

**Artifact Removal From Sleep-Disordered EEG by Wavelet Enhanced
Independent Component Analysis**

by

Xinyi Fan

A thesis submitted in partial fulfillment of the requirements for the degree of

Master of Science

in

Statistical Machine Learning

Department of Mathematical and Statistical Sciences
University of Alberta

© Xinyi Fan, 2024

Abstract

In the field of sleep research, the quantitative analysis of electroencephalography (EEG) data acquired during sleep offers invaluable insights. However, the presence of artifacts in such data can severely distort analytical outcomes. Therefore, this study aims to develop an innovative artifact detection and rejection tool to enhance the analysis of sleep EEG data from overnight polysomnography (PSG) in patients with sleep disorders.

While recent advancements have seen the trend of artifact removal using a hybrid method, these techniques typically require pre-labeled data for training machine learning models, introducing a dependency on prior knowledge. Addressing this limitation, our research introduces a novel unsupervised learning approach, utilizing hierarchical clustering to identify artifactual components within wavelet-enhanced independent component analysis (ICA)-separated data. We present a unique set of features for clustering, including kurtosis, zero-crossing count, skewness, Hjorth parameters, and Rényi entropy, tailored to discern artifacts in EEG recordings.

Our methodology affords the flexibility of fully automated artifact removal or semi-automated processes involving visual inspection of hierarchical dendrograms. Comparative analyses demonstrate that this new method not only refines EEG signal quality but also surpasses traditional manual cleaning techniques in performance. The findings underscore the potential of hierarchical clustering in the unsupervised learning landscape for artifact detection, heralding a significant step forward in the preprocessing of EEG data for sleep studies.

Preface

This thesis is an original work by Xinyi Fan. No part of this thesis has been previously published. The clinical study, of which the CF00N data presented and analyzed in this thesis is a part, received ethics approval by the Health Research Ethics Board of the University of Alberta Pro00057638.

Acknowledgements

I would like to thank my supervisor, Dr. Adam Kashlak, whose support and guidance have been instrumental throughout this thesis. Additionally, I would like to thank my co-supervisor, Dr. Giseon Heo, for her invaluable advice and assistance. Furthermore, I am deeply appreciative of the encouragement and support provided by my parents, Weixing Fan and Wenwei Liu. Their belief in me has been a constant source of motivation. Finally, I would like to express my sincere appreciation to my partner, Junhao Lu, whose love, companionship, and inspiration have supported me during this journey.

Table of Contents

1	Introduction	1
1.1	Background Overview	1
1.2	Polysomnography and Electroencephalography	1
1.3	Artifacts	4
1.3.1	Ocular Artifacts	5
1.3.2	Muscle Artifacts	6
1.3.3	Cardiac Artifacts	7
1.4	Overview	7
2	Previous Works	10
2.1	Single Artifact Removal Techniques	10
2.1.1	Independent component analysis	10
2.1.2	Wavelet transform	11
2.2	Hybrid Method	11
2.2.1	Wavelet enhanced Independent component analysis	12
2.3	Artifact Removal for Sleep EEG Data	13
3	Data	14
3.1	The Dataset	14
3.2	Data Pre-processing	15
4	Methodology	16
4.1	Independent Component Analysis	17
4.1.1	FastICA	18
4.2	Wavelet Transform	20
4.2.1	Continuous Wavelet Transform	21
4.2.2	Discrete Wavelet Transform	21
4.3	Feature Extraction	23
4.3.1	Kurtosis	23

4.3.2	Zero-Crossing Count	23
4.3.3	Skewness	25
4.3.4	Hjorth Features	25
4.3.5	Renyi's Entropy	26
4.4	Hierarchical Clustering	27
4.5	Inverse DWT	28
4.6	Inverse ICA	28
5	Outputs and Results	29
5.1	Outputs	29
5.2	Performance Evaluation	30
5.2.1	Accuracy	30
5.2.2	Results	31
6	Discussion	39
6.1	Results	39
6.2	Future Work	40
	Bibliography	42
	Appendix A: Method Comparison	48

List of Tables

3.1	EEG electrodes names and replacement.	15
3.2	Proportion of sleep states of CF050 patient.	15
5.1	Mean and standard deviation of accuracy for pre-processed only, automatic and semi-automatic clean methods across eight EEG channels over 50 model runs. Numbers in bold are the highest mean accuracy for each channel.	32

List of Figures

1.1	Labels for 10-20 electrode placement systems [5]	3
1.2	Signal of recorded EEG, pure EEG and three types of artifact, EOG, ECG and EMG [11].	5
4.1	Block diagram of the proposed method of wavelet enhanced ICA technique for artifact removal from EEG data.	17
4.2	Third level wavelet decomposition of an EEG signal [55].	22
4.3	Daubechies wavelet and scaling functions of different orders [51]. . . .	24
5.1	The raw eight-channel EEG signal is shown in (a) and the signal after pre-processing step is shown in (b). This step is conducted to retain only the frequency bands of interest (0.3 to 80 Hz).	34
5.2	Decomposed eight ICs from pre-processed EEG signal using ICA. . .	35
5.3	WCs from the first IC in Figure 5.2 by DWT. The wavelet coefficients from top to bottom are D1, D2, D3, D4 and A4 respectively.	35
5.4	Boxplots of features extracted from each WCs. The seven features presented are (a) kurtosis, (b) zero-crossing count, (c) skewness, (d) activity, (e) mobility, (f) complexity and (g) Renyi's entropy.	36
5.5	Dendrograms of each set of WCs. The x-axis represents the label of each epoch and each independent component. The WCs are shown in (a) D1, (b) D2, (c) D3, (d) D4 and (e) A4.	37
5.6	Cleaned EEG signals across eight channels. Figure (a) displays the EEG signal cleaned via the automatic method, and Figure (b) shows the EEG signals cleaned using the semi-automatic method.	38
A.1	Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel F4M1.	49
A.2	Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel F3M2.	49

- A.3 Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel C4M1. 50
- A.4 Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel C3M2. 50
- A.5 Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel T6M1. 51
- A.6 Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel T5M2. 51
- A.7 Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel O2M1. 52
- A.8 Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel O1M2. 52

Chapter 1

Introduction

1.1 Background Overview

Sleep is a fundamental biological process that is essential for all living organisms, constituting approximately one-third of human life. Inadequate or poor-quality sleep has been linked to dysfunction in numerous bodily systems [1]. As the importance of understanding sleep patterns has grown in recent decades, it is recommended to employ reliable subjective and objective measures for evaluation. Laboratory-based polysomnography (PSG) is considered the benchmark for objectively measuring sleep [2].

1.2 Polysomnography and Electroencephalography

Polysomnography (PSG) is a comprehensive procedure used to systematically collect physiological data during sleep. It involves the simultaneous monitoring of various parameters to assess and diagnose underlying causes of sleep disturbances. A typical polysomnogram includes the recording of electrocardiography (ECG) for monitoring heart activity, electromyography (EMG) for studying muscular contractions, electroencephalography (EEG) for observing brain activity, and electro-optigraphy (EOG) for assessing eye dipole fields. Additionally, pulse oximetry, airflow and respiratory effort are monitored to evaluate sleep-related breathing disorders, such as obstructive sleep apnea (OSA), central sleep apnea, and sleep-related hypoventila-

tion/hypoxia. PSG serves as the gold standard for diagnosing these conditions. Furthermore, it can be employed to assess other sleep disorders, including nocturnal seizures, narcolepsy, periodic limb movement disorder, and rapid eye movement sleep behavior disorder [3].

EEG is widely regarded as one of the most powerful techniques for studying the electrophysiological dynamics of the brain in a noninvasive manner [4]. It refers to the electrical activity of a fluctuating nature that is detected from the surface of the scalp through the use of metal electrodes and conductive substances. EEG primarily detects the electrical currents generated by the activation of brain cells, known as neurons. These currents arise during the synaptic excitations occurring in the dendrites of numerous pyramidal neurons located in the cerebral cortex. The differences in electrical potentials observed in EEG recordings are a result of the combined post-synaptic graded potentials originating from pyramidal cells. This creates electrical dipoles between the soma and the apical dendrites.

The electrical currents within the brain predominantly consist of ions such as Na^+ , K^+ , Ca^{++} , and Cl^- , which are transported across the neuronal membranes through specific channels. The movement of these ions is governed by the membrane potential, dictating the direction of ion flow [5].

EEG is recorded by placing metal electrodes on the scalp with a 10-20 electrode placement system, seen in Figure 1.1. This system establishes a standardized method for physically positioning and designating electrodes on the scalp. The head is divided into proportional distances from prominent landmarks on the skull, such as the nasion, preauricular points, and inion, in order to ensure comprehensive coverage of all brain regions. The 10-20 system utilizes proportional distances, expressed as percentages, between the ears and nose to determine the electrode placement points. Each electrode placement is labeled based on its adjacent brain area: F (frontal), C (central), T (temporal), P (posterior), and O (occipital). These labels are accompanied by odd numbers on the left side of the head and even numbers on the right side.

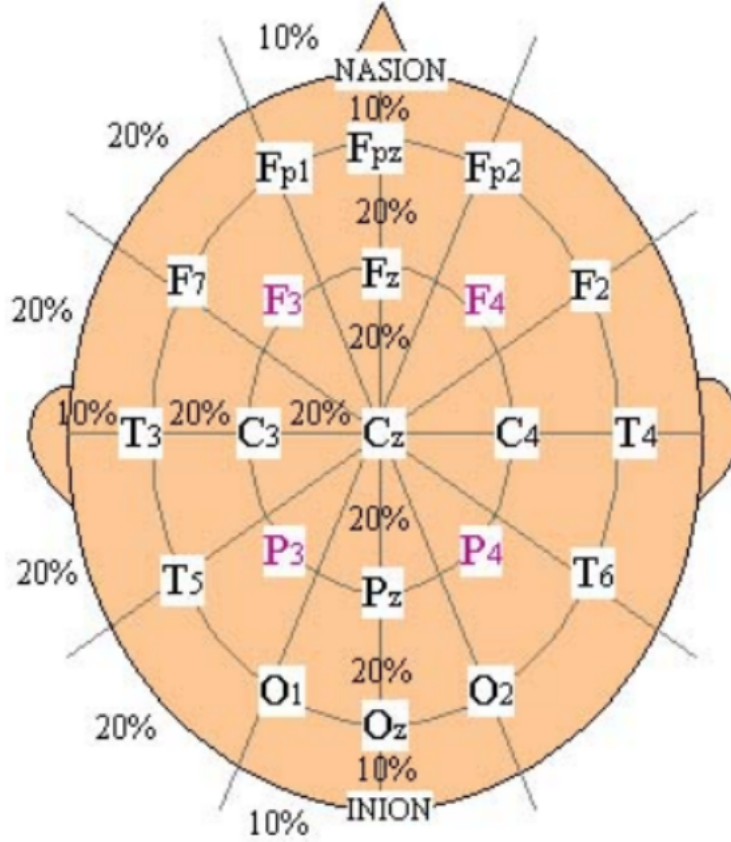


Figure 1.1: Labels for 10-20 electrode placement systems [5]

The determination of left and right sides follows the convention from the subject's perspective.

EEG captures a wide range of frequency components, but for clinical and physiological purposes, the focus lies between 0.3 and 30 Hz. This range can be roughly classified into various frequency bands, as outlined by Kellaway [6]:

Delta (<4 Hz): Delta rhythms primarily occur during deep sleep stages in healthy adults. If present outside of sleep, they may indicate pathological conditions.

Theta (4–8 Hz): Theta frequencies are observed in normal infants, children, and during drowsiness and sleep in adults. In awake adults, only a small amount of theta rhythms are typically present. High theta activity in awake adults suggests abnormal

and pathological conditions.

Alpha (8–14 Hz): Alpha rhythms are prominent in relaxed and mentally inactive states of wakefulness in healthy adults. They are most noticeable in the occipital area and typically have an amplitude of less than 50 V. Alpha rhythms are suppressed when the eyes are open (visual attention) or during mental exertion such as thinking.

Beta (14–30 Hz): Beta activity is predominantly observed in the frontocentral region and has lower amplitude compared to alpha rhythms. It is enhanced during states of anticipation and tension.

Gamma (>30 Hz): Gamma rhythms have a high frequency band and are usually not of clinical or physiological interest. Therefore, they are often filtered out in EEG recordings.

EEG provides valuable insights into the functional state of the brain and the mental condition of an individual. By analyzing the EEG, we can extract vital information that helps monitor a patient’s health, diagnose various brain conditions, and identify abnormalities in brain activity. The EEG serves as a powerful tool in understanding brain function, allowing healthcare professionals to assess cognitive processes, detect seizures, evaluate sleep patterns, and investigate various neurological disorders. Its non-invasive nature and ability to capture real-time brain activity make EEG a valuable tool in both clinical and research settings for understanding and monitoring brain health [7] [8].

1.3 Artifacts

EEG, with its high temporal resolution, allows for detailed analysis of brain activity over time. However, this advantage also brings the challenge of susceptibility to unwanted noise and the presence of artifacts in the EEG signals. These artifacts can significantly impact the quality and reliability of the recorded data [9]. Consequently, they have the potential to introduce substantial errors in measurement and diagnosis, diminishing the clinical usefulness of EEG signals.

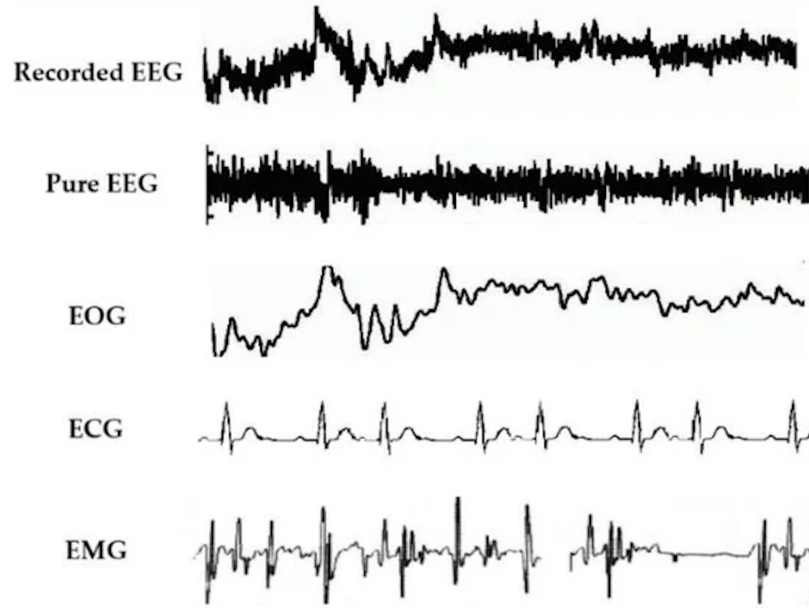


Figure 1.2: Signal of recorded EEG, pure EEG and three types of artifact, EOG, ECG and EMG [11].

Artifacts can be categorized into two main types: physiological artifacts and non-physiological artifacts. Physiological artifacts originate from the patient’s own physiological processes, while non-physiological artifacts arise from the patient’s external environment such as faulty electrodes, line noise, and high electrode impedance [10]. These artifacts can be mitigated through the use of more precise recording systems and strict recording procedures. By employing advanced equipment and following rigorous protocols during data acquisition, the occurrence of these artifacts can be minimized. On the other hand, removing physiological artifacts from EEG signals is a more complex task. These artifacts originate from the patient’s own physiological processes and necessitate specialized techniques and algorithms for accurate removal. There are three major physiological artifacts that have effects on EEG data.

