

UMAP for Dimensionality Reduction in Sleep Stage Classification Using EEG Data

Yangfan Deng*, Hamad Albidah*, Haoliang Cheng[†], Ahmed Dallal*, Jijun Yin*, and Zhi-Hong Mao*[‡]

*Department of Electrical and Computer Engineering, University of Pittsburgh

[†]Department of Computer Science, University of Pittsburgh

[‡]Department of Bioengineering, University of Pittsburgh

Pittsburgh, PA, USA

Emails: yad58@pitt.edu, hna4@pitt.edu, hac177@pitt.edu, ahd12@pitt.edu, davidyin96@gmail.com, zhm4@pitt.edu

Abstract—Sleep is vital for wellness, and electroencephalography (EEG) serves as an instrumental tool in the study of sleep. Sleep is classified into four stages: stages N1-N3, and rapid eye movement (REM). To acquire effective and robust EEG features for sleep detection and analysis, we explore the dimensionality reduction effects of Uniform Manifold Approximation and Projection (UMAP) on various features of the EEG signals. Compared with traditional band power analysis, UMAP demonstrated higher accuracy for sleep stage classification and better reliability. Using UMAP, we observed an average of 11% increase in accuracy and an average of 20% increase in Macro-F1 Score on the same dataset. Particularly, in the wakefulness stage, Macro-F1 Score increased by 23%. Moreover, the 2D visual analysis revealed the outstanding ability of UMAP to cluster EEG signals after significant dimensionality reduction of the data.

Index Terms—EEG signal, UMAP, dimensionality reduction, sleep stage classification

I. INTRODUCTION

Sleep has attracted tremendous interests in research as it is so important for human health. Sleep is a state of rest for the body and mind, during which consciousness is partially or completely lost and there is minimal physical activity. Important activities such as the body's recovery and regeneration, as well as the brain's conversion of short-term memories into long-term memories, all occur during sleep. Insufficient sleep can increase the risks of cardiovascular diseases, metabolic and endocrine disorders, obesity, etc. A full night's sleep typically includes 4-6 sleep cycles, each consisting of four stages, including rapid eye movement (REM) and stages N1-N3 [1]. Different sleep stages can be differentiated by the frequency contents of brain waves, which signify the depth and quality of sleep. Among the various methods for examining sleep patterns, the use of electroencephalograms (EEGs) is proficient in revealing the intricate activities of the brain. This paper focuses on the sleep stage classification and detection using EEG signals.

EEG signals are collected by the electrodes located on different parts of the scalp, which can record the electrical activities with millisecond-level temporal resolution. Therefore, the resulting EEG signals are a form of high-dimensional, sequential data. These data can be further processed by a statistical or machine learning model to detect the brain status. Using these high-dimensional EEG signals as input data di-

rectly can lead to poor interpretability and high complexity of the model, increased computational burden, and risk of model overfitting. Thus, transforming high-dimensional EEG data into a lower dimensional feature space, while still retaining key information within the data, is crucial for developing efficient machine learning algorithms for estimation of brain activities, specifically, classification of sleep stages in this study.

The existing dimension reduction algorithms applied to EEG signals are categorized into linear and nonlinear types. For linear algorithms, Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) have been proven to exhibit excellent performance in dimensionality reduction for EEG signals [2]. These algorithms are built upon the assumption that the data can be adequately represented by a linear structure, which may inadvertently neglect nonlinear relationships. This oversight has the potential to cause significant information loss in the compressed data. In the realm of nonlinear algorithms, t-Distributed Stochastic Neighbor Embedding (t-SNE) [3] and Isometric Mapping (ISOMAP) [4] are powerful tools, which can handle the complexities of EEG data that may not be well-represented in linear algorithms and can uncover the underlying structures in the data that traditional methods might miss. However, these nonlinear algorithms may not effectively balance the local structures and global structures in high-dimensional data and suffer from slow processing.

