

Towards High Generalization Performance on Electrocardiogram Classification

Hyeongrok Han^{†1}, Seongjae Park^{†2}, Seonwoo Min^{1,3}, Hyun-Soo Choi⁴, Eunji Kim¹, Hyunki Kim¹, Sangha Park¹, Jinkook Kim², Junsang Park², Junho An², Kwanglo Lee², Wonsun Jeong², Sangil Chon², Kwonwoo Ha², Myungkyu Han², Sungroh Yoon^{* 1,5}

¹Department of Electrical and Computer engineering, Seoul National University, Seoul, South Korea

²HUINNO Co., Ltd., Seoul, South Korea

³LG AI Research, Seoul, South Korea

⁴Department of Computer Science and Engineering, Kangwon National University, Chuncheon, South Korea

⁵Department of Biological Sciences, Interdisciplinary Program in Bioinformatics, Interdisciplinary Program in Artificial Intelligence, ASRI, INMC, and Institute of Engineering Research, Seoul National University, Seoul, South Korea

Abstract

Recently, many electrocardiogram (ECG) classification algorithms using deep learning have been proposed. The characteristics of ECG vary from dataset to dataset for various reasons (i.e., hospital, race, etc.). Therefore, it is important for a model to have high generalization performance consistently over all datasets. In this paper, as part of the PhysioNet / Computing in Cardiology Challenge 2021, we present a model developed to classify cardiac abnormalities from 12 lead and reduced-lead ECGs. In particular, to upgrade our previous model for improving generalization performance, we newly adopt constant-weighted cross-entropy loss, additional features, Mixup augmentation, and squeeze/excitation block, OneCycle learning rate scheduler, which are selected via evaluation of generalization performance using leave-one-dataset-out cross-validation setting. With the present model, our DSAIL_SNU team has received challenge scores of 0.55, 0.58, 0.58, 0.57 and 0.57 (ranked 2nd, 1st, 1st, 2nd, 2nd out of 39 teams) for the 12-lead, 6-lead, 4-lead, 3-lead, and 2-lead versions of the hidden test set, respectively. The present model achieves higher generalization performance over all versions of the hidden test set than the model submitted last year.

1. Introduction

Electrocardiogram (ECG) is an important tool for diagnosing cardiac abnormalities, and more than 300 mil-

[†]: equal contribution (Hyeongrok Han and Seongjae Park)

^{*}: corresponding author (Sungroh Yoon)

lion ECGs are obtained worldwide each year [1]. Standard ECGs, which are used to diagnose heart diseases, consist of 12 leads. However, it is not always possible to obtain all 12 leads due to the cost and limitations of measurement devices. Recently, it has been demonstrated that a subset of 12 leads also contains sufficiently meaningful information [2].

In last year, we developed a model to classify clinical cardiac abnormalities from 12-lead ECG classification [3]. Although our model showed a high challenge score on the validation dataset, it showed a much lower score on the hidden test dataset due to the lack of dataset-wise generalization performance. The characteristics of ECG vary from dataset to dataset for various reasons, i.e., hospital, race, etc. It is important to design a model to have generalization performance on dataset unseen during training. Therefore, it is necessary to check whether the proposed model shows high generalization performance consistently over various datasets.

In this paper, as part of the PhysioNet / Computing in Cardiology Challenge 2021, we present a model developed to classify cardiac abnormalities from 12 and reduced-lead ECGs [4–6]. To get the present model having high generalization performance, we have attempted various techniques and evaluated them using the leave-one-dataset-out cross-validation for model selection. The present model achieves 0.1 higher challenge score in average over all test set versions than the model we submitted last year [3].

2. Methods

2.1. Data

Table 1 shows the statistics of the data provided by the challenge with 26 scored SNOMED-CT labels [13] from

Dataset	Number of recordings	w/ Scored Labels	Average Length (second)
Ningbo[7]	34,905	34,485	10
PTB-XL[8]	21,837	21,604	10
Chapman[9]	10,247	9,710	10
G12EC[5]	10,344	9,458	9
CPSC[10]	6,877	5,279	15
CPSC-Extra[10]	3,453	1,278	15
PTB[11]	516	97	110
INCART[12]	74	33	1,800

Table 1: Data statistics

eight datasets [6]. Among them, PTB and INCART data are not used for training because of the long lengths and relatively small number of samples. We also do not use those without any positive scored labels for training. When training the model, the ratio of train and validation datasets is 9:1. In the leave-one-dataset-out cross-validation setting, one of the six datasets is used as the test dataset, and the remaining five datasets are used to train and validation datasets.

