the noises that were present at the beginning.

Our module narrowed down 37 peaks from a dataset of 100000 data points as shown in Fig.4. The original dataset is represented using yellow color, and the cleaned one is represented using black. Additionally, we have tested our module with an additional four datasets. In this case, for each of the datasets, we have compared two consecutive stages for all the phases. The stages are Initial, Cluster, 1st Wavelet, 1st Median, 2nd Wavelet, 2nd Median, and Final. We represented the data by superimposing the plots on top of one another. Each of the six graphs represents a comparison between the two phases of data denoising. For any graph, the comparisons are between (a) raw data and data after clustering, (b) data after clustering and data after first wavelet analysis, (c) data after first wavelet analysis and data after first filtering, (d) data after first filtering vs. data after second wavelet analysis, (e) data after second wavelet analysis and data after second filtering and (f) data after second filtering and final value. Also, unlike in our initial test where we have plotted a data points vs. fuel graph, we have drawn a graph where X-axis corresponds to date and time, and Y-axis corresponds to fuel value. Their graphical representations are presented from Fig.5 to Fig.8.

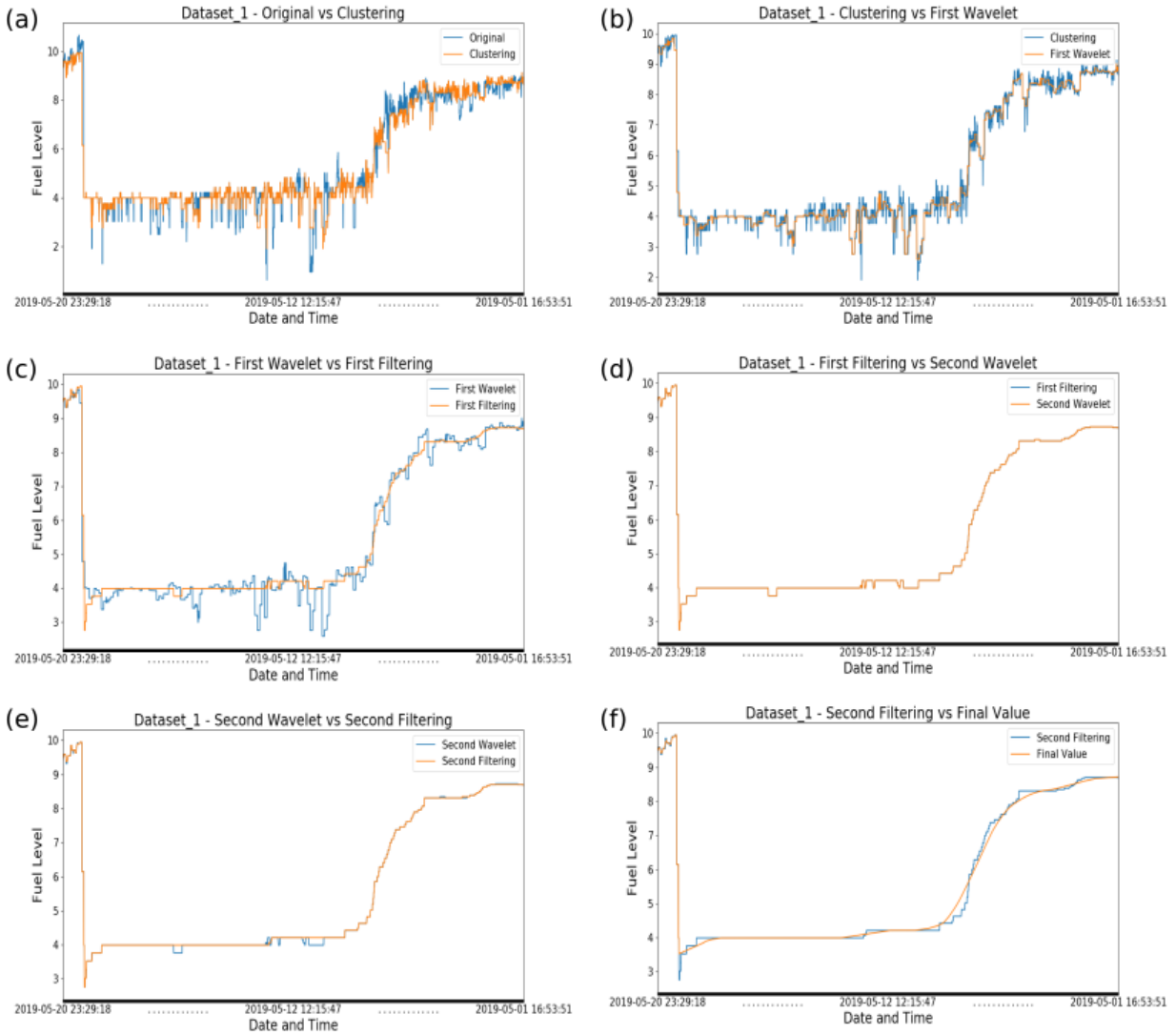


Fig.5 (a)-(c) One notable point here is most of the data is cleaned in the first three phases. (d) There seems to be no effect of the second wavelet after the first filtering is performed. (e) Second, filtering has caused small changes in the dataset. (f) Changes are significant as compared to immediately previous ones.

