

SCORE-IT: A Machine Learning-based Tool for Automatic Standardization of EEG Reports

Sam Rawal¹ and Yogatheesan Varatharajah^{2,3}

1. Carle Illinois College of Medicine, University of Illinois at Urbana Champaign.
2. Department of Bioengineering, University of Illinois at Urbana Champaign.
3. Department of Neurology, Mayo Clinic.
{scrawal2, varatha2}@illinois.edu

Abstract— Machine learning (ML)-based analysis of electroencephalograms (EEGs) is playing an important role in advancing neurological care. However, the difficulties in automatically extracting useful metadata from clinical records hinder the development of large-scale EEG-based ML models. EEG reports, which are the primary sources of metadata for EEG studies, suffer from lack of standardization. Here we propose a machine learning-based system that automatically extracts components from the SCORE specification from unstructured, natural-language EEG reports. Specifically, our system identifies (1) the type of seizure that was observed in the recording, per physician impression; (2) whether the session recording was *normal* or *abnormal* according to physician impression; (3) whether the patient was diagnosed with epilepsy or not. We performed an evaluation of our system using the publicly available TUH EEG corpus and report F1 scores of 0.92, 0.82, and 0.97 for the respective tasks.

I. INTRODUCTION

Analysis and interpretation of EEG data is important in the diagnosis of neurological conditions such as epilepsy. Recently, machine learning-based analyses are becoming the mainstay of quantitative EEG-based decision making [1]. EEG reports are the primary source of metadata for quantitative EEG studies because they contain rich information regarding the patients' condition at the time of EEG recording including seizures, interictal abnormalities, and background activity such as posterior dominant rhythm. However, due to a lack of standardization in the organization of information presented and terminology in EEG reports, automatic extraction of those information has proven to be a difficult task. As a result, manual information retrieval has been the primary way of generating metadata for EEG-based ML studies. Clearly, this is not scalable and warrants the development of natural language processing-based automated tools.

The standardized computer-based organized reporting of EEG (SCORE) guidelines seek to standardize reporting of EEG studies by specifying the content and terminology of characteristics that are described in an EEG report [2, 3]. As such, those guidelines provide an ideal target for automated information retrieval and validation. In this work, we present a semi-automated tool for extracting structured information from unstructured natural text EEG reports with an eye towards converting

past EEG reports to the standardized format stipulated in the SCORE guidelines. Furthermore, we also leverage the classification tasks made available through various labeled subsets of the TUH dataset as a valuable proxy to evaluate the capabilities of our system.

The development of such a system poses several challenges: 1) data sparsity: from an entire patient record consisting of dozens of sentences, often only a single phrase is relevant to making a classification decision; 2) lack of labeled data: the lack of clinical data available around many of the SCORE attributes; and 3) accounting for varied clinician practices. We took a two-step approach to address those challenges: first, we leverage a previously developed named entity recognition approach using BERT Transformer models trained on the National NLP Clinical Challenges (n2c2) dataset [4]; second, hand-crafted rules are applied to these extracted entities, considering factors such as the SCORE lexicon and numerical values to identify information relevant to SCORE attributes. This hybrid approach allows us to leverage a) existing well-trained ML models for relevant entity identification, and b) domain knowledge to further classify them into SCORE attributes.

In the current work, we focus on identifying three SCORE attributes: (1) the type of seizure if present (complex partial, simple partial, absence, myoclonic, generalized tonic-clonic, other/none); (2) whether the recording contained any abnormal events (e.g., seizures); and (3) whether the patient is being evaluated for epilepsy. In addition, we designed and validated our system on the Temple University Hospital (TUH) EEG dataset [5], containing over 16,000 EEG recordings and reports. We utilized subsets of the TUH EEG corpus consisting ground truth labels for the three tasks we focused on. We achieved weighted F1 scores of 0.92 (on 171 test records), 0.82 (on 561 records), and 0.97 (on 2727 records) for the seizure classification, normal/abnormal classification, and epilepsy classification tasks, respectively.

II. RELATED WORK

Existing work on classification from patient medical records includes Track 1 of the 2018 National NLP Clinical Challenges (n2c2), which involves making binary classification decisions about whether a patient

meets some criteria, such as alcoholism or history of myocardial infarction over past 6 months [6]. Existing strategies for this type of task include extraction of clinical entities, hand-crafted features that are used to train simple machine learning models, and creation of lexicons to feed into rules [7]. On the other hand, there are related classification tasks in the clinical domain for which there is ample labeled data available, e.g., clinical/biomedical named entity recognition (NER) [4]. Prior work on performing NER has included the use of classical Machine Learning models, such as Conditional Random Fields, as well as deep neural networks, including LSTM networks [8] and, more recently, fine-tuning transformer networks, such as BERT [9].

