# Focal and Generalized Seizures Distinction by Rebalancing Class Data and Random Forest Classification

**5 authors**, including:

Lina Abou-Abbas
McGill University
**14** PUBLICATIONS   **154** CITATIONS

SEE PROFILE

Imene Jemal
Institut National de la Recherche Scientifique
**5** PUBLICATIONS   **16** CITATIONS

SEE PROFILE

Khadidja Henni
**17** PUBLICATIONS   **74** CITATIONS

SEE PROFILE

Neila Mezghani
Télé-université
**128** PUBLICATIONS   **868** CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:

Pathways to better outcome View project

Improvement of knee kinematic analysis: Application to knee osteoarthritis and other musculoskeletal disorders View project

# Focal and Generalized Seizures Distinction by Rebalancing Class Data and Random Forest Classification

Lina Abou-Abbas[1,2*], Imene Jemal[2,3], Khadidja Henni[1,2], Amar Mitiche[2,3], and Neila Mezghani[1,2]

[1] Imaging and Orthopaedics Research Laboratory, The CHUM, Montreal, Canada
[2] Research Center LICEF, Teluq University, Montreal, Canada
[3] INRS- Centre Énergie, Matériaux et Télécommunications, Montréal, Canada

* E-mail: lina.abou-abbas.1@etsmtl.net (LAA)

***Abstract.*** Epileptic seizures are caused by abnormal electrical activity of brain cells, frequently accompanied by a short-lived loss of control or awareness. Epileptic seizures differ depending on their origin in the brain. They can be categorized as either focal or generalized in onset. The identification of seizure category is essential in brain surgery and in selecting medications that could help bring seizures under control. It is not always feasible to find out exactly if the seizure was generalized or focal without a thorough analysis of the continuous prolonged EEG waveforms. In this study, we propose an automatic classification method based on Hjorth parameters measured in electroencephalographic records (EEG). 1497 EEG signals from the Temple University Hospital Seizure Corpus (v.1.5.1) are used. Hjorth parameters (activity, complexity, and mobility) are extracted from these EEG records. To address class imbalance, data was rebalanced by Synthetic Minority Over Sampling (SMOTE). We also investigated the impact of changing the window length on the random forest classifier. For comparison, cost-sensitive learning has been applied by providing more weight to the minority class (generalized seizure) directly in the classifier. The performance of the proposed method was compared using accuracy, recall, and precision measures. Our method achieved a highest accuracy rate of 92.3% with a recall of 92.7% and precision of 91.8% using Hjorth parameters extracted from 10 seconds windows and rebalanced using SMOTE. A slight variation in performance measures occurred according to window size.

**Keywords:** Generalized and focal seizures, EEG, Hjorth parameters, SMOTE, Weighted random forest classification.

# 1    Introduction

Epileptic seizures are characterized by an intense sudden burst of electrical activity in the brain. They have various causes and treatments and can affect people of all ages. Epileptic seizures can be classified into two main categories: generalized and focal. Generalized seizures occur when an abnormal electrical activity involves concurrently both sides of the brain. Focal seizures, instead, are characterized by excessive electrical discharge in areas of a single brain side. Around seventy percent of seizures can be controlled by medications. Following  recent advances in machine learning based solutions for seizure detection [1, 2], the next challenge is the classification of seizures into focal or generalized. Recognition of epileptic seizure localization is crucial for drugs selection and surgery procedures. EEG signal interpretation remains the most effective and simple way for the localization of seizure origin. Because EEG visual inspection and interpretation is laborious, time consuming, and requires a trained expert, efficient automatic methods are necessary.

There have been several studies of automatic classification of seizures to characterize focal versus non-focal seizures: Sharma et al. [3] used a wavelet EEG representation to classify focal versus non-focal EEG signals, and reported an accuracy of 94.25%, whereas Bhattacharyya et al.[4] used rhythms extracted from empirical wavelet transforms to obtain an accuracy of 90%.  A novel method based on empirical mode decomposition and phase space reconstruction was proposed in [5] , results showed an accuracy of 96% in classifying focal EEG signals. In [6] Saputro et al. combined Mel Frequency Cepstral Coefficients, Hjorth components, and independent component analysis, reaching 91.4% recognition using a support vector machine. In [7] Roy et al. utilized a features extraction step based on calculating the eigenvalues by magnitude of the Fast Fourier Transform across all EEG channels and showed that a classification of seizure type is possible with an accuracy of 90.1% using the k-NN classifier. Das et al. in [8] discriminates between focal and non-focal signals, by using log-energy entropy derived from the combined empirical mode decomposition and discrete wavelet domain and reported a maximum accuracy of 89.4% with k-NN.