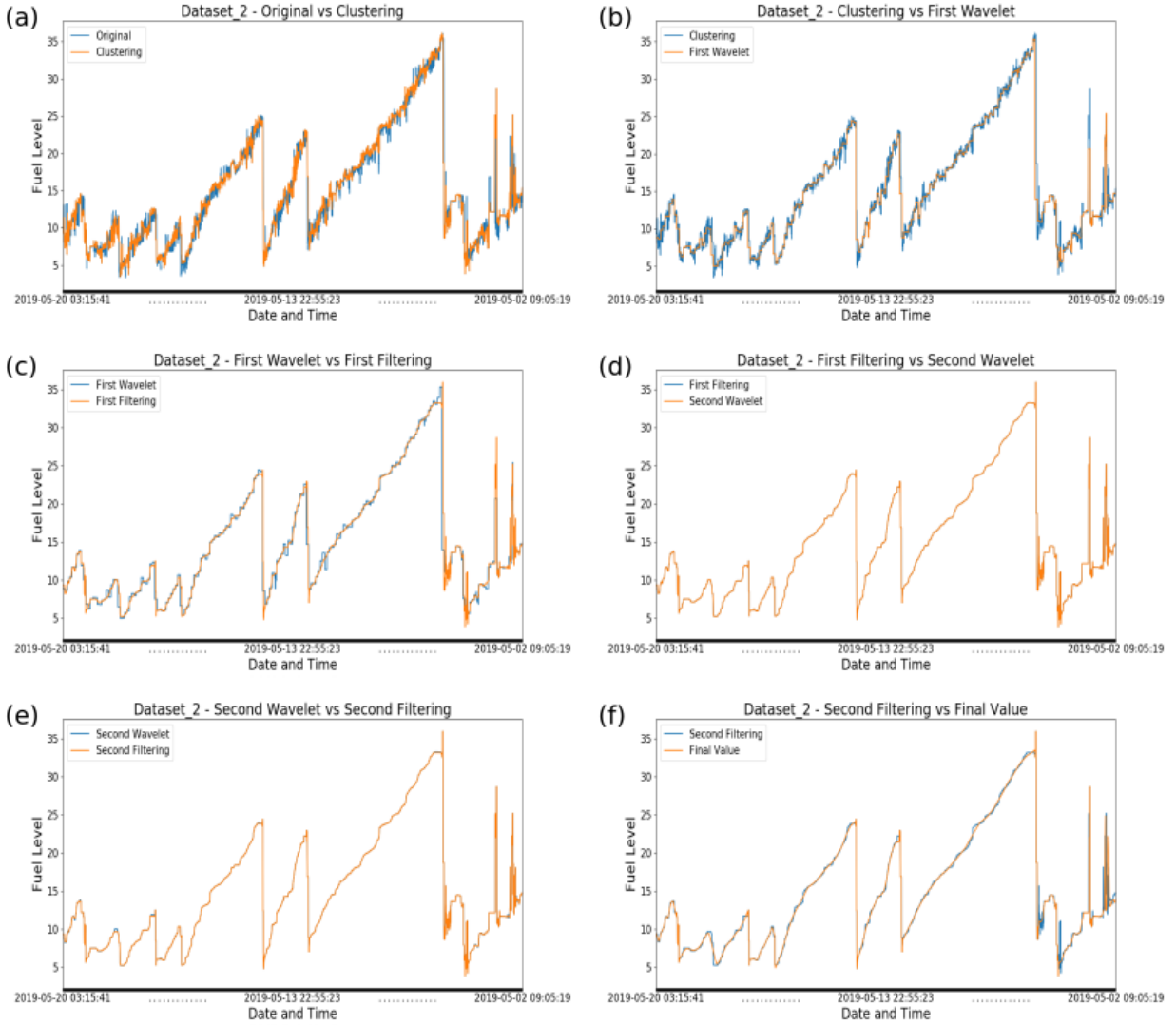


Fig.6 (a)-(c) Noise in prominent and heavy cleaning is performed. (d)-(e) Dissimilar to the previous graph, the state of the graph remains almost the same. Hence, there seems to be a small effect of second wavelet and second filtering after the first filtering is performed. (f) Final peak validation changes the graph by a small margin.