1.3.1 Ocular Artifacts

Ocular artifacts pose a significant challenge in EEG recordings as they can introduce notable disturbances. These artifacts are particularly prominent during specific sleep stages, such as awakeness and rapid eye movement (REM) sleep. Awakeness is char-

acterized by increased eye movements, including voluntary and reflexive eye blinks, which can propagate across the scalp and manifest as EEG activity. Similarly, during REM sleep, which is associated with vivid dreaming and rapid eye movements, EOG artifacts can be particularly prominent in the EEG signals. These stages of sleep pose additional challenges for accurate interpretation and analysis of EEG data, as the presence of EOG artifacts can interfere with the underlying brain activity [12].

Ocular artifacts in EEG recordings result from changes in the orientation of the retina and cornea dipole during eye movements, while blink artifacts arise from variations in ocular conductance caused by the contact between the cornea and eyelid. Moreover, due to the volume conduction effect, both ocular artifacts and EEG activity propagate to the surface of the head and are recorded by the electrodes. EOG recordings can capture these ocular signals, which typically exhibit amplitudes several times greater than EEG signals and share similar frequency characteristics [13][14].

It is important to recognize that the contamination between EEG and EOG is bidirectional [11]. This means that not only can EEG data be contaminated by EOG artifacts, but EOG signals themselves can also be affected by EEG activity. As a result, when attempting to remove EOG artifacts from the EEG recordings, bidirectional interference can occur, leading to potential errors in the artifact removal process.

1.3.2 Muscle Artifacts

Contamination of EEG data by muscle activity poses a significant challenge due to the involvement of various muscle groups [15]. These artifacts can arise from muscle contractions and stretches occurring in close proximity to the signal recording sites. Common activities such as talking, sniffing, swallowing, and other facial or body movements can contribute to muscle-related artifacts in EEG recordings. The amplitude and waveform of these artifacts are influenced by the degree of muscle contraction and stretch [16].

Compared to artifacts from eye movements, obtaining meaningful activity solely from a single-channel measurement in EMG artifacts is extremely difficult [17]. The complex nature of muscle activity makes it particularly challenging to effectively eliminate these artifacts from the EEG signals. The presence of EMG artifacts can obscure the underlying neural activity and introduce significant distortions, hindering accurate interpretation and analysis of the EEG data.

1.3.3 Cardiac Artifacts

Cardiac artifacts present a challenge in EEG recordings when the electrodes are placed in close proximity to a blood vessel, leading to artifacts caused by the movement of expansion and contraction associated with the heartbeat. These artifacts, known as pulse artifacts, can manifest in the EEG signals with a waveform similar to that of the cardiac activity. Due to the similarity in waveform, it becomes difficult to effectively remove these artifacts from the EEG recordings [15].

On the other hand, electrocardiogram (ECG) represents a distinct cardiac activity. Unlike pulse artifacts, ECG signals exhibit a characteristic regular pattern, separate from cerebral activity. This characteristic pattern makes it relatively easier to identify and remove ECG artifacts by utilizing a reference waveform [11].

1.4 Overview

Addressing noise and artifacts in EEG signals is a critical aspect of utilizing this technique in medical applications [18]. By employing appropriate artifact removal techniques, researchers and clinicians can enhance the reliability and clinical utility of EEG data, enabling more accurate analysis and interpretation of brain activity.

Researchers have been investigating EEG artifact removal methods for nearly 50 years, yet to this day, a consensus has not been reached regarding the optimal algorithm for specific applications. Consequently, researchers must conduct comprehensive studies to examine the advantages and disadvantages of each algorithm from

various perspectives. Factors such as automatic versus manual methods, online versus offline methods, and suitability for specific applications need to be thoroughly evaluated in order to determine the most suitable choice [19].

While the majority of proposed methodologies have proven effective in detecting and removing artifacts from short-term EEG data obtained from healthy individuals in an awake state, applying these methods to full-night EEG presents challenges. The presence of artifacts, stemming from various common sources such as repetitive respiratory events disrupting sleep and gross body movements, makes it difficult. Such artifacts in the EEG signal often necessitate manual over-read for identification and exclusion of contaminated data. This task is not only prohibitively time-consuming and laborious but also challenging to perform consistently [20].

Only a limited number of algorithms have been proposed specifically targeting the removal of various artifacts from full overnight multichannel sleep EEG data. Notably, D’Rozario et al. in 2015 [20] introduced and validated a method tailored for patients with sleep disorders. Their approach involved experienced sleep technologists manually reviewing and identifying EEG artifacts, which served as training data for the algorithm.

In contrast, our study takes a different approach by utilizing full overnight multichannel sleep EEG data from a sleep disorder patient without prior artifact labeling. The objective of this thesis is to introduce a method capable of automatically and semi-automatically detecting and removing artifacts from such data. This method leverages a combination of Independent Component Analysis (ICA) and Discrete Wavelet Transform (DWT) technologies, augmented by hierarchical clustering. The innovation lies in providing users with the flexibility to choose between automatic and semi-automatic modes for artifact detection, thereby enhancing versatility in EEG data analysis.

The rest of the thesis is organized as follows. Chapter 2 offers a comprehensive review of prior research on artifact removal methods employing ICA, wavelet trans-

form (WT), or hybrid approaches of these two techniques. And it reviews the artifact removal methods specifically designed for sleep EEG data. Chapter 3 presents the background information and preprocessing steps applied to the data utilized in this study. Chapter 4 outlines the artifact detection and removal process of the proposed method, beginning with an overview of the techniques discussed in Chapter 2. Additionally, it covers the extraction of features for clustering purposes. Chapter 5 presents the results and analysis of the performance evaluation conducted. Chapter 6 serves as the concluding chapter, providing a summary and comparison of the proposed methods against manual cleaning. Furthermore, it delves into potential areas for improvement and future research directions.

Chapter 2

Previous Works

2.1 Single Artifact Removal Techniques

The intricate nature of EEG artifacts overlaps with important signals across spectral, temporal, and sometimes spatial domains. This complexity challenges simple preprocessing techniques. Basic filtering or amplitude thresholding often leads to insufficient artifact removal and signal distortion. Consequently, numerous advanced methods and algorithms have been developed to more effectively detect and remove these artifacts from EEG signals [21]. Independent component analysis (ICA) and wavelet transform (WT) are considered as ones of most used single artifacts removal approaches [11][22].

2.1.1 Independent component analysis

ICA is primarily used to separate underlying sources in biomedical signal measurements. Vigaro et al. [23] demonstrated its effectiveness in separating linear mixtures and extracting ocular information from EOG signals. Jung et al. [24] effectively used extended ICA for removing eye activity artifacts from EEG, showing results comparable to regression algorithms. However these studies involved manual process which is one major limitation of ICA-based artifact detection and removal methods. They require manual intervention, as it is not inherently automatic. Typically, independent components (ICs) with visually detected artifacts need manual rejection after

decomposition. However, automation is achievable by labeling ICs through features that quantify the likelihood of being artifactual. This can be achieved by integrating ICA with complementary methods such as WT [21].

ICA is effective in reducing high-level additive noise and eliminating common component noises. Many methods involving ICA to detect and remove ocular artifacts have been explored in studies by Joyce et al. [25], Flexer et al. [26], Li et al. [27], Romero et al. [28], and Zhou and Gotman [29]. However, it falls short in effectively removing certain high-frequency noisy components when used in isolation [30].

2.1.2 Wavelet transform

Wavelet Transform (WT) is favored for processing non-stationary biomedical signals, including single-channel EEG, due to its time and frequency domain localization. WT has the capability to automatically remove artifacts by applying thresholding to its output. Studies by Zikov et al. [31], Ramanan et al. [32], Kumar et al. [33], and Khatun et al. [34] have shown the effectiveness of Stationary Wavelet Transform (SWT) with thresholding in automatic EOG artifact removal. Asaduzzaman et al. [35] also utilized Discrete Wavelet Transform (DWT) and thresholding for automated EOG and EMG artifact removal. Despite WT's effectiveness, its limitation in identifying overlapping artifacts has led to combined approaches with ICA [11].

2.2 Hybrid Method

Single-stage artifact removal methods often have limitations in effectively removing all types of artifacts from EEG data. These limitations can include incomplete artifact removal, requirement of reference channels, or the inability to handle certain types of artifacts effectively. As a result, researchers have proposed hybrid methods which combine two or more techniques or algorithms to overcome these limitations and enhance artifact removal performance. The ICA and WT hybrid method is regarded as the most popular hybrid method [22].

2.2.1 Wavelet enhanced Independent component analysis

Castellanos and Makarov [36] innovated in EEG artifact suppression with wavelet enhanced ICA (wICA), combining ICA and wavelet thresholding. Unlike traditional denoising, this method uniquely applies thresholding as a step towards refining demixed independent components. It has been proven effective on both actual and semi-simulated EEG data. Ghandeharion and Erfanian [37] developed a WT-based method to automatically detect and remove EOG artifacts in EEG data. This method assesses ICs by measuring correlation, kurtosis, and projection strength, and it stands out for its high accuracy, sensitivity, and specificity in ocular artifact detection. Notably, this approach uniquely automates the identification of artifacts without needing pre-calibration or predefined thresholds. Vázquez et al. [38] developed a method where they extracted various characteristics from wavelet coefficients on ICs, including statistical, frequency, spatial, and template correlation traits. This approach was used to train a supervised classification model, enabling automatic identification and removal of ocular, high-frequency muscle, and ECG artifacts from EEG data. Al-Qazzaz et al. [39] introduced a novel four-stage approach for EEG signal enhancement, particularly in post-stroke dementia studies. This method combines automatic independent component analysis (AICA) with WT. It starts by estimating independent components, then identifies artifacts using metrics like skewness, kurtosis, and sample entropy. Following this, it denoises artifacts using Discrete Wavelet Transform (DWT) and reconstructs clean EEG signals. In 2022, Maddirala and Veluvolu [40] introduced a comprehensive framework for eye blink artifact removal in EEG data, integrating ICA with continuous wavelet transform (CWT), k-means clustering, and singular spectrum analysis (SSA). This framework uses CWT and k-means to pinpoint the eye-blink artifact region, followed by an SVM-based classifier that automatically identifies the artifact-containing ICs. Wavelet-enhanced ICA has demonstrated effectiveness in detecting and removing various types of artifacts. Simultaneously, WT enables the

complete recovery of neural components from EEG channels corrupted by artifacts outside of the contaminated frequency range. Therefore, the enhancement provided by wavelet-enhanced ICA contributes to preserving cerebral activity.

2.3 Artifact Removal for Sleep EEG Data

The methods discussed above primarily utilize EEG data obtained from healthy adults in awake states. However, there are specific approaches that have been proposed for the analysis of sleep EEG data. These methods are specifically tailored to handle the unique characteristics and challenges associated with sleep recordings. By considering the specific dynamics and patterns present during sleep, these techniques aim to enhance the accuracy and reliability of EEG analysis in the context of sleep studies.

Schetinin and Schult [41] introduced a novel method that combines a polynomial neural network and decision tree for detecting artifacts in sleep EEG recordings of newborns. Betta et al. [42] and Dursun et al. [43] proposed techniques specifically designed for ocular artifact detection and removal in automated sleep analysis. Betta et al. [42] utilized WT and adaptive filtering, while Dursun et al. [43] employed correlation and wavelet-based rules to separate EOG artifacts. In 2022, Ranjan et al. [44] developed a method for automatic detection and removal of cardiac artifacts from single-channel EEG data. This method incorporated empirical wavelet transforms (EWTs), adaptive threshold-based nonlinear Teager-Kaiser energy operator (TEO), customized morphological filters, and modified ensemble average subtraction (MEAS).

Additionally, machine learning algorithms have been utilized for artifact detection. Saifutdinova et al. [45] employed a random forest classifier on features extracted from multi-channel sleep EEG data to detect various types of artifacts. In 2019, they proposed an unsupervised method based on Riemannian Geometry, which exhibited superior performance on multi-channel EEG data compared to other online automatic EEG artifact removal tools [45].

Chapter 3

Data

3.1 The Dataset

In the context of a comprehensive dataset encompassing 75 pediatric patients, including 58 males and 17 females, all-night sleep polysomnography (PSG) recordings were conducted. Patient ages ranged from 2-18 years with a mean of 8.81 years and standard deviation of 4.52 years. Each recording was accompanied by a clinician event file, providing detailed annotations for apnea-hypopnea, arousals, and sleep state labels for every epoch. The PSG recordings were integral to a clinical study (Pro00057638) approved by the University of Alberta, focusing on potential obstructive sleep apnea in pediatric patients. Exclusively EEG channels were employed for PSG recordings in this thesis. The reference electrodes, identified in the Table 3.1, are denoted interchangeably as M1 and M2, given their placement on the mastoid. PSG channels are designated by the combination of the recording site and reference electrodes. For instance, F4M1 represents the PSG channel capturing the electrical potential difference between electrodes F4 and M1.

The specific patient dataset utilized in this work pertains to one patient, labeled as CF050. This patient was selected from the larger cohort of 74 pediatric patients. The EEG signals were recorded at a sampling rate of 512 samples per second. These signals were segmented into a sequence of epochs, each lasting 30 seconds. CF050 contains 948 epochs, 15,360 observation records in each epoch. The 30-second segmentation

Channel Name	Placement
F4/F3	Right/Left Frontal
C4/C3	Right/Left Central
T6/T5	Right/Left Temporal
O2/O1	Right/Left Occipital

Table 3.1: EEG electrodes names and replacement.

aligned with the binary classification of sleep and awake states, simplifying the categorization into two distinct states. Sleep stages were labeled by a sleep technician, designating each epoch as either part of the sleep state or indicative of wakefulness. The proportion of sleep states in the CF050 dataset is shown in Table 3.2.

States	Asleep	Awake
# Epochs	695	253
% Epochs	73.312	26.688

Table 3.2: Proportion of sleep states of CF050 patient.

3.2 Data Pre-processing

In the CF050 EEG dataset, PSG recording captured the entire session from the initial calibration of the equipment to its deactivation at night’s end, covering a duration of 7.9 hours. This period included instances of muscle movement during sleep. To enhance signal quality and minimize interference, a 4th order Butterworth band-pass filter by eegfilter function in eegkit library in R was employed for eight EEG channels [46][47]. The filter, with a cutoff frequency range of 0.3 to 80 Hz, was chosen for its ability to isolate relevant brainwave frequencies and eliminate artifacts, including equipment noise and muscle movement, while ensuring minimal signal distortion. Then the pre-filtered EEG data was segmented into 30 second epochs for further analysis.

Chapter 4

Methodology

Figure 4.1 presents a block diagram outlining the process of artifact removal from EEG data using a wavelet enhanced ICA approach, starting with a Butterworth band-pass filtered EEG divided into 30-second epochs [47]. Each epoch of EEG data is first separated into eight independent components (ICs) via fastICA function in fastICA library in R [48]. Subsequently, each IC is further decomposed using DWT to the fourth level with 'db4' as the mother wavelet by dwf function in waveslim library in R [49].

