

Deep Learning Based Classification of True/False Arrhythmia Alarms in the Intensive Care Unit

Jack Boynton¹, Byung Suk Lee²

¹ Department of Electrical & Biomedical Engineering, University of Vermont, Burlington, VT, USA

² Department of Computer Science, University of Vermont, Burlington, VT, USA

Abstract

Once a cardiac alarm is triggered in the intensive care unit (ICU), accurately classifying whether the alarm is true or false is of critical importance. Incorrect classification may lead to patient's death if the alarm is true or to disruption in patient care if false. There has been a body of research, as signified by the 2015 PhysioNet/CinC Challenge; due accomplishments have been made in the relevant computational technology, and yet the highest accuracy known thus far is in the mid-80% range (85%). Our work achieved much higher accuracy and, additionally, very early classification almost at the onset of an arrhythmia alarm, by utilizing state of the art deep learning methods. The machine learning model used is a Residual Network (ResNet) and a Bi-directional Long Short Term Memory (BiLSTM) connected in tandem. Using the PhysioNet dataset of 750 recorded ECG segments published with the challenge, our method performed the classification with 96% accuracy in 0.52 seconds from the onset of an alarm on average over all test ECG segments.

1. Introduction

Research presented in this paper stems from the 2015 PhysioNet/CinC Challenge [1], particularly the “retroactive” classification test to determine within 10 seconds whether an arrhythmia alarm is true or false in the ICU. The accuracy of classification is undoubtedly important; misclassifying a true alarm as false may result in a patient’s death and misclassifying a false alarm as true may result in wasteful disruption and disturbance. Thus, the goal of the challenge was to achieve as high true/false alarm classification accuracy as possible. The result of the challenge [2] was limited to 85% accuracy for the top method, and there has been no further advancement since then. Our work adds novel advancements to this state of the art.

There are two advancements. First, our work enhanced the accuracy to a high 90% range (96%) by using a *deep learning* model as the computational method. None of the published papers on this problem (e.g., reported by Clif-

ford et al. [2]) used deep learning, which has proven to produce powerful models given adequate training data. Second, our work addresses *early* classification of true or false alarms, which is attributed to the excellent feature extraction ability of a deep learning model. Early classification is as important as accurate classification. A delayed classification of true alarm is potentially dangerous to the patient, and a delayed classification of false alarm deprives the opportunity to suppress it in time. Thus, this paper presents the computational methods and the results of using a deep learning model to determine whether an alarm is true or false accurately and early in the ICU.

The deep learning model we used was inspired by Zhou et al.’s work [3], where a model combining Residual Network (ResNet) and Bi-directional Long Short Term Memory (BiLSTM) in tandem has been used to achieve impressive ECG heartbeat classification accuracy. Prequential evaluation [4] was used to train the model interleaved with testing while augmenting training data set incrementally through the evaluation cycle. The training data set and test data set were obtained by splitting training data set published in the 2015 PhysioNet/CinC Challenge.

The deep learning model outputs the probability of an alarm being true. (Its 1’s complement is the probability of the alarm being false.) It was observed that the output probability is always either increasing or decreasing monotonously, and this observation enabled an early classification of true or false alarm as early as the first sign of the direction of change (namely the “polarity”), either positive or negative. Through a progressive study of decreasing the time interval of outputting the probability, a conclusion was made that the ECG sample interval achieved the earliest classification, amounting to 0.52 seconds (for 125 samples) on average while achieving higher classification accuracy (96%) than the published methods [2]. Source codes and data sets of the methods used in our work and select results are available in a GitHub repository [5].

The rest of this paper presents the computational methods in Section 2, the experiment results and discussion in Section 3, and the conclusion in Section 4.

2. Methods

2.1. Deep learning model

Experimenting with different deep learning models led to our choice of a model comprising ResNet and BiLSTM in tandem, based on the work by Yang Zhou et al. [3]. Figure 1 shows the model structure. ResNet is used to extract complex features from the ECG time series, and BiLSTM is used to build a prediction model based on the temporal order of features. The ResNet architecture we used is from the work by Ismail Fawaz [6], and consists of three residual convolutional blocks with filter sizes 64, 128, and 128, respectively, effectively extracting 128 features from a given ECG segment. The BiLSTM is a standard TensorFlow model ‘Bidirectional.’



