

Fuzzy categorization of news articles using improved neural network

Tibor Tajti^a, Tamás Balla^a, Csaba Biró^a

^aEszterházy Károly Catholic University
tajti.tibor@uni-eszterhazy.hu
balla.tamas@uni-eszterhazy.hu
biro.csaba@uni-eszterhazy.hu

Abstract

The use of binary classification in machine learning has been widely adopted for its simplicity and ease of implementation. However, in many real-world applications, the data and class labels are often ambiguous and may belong to multiple categories simultaneously. This has led to the exploration of alternative approaches, such as fuzzy classification, which can handle the uncertainty and vagueness inherent in such data. Fuzzy logic and fuzzy sets are known to be more powerful compared to binary logic, and for extreme values of 0 and 1, they are compatible with binary logic. [2, 7]

There are many publicly available datasets for natural language texts, such as reviews and news articles, where the samples are labeled into categories. It is usual for these labels to be handled as binary class membership values, but this can be inaccurate in some cases. For instance, a news article may be labeled as "sport", even if it can be considered as belonging to both "sport" and "business" categories.

In this paper we aim to research the effectiveness of updating the binary class membership values of the training data according to the actual knowledge of the model during the training. Similar experiments have been conducted using the MNIST dataset of handwritten digits, where single labels were given for the images, referring to one digit's class [3]. The paper introduces a simple algorithm for the fuzzification of the class membership values of the training data. The incorporation of fuzzy values in the training process of machine learning algorithms has the potential to improve the accuracy and robustness of the models, as well as provide a more nuanced understanding of the relationships between the features and class

labels.

In our current research we conduct experiments using similar fuzzification on a dataset of natural language texts. We also make improvements on the algorithm to better adapt the algorithm parameters to the actual dataset and to the learning model. We start with a low value of the fuzzification rate parameters at the beginning of the training process to be cautious with the fuzzification, and gradually increase them as the training progresses and the fitness of the model improves. This will allow the learners to gather useful knowledge before they start making significant changes to the class membership values.

The proposed algorithm can also be used as an ensemble method. Ensemble methods, including the committee machines, have proven their effectiveness in many fields [6]. Committee machine methods produce predictions and aggregate the findings of several instances of neural networks or other machine learning algorithms. They frequently produce better than average or even better than the best individual result. They can be used in many areas, including multiple measurements [5], and, in machine learning, using multiple learners [1, 4].

In this paper we evaluate our algorithm using individual learners and also using ensemble learners. We show that the fuzzification of the class membership values of the training data during training may improve the accuracy of the machine learning algorithm. We also show that beside the performance improvement the model gains a finer tuned output of the relationships between the features and class labels.

References

- [1] G. FUMERA, F. ROLI: *A theoretical and experimental analysis of linear combiners for multiple classifier systems*, en, IEEE Transactions on Pattern Analysis and Machine Intelligence 27.6 (2005), pp. 942–956.
- [2] J. ROUBOS: *Learning fuzzy classification rules from labeled data*, Information Sciences 150.1-2 (Mar. 2003), pp. 77–93, DOI: [10.1016/s0020-0255\(02\)00369-9](https://doi.org/10.1016/s0020-0255(02)00369-9), URL: [https://doi.org/10.1016/s0020-0255\(02\)00369-9](https://doi.org/10.1016/s0020-0255(02)00369-9).
- [3] T. TAJTI: *Fuzzification of training data class membership binary values for neural network algorithms*, Annales Mathematicae et Informaticae 52 (2020), DOI: [10.33039/ami.2020.10.001](https://doi.org/10.33039/ami.2020.10.001), URL: <https://doi.org/10.33039/ami.2020.10.001>.
- [4] T. TAJTI: *New voting functions for neural network algorithms*, Annales Mathematicae et Informaticae 52 (2020), DOI: [10.33039/ami.2020.10.003](https://doi.org/10.33039/ami.2020.10.003), URL: <https://doi.org/10.33039/ami.2020.10.003>.
- [5] T. TAJTI, N. BENEDEK: *Motion sensor data correction using multiple sensors and multiple measurements*, in: 2016 IEEE 14th International Symposium on Applied Machine Intelligence and Informatics (SAMI), IEEE, Jan. 2016, DOI: [10.1109/sami.2016.7423022](https://doi.org/10.1109/sami.2016.7423022), URL: <https://doi.org/10.1109/sami.2016.7423022>.
- [6] V. TRESP: *Committee machines*, en, in: Handbook for neural network signal processing, 2001, pp. 1–18.
- [7] L. ZADEH: *The concept of a linguistic variable and its application to approximate reasoning—I*, Information Sciences 8.3 (1975), pp. 199–249, DOI: [10.1016/0020-0255\(75\)90036-5](https://doi.org/10.1016/0020-0255(75)90036-5), URL: [https://doi.org/10.1016/0020-0255\(75\)90036-5](https://doi.org/10.1016/0020-0255(75)90036-5).