Wavelet coefficients (WCs) across all ICs and epochs that have the same frequency band are compiled together. Features such as kurtosis, zero-crossing count, skewness, activity, mobility, complexity and Renyi's entropy are extracted from each independent wavelet coefficients. Euclidean distance of features between epochs are calculated for hierarchical clustering with average-linkage using hclust function in R [50].

The method offers flexibility in artifact detection, enabling both automatic and semi-automatic identification. In automatic mode, the method designates the smaller cluster in the initial division as an artifact. For semi-automatic detection, the user can determine which cluster to label as an artifact by interpreting the dendrogram's output. The WCs marked as artifact are removed in each ICs and each epoch. Cleaned eight channels EEG data are reconstructed by reversing DWT and ICA processes.

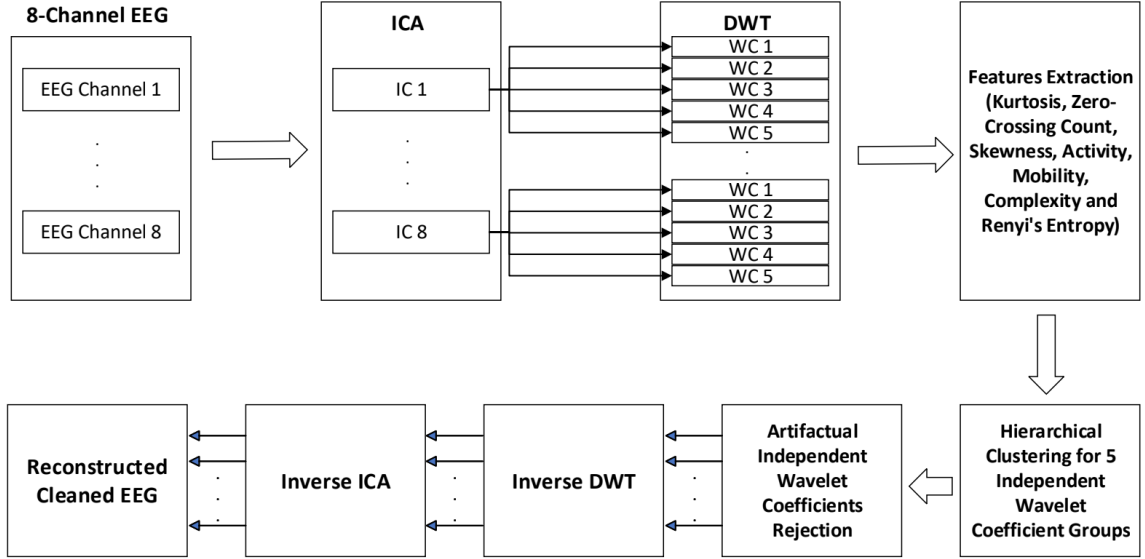


Figure 4.1: Block diagram of the proposed method of wavelet enhanced ICA technique for artifact removal from EEG data.

4.1 Independent Component Analysis

Since EEG measures the electrical activity generated by the brain through electrodes placed on the scalp, it often contains a mixture of signals originating from various sources within the brain. Independent Component Analysis (ICA) is a statistical signal processing technique that is commonly used to decompose the observed EEG signals into a set of statistically independent components, which correspond to different underlying brain processes or sources. This separation is achieved by assuming that the observed EEG signals are linear combinations of these independent sources.

ICA could be applied to successfully separate and analyze the underlying sources in the observed data based on several assumptions. It assumes that the observed EEG signals are a linear mixture of independent components and the mixing process is spatially stable, meaning that the mixing coefficients relating the sources to the recorded EEG signals remain constant over time. It also assumes that the independent components are statistically independent and the independent sources have non-Gaussian probability distributions, so that it allows ICA to exploit the non-Gaussianity of the

sources to separate them. In addition, ICA assumes that the number of independent sources contributing to the observed EEG signals is no greater than the number of recorded electrodes [36].

The observed signals are represented by the random vector $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_m)^T$ while the source components by the random vector $\mathbf{S} = (\mathbf{s}_1, \dots, \mathbf{s}_m)^T$. \mathbf{x}_n and \mathbf{s}_n are expressed as

$$\mathbf{x}_n = [x_n(1) \ x_n(2) \ \dots \ x_n(i)]^T, \text{ for } n = 1, 2, \dots, m$$

$$\mathbf{s}_n = [s_n(1) \ s_n(2) \ \dots \ s_n(i)]^T, \text{ for } n = 1, 2, \dots, m$$

where $x_n(i)$, $s_n(i)$ denote an observed signal and source component at a discrete time i , correspondingly. The relationship between the observed signal and the source component has the following expression:

$$\mathbf{X} = \mathbf{A}\mathbf{S}$$

where \mathbf{A} is the unknown mixing coefficients matrix which is to be estimated by ICA algorithms. The source components \mathbf{S} could be revealed with the un-mixing matrix $\mathbf{W} = \mathbf{A}^{-1}$ by the equation:

$$\mathbf{S} = \mathbf{X}\mathbf{W}$$

ICA relies on the assumption that the independent components have a non-Gaussian distribution, allowing for their statistical independence to be exploited. Several algorithms, such as Infomax, FastICA and JointICA, have been developed to estimate the mixing matrix \mathbf{A} and the independent components \mathbf{S} based on this assumption.

4.1.1 FastICA

The rows of un-mixing matrix \mathbf{W} is denoted by \mathbf{w}_i^T , $i = 1, \dots, n$. The FastICA algorithm is to maximize the non-Gaussianity of the estimated components to achieve

separation. Non-Gaussianity is quantified with the approximation of negentropy J_G with expression:

$$J_G = \sum_{i=1}^n E[G(\mathbf{w}_i^T \mathbf{Z})]$$

Where G is a non-quadratic function with assumptions of even and symmetric and $\log \cosh$ is used as the non-quadratic function in this thesis. The symbol E represents the mean value over the whitened matrix \mathbf{Z} . FastICA follows several steps to find the local maxima of the cost function:

1. Preprocessing:

- (a) Centering: Center the observed signal \mathbf{X} by subtracting the mean so that \mathbf{X} could be assumed to have the zero mean. The centering is expressed as:

$$\mathbf{x}_n = \mathbf{x}_n - E[\mathbf{x}_n], \text{ for } n = 1, 2, \dots, m$$

where E is the mean value.

- (b) Whitening: Linear transform the centered observed signal into uncorrelated variables with unit variances. Whitening helps to simplify the estimation process. The observed signal \mathbf{X} is transformed linearly and the whitened matrix \mathbf{Z} is obtained, which

$$E[\mathbf{Z}\mathbf{Z}^T] = \mathbf{I}$$

Eigenvalue decomposition of a matrix is used to decompose the matrix \mathbf{X} , $E[\mathbf{X}\mathbf{X}^T] = \mathbf{E}\mathbf{D}\mathbf{E}^T$, where \mathbf{E} is the eigenvector matrix and \mathbf{D} denotes the diagonal matrix of eigenvalues. The observed signal matrix can be whitened by

$$\mathbf{X} = \mathbf{D}^{-1/2}\mathbf{E}^T\mathbf{X}$$

2. Initialization: Initialize a random vector \mathbf{w}_i as the initial estimate of the mixing matrix.

3. Iterative Update:

(a)

$$\mathbf{w}_i^* = E[\mathbf{Z}g(\mathbf{w}_i^T \mathbf{Z})^T] - E[g'(\mathbf{w}_i^T \mathbf{Z})]\mathbf{w}_i$$

where g is the derivative of the nonlinear function G , \mathbf{Z} denotes the whitened observed signal.

(b)

$$\mathbf{w}_i = \mathbf{w}_i - \sum_{j=1}^{i-1} (\mathbf{w}_i^T \mathbf{w}_j) \mathbf{w}_j$$

(c)

$$\mathbf{w} = \frac{\mathbf{w}_i^*}{\|\mathbf{w}_i^*\|}$$

(d) Repeat step a to c if not converged, otherwise back to step 2 with $i = i + 1$ until all independent components are extracted.

4.2 Wavelet Transform

Wavelet transform is a mathematical technique used for analyzing signals. It decomposes a given signal into a set of wavelets, which are small wave-like functions. Unlike other transforms such as the Fourier transform, which uses sinusoidal basis functions, wavelet transform uses wavelets that are typically localized in both time and frequency [51].

The wavelet transform allows you to examine different frequency components of a signal at different resolutions. It provides a time-frequency representation of the signal, capturing both time and frequency information simultaneously. This makes it particularly useful for analyzing signals with non-stationary characteristics, where the properties change over time and the EEG signal is considered non-stationary [52].

Wavelet transform is usually applied in two ways, continuous wavelet transform (CWT) and discrete wavelet transform (DWT).

4.2.1 Continuous Wavelet Transform

The process of wavelet transform involves convolving the signal with a set of wavelet functions, known as the mother wavelet, at different scales and positions. The resulting coefficients represent the strength of the correlation between the wavelet and the signal at different time points and scales. By analyzing these coefficients, the presence of specific frequencies in the signal and their temporal localization could be identified [53]. The continuous wavelet transform can be defined by

$$CWT\{x(t); a, b\} = \frac{1}{\sqrt{|a|}} \int_{-\infty}^{\infty} f(t) \psi\left(\frac{t-b}{a}\right) dt$$

Where f is the signal, ψ is known as the basic wavelet, a is the scaling parameter and b is position parameter. Time t and the parameters $a, b \in \mathbb{R}$ are real continuous variables. $|a|^{-1/2}$ is used to normalize the energy.

The CWT offers high resolution in time for high-frequency components and high resolution in frequency for low-frequency components. However, it requires continuous sampling of the signal, making it computationally intensive for large datasets.

4.2.2 Discrete Wavelet Transform

The discrete wavelet transform is a variation of the wavelet transform that operates on discrete-time signals. It decomposes a signal into a series of discrete wavelets at different scales. The DWT provides a discrete-time representation of a signal's frequency content and is particularly suited for digital signal processing applications.

The DWT uses a set of discrete wavelet functions, called the wavelet basis or filter banks. These wavelets are scaled and translated to analyze the signal at different scales and positions. In this scenario, the parameters a and b are frequently derived from powers of two, known as dyadic scales and translations.

$$a = 2^j, b = k2^j$$

j is for level, k is for location $j, k \in \mathbb{Z}$

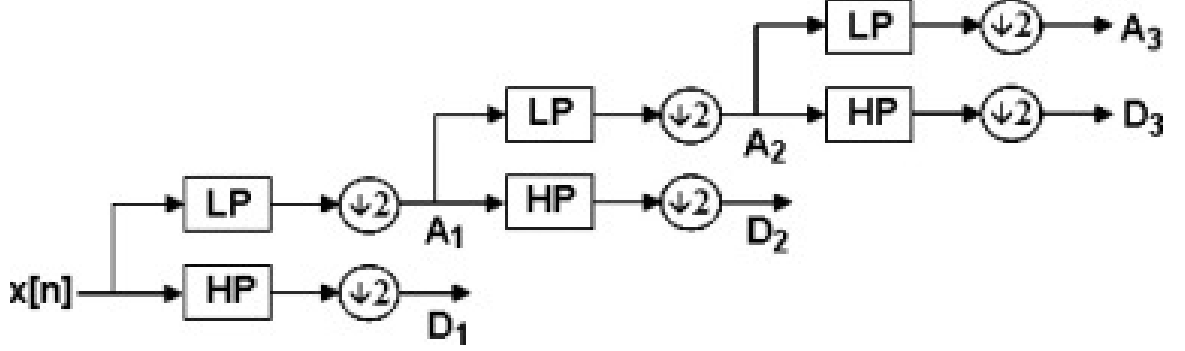


Figure 4.2: Third level wavelet decomposition of an EEG signal [55].

Mallet [54] developed a series of low-pass (LP) and high-pass (HP) filters applied to the signal, separating it into approximation coefficients which represents the low-frequency components and detail coefficients which represents the high-frequency components at each level of decomposition.

According to Figure 4.2, in the initial step of the DWT, the signal undergoes simultaneous filtering using a LP and HP filter, and the signals are down-sampled with the factor of 2. The resulting outputs from the LP and HP filters are referred to as the approximation (A_1) and detail (D_1) coefficients of the first level, respectively. These output signals have half the frequency bandwidth of the original signal and can be downsampled by a factor of two, following the Nyquist rule. The same process is repeated for the first level approximation and detail coefficients, generating the second level coefficients. This decomposition procedure doubles the frequency resolution through filtering and halves the time resolution through downsampling at each step.

By recursively applying the DWT on the approximation coefficients, the signal can be decomposed into multiple levels or scales. This multilevel decomposition allows for a hierarchical representation of the signal's frequency content. The DWT also offers efficient implementation and is well-suited for analyzing discrete-time signals, such as digital audio or image data.

Daubechies wavelets

The Daubechies' family of wavelets [56] is widely recognized as one of the most frequently employed sets of orthogonal wavelets that satisfy the necessary admissibility conditions. This property enables the reconstruction of the original signal from its wavelet coefficients. Figure 4.3 provides illustrations of wavelet and scaling functions for the Daubechies' family of orthogonal wavelets. The research on bio-signal classification using wavelet techniques has predominantly focused on the utilization of the Daubechies family of wavelets, specifically of order 2 or 4 [57].

4.3 Feature Extraction

4.3.1 Kurtosis

Kurtosis is a commonly used measure for quantifying the non-Gaussianity of a distribution. It serves as a classical method to assess the departure from Gaussianity. When data is preprocessed to have a unit variance, the kurtosis corresponds to the fourth moment of the data.