To address the challenges of the algorithms mentioned above, we adopt Uniform Manifold Approximation and Projection (UMAP) [5] method for dimensionality reduction of EEG signals. UMAP is commonly employed for visualization and preprocessing in machine learning algorithms. It has exhibited the state-of-the-art effectiveness and accuracy in reducing the dimensionality of data. Additionally, this algorithm maintains a relatively fast computation speed while effectively capturing both local and global structures. The objective of this paper is to apply UMAP to identify robust, reduced-dimensional features of EEG signals for sleep stage classification. We explore the geometric properties of EEG data in the low dimensional feature space obtained by UMAP, aiming to understand the effectiveness of EEG data grouping within the feature space. This approach not only facilitates the development of classifiers based on these condensed features

but also enables us to integrate the clustering patterns of unlabeled EEG data in the classifier training. In this paper, we utilize various geometric properties from EEG signals and combine them linearly as input data for UMAP to perform dimensionality reduction. By comparing various metrics and visualizing the distribution in a 2D feature space, we demonstrate the superiority of UMAP in EEG feature extraction.

The subsequent sections of the paper are organized as follows: Section II introduces the dimensionality reduction methods we have employed, Section III presents various metrics and visual representations for assessing the effectiveness of the dimensionality reduction, and Section IV concludes the paper.

II. METHODS

A. Filtering Procedure

We utilize data from PhysioNet [6], which comprises EEG signals sourced from Caucasian males and females aged 21-35 years not on any medication. Each collected raw EEG signal was sampled at the frequency of 100 Hz in 30 sec intervals. Through the use of a band-pass filter, we obtain signals at four different frequency ranges: delta wave (1-4 Hz), theta wave (5-8 Hz), alpha wave (9-12 Hz), and beta wave (13-25 Hz). These four distinct frequency bands of EEG signals show various power intensities during different sleep stages: deep sleep or slow-wave sleep (N3), shallow sleep (N1 and N2), rapid eye movement stage (REM), and wakefulness (W). We calculate several types of geometric properties for these waveforms: frequency band powers, harmonic parameters, and relative percent spectral energy bands (RPEB). By combining these properties from high-dimensional EEG data and using UMAP for dimensionality reduction, we can obtain their characteristics in a lower-dimensional feature space.

B. UMAP

UMAP is a nonlinear algorithm used for reducing dimensionality of high dimensional data and visualizing the low dimensional representation of the data [5]. This algorithm is based on the concept of manifold learning and aims to map high-dimensional data to a lower-dimensional feature space while preserving both local and global structures of the data. UMAP includes four steps: (i) nearest neighbor graph construction, (ii) weight calculation, (iii) layout optimization, and (iv) iterative optimization. Following the procedure in [5], we implement UMAP as described in Algorithm 1 below. The construction of the graph employs a local metric to calculate the distances between neighboring data points, thereby preserving the local structure. When reducing data from high-dimensional graph representations to lower dimension, the concept of a fuzzy simplicial complex in topology ensures that UMAP adequately retains global features. During the dimensionality reduction process, UMAP employs stochastic gradient descent to align the low-dimensional representation closely with the high-dimensional data.

Various geometric properties are calculated corresponding to different frequency contents of the raw EEG signals and linearly stacked together as input data for UMAP. Subsequently, the projections processed by UMAP in the 2D feature space are utilized for visualization and further classification of sleep stages. Here we use a “vanilla” classifier, K-Nearest Neighbors (KNN) classifier, for verifying the effectiveness of the proposed UMAP procedure compared with the other feature extraction methods. We identify the optimal hyperparameters for UMAP and KNN that achieve the best classification performance through an iterative process. This iterative process is repeated for the data of each participant across the entire dataset.

Algorithm 1: UMAP dimension reduction

Input: Data $X = \{x_1, x_2, \dots, x_n\}$ and UMAP hyperparameters $n_neighbors, n_components, min_dist, n_epochs$ 1: Construct the high-dimensional graph G from X using $n_neighbors$. 2: Perform spectral embedding on G to obtain initial low-dimensional representation Y . 3: for epoch in 1 to n_epochs do 4: for each y_i in Y do 5: Define Σ_attr as the sum of attractive forces on y_i from its neighbors. 6: Define Σ_rep as the sum of repulsive forces on y_i from all other points. 7: Update the position of y_i based on Σ_attr, Σ_rep , and min_dist . 8: end for 9: end for 10: Return the optimized low-dimensional representation Y .

