

Label noise handling on MNIST with neural networks ¹

István Fazekas^a, Attila Barta^b, László Fórián^b

^aUniversity of Debrecen, Department of Applied Mathematics and Probability Theory

^bUniversity of Debrecen, Doctoral School of Informatics

fazekas.istvan@inf.unideb.hu, barta.attila@inf.unideb.hu,
forian.laszlo@inf.unideb.hu

In recent years, deep neural networks have reached very impressive performance in the task of image classification. However, these models require very large datasets with labeled training examples, and such datasets are not always available. The labeling process is often very expensive, or it is very difficult even for experts in a particular field. That is what can lead to the use of databases with label noise, which contain incorrectly labeled instances. Therefore, it is important to examine training on this kind of datasets. According to a widely accepted assumption, deep networks learn consistent, simple patterns in the beginning, and then it is followed by the learning of the harder examples with possibly incorrect labels. So treating the label noise in the train set can lead to a better generalization ability instead of overfitting to the wrong examples.

In this work, we investigate the possibilities of improving a classifier (which is an ensemble of deep neural networks) by handling the label noise in the training dataset. We classify with an ensemble of convolutional neural networks (CNNs). At the start, we train that ensemble on the original training dataset. Then we are going to apply a label noise cleansing technique on that data. Finally, we take a CNN ensemble with the same structure as our original CNN ensemble, and train it on the new dataset gained after treating the label noise. We evaluate and compare the performance of the ensemble classifiers and draw conclusions. Our goal is to study label correcting neural networks for preprocessing purpose. Preprocessing can be either relabeling or deleting items detected to have noisy labels. After preprocessing, usual CNNs are applied for the data. With preprocessing, the performance of very accurate convolutional networks can be further improved.

We conduct experiments on the MNIST dataset [2], which contains handwritten digits. It consists of images with 28×28 grayscale pixels. The size of the training set is 60 000 examples and the test set has 10 000 samples. However, it contains some misleading items. We shall consider these misleading instances as inaccurately labeled ones so we can apply some known methods elaborated to handle noisy labels.

References

- [1] I. Fazekas, A. Barta, L. Fórián, Ensemble noisy label detection on MNIST, *Annales Mathematicae et Informaticae* **53** (2021): Selected papers of the 1st Conference on Information Technology and Data Science, 125-137.
- [2] Y. LeCun, C. Cortes, C. Burges, MNIST handwritten digit database, ATT Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist>
- [3] D. Tanaka, D. Ikami, T. Yamasaki, K. Aizawa, Joint Optimization Framework for Learning With Noisy Labels, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, 5552-5560.
- [4] K. Yi, J. Wu, Probabilistic End-To-End Noise Correction for Learning With Noisy Labels, *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, 7017-7025.

¹This work was supported by the construction EFOP-3.6.3-VEKOP-16-2017-00002. The project was supported by the European Union, co-financed by the European Social Fund.