We apply the following data pre-processing procedures. First, we upsample or downsample ECGs into 500Hz. Then, we apply a Finite Impulse Response bandpass filter with a bandwidth of 3 to 45Hz. Normalization is applied using the minimum and maximum values of each sample. Finally, for any recording with a data length longer than 7,500, we use a randomly selected segment with a length of 7,500 as input. If the length is shorter than 7,500, we use zero-padding to make the length to be 7,500. For reduced-lead model training, pre-defined leads are extracted from the 12-lead sample [6].

2.2. Model Architecture

For the baseline model, we use our previous work [3]. We use the WRN model architecture with 14 convolution/dense layers and widening factor 1 [14]. The overall structure of the model is shown in Figure 1. The additional parts from the baseline are depicted in purple. The baseline model consists of the basic residual block, but we use the Squeeze and Excitation (SE) block to let the model learn interdependency between channels [15]. For the model to consider the demographic information, we add additional features to the dense layer of the output stem.

2.3. Training

First, we describe the experiment settings. Each model is trained for 100 epochs using Pytorch with an NVIDIA GeForce RTX 3080 [16]. We use Adam optimizer, L2 weight decay of 0.0005, a dropout rate of 0.3, a batch size of 128, and a learning rate of 0.001 through hyper-

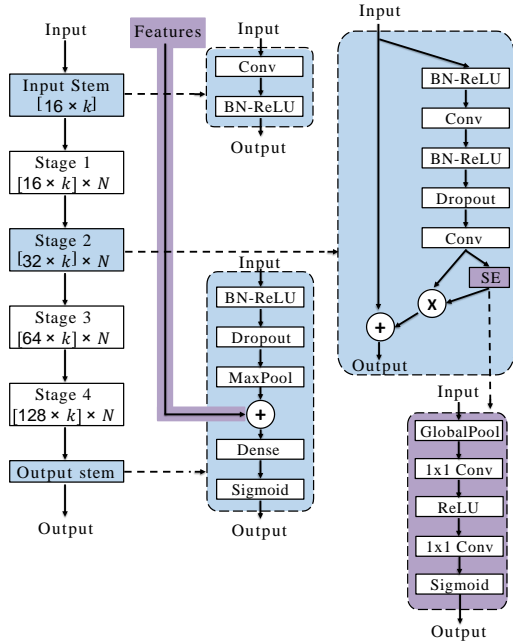


Figure 1: Model overview.

parameter search. In the next part, we explain the training refinements to improve dataset-wise generalization.

Constant-weighted binary cross-entropy loss

In last year, we adopted confusion-weighted binary-cross-entropy (CoW-BCE) [3] loss designed to resemble an evaluation metric called challenge score [6]. Although the model trained via CoW-BCE loss showed a high challenge score on the validation dataset, it showed a much lower score on the hidden test dataset.

In this work, we use constant-weighted binary-cross-entropy inspired via asymmetric loss (ASL) [17]. To overcome the inherent positive-negative imbalance in typical multi-label classification problems, the ASL uses asymmetric focusing and asymmetric probability shifting as follows:

$$\text{ASL} = \begin{cases} -(1-p)^{\gamma^+} \log(p), & \text{if } y \text{ is } 1 \\ -(p_m)^{\gamma^-} \log(1-p_m), & \text{otherwise} \end{cases} \quad (1)$$

where p is the output probability of the model, p_m is the shifted probability, and γ^+ , γ^- are positive and negative focusing parameters, respectively.

For ease of implementation, we assume the positive focusing parameter γ^+ to be 0. We investigate the constant value of the negative coefficient, which depends on the optimal negative focusing parameters γ^- and shifted probability p_m . Experimentally, we set the negative coefficient to be 0.1, which is approximately the ratio of positive to negative classes in the whole dataset.

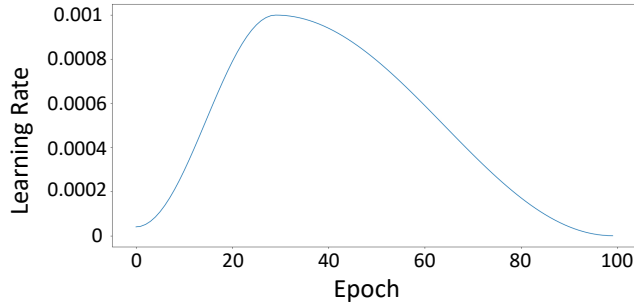


Figure 2: OneCycle Learning rate Scheduler.