Figure 1: Deep learning model structure (source: [3]).

2.2. Prequential evaluation

Time series data are arriving as chronologically ordered samples. Therefore, we used prequential evaluation [4] (as opposed to the conventional cross-validation) to consider the effect of such temporal ordering. Basically, each new batch of data is first used as test data and then appended to the existing training data. Thus, the training data size keeps increasing, and so does the training time as the total training size increases (see Figure 2). Our prequential evaluation is a ‘growing window’ version adopted from the empirical study done by Cerqueira et al. [7]

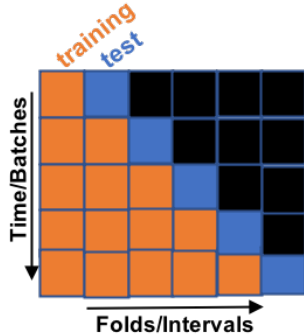


Figure 2: Prequential evaluation (source: [7]).

2.3. Data sets

The 2015 PhysioNet/CinC Challenge public dataset contains 750 five-minute ECG recordings from four hospitals with known and categorized life-threatening arrhythmia ICU alarms that occur at the beginning of the last 10 seconds of each segment. ECG signals in the public dataset were re-sampled to 250 Hz, and passed through a band-pass filter of 0.05 to 40 Hz to reduce baseline drift and

noise. Two ECG leads (II and VII) and one arterial blood pressure lead are included in the dataset, and the ECG lead I was used in this work.

The training dataset provided in the 2015 PhysioNet/CinC Challenge was split to training and test data sets at the ratio of 80% to 20% for the purpose of our evaluation, as the test data set used in the challenge was not available.

3. Results and Discussion

3.1. True/false alarm classification performance

The performance has two factors: classification accuracy (i.e., how accurate the classification output is) and classification time (i.e., how quickly the classification can be made). The classification accuracy was measured as in Equation 1 following the 2015 PhysioNet/CinC Challenge, and the classification time was measured as wall-clock time (or, equivalently, the number of ECG samples).

$$\text{Accuracy score} = \frac{TP + TN}{TP + TN + FP + 5 \times FN} \quad (1)$$

We first tried a threshold-based approach, that is, wait until the probability reaches a certain threshold value before outputting the classification result. Specifically, the model training and testing were done at different batch intervals of ECG samples, progressively decreased from 10 seconds down to 2, 1, 0.5 seconds and then further down to one sample interval (4 milliseconds). The result showed an increase in the classification time without any gain in the classification accuracy when the interval was reduced. That is, the time until the probability reaches a threshold was the same (0.96 seconds) regardless of the batch interval, but a longer interval delayed the classification time due to the latency to reach the end of the interval (see Figure 3).

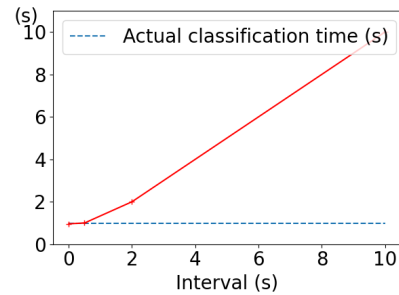


Figure 3: Mean classification time for varying interval.

Careful inspection of this phenomenon revealed that, for all test ECG segments, the probability changed rapidly and monotonously, either positive (true alarm case) or negative (false alarm case)—see Figure 4. Based on this observation, the notion of threshold was dismissed, and the earliest possible sample point of detecting the polarity of the

probability change was examined. We would call this a “polarity” approach. The resulting classification time was about 125 samples, amounting to 0.52 seconds, on average, which was far shorter than 1.88 seconds on average for the threshold-based approach (see Figure 5). Notably, the classification time in the polarity approach was also more consistent across different test ECG segments; moreover, earlier classification did not compromise the accuracy at all, and actually raised it a bit to 96.23% compared with 95.00% for the threshold approach.