The automated ability to differentiate seizure types such as focal vs generalized remains a largely neglected topic due to both a lack of clinical datasets and annotations complexity. The TUH EEG corpus [9] has recently become the largest publicly available dataset to support epilepsy research, offering the opportunity to develop automatic prediction, detection and classification systems for epileptic seizures. To date, only a limited number of studies have used this challenging database for the task of seizure classification [6, 7]. While these previous works showed promising results in classifying seizure types by analyzing EEG signals in time and frequency domains, computing complexity remains a major issue.

The Hjorth descriptor is a set of nonlinear features providing spectral properties of the EEG signals in the time domain [10]. It consists of three parameters: activity, mobility, and complexity (see Table 1). The activity represents the mean power of the signal and mobility its mean frequency. Complexity is the estimate of the signal bandwidth [10]. Hjorth parameters can be computed fast, and implementation is straightforward because their calculation is based on the signal variance and its derivatives. The computational

cost is generally considered low compared to other methods. The Hjorth descriptor was shown by several studies to be useful to analyze nonstationary EEG signals and was successfully used in different applications, such as emotion recognition, mental-task discrimination, epilepsy prediction, and focal EEG signals classification [6, 11–13].

Over the last decade, Random forest classifier (RF) has received growing attention due to its robust performance across a wide range of medical applications such as early seizure detection, automated sleep stage identification and recognition of Alzheimer's disease [14, 15]. The RF classifier, an ensemble learning method, uses a bagging scheme where classification is determined by majority voting [16]. Since medical data is often subject to class imbalance, which means the different classification categories are not equally represented, classifiers generally tend to be biased in favor of the majority class when equal weights are assigned to classes. Therefore, two techniques were introduced to address data imbalance: (1) cost-sensitive learning (giving each class a mis-classification cost or weight according to its distribution in the whole training dataset) and, (2) data resampling (under-sampling or over-sampling). The Synthetic Minority Oversampling Technique (SMOTE) is based on generating synthetic minority examples by interpolation to oversample the minority class in the original training set [17]. In this study, we propose and investigate a method based on Hjorth parameters representation and random forest classification to distinguish focal from generalized epileptic seizures. We study the effect of varying the processing EEG window size. To address class imbalance, which is due to the uneven representation of generalized and focal seizure classes, we implemented the weighted sampler function [18] into the classifier and compared its results to a conventional random forest classifier preceded with a Synthetic Minority Oversampling Technique (SMOTE) [17]. We show that the latter method enhanced the overall performance of the classifier.

**Table 1.** Hjorth Parameters- y(t) is the signal and y'(t) is its derivative and var is the variance

| Parameter | Equation |
|---|---|
| **Activity** | $var(y(t))$ |
| **Mobility** | $\sqrt{\dfrac{var\,(y'(t))}{var\,(y(t))}}$ |
| **Complexity** | $\dfrac{mobility\,(y'(t))}{mobility\,(y(t))}$ |

## 2 Materials and Methods

### 2.1 Database

Data is from the TUH EEG Seizure Corpus (TUSZ) v1.5.1 [19]. It was recorded in a real-time clinical environment using the International standard 10/20 system with 24 to 36 channels. The standard 19 EEG channels were used in this study. EEG segments have been labeled by experts. All uninteresting portions of the data, including eye blinks, artifacts, and noise were eliminated from the EEG records. The seizure segments annotated

as generalized and focal were considered by this study. The dataset consists of EEG signals collected from 115 patients of which 61 are females. The data contains 218 sessions that were broken to 1497 files, of which 1069 contain focal seizures. The sampling frequency varies between 250 and 500 Hz. Table 1 summarizes the database of this study. Figure 1 shows plots of generalized and focal seizure EEG records . More details about the dataset can be found in [9, 19], for instance.