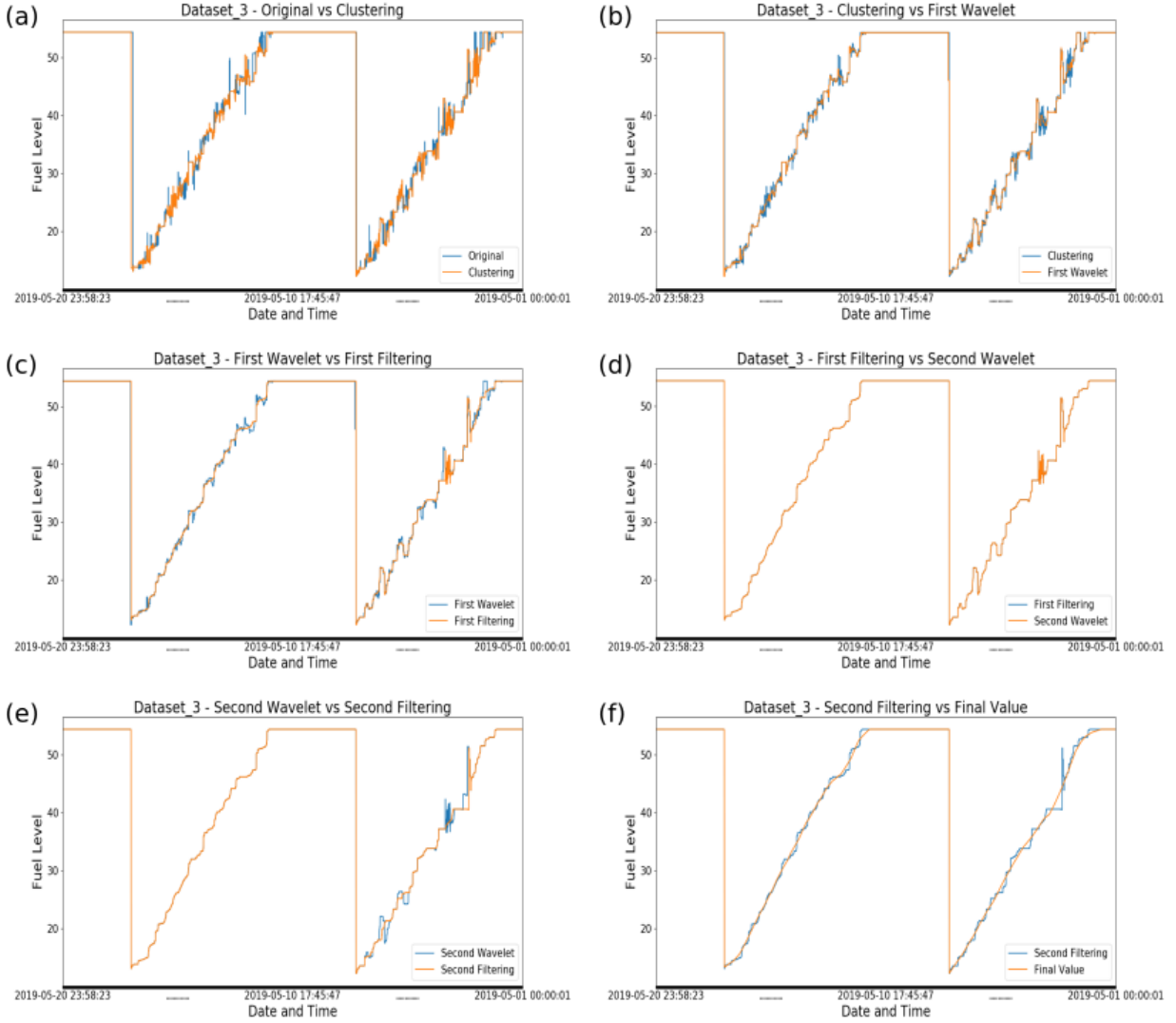


Fig.7 (a)-(c) Another notable point is that noise is more condensed along with the consumption rate, which is shown in the graph by a downward trend. (d) Similar to Fig.6, the difference is minimal after the implementation of wavelet for the second time. (e)-(f) Cleaning effect is visible in the downward trend.