However, our work on automatically extracting useful metadata such as seizure type, EEG classification, and diagnostic information from EEG reports, to our knowledge, is the first attempt at developing automated information retrieval approaches for EEG reports. Considering that the majority of the (past and present) EEG reports are written in natural text format, our work offers the potential to standardize EEG reports and to generate useful metadata for subsequent ML-based analyses.

III. DATA

The TUH EEG Corpus (TUEG) is a collection of over 30,000 clinical EEG records collected and made available by Temple University Hospital (TUH) [5]. Corresponding patient medical reports, in plaintext format, are also made available alongside the EEG recordings. There are also subsets of the corpus, containing labeled data for several tasks. In this work, we utilize the TUH EEG Epilepsy Corpus (TUEP), TUH Abnormal EEG Corpus (TUAB), and TUH EEG Seizure Corpus (TUSZ). Table 1 describes the number of labeled samples for each of the subsets. (As the Epilepsy subset had no train and test partition, unlabeled records from outside the subset were used to develop rules.)

Table 1. TUH Seizure Dataset Support

Dataset	Class	Train Support	Test Support
Seizure	Absence	10	6
	Complex Partial	45	13
	Myoclonic	1	12
	Simple Partial	2	0
	Tonic-Clonic	12	4
	None	913	97
Epilepsy	Epilepsy	428	428
	No Epilepsy	133	133
Abnormal	Normal	1371	150
	Abnormal	1346	126

IV. METHODOLOGY

In this work, we present a system that, given a natural-language patient report, predicts, 1) the type of seizure that was observed in the recording, per physician impression; 2) whether the session recording was *normal*

or *abnormal* according to physician impression; and 3) whether a patient has epilepsy or not.

To construct such a system, we provide a framework for classifying EEG records that divides the process into two steps: a *broad parsing* step, for which tasks are well-defined and considerable training data exists; and a *narrow parsing/classification* step built on top of the broad parsing, for which tasks are highly domain-specific and little training data is available. In this manner, we address issues of data sparsity and support while simultaneously leveraging effective data-driven, deep learning techniques that can be adapted to this specific domain.

A high-level architecture of the system is detailed in Figure 1. In the following subsections, we discuss each component.

IV-A. Broad Parsing

The overall goal of this step is to reduce the data sparsity by only extracting information from the record that will be potentially useful for classification, while leveraging existing well-trained machine learning models and methods that utilize the structure of the medical records.

First, section headers are identified, and each sentence in the report is matched with its corresponding header. The headers are extracted using a regular expression that matches the format of the medical records present in the TUH corpus. Second, we perform Named Entity Recognition using BERT Transformer models [9] trained on datasets released by national NLP clinical challenges (n2c2). Specifically, we extract clinical entities, such as medical problems, labs, and treatments, and medication entities, such as medication name, dose, frequency, duration, and reason. The information extracted in this step is then passed to the *narrow parsing* step.

IV-B. Narrow Parsing

The *narrow parsing* step that is built on top of the broad parsing step. This step is motivated by the fact that the classification tasks we address have a small support of training data; therefore, it is not possible to train a system only from the supervised data. Consequently, this phase consists of a series of hand-crafted rules built around the extracted entities and sections in the prior step, in effect performing domain adaptation from models trained on related tasks in the broad parsing step to these three tasks with limited data.

Classification of each of the three tasks (epilepsy, normal/abnormal, seizure type) is performed using hand-crafted rules that are constructed using the outputs of the section and entity extraction from the previous step. Thus, the broad parsing step dramatically compresses the input dimensionality from an unstructured, free-text medical report to a collection of named entities and the