The kurtosis of an independent wavelet coefficient x with the assumption of zero mean is denoted by $kurt(x)$ with the expression:

$$kurt(x) = E[x^4] - 3(E[x^2])^2$$

4.3.2 Zero-Crossing Count

Zero-Crossing Count is a time domain feature that quantifies how frequently the independent wavelet coefficient crosses the zero axis [58]. The expression of zero crossing count is as below

$$ZC(i) = \sum_{n=0}^{N-1} |\text{sign}[x_i(n)] - \text{sign}[x_i(n-1)]|$$

where $\text{sign}[x_i(n)]$ is defined as $\text{sign}[x_i(n)] = \begin{cases} 1, & \text{if } x_i(n) > 0 \\ -1, & \text{if } x_i(n) < 0 \end{cases}$

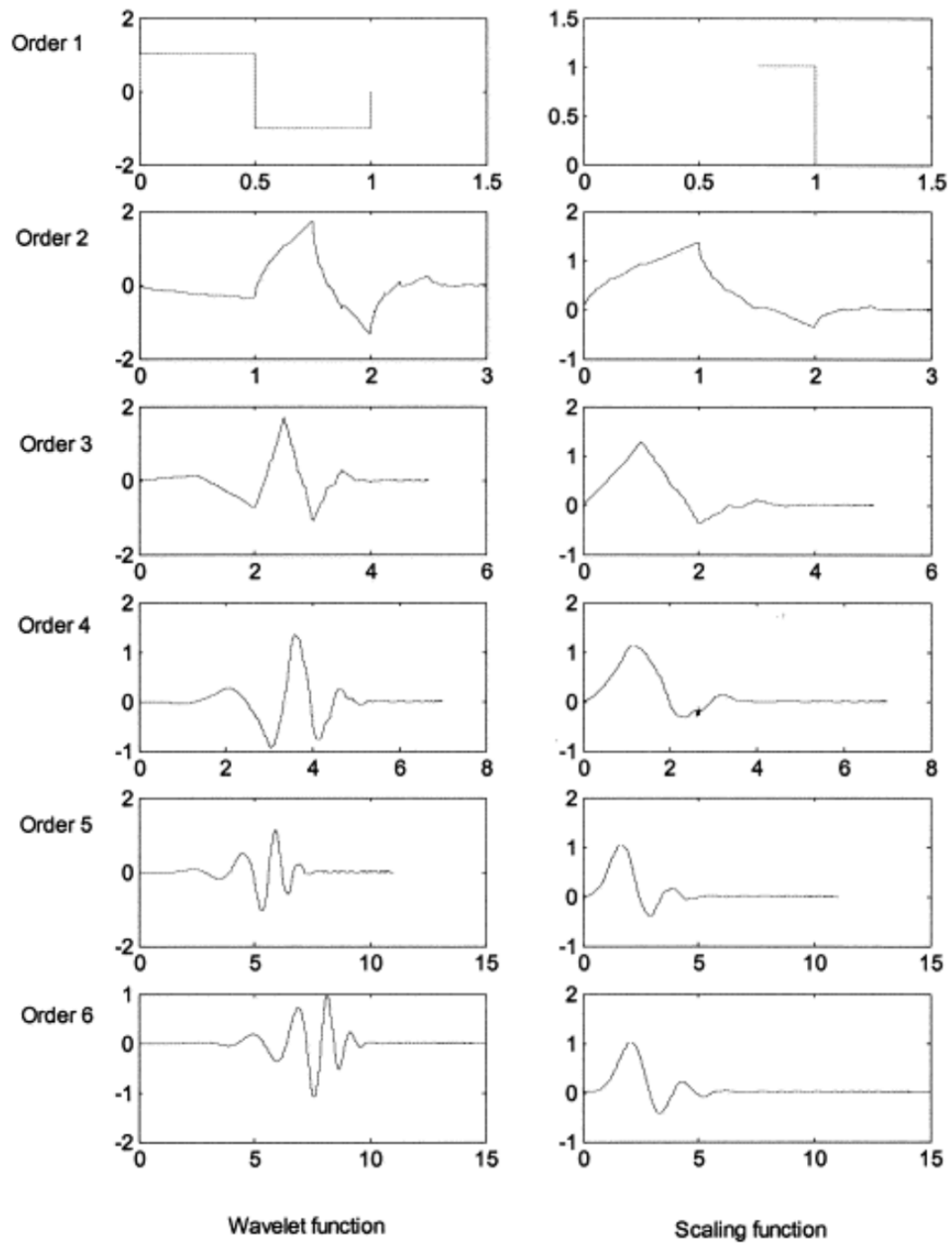


Figure 4.3: Daubechies wavelet and scaling functions of different orders [51].

4.3.3 Skewness

Skewness measures the asymmetry in the distribution of the independent wavelet coefficient amplitudes.

$$\text{Skewness} = \frac{1}{N} \sum_{i=1}^N \left(\frac{x_i - \mu}{\sigma} \right)^3$$

where N is the number of data points in each epoch, x_i are the individual data points, μ and σ are the mean and standard deviation of each epoch of each channel.

4.3.4 Hjorth Features

Hjorth parameters are statistical descriptors in the time domain used to characterize the independent wavelet coefficient, encompassing three distinct types: Activity, which measures signal variance; Mobility, indicative of the signal's mean frequency; and Complexity, reflecting the signal's frequency variation [59].

Activity

Activity indicates the variance of the EEG data which is a statistical measure used to quantify the amount of variation or dispersion. It is defined as the average of the squared deviations of each data point from the corresponding mean of the channel [60].

The Activity of EEG data could be expressed as

$$\sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

where N is the number of data points in each epoch, x_i are the individual data points, μ is the mean of each epoch of each channel.

Mobility

Mobility represents the mean frequency of the EEG data. It is calculated based on the variance and the variance of the first derivative of the data. It is calculated with

expression

$$\text{Mobility}(t) = \sqrt{\frac{\text{var}\left(\frac{dx(t)}{dt}\right)}{\text{var}(x(t))}}$$

where $\text{var}(x(t))$ is the variance of the EEG data, $\text{var}\left(\frac{dx(t)}{dt}\right)$ is the variance of the first derivative of the data.

Complexity

Complexity is a measure that compares the signal's similarity to a pure sine wave. A signal with complexity closer to 1 is more sinusoidal [60]. It is expressed based on the mobility of the EEG data and the first derivative of the data with expression

$$\text{Complexity}(t) = \frac{\text{Mobility}\left(\frac{dx(t)}{dt}\right)}{\text{Mobility}(x(t))}$$

4.3.5 Renyi's Entropy

Entropy, in the context of EEG signals, serves as a statistical descriptor quantifying the unpredictability or variability within the signal. It assesses the randomness inherent in the EEG data, reflecting the degree of disorder and the distribution of different frequency components in the signal's time series. It is defined mathematically as

$$\text{Entropy} = -\sum_{i=1}^n p(x_i) \log p(x_i)$$

where $p(x_i)$ is the probability of the i -th outcome of the random variable, n is the total number of possible events.

Renyi's entropy is a generalization of the entropy and it represents a whole family of entropy measures. It is calculated as

$$\text{Renyi's Entropy} = \frac{1}{1-\alpha} \log \left(\sum_{i=1}^n p(x_i)^\alpha \right)$$

where $p(x_i)$ and n have the same meanings as entropy and α is the order of Renyi's entropy. In this study, α was set as 2 and the histogram with $n=500$ possible events was used to estimate the probability density of the wavelet coefficients.

4.4 Hierarchical Clustering

In EEG artifact detection, an unsupervised learning approach like hierarchical clustering is valuable due to the lack of labeled data indicating whether an epoch is an artifact. However, the other unsupervised learning algorithms such as K-means and fuzzy C-means require the number of clusters as a priori. Hierarchical clustering is chosen for its ability to reveal the overall structure and relationships among clusters. It generates a hierarchical tree, providing a nested set of partitions visualized through a dendrogram. This dendrogram illustrates the strength of inter-electrode similarities, offering insights into the EEG data structure without predefined labels, making it an effective tool for identifying artifacts in EEG signals.

Hierarchical clustering in this study adopts an agglomerative approach, which groups smaller clusters into larger ones, as opposed to the divisive method that splits larger clusters into smaller ones. The work follows the methodology used by Mert and Akan [61] and Mammone [62]. Among the various hierarchical methods, single-linkage, complete-linkage, average-linkage, and Ward’s method, average-linkage was chosen after testing for its superior results in this context. This method effectively splits clusters with artifact and without in EEG data by progressively merging clusters based on average distances.

Average linkage is the method where the distance between two clusters is defined as the average distance between all pairs of objects, where one member of the pair is from each cluster. It is represented by

$$D(A, B) = \frac{1}{|A| \times |B|} \sum_{a \in A, b \in B} d(a, b)$$

where $D(A, B)$ is the distance between two clusters A and B , $|A|$ and $|B|$ are the number of elements in clusters A and B respectively. $d(a, b)$ is the distance between elements a and b . In this work, Euclidean distance was used and calculated among elements

$$d(P, Q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

where $d(P, Q)$ represents the Euclidean distance between P and Q , $P = (p_1, p_2, \dots, p_7)$, $Q = (q_1, q_2, \dots, q_7)$ are features extracted from independent component wavelet coefficients of each epoch EEG data.

4.5 Inverse DWT

Inverse Discrete Wavelet Transform (IDWT) reconstructs the ICs from WCs. After a signal is decomposed using DWT into components at different scales or resolutions, IDWT combines these components back into the artifactual components removed ICs. The process involves summing up the wavelet series, typically represented by:

$$s(t) = \sum_j \sum_k c_{j,k} \psi_{j,k}(t)$$

where $s(t)$ represents the reconstructed IC, $c_{j,k}$ are the wavelet coefficients, $\psi_{j,k}(t)$ are the wavelet basis functions at scale j and position k .

4.6 Inverse ICA

Inverse ICA refers to the process of reconstructing the multi-channel EEG signal from the ICs that were separated by ICA. After artifact components are detected and removed, the remaining components are recombined using inverse ICA to reconstruct the artifact-free EEG data by applying the mixing matrix \mathbf{A} , the inverse matrix of unmixing matrix \mathbf{W} , to ICs.

$$\mathbf{X}' = \mathbf{A}\mathbf{S}'$$

where \mathbf{X}' is the reconstructed artifact free multi-channel EEG data, and \mathbf{S}' is the ICs with artifactual components removed.

Chapter 5

Outputs and Results

5.1 Outputs

We evaluated our proposed method using raw, full-night sleep EEG data and conducted a comparative analysis with EEG signals that were either only preprocessed or manually cleaned. The proposed system demonstrated effective removal of artifactual components while preserving the cerebral activities of interest.

Given that the dataset encompasses a full night’s sleep EEG signal, consisting of 948 epochs, we specifically illustrate the proposed method by examining a random epoch among 948 epochs, epoch number 700, as a representative example.

The preprocessing and methodology for recording EEG signals is outlined in Chapter 3 and Chapter 4. Initially, the captured signal undergoes processing through a 4th order Butterworth band-pass filter, which isolates the desired frequency bands. Illustrations of the raw eight-channel EEG signal, alongside the filtered signal, are presented in Figure 5.1. It is observed that artifacts not within the target frequency bands are removed by this filtering process, preserving the integrity of both the time and frequency resolution of the EEG signal.

Nevertheless, artifacts that overlap with the EEG signal’s frequency resolution were largely not eliminated. To address this, ICA was applied to the eight-channel signal depicted in Figure 5.1b, aiming to break down the signal into its statistically independent and non-Gaussian components. The decomposition process yielded eight

independent components (ICs) as demonstrated in Figure 5.2. The decomposition of each IC into wavelet coefficients (WCs) was further carried out using the fourth level DWT. Figure 5.3 illustrates the wavelet coefficients from the first IC of epoch 700.

Values of kurtosis, zero-crossing count, skewness, Hjorth parameters and Rényi entropy are extracted from each WC and each epoch are shown by boxplots in Figure 5.4.

For each set of WCs, dendrograms were generated. Figure 5.5 shows the dendrograms based on the Euclidean distance of features between epochs with average-linkage.

In the automatic cleaning method, dendrograms are automatically segmented into two clusters, with the smaller cluster identified as artifacts and the larger as signals. Conversely, the semi-automatic cleaning method necessitates visual inspection to determine which clusters are artifactual components. Upon inspection, D1 is categorized into two clusters, while D2, D3, and A4 are each divided into four clusters, and D4 into five clusters. Clusters containing the majority of data points are classified as signal, with the remaining designated as artifacts and subsequently removed. After excising artifactual components and reconstructing the EEG signal using inverse DWT and inverse ICA, Figure 5.6 showcases the cleaned EEG data achieved through both automatic and semi-automatic methods.

5.2 Performance Evaluation

5.2.1 Accuracy

Accuracy is a measure used to evaluate the performance of a classification model. It represents the ratio of correctly predicted instances to the total instances in a dataset. In the context of binary classification, accuracy is calculated as:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

5.2.2 Results

Due to the absence of labeled artifactual components in the raw dataset, assessing the accuracy for the proposed method is challenging. To evaluate performance, we employed a sleep stage classification algorithm for two sleep stage classification, specifically the Hidden Markov Model (HMM) proposed by Kashlak [63]. Using both manually cleaned EEG data and EEG data cleaned by our proposed method, we applied the same sleep stage classification algorithm. Comparison of the outcomes enables an assessment of whether the proposed method enhances or achieves comparable results in sleep stage classification compared to manually cleaned EEG data.

Only channel C4M1 underwent manual cleaning, and its performance is contrasted with all eight channels cleaned by our proposed method to identify potential improvements. To assess potential improvements over the raw EEG dataset, we applied the sleep stage classification model solely to pre-processed EEG data as a control group.

To facilitate the application of the HMM algorithm, EEG data underwent transformation into Power Spectral Density (PSD) using Welch’s method. Given the randomness start in the HMM model, we executed the model for each channel of each method 50 times to evaluate performance and stability. Accuracy metrics were computed from individual confusion matrices, with mean and standard deviation calculated from the accuracy obtained over the 50 model runs.

According to Table 5.1, the Semi-automatic method consistently outperforms the other two in terms of mean accuracy, with its superiority most evident in channel C4M1. And the Pre-processed method obtains a lower mean accuracy compared with the Automatic and Semi-automatic method. However, the Semi-automatic method does not maintain this lead in channel F3M2, where it falls behind the Pre-Processed method, and the Automatic method also sees its lowest mean accuracy in this same channel. Exceptional performance is observed from all methods in channels C3M2, O2M1, and O1M2, with the Semi-automatic method demonstrating perfect consis-

Method/Channel	Pre-Processed		Automatic		Semi-automatic	
	Mean	Std Dev	Mean	Std Dev	Mean	Std Dev
F4M1	0.7484	0.1393	0.8080	0.1583	0.8062	0.1590
F3M2	0.5667	0.0907	0.5408	0.0670	0.5636	0.0973
C4M1	0.7432	0.1275	0.8630	0.1474	0.9391	0.0698
C3M2	0.9473	0.0000	0.9387	0.0422	0.9473	0.0000
T6M1	0.6939	0.0609	0.7389	0.0402	0.7565	0.0289
T5M2	0.5761	0.0701	0.5865	0.0785	0.5777	0.0769
O2M1	0.9092	0.0988	0.9641	0.0000	0.9641	0.0000
O1M2	0.8676	0.1146	0.9641	0.0000	0.9641	0.0000

Table 5.1: Mean and standard deviation of accuracy for pre-processed only, automatic and semi-automatic clean methods across eight EEG channels over 50 model runs. Numbers in bold are the highest mean accuracy for each channel.

tency, as evidenced by a standard deviation of zero in these channels. Moreover, it shows enhanced stability compared to the Automatic method, particularly in channel C4M1 where it boasts a reduced standard deviation.

Additionally, we conducted two paired t-tests to assess differences between the Pre-processed method and both the Automatic and Semi-automatic methods, focusing specifically on the C4M1 channel. Our hypothesis is that both the Automatic and Semi-automatic methods would yield the same accuracies as the Pre-processed method for C4M1 channel. The paired t-test comparing the Pre-processed only method with the automatic method resulted in a p-value of less than 0.001, with a t-statistic of -8.0933. Similarly, the paired t-test comparing the Pre-processed only method with the Semi-automatic method yielded a p-value of less than 0.001, with a t-statistic of -12.4886. As both p-values are less than 0.001, null hypothesis is rejected.

Furthermore, we calculated 95% confidence intervals to further assess the accuracy of the Automatic and Semi-automatic methods for the C4M1 channel. The 95% con-

fidence interval for the Automatic method was (0.8221, 0.9039), indicating that the accuracy obtained from manual cleaning the C4M1 channel (0.9146) falls outside this interval, suggesting that the confidence interval is below the accuracy obtained from manual cleaning. Conversely, the 95% confidence interval for the Semi-automatic method was (0.9198, 0.9584), indicating that the accuracy obtained from manual cleaning also falls outside this interval, but the confidence interval is above the accuracy obtained from manual cleaning.