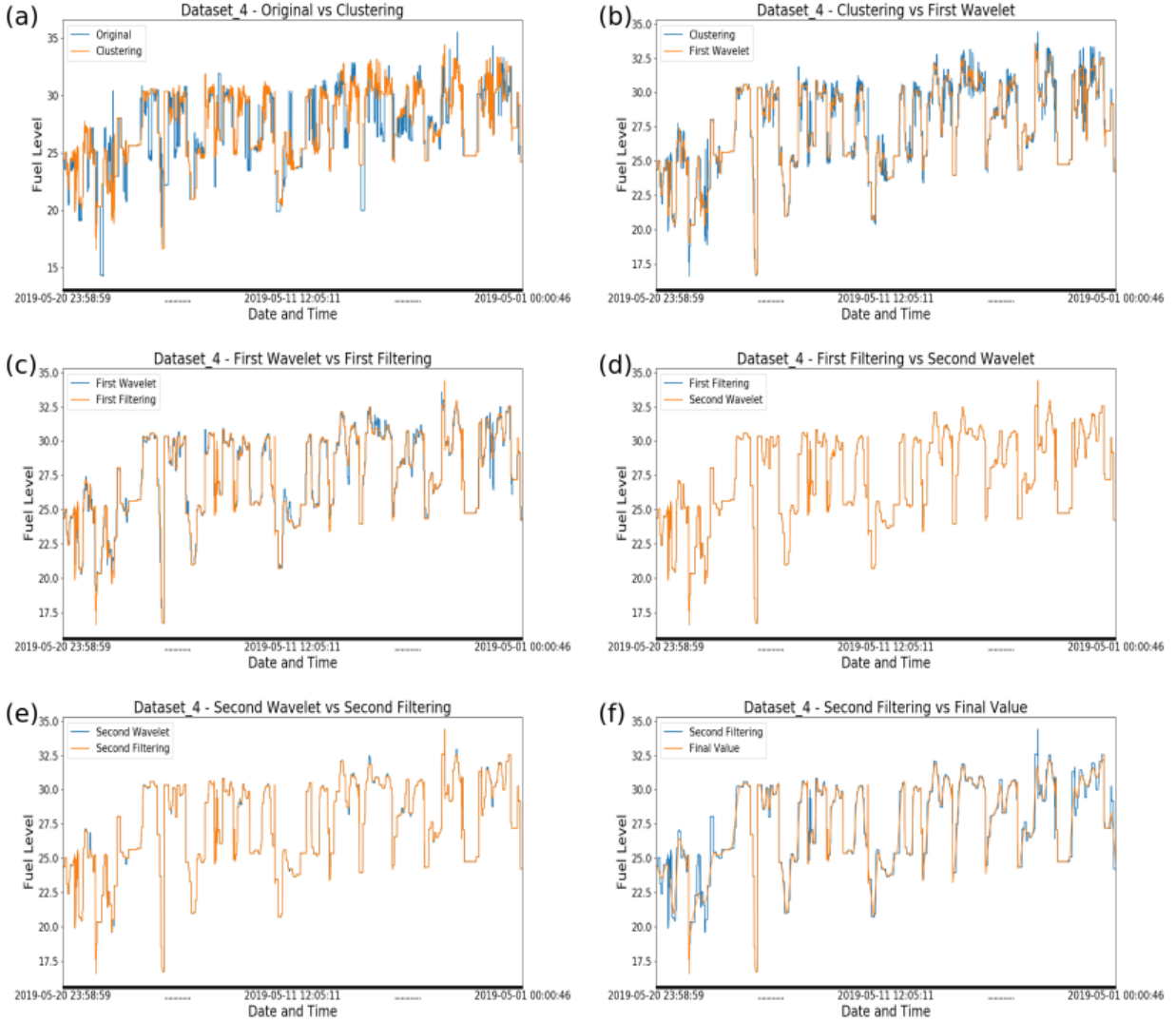


Fig.8 (a)-(c) Fluctuations in fuel level is much prominent as compared to the previous dataset. (d)-(e) Although, state of the graph remains almost the same, (f) in the last stage, due to the peak validation method, we see changes in the final dataset as compared to the previous adjacent graph.

4. Conclusion

In conclusion, the real dataset (on which we have worked) has contained severe noise, which may have occurred due to poor street condition of some roads of Bangladesh, noisy dynamics, mobile environment and so on. The proposed module in this work, regardless of the intensity of data noise, is capable of detecting peaks without removing the noise in the first place. The reason behind the capability of doing that: different methods have been arranged in such a way that it can store the data before sending it to the next phase or process. For instance, data generated from clustering is stored first before sending it to the following phase to perform wavelet transform. This strategy allows the system to compare the values of the datasets retrieved from two adjacent processes. Before removing the noise, this entire task is conducted using multiple peak validity and peak detection methods. Consequently, the system works more efficiently as it is independent of the noise percentage. So regardless of the percentage of noise present in the datasets, peaks have been detected. However, after successful peak detection, the noise has been gotten removed and the correct value of the consumption rate has also been generated. So, this system can isolate noise and peak first; and then remove noises from the rest of the datasets, which provides not only correct peak values but also consumption values as well. Since this is a two layer-based noise removal system, at first, peaks have been separated from the noisy dataset. Then, the consumption rate is deduced using these separated peaks. This approach helps to maintain the integrity of the peaks. This type of design allows us to minimize data loss during noise removal, which has not been reported previously in the presence of bizarre data patterns i.e. constant values, accuracy subsides. Without such extreme cases, accuracy may increase a lot. Our work can be implemented in any of the industrial fields, where exist the complicated issues of fuel measurement i.e. vehicle industry, fuel station, power plant and so on. As long as any time-series data distribution resembles our dataset distribution, the module can work efficiently to reduce noise and produce correct values.

Acknowledgement

M.R.C. Mahdy acknowledges Dr. K.M. Masum Habib and Mr. Mahmudul Hasan Sohag of Pi-Labs Bangladesh Limited for several important discussions. M.R.C. Mahdy also acknowledges Mr. Saikat Chandra Das for several important discussions.

Conflict of Interest

Authors declare no competing financial interest.