III. RESULTS AND DISCUSSION

A. Classified Performance

Inspired by DeepSleepNet [7], we chose the following metrics for measuring the effectiveness of the proposed approach: Average Overall Accuracy (ACC), Macro-F1 Score (MF1), Cohen’s Kappa (κ), and the per-class F1-Score for each distinct sleep stage. We calculated the mean and variance of these metrics based on each tester’s EEG data in the dataset. We tested two types of input data for UMAP. The first type is the band powers of EEG data, and the second type is the band powers with harmonic parameters (the center frequency, band width, and power value at the center frequency). Furthermore, band power analysis was selected as a basis of comparison for UMAP. Results from Mao et al. [2] confirm that band power analysis outperforms traditional linear algorithms such as PCA in terms of classification accuracy and clustering effect. Through comparison among metrics from different algorithms, we observed that UMAP exhibits outstanding performance in the dimensionality reduction of EEG data.

Comparing the values of ‘band power*’ and ‘band power’ in Table I, we observed that using UMAP for dimensionality reduction not only enhanced the accuracy of sleep stage classification but also improved the consistency and robustness of the prediction model. Furthermore, these improvements were amplified upon incorporating harmonic parameters as input data. Table II shows a comprehensive enhancement in classification performance for each class when using UMAP to perform dimensionality reduction on band powers compared with traditional band power analysis. With the addition of harmonic parameters, the method utilizing UMAP managed to

TABLE I
COMPARISON AMONG DIFFERENT EEG GEOMETRIC PROPERTIES ACROSS
OVERALL ACCURACY (ACC), MACRO-F1 SCORE (MF1), AND COHEN'S
KAPPA (κ)

Properties	Band power	Band power*	Band power with harmonic parameters*
Accuracy	71.2% ($\pm 0.4\%$)	78.8% ($\pm 0.3\%$)	79.2% ($\pm 0.3\%$)
MF1	0.565 (± 0.009)	0.671 (± 0.006)	0.639 (± 0.008)
κ	0.554 (± 0.007)	0.639 (± 0.008)	0.691 (± 0.006)

* Using UMAP for dimensionality reduction.

TABLE II
COMPARISON AMONG DIFFERENT EEG GEOMETRIC PROPERTIES ACROSS
PER-CLASS F1-SCORE (F1)

Properties	Band power	Band power*	Band power with harmonic parameters*
W	0.698 (± 0.023)	0.838 (± 0.010)	0.868 (± 0.005)
N1	0.272 (± 0.035)	0.296 (± 0.047)	0.298 (± 0.042)
N2	0.741 (± 0.011)	0.807 (± 0.007)	0.821 (± 0.007)
N3	0.466 (± 0.106)	0.537 (± 0.130)	0.476 (± 0.134)
REM	0.586 (± 0.032)	0.686 (± 0.019)	0.665 (± 0.029)

* Using UMAP for dimensionality reduction.

further increase the classification accuracy for the W, N1, and N2 stages without significantly decreasing the classification accuracy of N3 and REM stages.

B. Visual Representation

Figs. 1-3 demonstrate the dimensionality reduction effect of UMAP. Each comprises a scatter plot depicting the distribution of data in a two-dimensional space and a Kernel Density Estimation (KDE) plot derived from the scatter plot. Fig. 1 displays the clustering effect when using only the band powers of delta and alpha waves as the coordinates. According to the results of Mao et al. [2], the combination of delta and alpha waves offers the best clustering result among all possible combinations of frequency bands, and this result is comparable to the best results obtained from the methods of PCA and autoencoder in [2]. Fig. 2 illustrates the effect of dimensionality reduction on a 2D feature space using UMAP, where the band powers of all four frequency bands are combined. In Fig. 3, in addition to the band powers, the harmonic parameters of the four frequency bands are also integrated with band powers, employing UMAP for reduction to a two-dimensional space.