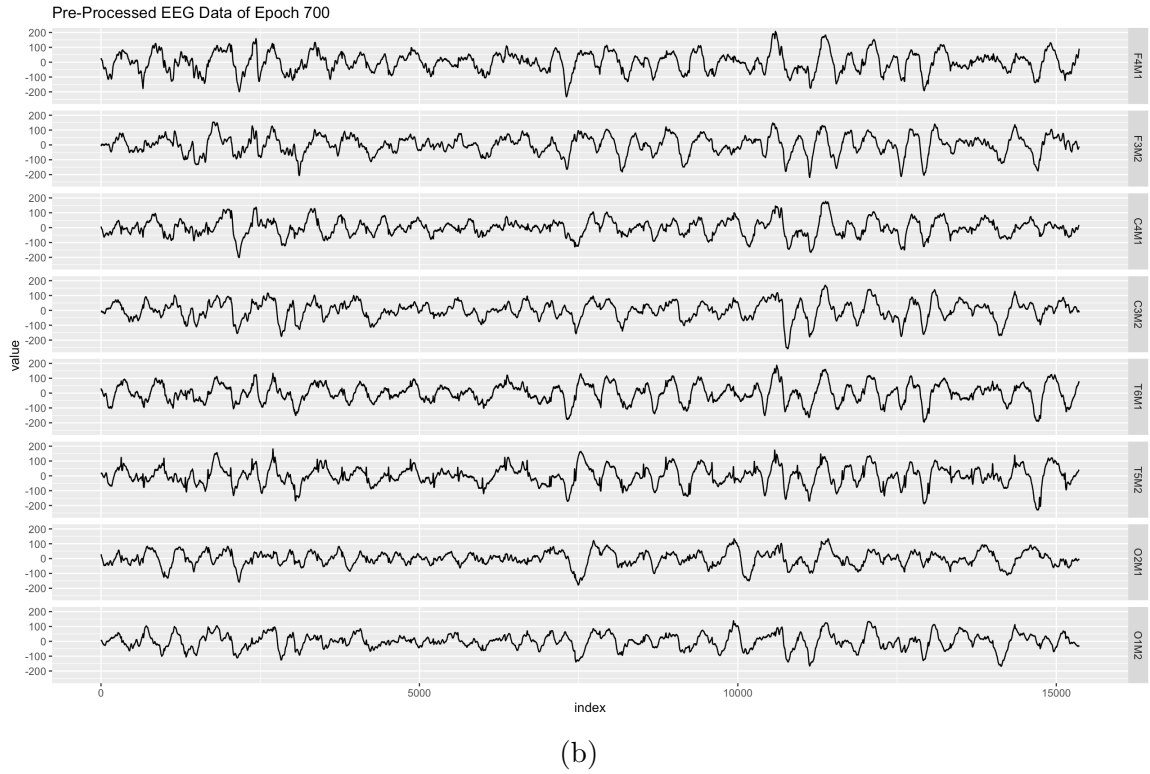
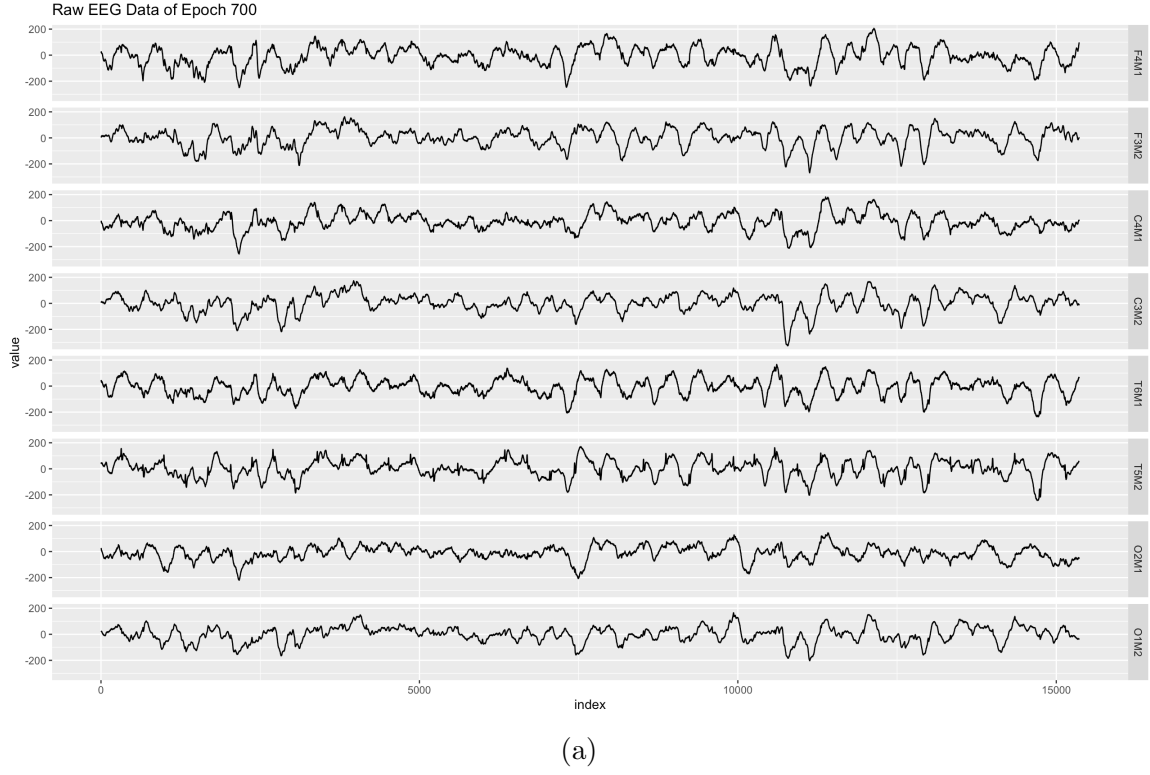


Figure 5.1: The raw eight-channel EEG signal is shown in (a) and the signal after pre-processing step is shown in (b). This step is conducted to retain only the frequency bands of interest (0.3 to 80 Hz).

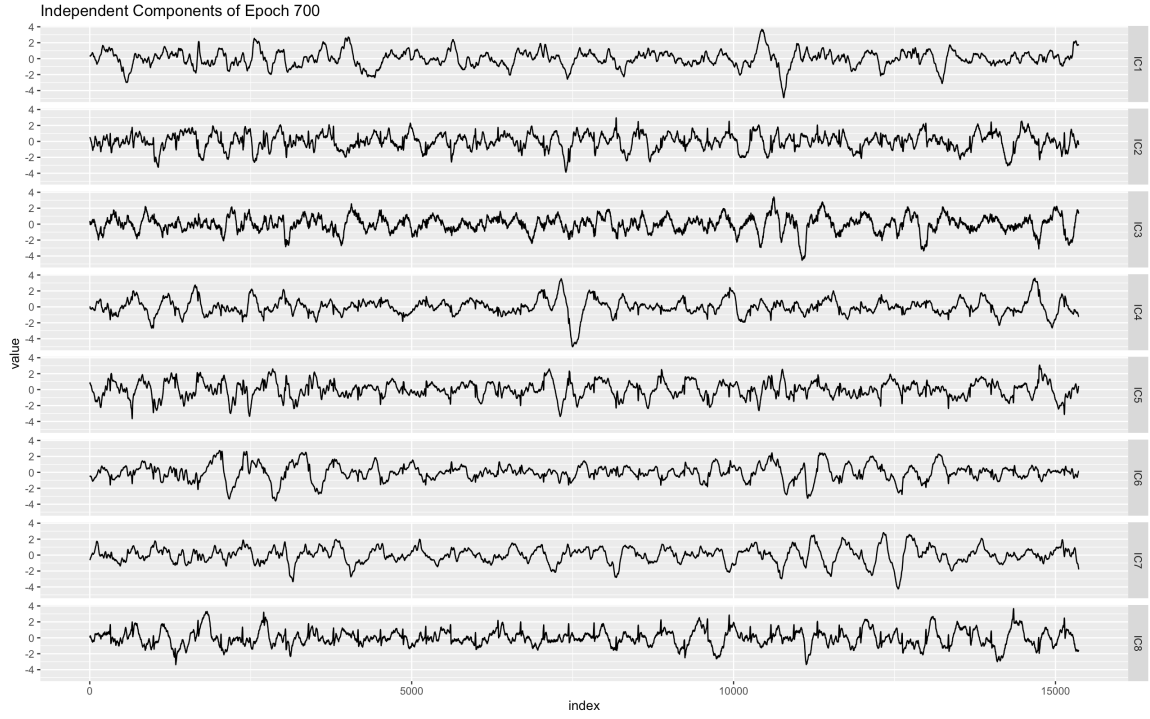


Figure 5.2: Decomposed eight ICs from pre-processed EEG signal using ICA.

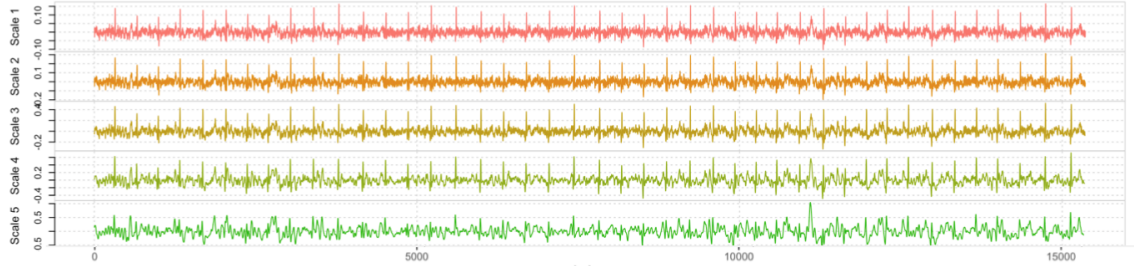


Figure 5.3: WCs from the first IC in Figure 5.2 by DWT. The wavelet coefficients from top to bottom are D1, D2, D3, D4 and A4 respectively.

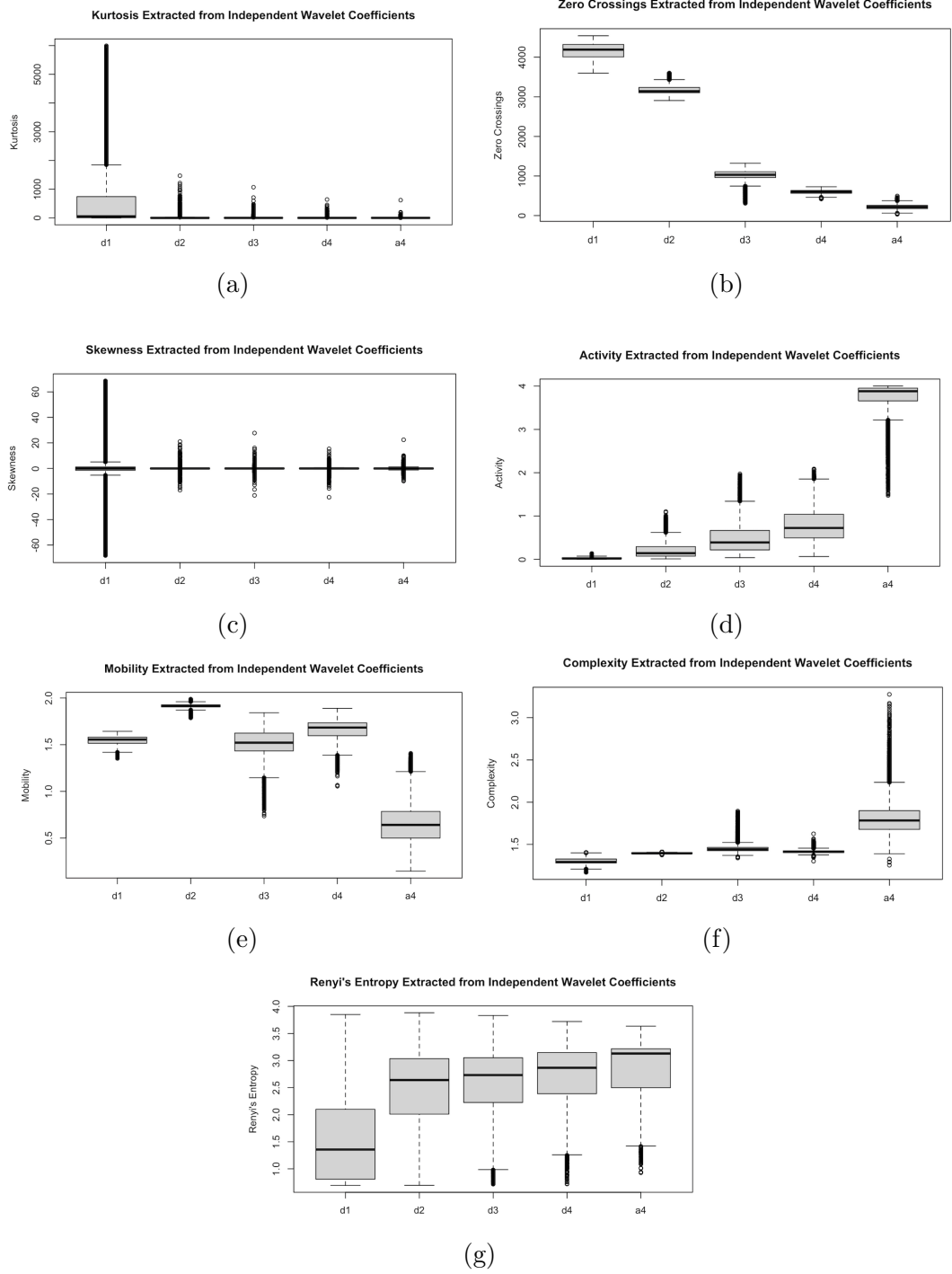


Figure 5.4: Boxplots of features extracted from each WCs. The seven features presented are (a) kurtosis, (b) zero-crossing count, (c) skewness, (d) activity, (e) mobility, (f) complexity and (g) Renyi's entropy.