References

1. Kollem, S., Reddy, K. R. L., & Rao, D. S. (2019). A Review of Image Denoising and Segmentation Methods Based on Medical Images. *International Journal of Machine Learning and Computing*, 9(3).
2. Gondara, L. (2016, December). Medical image denoising using convolutional denoising autoencoders. In *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)* (pp. 241-246). IEEE.
3. Singh, A., & Verma, R. S. (2017). An efficient non-local approach for noise reduction in natural images. *International Journal of Advanced Research in Computer Science*, 8(5).
4. Kaur, C., & Bansal, N. (2016). De-Noising Medical Images Using Low Rank Matrix Decomposition NN and SVM. (*IJCSIT*) *International Journal of Computer Science and Information Technologies*, 7(1), 45-48.
5. Saihood, A. A. (2014). Image Denoising using Neural Network with SVM (Support Vector Machine) and LDA (Linear Discriminant Analysis). *Int. J. Comput. Sci. Inf. Technol.*, 5(4), 5363-5367.
6. Wu, J., Huang, D. Y., Xie, L., & Li, H. (2017, August). Denoising Recurrent Neural Network for Deep Bidirectional LSTM Based Voice Conversion. In *INTERSPEECH* (pp. 3379-3383).
7. Yang, Q., Yan, P., Zhang, Y., Yu, H., Shi, Y., Mou, X., & Wang, G. (2018). Low-dose CT image denoising using a generative adversarial network with Wasserstein distance and perceptual loss. *IEEE transactions on medical imaging*, 37(6), 1348-1357.
8. Miranda, A. L., Garcia, L. P. F., Carvalho, A. C., & Lorena, A. C. (2009, June). Use of classification algorithms in noise detection and elimination. In *International Conference on Hybrid Artificial Intelligence Systems* (pp. 417-424). Springer, Berlin, Heidelberg.
9. Germain, F. G., Chen, Q., & Koltun, V. (2018). Speech denoising with deep feature losses. *arXiv preprint arXiv:1806.10522*.

10. Singla, E. M., & Singh, M. H. (2015). Paper on frequency based audio noise reduction using butterworth, Chebyshev & Elliptical filters. *Int. J. Recent Innov. Trends Comput. Commun*, 3, 5989-5995.

11. Haque, M., & Bhattacharyya, K. (2018). Speech Background Noise Removal Using Different Linear Filtering Techniques. In *Advanced Computational and Communication Paradigms* (pp. 297-307). Springer, Singapore.

12. Welk M., Bergmeister A., Weickert J. (2005) Denoising of Audio Data by Nonlinear Diffusion. In: Kimmel R., Sochen N.A., Weickert J. (eds) *Scale Space and PDE Methods in Computer Vision. Scale-Space 2005. Lecture Notes in Computer Science*.

13. Magsi, H., Sodhro, A. H., Chachar, F. A., & Abro, S. A. K. (2018, March). Analysis of signal noise reduction by using filters. In *2018 International Conference on Computing, Mathematics and Engineering Technologies (iCoMET)* (pp. 1-6). IEEE.

14. Magondou, S. (2016). *Using Neural Networks to reduce noise in internet of things data streams* (Doctoral dissertation, University of Nairobi).

15. Agresti, G., Schaefer, H., Sartor, P., & Zanuttigh, P. (2019). Unsupervised domain adaptation for ToF data denoising with adversarial learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 5584-5593).

16. He, Q., Wang, X., & Zhou, Q. (2014). Vibration sensor data denoising using a time-frequency manifold for machinery fault diagnosis. *Sensors*, 14(1), 382-402.

17. Manmohan, Gupta, S., & Kuldeep. (2015). An Novel Approach for image denoising using Wavelet Transforms. *International Journal of Engineering Research and Management (IJERM)*, 02(1), 142-146.

18. Mythili, C., & Kavitha, V. (2011). Efficient technique for color image noise reduction. *The research bulletin of Jordan, ACM*, 1(11), 41-44.