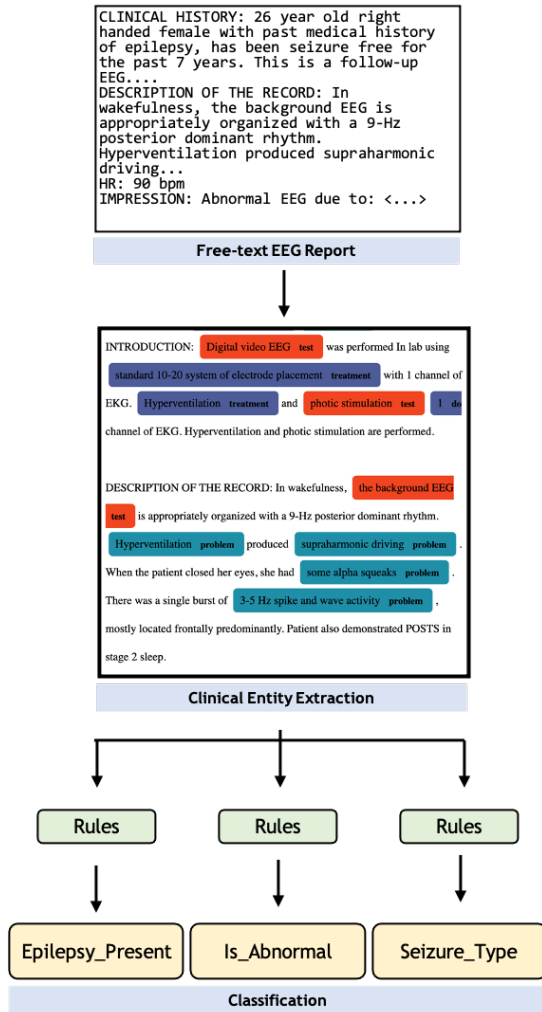


Figure 1. System Architecture.

section of the note they are identified in, thus making it possible to create simpler and more generalizable rules.

For example, the rule used for binary classification in the Epilepsy identification task consists of extracting all *problem* entities mentioned in the “CLINICAL HISTORY” section of the patient record, then outputting a positive classification if any of the entities match *epilepsy* (or synonyms and related phrases).

V. EXPERIMENTAL DESIGN

The system was evaluated on three specific tasks that are subsets of the TUH EEG dataset, and for which gold annotations are provided:

- 1) Seizure type classification: given a record, determine the type of seizure noted in the EEG record (complex partial, simple partial, absence, myoclonic, tonic-clonic, other/none).

- 2) Epilepsy classification: given a record, determine whether the patient is being evaluated for epilepsy or not.
- 3) Normal/abnormal classification: determine if the clinical impression in a given record is *normal* or *abnormal*.

A detailed information for each of these tasks is provided in the following subsections.

V-A. Seizure Classification

The Seizure subset of the TUH dataset contains 1010 unique patient records classified by the type of seizure that is presented in the corresponding EEG recording. One important aspect to note is that, in this work, we focused on specific types of seizures or NONE. Consequently, records labeled “GNSZ” or “FNSZ” (for generalized or focal seizures respectively, without any further specifics) were dropped from the dataset. The breakdown of support for each type of seizures is described in Table 1.

V-B. Epilepsy Classification

The Epilepsy subset of the TUH dataset contains 428 patient records labeled *epilepsy* and 133 records labeled *no epilepsy*. One complication of this dataset is that the labels are associated with the EEG recording findings, rather than what is detailed in the patient records; in a manual review of the dataset, it appears there are numerous instances of the gold label not being substantiated within the clinical record alone.

V-C. Normal/Abnormal Classification

The Abnormal subset of the TUH corpus consists of notes that are labeled *normal* or *abnormal*. This classification is made based off physician comments in the “CLINICAL IMPRESSIONS” section of medical records. There are 126 notes with the *abnormal* label and 150 notes with the *normal* label in this subset.

VI. RESULTS

On the seizure classification task, our system performed with a weighted F1 score of 0.92 on 171 test records. On the epilepsy classification task, our system performed with a weighted F1 score of 0.82 on 561 records. On the abnormal classification task, our system performed with a weighted F1 of 0.97 on 2727 records. Notably, the system performed poorer on classes with a very small support, due to difficulty in creating generalizable hand-crafted rules from a very small dataset. More detailed results are provided in Table 2.

VII. DISCUSSION

We proposed an automated approach for extracting clinically useful metadata from natural text EEG reports.

Table 2. Classification Results – Weighted Average

Task	Precision	Recall	F1-score	Support
Seizure	0.93	0.93	0.93	121
Abnormal	0.98	0.97	0.97	276
Epilepsy	0.82	0.82	0.82	561

Major challenges addressed in our work include 1) designing strategies for integration of clinical domain knowledge, 2) accounting for varied clinician practices, 3) establishing ground truth standards, and 4) robust verification of the system. Applications of this system include enabling better search and filtering of clinical reports, as well as auto-labeling unstructured datasets for research purposes.