**Table 2.** Overview of the subset of the tusz eeg corpus used in our study for seizure type classification

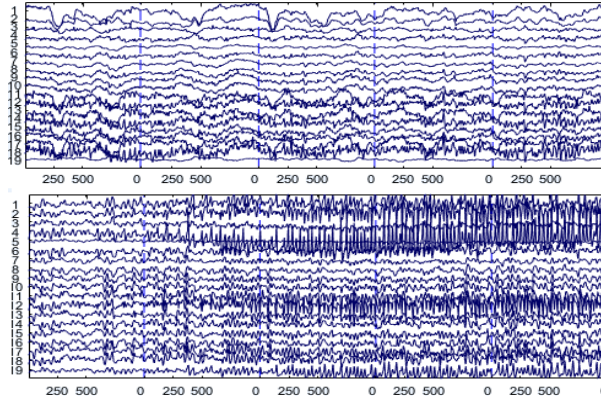| | |
|---|---|
| Nb of Patients (F) | 115 (61) |
| Nb of Sessions | 218 |
| Nb of Files | 1497 |
| Nb of Focal Seizures | 1069 |
| Nb of Generalized Seizures | 428 |
| Duration of Focal Seizure in hrs | 147.49 |
| Total duration in hrs | 287.79 |



**Fig. 1.** A) An example of raw EEG illustrating focal seizure epochs. B) An example of raw EEG illustrating generalized seizure epochs

## 2.2 The proposed method

A bandpass filter with cutoff frequencies (0.5:75) Hz has been used, followed by a 60Hz notch filter. The data has been re-referenced to the average of all electrodes, followed by re-sampling to 256Hz. Window sizes of 5-, 10-, 15- and 20-seconds were considered. Hjorth activity, mobility, and complexity parameters were then extracted from each channel of the pre-processed EEG signals and for each window size separately. A feature vector of dimension 57 was considered at each run. The supervised classification was carried out using the RF Classifier. The classifier was chosen based on its successful use in previous works [14, 15, 20]. 10-fold cross validation was employed for training and testing to avoid overfitting and to ensure stable and reliable results, where EEG signals are partitioned randomly into 10 subsets, where nine are for training and the remainder

for testing. In a first set of experiments, the weighted random forest classifier was considered to balance class weights. In a second set of experiments, data was rebalanced using the SMOTE prior to being used as input to classification. The performance of the classifier was evaluated using accuracy, recall and precision measures.

## 3      Experimental Results

The analysis described in this work was carried out using Matlab R2020b and Python. The study explored the use of four window sizes. Results of classification of focal vs generalized seizure EEG records are summarized in Tables 3 and Table 4. The performance of classification was calculated by averaging the accuracy, recall, and precision obtained using the test data in each of the 10 iterations. In both Tables 3 and Table 4, windows of 5, 10, 15, and 20s, were compared. The first evaluation uses the weighted random forest classifier. Results in Table 3 show a maximum accuracy of 87.3% obtained using the 10s window with a recall of 58.9% and precision of 90.6%. The second evaluation uses SMOTE followed by a conventional random forest. Results in Table 4 show a maximum accuracy of 92.3% with a recall of 91.9%, and precision of 92.6%, obtained using the 10s window size. An increase in window size corresponds to a slight decrease in performance.

**Table 3.** classification performance of the weighted random forest classifier using four window sizes

|  | Window Size in second | | | |
|---|---|---|---|---|
|  | 5s | 10s | 15s | 20s |
| Accuracy | 0.867 | **0.873** | 0.867 | 0.862 |
| Recall | 0.619 | 0.589 | 0.562 | 0.485 |
| Precision | 0.896 | 0.906 | 0.891 | 0.894 |

Performance was also investigated using the receiver operating characteristic (ROC) curve analysis and area under the curve (AUC) metric. The ROC curve represents the cut-off values between the true positive and false positive rates. Figure 2 displays ROC curves when weighted RF is used. Figure 3 displays ROC curves when SMOTE is used followed by a conventional RF. Figures 3 and 4 give the performance with 10-fold cross validation. AUC for each fold is shown in each figure, in addition to the mean of AUC.