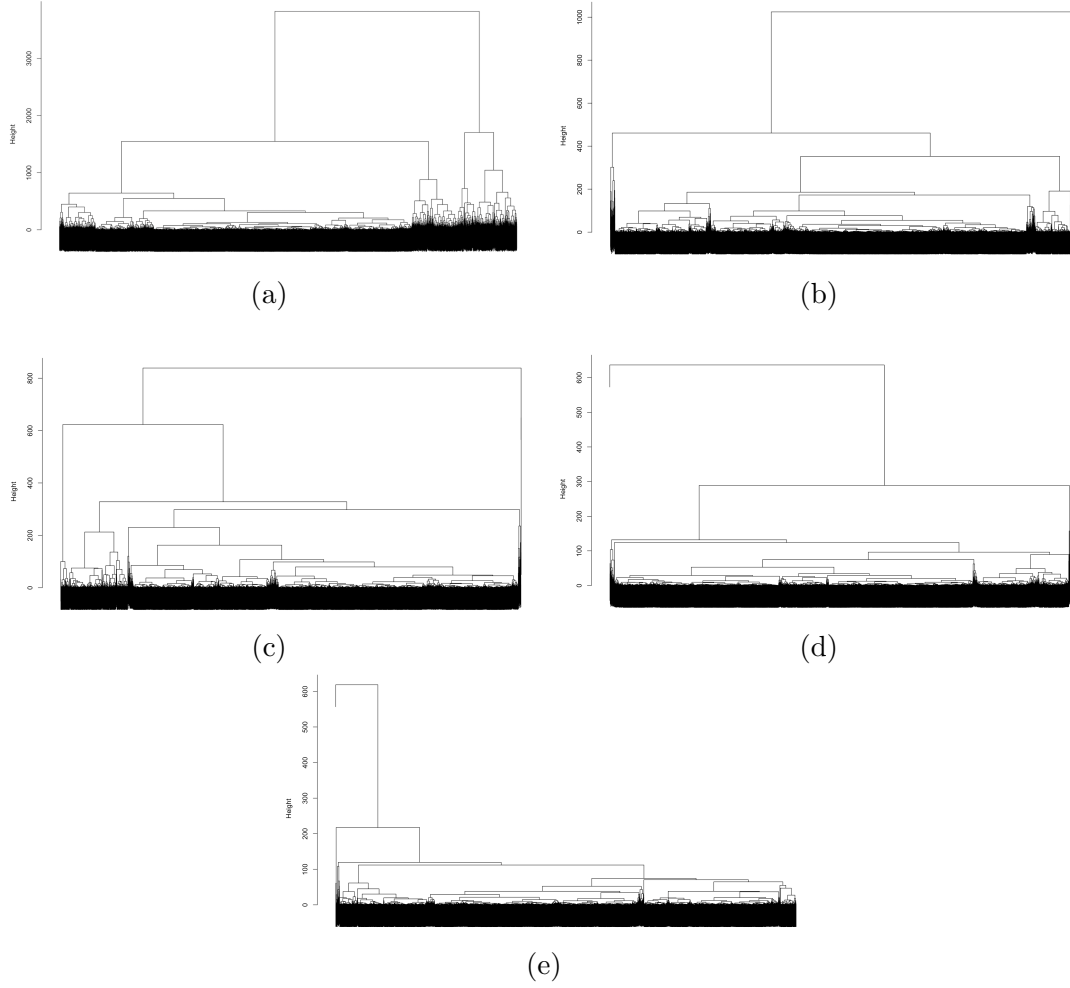


Figure 5.5: Dendrograms of each set of WCs. The x-axis represents the label of each epoch and each independent component. The WCs are shown in (a) D1, (b) D2, (c) D3, (d) D4 and (e) A4.

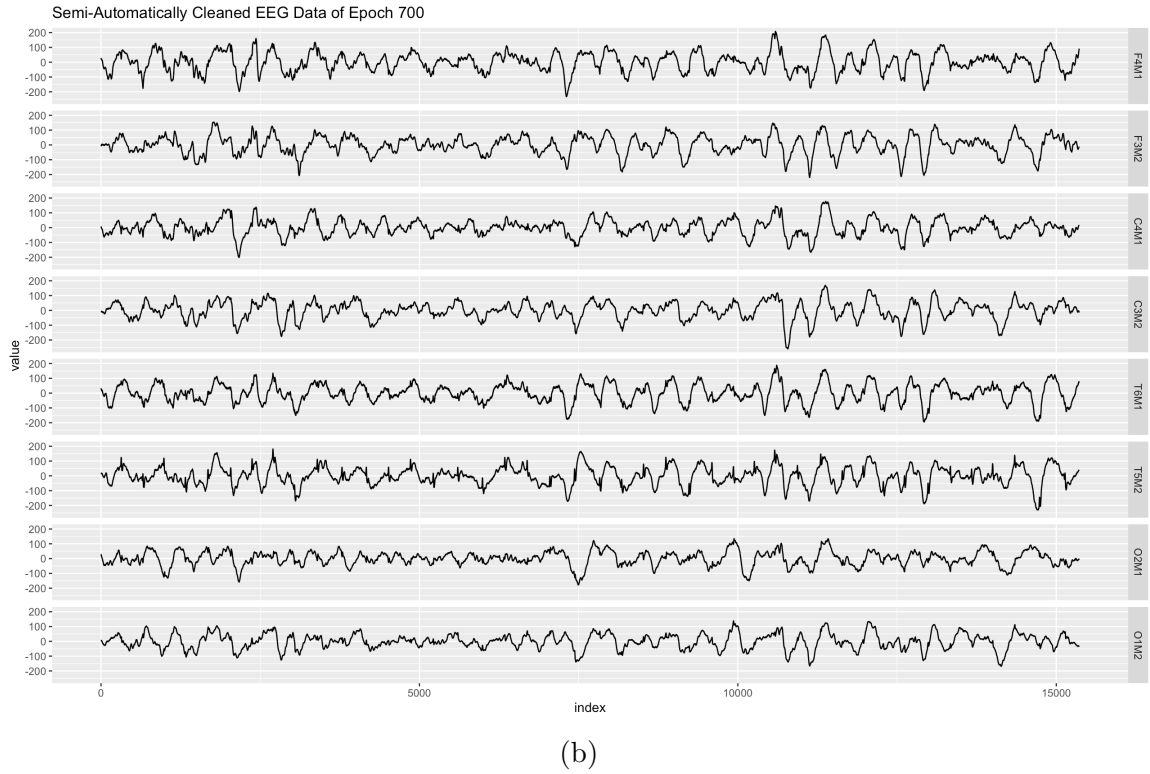
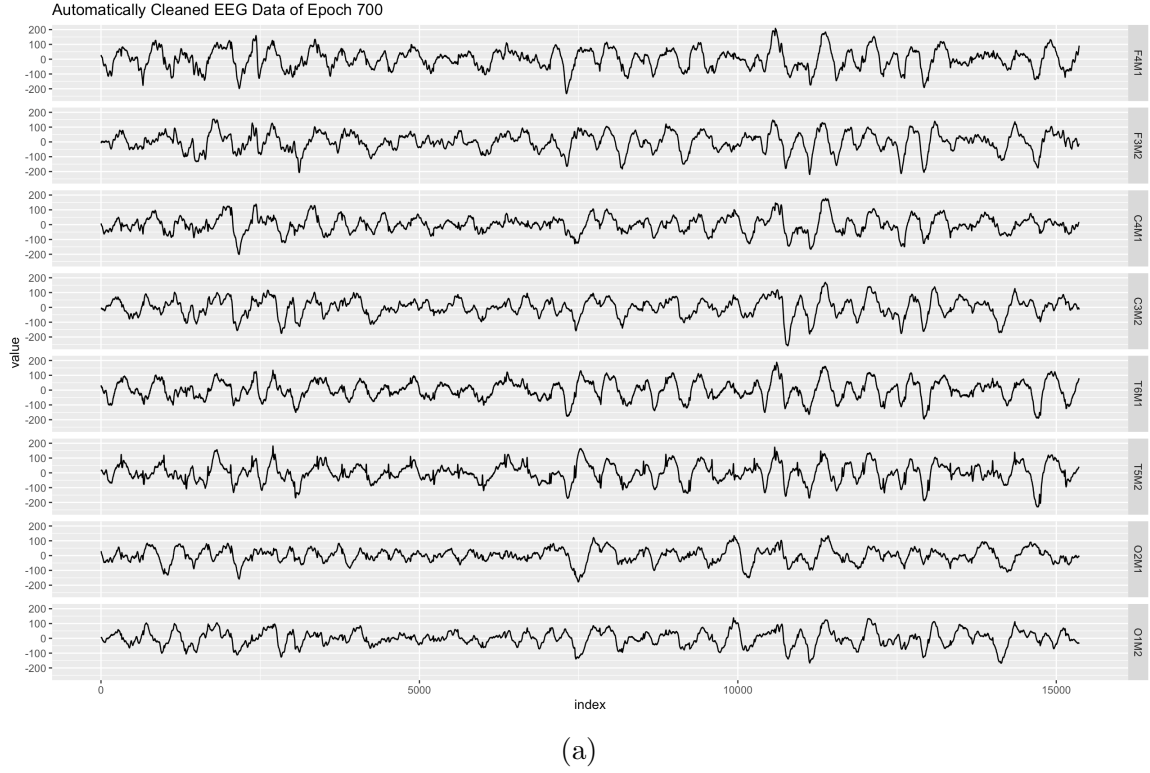


Figure 5.6: Cleaned EEG signals across eight channels. Figure (a) displays the EEG signal cleaned via the automatic method, and Figure (b) shows the EEG signals cleaned using the semi-automatic method.

Chapter 6

Discussion

6.1 Results

This study evaluates three EEG data cleaning approaches: baseline pre-processing, which acts as the control, and two proposed methods, automatic and semi-automatic cleaning, to determine their effectiveness in improving data quality.

The investigation reveals an increase in sleep stage classification accuracy from the baseline pre-processing, with an accuracy rate of 75.66%, to the automatic cleaning method, which achieves an accuracy of 80.05%. This improvement suggests that even with the automatic method's relatively straightforward approach of classifying wavelet coefficients into two categories and marking a small portion of EEG data as artifacts, it still manages to enhance data quality over just pre-processing.

In addition, the semi-automated cleaning method demonstrates better performance compared to both the pre-processed and automatically cleaned EEG data, achieving an average accuracy of 81.48%.

Moreover, for channels that already exhibited high performance (C3M2, O2M1 and O1M2), the semi-automated and automated cleaning methods did not negatively impact their performance, maintaining their high accuracy levels.

There is also sufficient evidence to show that both automatically cleaned and semi-automatically cleaned C4M1 channel have shown an improved performance compared with pre-processed only C4M1 data with both p values less than 0.001. When these

methods are compared against manually cleaned EEG data, which recorded an accuracy of 91.46%, both the pre-processed only and automatically cleaned EEG data demonstrate inferior performance. However, the semi-automatic cleaning method not only outperforms these two methods but also surpasses manual cleaning with an accuracy of 93.91%. This indicates that the semi-automatic approach, while leveraging the benefits of automation, also incorporates critical elements of manual oversight, resulting in enhanced accuracy and reliability.

Overall, the findings suggest that while automated processes offer some improvement over basic pre-processing, the integration of semi-automated techniques provides a more effective balance between automated efficiency and the nuanced accuracy of manual intervention. This approach appears particularly beneficial in channels where EEG data quality can be significantly improved, thus making it a promising method for EEG data cleaning and analysis.

6.2 Future Work

Based on the insights gained from the current analysis, future work on EEG artifact cleaning methods can be significantly enhanced by addressing certain key areas. The present automated cleaning approach, which relies on dividing independent wavelet coefficients into two clusters without adequately considering the size disparities between these clusters, has shown limitations. This method tends to inadequately capture all potential artifactual components, leading to suboptimal cleaning.

To refine this approach, future efforts should involve collaboration with EEG sleep data experts. This collaboration can provide valuable insights into the typical proportion of artifact components present in raw EEG data. Armed with this knowledge, the methodology for dividing signal and artifact clusters can be adjusted to achieve a more effective separation. This refined division method is expected to enhance the accuracy of artifact detection and removal.

Additionally, an interesting insight observed in the current study is that channels

exhibiting better performance are located on the same side of the head. This observation suggests that spatial features of EEG data are influential and should be incorporated into the clustering process. By considering these spatial characteristics, the artifact cleaning method could become more precise and tailored to the unique topographical patterns of EEG signals.

Furthermore, while the current method has been tested on EEG data related to sleep disorders, its applicability could be broadened to other types of EEG data. By extending the testing to various EEG datasets, the robustness and versatility of the artifact cleaning method can be thoroughly evaluated. This expansion could lead to a more universally applicable EEG cleaning approach, beneficial across different EEG applications and conditions.

Bibliography

- [1] M. K. Pavlova and V. Latreille, “Sleep disorders,” *The American journal of medicine*, vol. 132, no. 3, pp. 292–299, 2019.
- [2] A. T. Van de Water, A. Holmes, and D. A. Hurley, “Objective measurements of sleep for non-laboratory settings as alternatives to polysomnography—a systematic review,” *Journal of sleep research*, vol. 20, no. 1pt2, pp. 183–200, 2011.
- [3] M. J. Aminoff, F. Boller, G. Bruyn, H. L. Klawans, D. F. Swaab, and P. Vinken, *Handbook of clinical neurology*. North-Holland Publishing Company, 1968.
- [4] M. X. Cohen, “Where does eeg come from and what does it mean?” *Trends in neurosciences*, vol. 40, no. 4, pp. 208–218, 2017.
- [5] M. Teplan *et al.*, “Fundamentals of eeg measurement,” *Measurement science review*, vol. 2, no. 2, pp. 1–11, 2002.
- [6] P. Kellaway, “An orderly approach to visual analysis: Characteristics of the normal eeg of adults and children,” *Clinical practice of clinical electroencephalography*, 1997.
- [7] P. L. Nunez and R. Srinivasan, *Electric fields of the brain: the neurophysics of EEG*. Oxford University Press, USA, 2006.
- [8] L. J. Hirsch and R. P. Brenner, *Atlas of EEG in critical care*. John Wiley & Sons, 2011.
- [9] T. Radüntz, J. Scouten, O. Hochmuth, and B. Meffert, “Automated eeg artifact elimination by applying machine learning algorithms to ica-based features,” *Journal of neural engineering*, vol. 14, no. 4, p. 046 004, 2017.
- [10] M. Sazgar and M. G. Young, *Absolute epilepsy and EEG rotation review: Essentials for trainees*. Springer, 2019.
- [11] X. Jiang, G.-B. Bian, and Z. Tian, “Removal of artifacts from eeg signals: A review,” *Sensors*, vol. 19, no. 5, p. 987, 2019.
- [12] S Romero, M. Mananas, S Clos, S Gimenez, and M. Barbanoj, “Reduction of eeg artifacts by ica in different sleep stages,” in *Proceedings of the 25th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (IEEE Cat. No. 03CH37439)*, IEEE, vol. 3, 2003, pp. 2675–2678.

- [13] G. L. Wallstrom, R. E. Kass, A. Miller, J. F. Cohn, and N. A. Fox, "Automatic correction of ocular artifacts in the eeg: A comparison of regression-based and component-based methods," *International journal of psychophysiology*, vol. 53, no. 2, pp. 105–119, 2004.
- [14] A. Schlögl, C. Keinrath, D. Zimmermann, R. Scherer, R. Leeb, and G. Pfurtscheller, "A fully automated correction method of eeg artifacts in eeg recordings," *Clinical neurophysiology*, vol. 118, no. 1, pp. 98–104, 2007.
- [15] A. Q. Hamal and A. W. bin Abdul Rehman, "Artifact processing of epileptic eeg signals: An overview of different types of artifacts," in *2013 International Conference on Advanced Computer Science Applications and Technologies*, IEEE, 2013, pp. 358–361.
- [16] S. Devuyst, T. Dutoit, P. Stenuit, M. Kerkhofs, and E. Stanus, "Removal of eeg artifacts from eeg using a modified independent component analysis approach," in *2008 30th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, IEEE, 2008, pp. 5204–5207.
- [17] X. Chen, A. Liu, J. Chiang, Z. J. Wang, M. J. McKeown, and R. K. Ward, "Removing muscle artifacts from eeg data: Multichannel or single-channel techniques?" *IEEE Sensors Journal*, vol. 16, no. 7, pp. 1986–1997, 2015.
- [18] S. Barua and S. Begum, "A review on machine learning algorithms in handling eeg artifacts," in *The Swedish AI Society (SAIS) Workshop SAIS, 14, 22-23 May 2014, Stockholm, Sweden*, 2014.
- [19] W. Mumtaz, S. Rasheed, and A. Irfan, "Review of challenges associated with the eeg artifact removal methods," *Biomedical Signal Processing and Control*, vol. 68, p. 102741, 2021.
- [20] A. L. D’Rozario *et al.*, "An automated algorithm to identify and reject artefacts for quantitative eeg analysis during sleep in patients with sleep-disordered breathing," *Sleep and Breathing*, vol. 19, pp. 607–615, 2015.
- [21] M. K. Islam, A. Rastegarnia, and Z. Yang, "Methods for artifact detection and removal from scalp eeg: A review," *Neurophysiologie Clinique/Clinical Neurophysiology*, vol. 46, no. 4-5, pp. 287–305, 2016.
- [22] C Rashmi and C Shantala, "Eeg artifacts detection and removal techniques for brain computer interface applications: A systematic review," *Int. J. Adv. Technol. Eng. Explor*, vol. 9, p. 354, 2022.
- [23] R. N. Vigário, "Extraction of ocular artefacts from eeg using independent component analysis," *Electroencephalography and clinical neurophysiology*, vol. 103, no. 3, pp. 395–404, 1997.
- [24] T.-P. Jung, S. Makeig, M. Westerfield, J. Townsend, E. Courchesne, and T. J. Sejnowski, "Removal of eye activity artifacts from visual event-related potentials in normal and clinical subjects," *Clinical neurophysiology*, vol. 111, no. 10, pp. 1745–1758, 2000.