As mentioned previously, one complication with interpreting these results is that for certain tasks, such as epilepsy classification, the labels are derived from the EEG findings and not necessarily explicitly mentioned in the patient medical record, making it difficult to fully accurately judge the performance of the system.

Although the broad and narrow parsing architecture of the system helped in the development of more general rules, the system still performed better on classes with greater support. This can be attributed to the fact that having access to a wider range of training samples allowed rules to be developed that could capture more of the variance. On the other hand, the system performed poorer on classes with few samples, as the rules developed with under 5 training samples were not able to be validated to ensure sufficient generalizability.

VIII. FUTURE WORK

The current iteration of the system still performs relatively poorly in classification circumstances with low data support. Consequently, one important area of focus is reducing the amount of hand-crafted rules needed, while increasing the generalizability of the system, in order to be able to apply the system across a wider range of use-cases. To that effect, a prominent area of future focus remains using existing NLP tasks, such as semantic textual similarity [10], as a means to perform zero-shot classification in place of hand-crafting rules. Another area for future work is implementing and validating our approach on additional datasets.

IX. CONCLUSION

We described a semi-automated ML-based system to extract standardized components from unstructured EEG reports. Our system can facilitate better indexing, searching and organization of existing EEG report collections by providing an semi-automated methodology of standardizing important data within free-text EEG reports. Furthermore, our system can be used to generate labels for large, unlabeled EEG corpora like the TUH

dataset to motivate future research.

ACKNOWLEDGEMENTS

This research was supported by National Science Foundation’s Computer and Information Science and Engineering Research Initiation Initiative Award SCH-2105233. We would like to thank Dr. Susan Herman and her team from the Barrow Neurological Institute, as well as Neeraj Wagh from the University of Illinois at Urbana Champaign, for contributions to discussions on where this work can be applied.

REFERENCES

- [1] M. Golmohammadi, A. H. Harati Nejad Torbati, S. Lopez de Diego, I. Obeid, and J. Picone, “Automatic analysis of eegs using big data and hybrid deep learning architectures,” *Frontiers in human neuroscience*, vol. 13, p. 76, 2019.
- [2] S. Beniczky, H. Aurlien, J. C. Brögger, A. Fuglsang-Frederiksen, A. Martins-da Silva, E. Trinka, G. Visser, G. Rubboli, H. Hjalgrim, H. Stefan *et al.*, “Standardized computer-based organized reporting of eeg: Score,” *Epilepsia*, vol. 54, no. 6, pp. 1112–1124, 2013.
- [3] S. Beniczky, H. Aurlien, J. C. Brögger, L. J. Hirsch, D. L. Schomer, E. Trinka, R. M. Pressler, R. Wennberg, G. H. Visser, M. Eisermann *et al.*, “Standardized computer-based organized reporting of eeg: Score—second version,” *Clinical Neurophysiology*, vol. 128, no. 11, pp. 2334–2346, 2017.
- [4] J. Patrick and M. Li, “High accuracy information extraction of medication information from clinical notes: 2009 i2b2 medication extraction challenge,” *Journal of the American Medical Informatics Association*, vol. 17, no. 5, pp. 524–527, 2010.
- [5] A. Harati, S. Lopez, I. Obeid, J. Picone, M. Jacobson, and S. Tobochnik, “The tuh eeg corpus: A big data resource for automated eeg interpretation,” *2014 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE, 2014, pp. 1–5.
- [6] A. Stubbs, M. Filannino, E. Soysal, S. Henry, and Ö. Uzuner, “Cohort selection for clinical trials: n2c2 2018 shared task track 1,” *Journal of the American Medical Informatics Association*, vol. 26, no. 11, pp. 1163–1171, 2019.
- [7] S. Rawal, A. Prakash, S. Adhya, S. Kulkarni, S. Anwar, C. Baral, and M. Devarakonda, “Semi-automated clinical lexicon induction and its use in cohort selection from clinical notes,” *2020 IEEE International Conference on Healthcare Informatics (ICHI)*. IEEE, 2020, pp. 1–2.
- [8] S. C. Rawal, “Prescription information extraction from electronic health records using bilstm-crf and word embeddings,” 2018.
- [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [10] Y. Wang, S. Fu, F. Shen, S. Henry, O. Uzuner, and H. Liu, “The 2019 n2c2/ohnlp track on clinical semantic textual similarity: overview,” *JMIR Medical Informatics*, vol. 8, no. 11, p. e23375, 2020.