Demographic features

For the model to consider demographic information, we additionally use two kinds of features, *i.e.*, age and gender. The demographic feature vector consists of 5 values for age, one-hot encoded sex, and two flags for missing values. If there are age and gender values in the header, the values are used directly, and missing flags are set to 0. Otherwise, pre-defined default values are used, and the missing flags are set to 1. The default age value is 60.37, and the default gender value (female/male ratio) is 0.471/0.519. As shown in the purple path in Figure 1, the feature vector is concatenated with the feature extracted by Deep Neural Network(DNN) before the last dense layer.

Mixup

Mixup is one of the data augmentation techniques for better generalization [18]. It makes the decision boundary smoother by regularizing the model. Assuming that two arbitrary input signals in the batch are x_1, x_2 , the features of the samples are f_1, f_2 , and the labels are l_1, l_2 , the mixup samples x', f', l' are created as follows.

$$x' = \lambda x_1 + (1 - \lambda)x_2 \quad (2)$$

$$f' = \lambda f_1 + (1 - \lambda)f_2 \quad (3)$$

$$l' = \lambda l_1 + (1 - \lambda)l_2 \quad (4)$$

As used in the original mixup paper, mixing coefficient λ is sampled from a Beta(0.2,0.2) distribution. The model is trained using the generated x', f' , and l' .

Learning rate scheduler

We use the OneCycle learning rate scheduler [19]. It is known as a method for effective training by “super-convergence” of residual blocks. At the beginning of training, the learning rate is set to a small value, and it is gradually increased and then decreased again after reaching the pre-defined maximum value. The learning rate values per epoch are shown in Figure 2. The maximum learning rate value is set to 0.001. The model is trained for a total of 100 epochs using a cosine annealing strategy.

Leads	Training	Validation	Test	Ranking
12	0.654	0.610	0.550	2nd
6	0.680	0.580	0.580	1st
4	0.691	0.600	0.580	1st
3	0.689	0.590	0.570	2nd
2	0.673	0.590	0.570	2nd

Table 2: Challenge scores for our model using whole six datasets.

3. Experiments results

The experiment results of the present model trained using the whole six datasets are shown in Table 2. We report the training, validation, and test challenge score, and team ranking for our proposed 12-lead, 6-lead, 4-lead, 3-lead, and 2-lead models. The average validation challenge score is 0.594, and the average test challenge score is 0.570. As a result of adding up the scores up to the test set, our model records 2nd, 1st, 1st, 2nd, and 2nd place on 12-lead, 6-lead, 4-lead, 3-lead, and 2-lead, respectively.

In Table 3, we compare the performance of the baseline and our 12-lead model. We show the results from the leave-one-dataset-out cross-validation setting (3-9 column) and using the all six datasets (10 column). We report the challenge scores when the dataset in the first row is used as a test dataset. The last column shows the challenge scores of 12-lead models trained and tested using six whole datasets with 9:1 ratio. The challenge score from using six whole datasets obtained by our model is 0.654, which is 0.08 lower than the baseline. Although the present model obtains a lower challenge score when trained using the whole six datasets, the dataset-wise generalization performance becomes better compared to the baseline. The average dataset-wise challenge score of our proposed model is 0.483, which is 0.1 higher than the baseline. The usage of a constant-weighted binary cross-entropy loss instead of CoW-BCE loss function makes the most of the improvement in the dataset-wise generalization performance. In particular, the changed loss function improves the generalization performance for the PTB-XL dataset.

4. Concluding Remarks

In this paper, as a participating team in the PhysioNet Challenge 2021, we proposed 12 and reduced-lead models for automatically classifying cardiac abnormalities from ECGs. We focused on building the classification model to have a high dataset-wise generalization performance. We used the SE-WRN-14-1 network with constant binary cross-entropy loss, feature extraction, mixup, and OneCycle learning rate scheduler by evaluating with the leave-one-dataset-out cross-validation setting. The average score

Test dataset	Ningbo	PTB-XL	Chapman	G12EC	CPSC	CPSC-Extra	Average	All	
Model	Baseline	0.545	-0.101	0.659	0.428	0.463	0.310	0.384	0.732
	Our model	0.626	0.200	0.723	0.519	0.506	0.424	0.483	0.654

Table 3: Comparison of the performance between the baseline and our model in leave-one-dataset-out cross-validation setting (3-9 columns) and using all six dataset (10 column))

of our proposed model was 0.1 higher than the baseline. Again this year, although the ranking of the validation set was not very good (15th), but the ranking on the hidden test set was high (2nd) due to the increase in the dataset-wise generalization performance.