Fig. 1 reveals that although the band powers of delta and alpha waves can differentiate between various sleep stages to a certain extent, the distinction between different sleep stages is not pronounced. The proximity between stages is relatively small, easily leading to a potential misclassification while employing distance-based classifiers such as KNN. In contrast, Fig. 2 employs UMAP for dimensionality reduction, which effectively addressed this issue. There is a greater degree of separation between different sleep stages, with particularly superior clustering effects for stages N2 and N3. Nonetheless, for the other sleep stages, there are still some degrees of data point overlapping, indicating that they are not distinctly

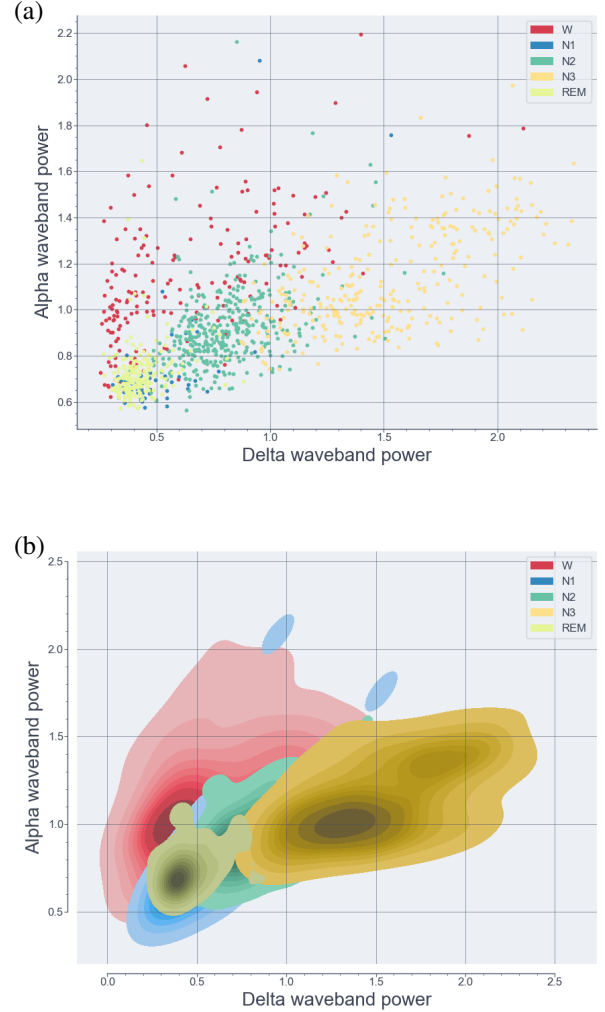


Fig. 1. Clustering results based on delta and alpha wavebands: (a) scatter plot and (b) KDE plot.

segregated, which can be observed in the KDE plot in Fig. 2(b). This issue was resolved when we used both harmonic parameters and band powers as input data for UMAP. In Fig. 3, it is apparent that aside from some coupling between stages N1 and REM, we achieved excellent clustering results for the other three sleep stages, including wakefulness (W).

IV. CONCLUSION

This study demonstrates that the method of UMAP for dimensionality reduction significantly enhances the accuracy and robustness of sleep stage classification using EEG data. This improvement is further amplified when incorporating harmonic parameters along with band powers as input features. Visual analyses reinforce that UMAP effectively clusters and separates different sleep stages, highlighting its efficiency and superiority over traditional band power analysis. Overall, UMAP emerges as a highly effective tool for EEG data analysis in the field of sleep research.

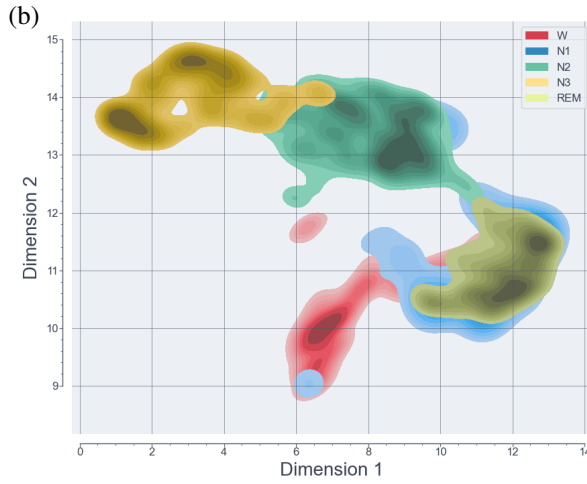
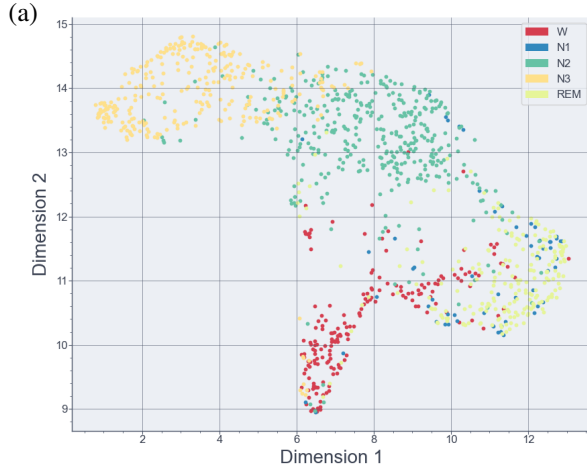