**Table 4.** classification performance of the our approach using smote followed by random forest classifier

|  | Window Size in second | | | |
|---|---|---|---|---|
|  | 5s | 10s | 15s | 20s |
| Accuracy | 0.918 | **0.923** | 0.922 | 0.921 |
| Recall | 0.903 | 0.919 | 0.927 | 0.928 |
| Precision | 0.931 | 0.926 | 0.918 | 0.916 |

Results indicate that performance is good for all folds, giving an average of AUC equal to 92% for weighted RF. A high value of 98% AUC is obtained with RF preceded with SMOTE. In summary, RF classification preceded by SMOTE yields the best performance for each of the 10-folds compared to weighted RF.
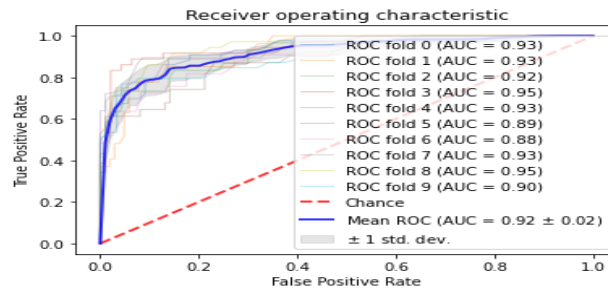


**Fig. 2.** Receiver Operating Characteristic (ROC) curve for the weighted RF classifier. Each curve denotes the ROC of one-fold of the 10-fold cross validation: AUC is displayed for each fold and mean AUC for the 10-fold.
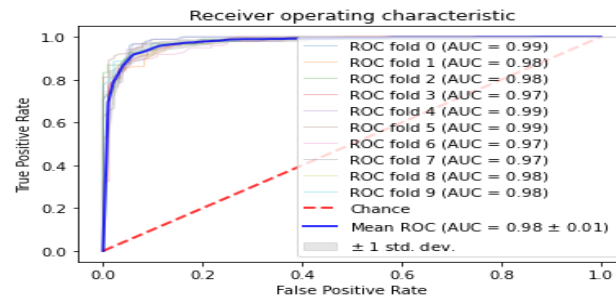


**Fig. 3.** Receiver Operating Characteristic (ROC) curve for RF classifier preceded with SMOTE. Each curve denotes the ROC of one-fold of the 10-fold cross validation: AUC is displayed for each fold and mean AUC for the 10-fold.

As can be seen from results in Tables 3 and 4, the proposed hybrid method (combination of Hjorth descriptor, SMOTE and RF) outperforms the weighted RF with consistent improvement of 5% in classification accuracy. The method achieved an accuracy of 92.3% for the focal and generalized seizures classification. The precision and recall were also good. These results demonstrate the advantage in balancing the input feature space using SMOTE before classification, instead of giving each class a mis-classification cost. With regard to the window size, the results show a slight difference in accuracy, recall and precision as the window size varies between 5 and 20 seconds. The best results were obtained for a 10s window size in both experiments.

## 4 Conclusion

In this study, we investigated focal vs generalized seizure classification using a sub-set of the TUSZ corpus. We investigated two ways to address imbalanced data, Random Forest supervised classification, and Hjorth data description extracted from 19 EEG

channels. Weighted Random Forest classification was compared to a conventional Random Forest classifier preceded by SMOTE oversampling. Results indicate a good separation between focal and generalized seizure using SMOTE applied on EEG segments of 10s window size. A system for classifying seizure types could be a clinically relevant tool for experts to have better diagnosis. In future work, we intend to expand our study to include classification of specific types of generalized and focal seizures such as tonic, clonic, tonic-clonic, complex partial, and simple partial seizures.

## Acknowledgment

## References

1. Shoeb, A.H.: Application of Machine Learning to Epileptic Seizure Onset Detection and Treatment MASS NSl OF TECHNOLOGY. Massachusetts Institute of Technology (2009).
2. Song, Y., Crowcroft, J., Zhang, J.: Automatic epileptic seizure detection in EEGs based on optimized sample entropy and extreme learning machine. Journal of Neuroscience Methods. 210, 132–146 (2012). https://doi.org/10.1016/j.jneumeth.2012.07.003.
3. Sharma, M., Dhere, A., Pachori, R.B., Acharya, U.R.: An automatic detection of focal EEG signals using new class of time–frequency localized orthogonal wavelet filter banks. Knowledge-Based Systems. 118, 217–227 (2017). https://doi.org/10.1016/j.knosys.2016.11.024.
4. Bhattacharyya, A., Sharma, M., Pachori, R.B., Sircar, P., Acharya, U.R.: A novel approach for automated detection of focal EEG signals using empirical wavelet transform. Neural Computing and Applications. 29, 47–57 (2018). https://doi.org/10.1007/s00521-016-2646-4.
5. Zeng, W., Li, M., Yuan, C., Wang, Q., Liu, F., Wang, Y.: Classification of focal and non focal EEG signals using empirical mode decomposition (EMD), phase space reconstruction (PSR) and neural networks. Artificial Intelligence Review. 52, 625–647 (2019). https://doi.org/10.1007/s10462-019-09698-4.
6. Dwi Saputro, I.R., Maryati, N.D., Solihati, S.R., Wijayanto, I., Hadiyoso, S., Patmasari, R.: Seizure Type Classification on EEG Signal using Support Vector Machine. In: Journal of Physics: Conference Series. p. 12065. IOP Publishing (2019). https://doi.org/10.1088/1742-6596/1201/1/012065.
7. Roy, S., Asif, U., Tang, J., Harrer, S.: Seizure Type Classification Using EEG Signals and Machine Learning: Setting a Benchmark. In: 2020 IEEE Signal Processing in Medicine and Biology Symposium, SPMB 2020 - Proceedings