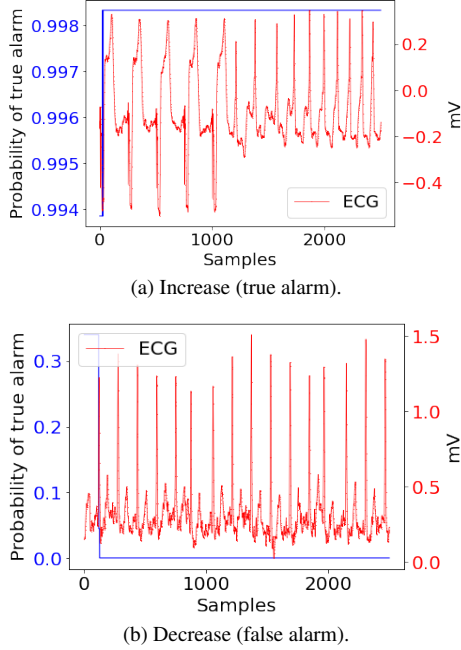


Figure 4: Change of model’s output probability over time.

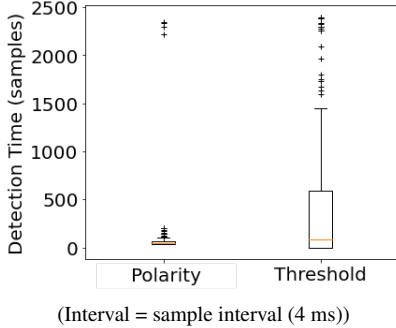


Figure 5: Classification times of the two approaches.

The remarkably early classification with such a high accuracy drew our suspicion at first and triggered a thorough investigation. The prequential evaluation was part of the investigation effort, to prevent overfit of the model for time series. Visual inspection of a significant number of the ECG segments used suggested that the key to such early classification is the substantial regularity inherent in the ECG waveform morphology. This regularity enables the deep learning model to quickly capture the signature features of samples that predict the polarity of the output

probability (i.e., increase or decrease). In our work, it was typically after “seeing” the first wavelet (e.g., P-, QRS, T-wave) in an ECG beat. A further, larger-scale investigation is warranted involving more diverse ECG segments reflecting different patient cohorts (e.g., gender, age, body mass index).

3.2. Ablation study of the model

To assess the merit of using the combined model of ResNet and BiLSTM in tandem as opposed to using either one of them, an ablation study was done to compare with ResNet only and BiLSTM only. The result was that ResNet+BiLSTM achieved 5.7% and 34.51% higher accuracy than ResNet-only and BiLSTM-only, respectively (see Table 1). Note that Table 1 also shows accuracy measured without bias against FN (see Equation 2). The result is even higher accuracy, which is encouraging as detecting true alarms is as important as, or more important than, detecting false alarms.

$$\text{Accuracy score} = \frac{TP + TN}{TP + TN + FP + FN} \quad (2)$$

Table 1: Classification accuracy among the three model structures (interval = sample interval (4 ms)).

Measure	ResNet+BiLSTM	ResNet	BiLSTM
Equation 1	96.23%	91.05%	71.54%
Equation 2	98.71%	93.43%	78.93%

3.3. Sanity check on incorrect classification

Despite the high accuracy of 96.23%, we made an effort to identify the cause of the 3.77% incorrect classifications. We first checked if the diagnostic type of ECG anomaly was the cause, but found no correlation with the five diagnostic types of ECG in the data set (i.e., asystole, extreme bradycardia, extreme tachycardia, ventricular tachycardia, ventricular flutter/fibrillation [1]; Table 1). Then, upon examination of visualized ECG signals of the 3.77% concerning segments, it was revealed that all incorrect classification was attributed to extreme noise in the ECG signal (see Figure 6) independently of the diagnostic type.

3.4. Comparison with prior art

Table 2 summarizes the top four methods that achieved higher than 80% accuracy (according to Equation 1) in the retroactive classification test among the contestants [2]; the ‘voting algorithm’ was not a contestant but added by the challenge organizer to always pick the best result via majority vote of the select top 13 methods. Note that none of them used deep learning. Note as well that the aspect of “early classification” is unique to our work, so is not applicable to their work.

It should be noted that the accuracy results of the compared work are based on a test dataset used by organizer of the 2015 PhysioNet/CinC Challenge, whereas the accuracy

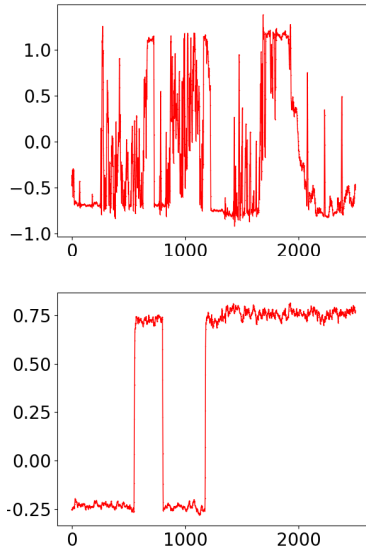


Figure 6: Example noisy post-alarm segments.