Acknowledgments

This work was supported by the National Research Foundation of Korea(NRF) grant funded by the Korea government (Ministry of Science and ICT, 2018R1A2B3001628), the BK21 FOUR program of the Education and Research Program for Future ICT Pioneers, Seoul National University in 2021, and Samsung Electronics(DS and Foundry).

References

- [1] Holst H, Ohlsson M, Peterson C, Edenbrandt L. A confident decision support system for interpreting electrocardiograms. *Clinical Physiology* 1999;19(5):410–418.
- [2] Drew BJ, Pelter MM, Brodnick DE, Yadav AV, Dempel D, Adams MG. Comparison of a new reduced lead set ecg with the standard ecg for diagnosing cardiac arrhythmias and myocardial ischemia. *Journal of electrocardiology* 2002; 35(4):13–21.
- [3] Min S, Choi HS, Han H, Seo M, Kim JK, Park J, et al. Bag of tricks for electrocardiogram classification with deep neural networks. In *2020 Computing in Cardiology. IEEE*, 2020; 1–4.
- [4] Goldberger AL, Amaral LA, Glass L, Hausdorff JM, Ivanov PC, Mark RG, et al. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *Circulation* 2000;101(23):e215–e220.
- [5] Perez Alday EA, Gu A, Shah A, Robichaux C, Wong AKI, Liu C, et al. Classification of 12-lead ECGs: the PhysioNet/Computing in Cardiology challenge 2020. *Physiological Measurement* 2020;41.
- [6] Reyna MA, Sadr N, Perez Alday EA, Gu A, Shah A, Robichaux C, et al. Will Two Do? Varying dimensions in electrocardiography: the PhysioNet/Computing in Cardiology challenge 2021. *Computing in Cardiology* 2021;48:1–4.
- [7] Zheng J, Cui H, Struppa D, Zhang J, Yacoub SM, El-Askary H, et al. Optimal multi-stage arrhythmia classification approach. *Scientific Data* 2020;10(2898):1–17.
- [8] Wagner P, Strodthoff N, Bousseljot RD, Kreiseler D, Lunze FI, Samek W, et al. PTB-XL, a large publicly available electrocardiography dataset. *Scientific Data* 2020;7(1):1–15.
- [9] Zheng J, Zhang J, Danioko S, Yao H, Guo H, Rakovski C. A 12-lead electrocardiogram database for arrhythmia research covering more than 10,000 patients. *Scientific Data* 2020; 7(48):1–8.
- [10] Liu F, Liu C, Zhao L, Zhang X, Wu X, Xu X, et al. An open access database for evaluating the algorithms of electrocardiogram rhythm and morphology abnormality detection. *Journal of Medical Imaging and Health Informatics* 2018;8(7):1368–1373.
- [11] Bousseljot R, Kreiseler D, Schnabel A. Nutzung der ekg-signaldatenbank CARDIODAT der PTB über das Internet. *Biomedizinische Technik* 1995;40(S1):317–318.
- [12] Tihonenko V, Khaustov A, Ivanov S, Rivin A, Yakushenko E. St Petersburg INCART 12-lead arrhythmia database. *PhysioBank PhysioToolkit and PhysioNet* 2008;Doi: 10.13026/C2V88N.
- [13] Shahpori R, Doig C. Systematized nomenclature of medicine—clinical terms direction and its implications on critical care. *Journal of critical care* 2010;25(2):364–e1.
- [14] Zagoruyko S, Komodakis N. Wide residual networks. *arXiv preprint arXiv160507146* 2016;.
- [15] Hu J, Shen L, Sun G. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018; 7132–7141.
- [16] Paszke A, Gross S, Chintala S, et al. Automatic differentiation in PyTorch. *NIPS Autodiff Workshop* 2017;.
- [17] Ben-Baruch E, Ridnik T, Zamir N, Noy A, Friedman I, Protter M, et al. Asymmetric loss for multi-label classification, 2020.
- [18] Zhang H, Cisse M, Dauphin YN, Lopez-Paz D. Mixup: beyond empirical risk minimization. *arXiv preprint arXiv171009412* 2017;.
- [19] Smith LN, Topin N. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006. International Society for Optics and Photonics, 2019; 1100612.

Address for correspondence:

Sungroh Yoon
Rm. 908, Bldg. 301, 1 Gwanak-ro, Gwanak-gu, Seoul 08826, South Korea
sryoon@snu.ac.kr