- [25] C. A. Joyce, I. F. Gorodnitsky, and M. Kutas, "Automatic removal of eye movement and blink artifacts from eeg data using blind component separation," *Psychophysiology*, vol. 41, no. 2, pp. 313–325, 2004.
- [26] A. Flexer, H. Bauer, J. Pripfl, and G. Dorffner, "Using ica for removal of ocular artifacts in eeg recorded from blind subjects," *Neural Networks*, vol. 18, no. 7, pp. 998–1005, 2005.
- [27] Y. Li, Z. Ma, W. Lu, and Y. Li, "Automatic removal of the eye blink artifact from eeg using an ica-based template matching approach," *Physiological measurement*, vol. 27, no. 4, p. 425, 2006.
- [28] S Romero, M. Mañanas, and M. J. Barbanoj, "Ocular reduction in eeg signals based on adaptive filtering, regression and blind source separation," *Annals of biomedical engineering*, vol. 37, pp. 176–191, 2009.
- [29] W. Zhou and J. Gotman, "Automatic removal of eye movement artifacts from the eeg using ica and the dipole model," *Progress in Natural Science*, vol. 19, no. 9, pp. 1165–1170, 2009.
- [30] C. Q. Lai, H. Ibrahim, M. Z. Abdullah, J. M. Abdullah, S. A. Suandi, and A. Azman, "Artifacts and noise removal for electroencephalogram (eeg): A literature review," in *2018 IEEE Symposium on Computer Applications & Industrial Electronics (ISCAIE)*, IEEE, 2018, pp. 326–332.
- [31] T. Zikov, S. Bibian, G. A. Dumont, M. Huzmezan, and C. R. Ries, "A wavelet based de-noising technique for ocular artifact correction of the electroencephalogram," in *Proceedings of the Second Joint 24th Annual Conference and the Annual Fall Meeting of the Biomedical Engineering Society*[[*Engineering in Medicine and Biology*], IEEE, vol. 1, 2002, pp. 98–105.
- [32] S. V. Ramanan, N. Kalpakam, and J. Sahambi, "A novel wavelet based technique for detection and de-noising of ocular artifact in normal and epileptic electroencephalogram," 2004.
- [33] P. S. Kumar, R. Arumuganathan, K. Sivakumar, and C Vimal, "Removal of ocular artifacts in the eeg through wavelet transform without using an eeg reference channel," *Int. J. Open Problems Compt. Math*, vol. 1, no. 3, pp. 188–200, 2008.
- [34] S. Khatun, R. Mahajan, and B. I. Morshed, "Comparative study of wavelet-based unsupervised ocular artifact removal techniques for single-channel eeg data," *IEEE journal of translational engineering in health and medicine*, vol. 4, pp. 1–8, 2016.
- [35] K Asaduzzaman, M. Reaz, F Mohd-Yasin, K. Sim, and M. Hussain, "A study on discrete wavelet-based noise removal from eeg signals," in *advances in computational biology*, Springer, 2010, pp. 593–599.
- [36] N. P. Castellanos and V. A. Makarov, "Recovering eeg brain signals: Artifact suppression with wavelet enhanced independent component analysis," *Journal of neuroscience methods*, vol. 158, no. 2, pp. 300–312, 2006.

- [37] H. Ghandeharion and A. Erfanian, “A fully automatic ocular artifact suppression from eeg data using higher order statistics: Improved performance by wavelet analysis,” *Medical engineering & physics*, vol. 32, no. 7, pp. 720–729, 2010.
- [38] R. R. Vázquez, H. Velez-Perez, R. Ranta, V. L. Dorr, D. Maquin, and L. Mailard, “Blind source separation, wavelet denoising and discriminant analysis for eeg artefacts and noise cancelling,” *Biomedical signal processing and control*, vol. 7, no. 4, pp. 389–400, 2012.
- [39] N. K. Al-Qazzaz, S. Hamid Bin Mohd Ali, S. A. Ahmad, M. S. Islam, and J. Escudero, “Automatic artifact removal in eeg of normal and demented individuals using ica-wt during working memory tasks,” *Sensors*, vol. 17, no. 6, p. 1326, 2017.
- [40] A. K. Maddirala and K. C. Veluvolu, “Ica with cwt and k-means for eye-blink artifact removal from fewer channel eeg,” *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 30, pp. 1361–1373, 2022.
- [41] V. Schetinin and J. Schult, “The combined technique for detection of artifacts in clinical electroencephalograms of sleeping newborns,” *IEEE Transactions on Information Technology in Biomedicine*, vol. 8, no. 1, pp. 28–35, 2004.
- [42] M. Betta, A. Gemignani, A. Landi, M. Laurino, P. Piaggi, and D. Menicucci, “Detection and removal of ocular artifacts from eeg signals for an automated rem sleep analysis,” in *2013 35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, IEEE, 2013, pp. 5079–5082.
- [43] M. Dursun *et al.*, “A new approach to eliminating eeg artifacts from the sleep eeg signals for the automatic sleep stage classification,” *Neural Computing and Applications*, vol. 28, pp. 3095–3112, 2017.
- [44] R. Ranjan, B. C. Sahana, and A. K. Bhandari, “Cardiac artifact noise removal from sleep eeg signals using hybrid denoising model,” *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–10, 2022.
- [45] E. Saifutdinova, M. Congedo, D. Dudysova, L. Lhotska, J. Koprivova, and V. Gerla, “An unsupervised multichannel artifact detection method for sleep eeg based on riemannian geometry,” *Sensors*, vol. 19, no. 3, p. 602, 2019.
- [46] I. W. Selesnick and C. S. Burrus, “Generalized digital butterworth filter design,” *IEEE Transactions on signal processing*, vol. 46, no. 6, pp. 1688–1694, 1998.
- [47] N. E. Helwig, *Eegkit: Toolkit for electroencephalography data*, R package version 1.0-4, 2018. [Online]. Available: <https://CRAN.R-project.org/package=eegkit>.
- [48] J. L. Marchini, C. Heaton, and B. D. Ripley, *Fastica: Fastica algorithms to perform ica and projection pursuit*, R package version 1.2-3, 2021. [Online]. Available: <https://CRAN.R-project.org/package=fastICA>.

- [49] B. Whitcher, *Waveslim: Basic wavelet routines for one-, two-, and three-dimensional signal processing*, R package version 1.8.4, 2022. [Online]. Available: <https://CRAN.R-project.org/package=waveslim>.
- [50] R Core Team, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria, 2022. [Online]. Available: <https://www.R-project.org/>.
- [51] H. Adeli, Z. Zhou, and N. Dadmehr, “Analysis of eeg records in an epileptic patient using wavelet transform,” *Journal of neuroscience methods*, vol. 123, no. 1, pp. 69–87, 2003.
- [52] N. Bajaj, “Wavelets for eeg analysis,” *Wavelet Theory*, pp. 1–16, 2020.
- [53] O. Rioul and P. Duhamel, “Fast algorithms for discrete and continuous wavelet transforms,” *IEEE transactions on information theory*, vol. 38, no. 2, pp. 569–586, 1992.
- [54] S. G. Mallat, “A theory for multiresolution signal decomposition: The wavelet representation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 11, no. 7, pp. 674–693, 1989.
- [55] H. Ocak, “Automatic detection of epileptic seizures in eeg using discrete wavelet transform and approximate entropy,” *Expert Systems with Applications*, vol. 36, no. 2, pp. 2027–2036, 2009.
- [56] I. Daubechies, “Orthonormal bases of compactly supported wavelets,” *Communications on pure and applied mathematics*, vol. 41, no. 7, pp. 909–996, 1988.
- [57] T. Gandhi, B. K. Panigrahi, and S. Anand, “A comparative study of wavelet families for eeg signal classification,” *Neurocomputing*, vol. 74, no. 17, pp. 3051–3057, 2011.
- [58] T. Staudinger and R. Polikar, “Analysis of complexity based eeg features for the diagnosis of alzheimer’s disease,” in *2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, IEEE, 2011, pp. 2033–2036.
- [59] S.-H. Oh, Y.-R. Lee, and H.-N. Kim, “A novel eeg feature extraction method using hjorth parameter,” *International Journal of Electronics and Electrical Engineering*, vol. 2, no. 2, pp. 106–110, 2014.
- [60] R. Jenke, A. Peer, and M. Buss, “Feature extraction and selection for emotion recognition from eeg,” *IEEE Transactions on Affective computing*, vol. 5, no. 3, pp. 327–339, 2014.
- [61] A. Mert and A. Akan, “Hilbert-huang transform based hierarchical clustering for eeg denoising,” in *21st European Signal Processing Conference (EUSIPCO 2013)*, IEEE, 2013, pp. 1–5.
- [62] N. Mammone *et al.*, “Hierarchical clustering of the electroencephalogram spectral coherence to study the changes in brain connectivity in alzheimer’s disease,” in *2016 IEEE Congress on Evolutionary Computation (CEC)*, IEEE, 2016, pp. 1241–1248.

- [63] A. B. Kashlak, P. Loliencar, and G. Heo, “Topological hidden markov models,” *Journal of Machine Learning Research*, vol. 24, no. 340, pp. 1–49, 2023.

Appendix A: Method Comparison

We have additionally transformed the outputs into an alternative format to facilitate comparison across various methodologies. Specifically, we arranged the raw, preprocessed, automatically cleaned, and semi-automatically cleaned EEG signals for eight channels of the randomly selected epoch (epoch 700), as shown in the following figures.

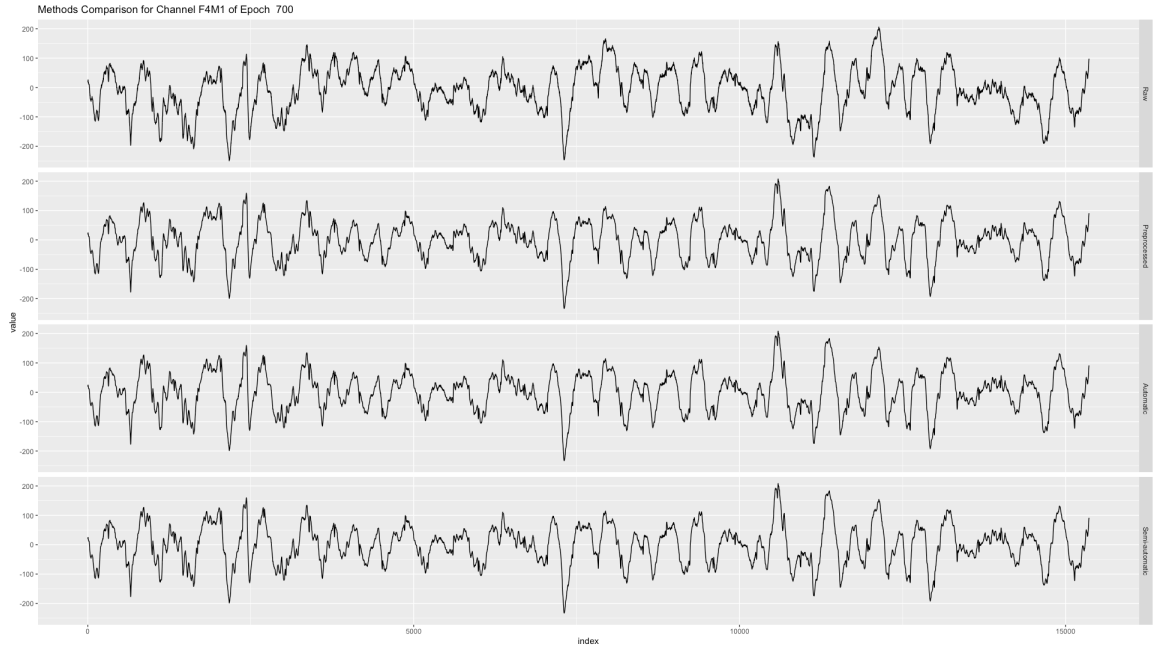


Figure A.1: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel F4M1.

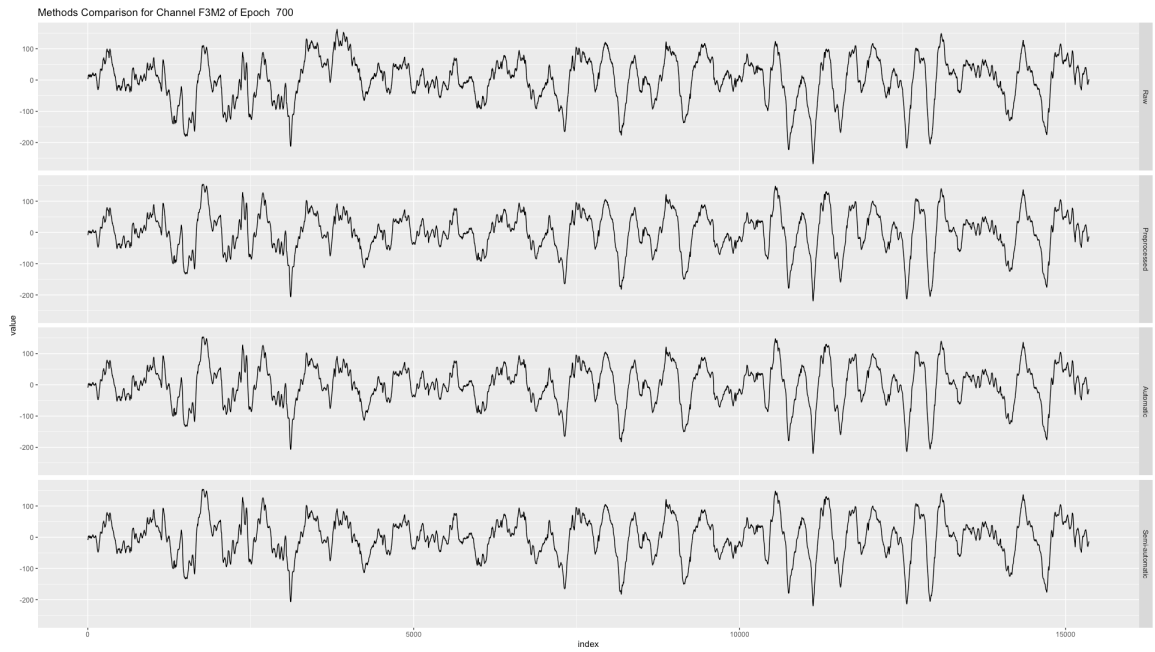


Figure A.2: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel F3M2.