Fig. 2. Clustering results of UMAP using the band powers of the four wavebands as input: (a) scatter plot and (b) KDE plot.

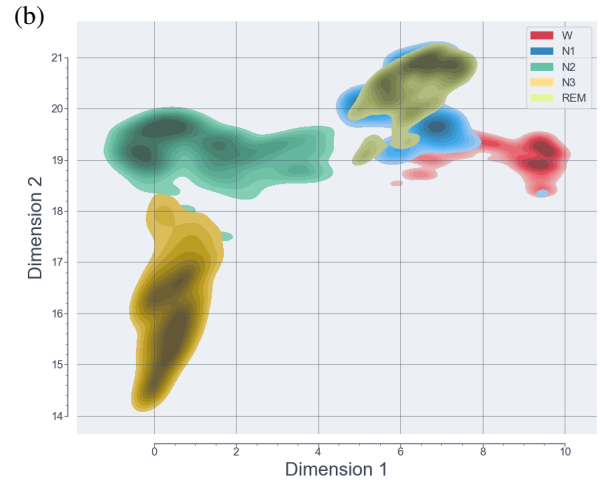
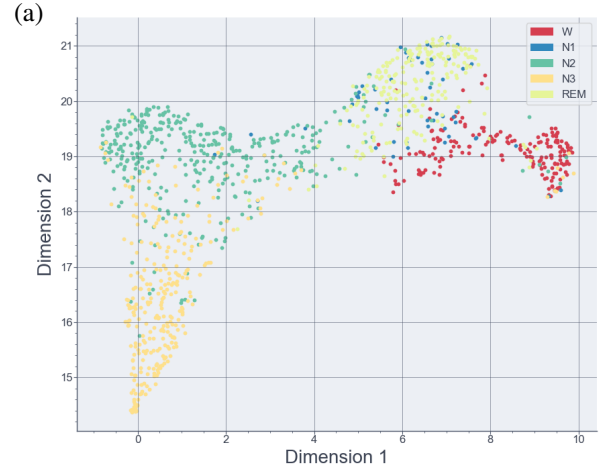


Fig. 3. Clustering results of UMAP using the band powers and harmonic parameters of the four wavebands as input: (a) scatter plot and (b) KDE plot.

REFERENCES

- [1] V. Brodbeck, A. Kuhn, F. von Wegner, A. Morzelewski, E. Tagliazucchi, S. Borisov, C. M. Michel, and H. Laufs, "EEG microstates of wakefulness and NREM sleep," *NeuroImage*, vol. 62, no. 3, pp. 2129–2139, 2012.
- [2] H. X. Mao, J. Widjaja, Y. Guo, J. Yin, and R. Vinjamuri, "Finding robust low dimensional features for sleep detection using EEG data," in *2022 IEEE 2nd International Conference on Data Science and Computer Application (ICDSCA)*. IEEE, 2022, pp. 42–45.
- [3] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, no. 11, pp. 2579–2625, 2008.
- [4] T. M. Rassias, "Properties of isometric mappings," *Journal of Mathematical Analysis and Applications*, vol. 235, no. 1, pp. 108–121, 1999.
- [5] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018.
- [6] PhysioNet. (2020) Sleep-edf database expanded (sleep-edfx). [Online]. Available: <https://physionet.org/content/sleep-edfx/1.0.0/>
- [7] A. Supratak, H. Dong, C. Wu, and Y. Guo, "DeepSleepNet: A model for automatic sleep stage scoring based on raw single-channel EEG," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 11, pp. 1998–2008, 2017.