19. Koskinen, K., Auvinen, P., Björkroth, K. J., & Hultman, J. (2015). Inconsistent denoising and clustering algorithms for amplicon sequence data. *Journal of Computational Biology*, 22(8), 743-751.
20. Chatterjee, P., & Milanfar, P. (2009). Clustering-based denoising with locally learned dictionaries. *IEEE transactions on Image Processing*, 18(7), 1438-1451.
21. Mahdy, M. R. C., Rivy, H. M., Jony, Z. R., Alam, N. B., Masud, N., Al Quaderi, G. D., & Rahman, M. S. (2020). Dielectric or plasmonic Mie object at air–liquid interface: The transferred and the traveling momenta of photon. *Chinese Physics B*, 29(1), 014211.
22. Chowdhury, R., Mahdy, M. R. C., Alam, T. N., Al Quaderi, G. D., & Rahman, M. A. (2020). Predicting the stock price of frontier markets using modified Black–Scholes Option pricing model and machine learning. *Physica A: Statistical Mechanics and its Applications*, 124444.
23. Chowdhury, R., Rahman, M. A., Rahman, M. S., & Mahdy, M. R. C. (2020). An approach to predict and forecast the price of constituents and index of cryptocurrency using machine learning. *Physica A: Statistical Mechanics and its Applications*, 124569.
24. Ahmed, Nahian, M. Nazmul Islam, Ahmad Saraf Tuba, M. R. C. Mahdy, and Mohammad Sujauddin. "Solving visual pollution with deep learning: A new nexus in environmental management." *Journal of environmental management* 248 (2019): 109253.
25. Hambal, A. M., Pei, Z., & Ishabailu, F. L. (2017). Image noise reduction and filtering techniques. *International Journal of Science and Research (IJSR)*, 6(3), 2033-2038.
26. Ding, S., Zhang, L., & Zhang, Y. (2010, April). Research on spectral clustering algorithms and prospects. In *2010 2nd International Conference on Computer Engineering and Technology* (Vol. 6, pp. V6-149). IEEE.
27. Murtagh, F. (2011). Hierarchical Clustering. *Computing Research Repository* 633-635.

28. Sudha, S., Suresh, G. R., & Sukanesh, R. (2009). Speckle noise reduction in ultrasound images by wavelet thresholding based on weighted variance. *International journal of computer theory and engineering*, 1(1), 7.
29. Aggarwal, R., Singh, J. K., Gupta, V. K., Rathore, S., Tiwari, M., & Khare, A. (2011). Noise reduction of speech signal using wavelet transform with modified universal threshold. *International Journal of Computer Applications*, 20(5), 14-19.
30. Üstündağ, M., Şengür, A., Gökbulut, M., & Ata, F. (2013). Performance comparison of wavelet thresholding techniques on weak ECG signal denoising. *Przegląd Elektrotechniczny*, 89(5), 63-66.
31. Boateng, K. O., Asubam, B. W., & Laar, D. S. (2012). Improving the effectiveness of the median filter.
32. Gao, Jianbo & Sultan, Hussain & Hu, Jing & Tung, Wen-Wen. (2010). Denoising Nonlinear Time Series by Adaptive Filtering and Wavelet Shrinkage: A Comparison. *Signal Processing Letters, IEEE*. 17. 237 - 240. 10.1109/LSP.2009.2037773.
33. Boto Giralda, Daniel & Díaz-Pernas, Francisco & González-Ortega, David & Díez, Jose & Antón-Rodríguez, Miriam & Martínez Zarzuela, Mario & De la Torre Díez, Isabel. (2010). Wavelet-Based Denoising for Traffic Volume Time Series Forecasting with Self-Organizing Neural Networks. *Comp.-Aided Civil and Infrastruct. Engineering*. 25. 530-545. 10.1111/j.1467-8667.2010.00668.x.
34. Chou, Chien-Ming. (2011). A Threshold Based Wavelet Denoising Method for Hydrological Data Modelling. *Water Resources Management*. 25. 1809-1830. 10.1007/s11269-011-9776-3.
35. Pandey, Brij & Tiwari, Harinarayan & Khare, Deepak. (2017). Trend analysis using discrete wavelet transform (DWT) for long-term precipitation (1851–2006) over India. *Hydrological Sciences Journal/Journal des Sciences Hydrologiques*. 62. 10.1080/02626667.2017.1371849.
36. Musial, Jan & Verstraete, Michel & Gobron, Nadine. (2011). Comparing the effectiveness of recent algorithms to fill and smooth incomplete and noisy time series. *Atmospheric Chemistry and Physics Discussions*. 11. 10.5194/acpd-11-14259-2011.

37. Du, Pan & Kibbe, Warren & Lin, Simon. (2006). Improved peak detection in mass spectrum by incorporating continuous wavelet transform-based pattern matching. *Bioinformatics* (Oxford, England), 22. 2059-65. 10.1093/bioinformatics/btl355.
38. Palshikar, Girish. (2009). Simple Algorithms for Peak Detection in Time-Series.