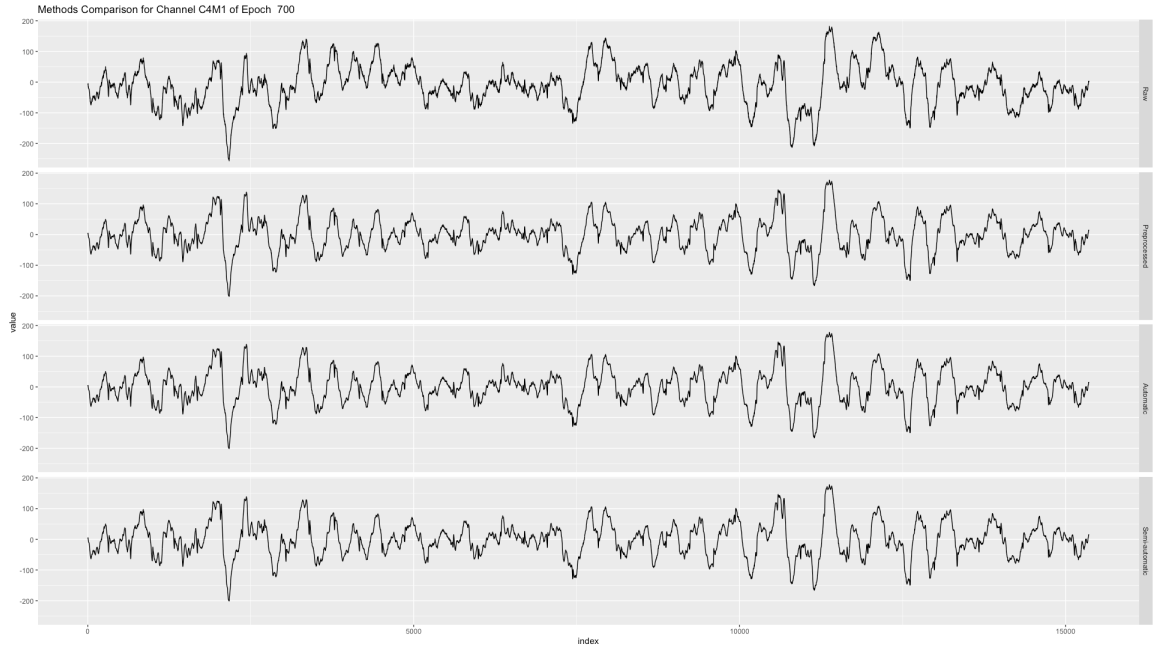


Figure A.3: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel C4M1.

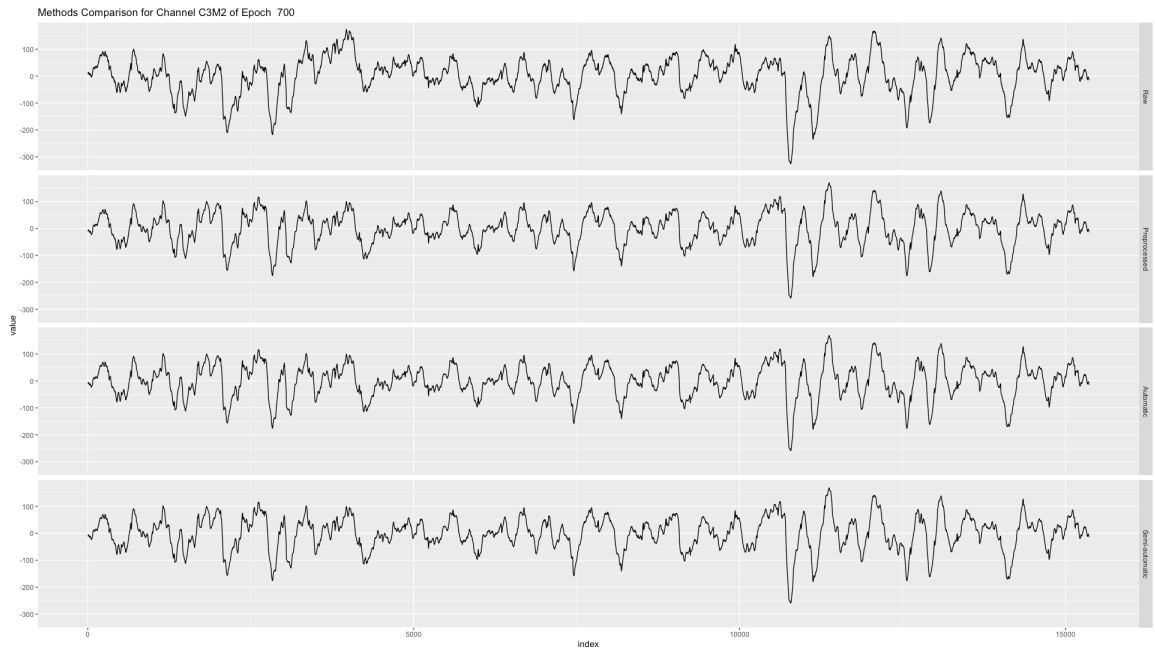


Figure A.4: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel C3M2.

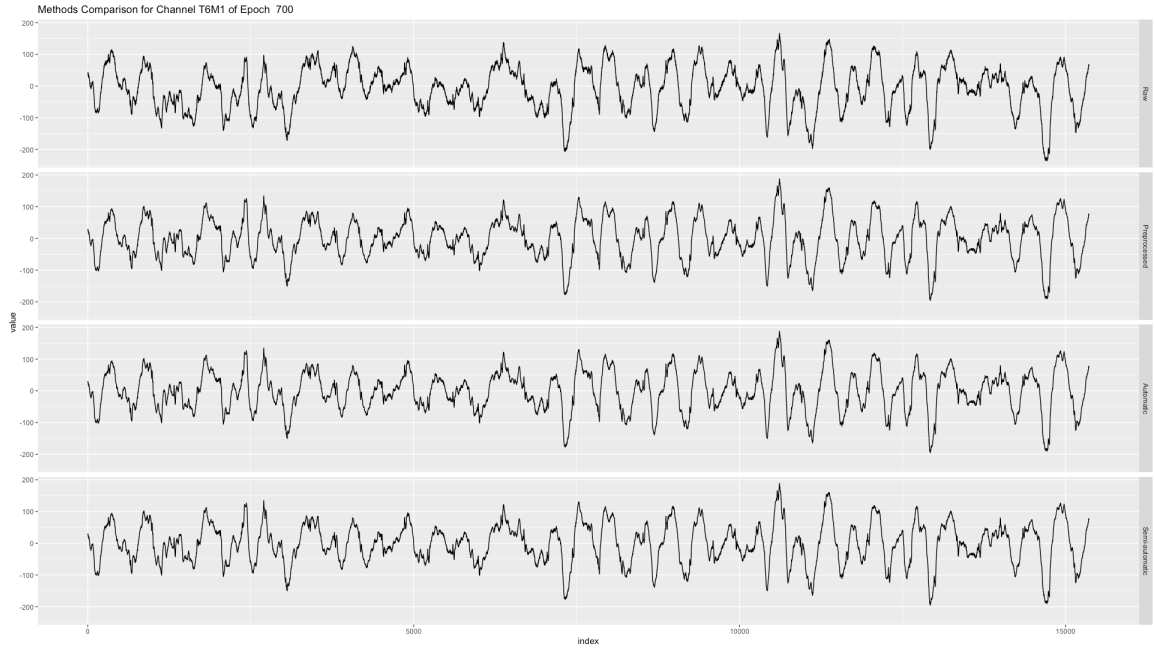


Figure A.5: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel T6M1.

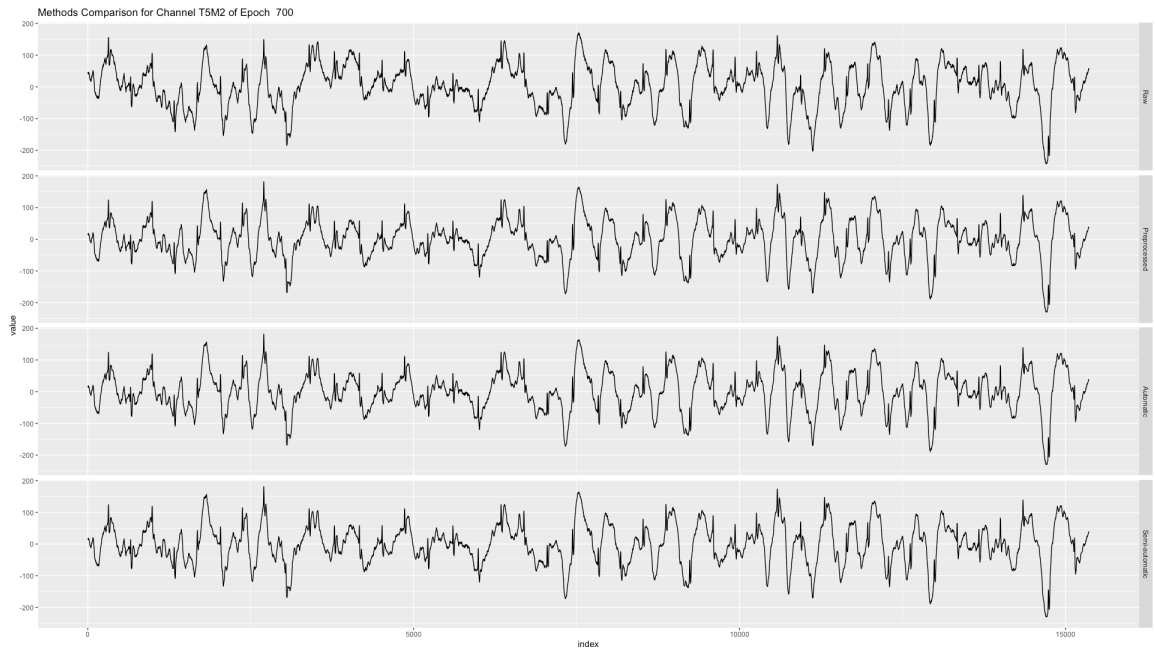


Figure A.6: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel T5M2.

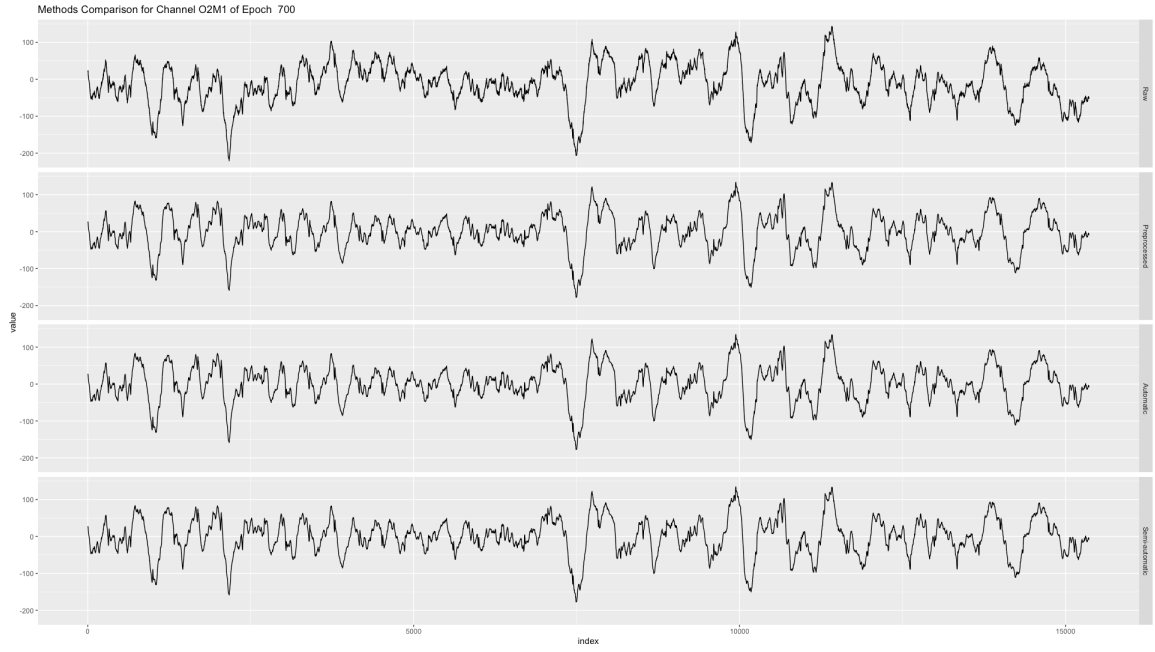


Figure A.7: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel O2M1.

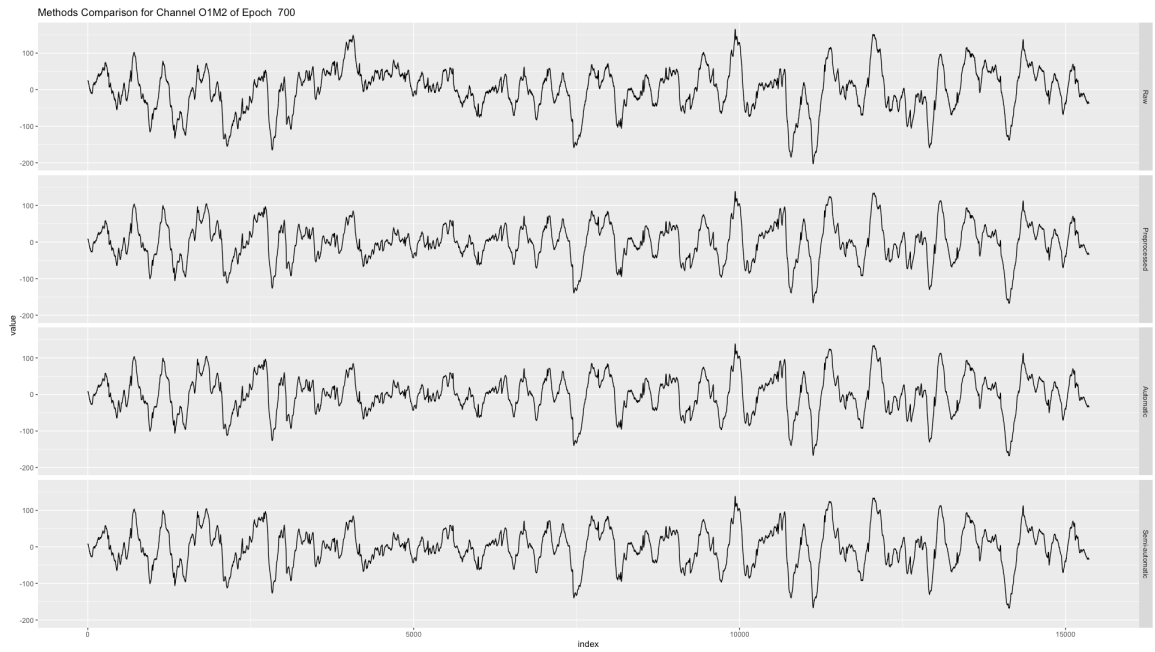


Figure A.8: Methods comparison among raw, preprocessed, automatically cleaned and semi-automatically cleaned EEG signal of epoch 700 Channel O1M2.