Table 2: Comparison with other work.

Publication	Method	Accuracy
This paper	ResNet+BiLSTM	96.23%
Voting algorithm (top 13) [2]	Majority vote	87.04%
Fallet et al. [8]	Rule-based	85.04%
Plesinger et al. [9]	Rule-based	84.96%
Kalidas et al. [10]	Support vector machine	81.85%

result of our work is based on a random 20% of the PhysioNet training dataset that was used in all the compared work. We used the remaining 80% of PhysioNet dataset as our own training dataset without any overlap to the 20% used as our test dataset. We believe this has no or little impact on the validity of the accuracy results.

4. Conclusion

This paper presented a novel work that achieved accurate and early classification of true or false arrhythmia alarm in the ICU, surpassing the state of the art. The enabling computational method was deep learning. The deep learning model, the prequential evaluation, and the experiments were presented. Our immediate further work is to adopt the method into the “real-time test” case of the PhysioNet/CinC 2015 Challenge [2], to predict a true alarm early and accurately before an arrhythmia alarm is triggered. The method can be also applied more broadly to other time-critical arrhythmia monitoring settings as well, like remote cardiac care via an implanted monitor [11].

Acknowledgments

This research was funded by the University of Vermont Summer Undergraduate Research Fellowship (SURF) and College of Engineering and Mathematical Sciences Re-

search Experience for Undergraduates (REU) award.

References

- [1] Clifford G, Silva I, Moody B, Mark R. Reducing false arrhythmia alarms in the ICU - the PhysioNet computing in cardiology challenge 2015. <https://www.physionet.org/content/challenge-2015/1.0.0/>, February 2015.
- [2] Clifford G, Silva I, Moody B, Li Q, Kella D, Shahin A, Kooistra T, Perry D, Mark R. The PhysioNet/Computing in Cardiology challenge 2015: Reducing false arrhythmia alarms in the ICU. In Proceedings of the 42nd Annual Conference on Computing in Cardiology. September 2015; 273–276.
- [3] Zhou Y, Zhang H, Li Y, Ning G. ECG heartbeat classification based on ResNet and Bi-LSTM. IOP Conference Series Earth and Environmental Science January 2020; 428:012014.
- [4] Gama J, Sebastião R, Rodrigues P. On evaluating stream learning algorithms. Machine Learning October 2013; 90:317–346.
- [5] Boynton J, Lee BS. Classification of true or false cardiac alarms in the ICU. https://github.com/JackWBoynton/ECG_Alarm_Classification_ICU/, April 2021.
- [6] Ismail Fawaz H. Deep learning for time series classification. <https://github.com/hfawaz/dl-4-tsc>, April 2020.
- [7] Cerqueira V, Torgo L, Mozetič I. Evaluating time series forecasting models: An empirical study on performance estimation methods. Machine Learning November 2020; 109:1997–2028.
- [8] Fallet S, Yazdani S, Vesin J. A multimodal approach to reduce false arrhythmia alarms in the intensive care unit. In Proceedings of the 42nd Annual Conference on Computing in Cardiology. September 2015; 277–280.
- [9] Plesinger F, Klimes P, Halamek J, Jurak P. False alarms in intensive care unit monitors: Detection of life-threatening arrhythmias using elementary algebra, descriptive statistics and fuzzy logic. In Proceedings of the 42nd Annual Conference on Computing in Cardiology. September 2015; 281–284.
- [10] Kalidas V, Tamil LS. Enhancing accuracy of arrhythmia classification by combining logical and machine learning techniques. In Proceedings of the 42nd Annual Conference on Computing in Cardiology. September 2015; 733–736.
- [11] Cheung CC, Deyell MW. Remote monitoring of cardiac implantable electronic devices. Canadian Journal of Cardiology July 2018;34(7):941–944.

Address for correspondence:

Byung Suk Lee

E409 Innovation Hall, 82 University Place, Burlington, VT 05405, U.S.A.

bslee@uvm.edu