(2020). https://doi.org/10.1109/SPMB50085.2020.9353642.

8. Das, A.B., Bhuiyan, M.I.H.: Discrimination and classification of focal and non-focal EEG signals using entropy-based features in the EMD-DWT domain. Biomedical Signal Processing and Control. 29, 11–21 (2016). https://doi.org/10.1016/j.bspc.2016.05.004.

9. Obeid, I., Picone, J.: The Temple University Hospital EEG Data Corpus. Frontiers in Neuroscience. 10, 196 (2016). https://doi.org/10.3389/fnins.2016.00196.

10. Hjorth, B.: EEG analysis based on time domain properties. Electroencephalography and Clinical Neurophysiology. 29, 306–310 (1970). https://doi.org/10.1016/0013-4694(70)90143-4.

11. Oh, S.-H., Lee, Y.-R., Kim, H.-N.: A Novel EEG Feature Extraction Method Using Hjorth Parameter. International Journal of Electronics and Electrical Engineering. 106–110 (2014). https://doi.org/10.12720/ijeee.2.2.106-110.

12. Vourkas, M., Papadourakis, G., Micheloyannis, S.: Use of ANN and Hjorth parameters in mental-task discrimination. IEE Conference Publication. 327–332 (2000). https://doi.org/10.1049/cp:20000356.

13. Jemal, I., Mitiche, A., Mezghani, N.: A Study of EEG Feature Complexity in Epileptic Seizure Prediction. Applied Sciences. 11, 1579 (2021). https://doi.org/10.3390/app11041579.

14. Fraiwan, L., Lweesy, K., Khasawneh, N., Wenz, H., Dickhaus, H.: Automated sleep stage identification system based on time-frequency analysis of a single EEG channel and random forest classifier. Computer Methods and Programs in Biomedicine. 108, 10–19 (2012). https://doi.org/10.1016/j.cmpb.2011.11.005.

15. Lehmann, C., Koenig, T., Jelic, V., Prichep, L., John, R.E., Wahlund, L.O., Dodge, Y., Dierks, T.: Application and comparison of classification algorithms for recognition of Alzheimer's disease in electrical brain activity (EEG). Journal of Neuroscience Methods. 161, 342–350 (2007). https://doi.org/10.1016/j.jneumeth.2006.10.023.

16. Breiman, L.: Random forests. Machine Learning. 45, 5–32 (2001). https://doi.org/10.1023/A:1010933404324.

17. Chawla, N. V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: SMOTE: Synthetic minority over-sampling technique. Journal of Artificial Intelligence Research. 16, 321–357 (2002). https://doi.org/10.1613/jair.953.

18. Chen, C., Liaw, A., Breiman, L.: Using Random Forest to Learn Imbalanced Data | Department of Statistics. (2004).

19. Shah, V., von Weltin, E., Lopez, S., McHugh, J.R., Veloso, L., Golmohammadi, M., Obeid, I., Picone, J.: The Temple University Hospital Seizure Detection Corpus. Frontiers in Neuroinformatics. 12, 83 (2018). https://doi.org/10.3389/fninf.2018.00083.

20. Donos, C., Dümpelmann, M., Schulze-Bonhage, A.: Early Seizure Detection Algorithm Based on Intracranial EEG and Random Forest Classification. International Journal of Neural Systems. 25, (2015). https://doi.org/10.1142/S0129065715500239.