

# Counteracting Learnability Imbalance in Multiclass Classification

Theodor Peifer

email: thp7219@thi.de

Technische Hochschule Ingolstadt

**Abstract**—Neural Networks have proven themselves to be powerful classification models by solving problems in a range of domains with high accuracy. Yet, this accuracy is never evenly distributed across all classes, which means that the true positive rates of each class differ. Although this problem is most common and severe when working with class imbalanced datasets, it can also occur with balanced datasets. This is because there are always classes that are more difficult for the model to learn than others (this phenomenon is further referred to as *learnability-imbalance*). This research will examine whether and how the methods to counteract class imbalance can be adapted to the more general problem of learnability imbalance. Since those methods use the different class proportions as a measure of how much they should affect each class, this research will propose to use the true positive rates of a model previously trained on the same dataset instead.

## I. INTRODUCTION

A frequent problem in classification appears when working with datasets that have an unequal number of samples per class and are therefore called a *class imbalanced* datasets [1]. Since there will be some classes, that have fewer elements for the model to learn from, their features will be harder to extract what finally will result in a lower true positive rate, i.e. a per-class accuracy. Thus, the consequence of having different class sizes can be described as having a *learnability imbalance* in the dataset, since some classes are more difficult to learn than others.

However, such learnability differences can appear also in balanced datasets for a variety of reasons, e.g. when the quality of the data of a class is lower than the rest of the data. A second reason, that this research will be focusing on, is that when some classes are similar, the model can confuse their samples with each other more easily which will often result in a lower accuracy of those classes.

Even though this issue is an inevitable product of every normal classification, in most cases the learnability difference of the classes is either low or not of great interest. But there are more extreme cases in which a model needs to produce fair and unbiased results, even though it has to be trained on a dataset that has an obvious learnability imbalance. An example for that is *name-ethnicity classification*, where a model predicts the ethnicity of a name only by its letters [2]. Nationalities that use the same language and therefore have similar names (e.g. *British* and *American*) result in a lower accuracy (see figure 1). Thus, when a model trained on such a dataset is used, e.g. for social science experiments, it can lead to unfair results and wrong interpretations.

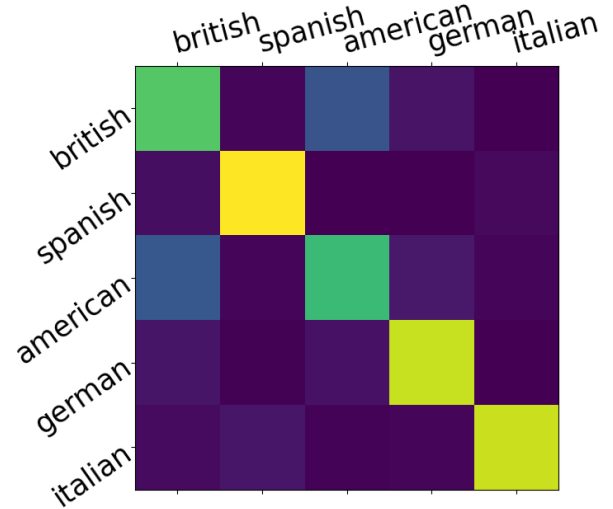


Fig. 1: A confusion matrix representing the true positive distribution produced by a recurrent neural network that was trained on a name-nationality dataset. It shows a strong confusion of British and American names and therefore a lower accuracy of both classes.

Another dataset that is a good showcase for the learnability imbalance is the CIFAR-10 [3] dataset (32×32 RGB images of ten different classes) because it also contains similar classes such as *dog* and *cat*. The CIFAR-10, a name-nationality and a further tabular dataset will be used for experiments in order to find methods for minimizing the variance in the evaluation scores.

To develop such methods, one should first resort to the solutions of the class imbalance problem. One of them is to weight the error function [4] according to the size of each class, i.e. the number of samples it contains. Therefore, there is a weight assigned to every class which is greater the fewer elements it contains and which is multiplied by the loss produced by its samples. This has the effect that samples from a smaller class will produce a higher error and, since the aim of a neural network is to minimize the error [5], will receive a higher focus in the learning process. Other methods will be discussed in the literature review.

Almost all the approaches to counteract class imbalance rely on the different proportions of the classes in the dataset. This raises the question if and how they are applicable on datasets that, even though they are not affected by class imbalance, still show a big learnability imbalance. When working with such

datasets the learnability differences of the individual classes are mostly only identifiable after the model has been trained normally, i.e. without applying any methods. The calculation of the true positive rates and a visualization (e.g. confusion matrix) can then reveal what classes were learned best and which should have received a higher weight. These evaluation scores can then be used, for example, to determine which loss weights should be used for each class. Then the model can be trained again using a weighted loss (as shown in *algorithm 1*).

---

**Algorithm 1** creating loss weights for a balanced dataset

---

$C \triangleq$  amount of classes

*train(weights)* describes the initialization and training process of a classification model in which  $weights_i$  will be multiplied to every loss generated by a sample of class  $i$ .

*evaluate()* creates the set  $s$  with  $|s| = C$  and  $s_i \in [0; 1]$  which contains the true positive scores of all classes.

$W(s)$  is a function that creates a set of loss-weights  $w$  with  $|w| = C$  and  $w_i \in \mathbb{R}^+$  using a set of true positive scores  $s$ :  $W : [0; 1] \rightarrow \mathbb{R}^+ ; \{s_1, \dots, s_C\} \rightarrow \{w_1, \dots, w_C\}$ .

**process:**

- 1:  $train(weights=\{w_1=1, w_2=1, \dots, w_C=1\})$
  - 2:  $s = evaluate()$
  - 3:  $w = W(s)$
  - 4:  $train(weights=w)$
  - 5:  $s' = evaluate()$
  - 6:  $compare(s, s')$
- 

In summary, this research will focus on adapting the methods that are used to counteract class imbalance to the learnability imbalance problem of balanced datasets. This will be done using the true positive rates produced by the evaluation of an independently, previously trained model. Instead of class sizes, these rates will be used as a measure for how much a method should affect each class.

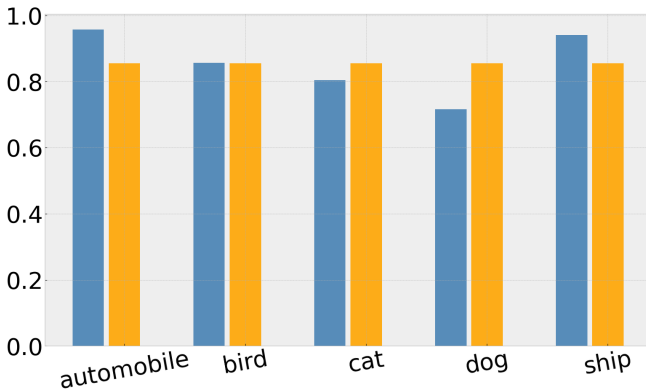


Fig. 2: True positive rates produced by a model that was trained normally (blue) on five CIFAR-10 classes compared with the desired true positive rates (orange).

The goal of this research is to find approaches that bring all per-class accuracies to nearly the same value by increasing those of the harder-to-learn classes and decreasing those of the easier-to-learn classes without causing a reduction of the overall

accuracy (see figure 2). Further, it will address limitations and risks that could potentially arise from this process.

## II. LITERATURE REVIEW

As mentioned before, the learnability imbalance has to be handled the same way as the class imbalanced dataset problem, except the true positive rates are used instead of the class proportions of the dataset. Therefore, the proposed solutions to class imbalance are fundamental to this research and will be reviewed in the following.

Since it was used as a main example in the introduction, the first method to examine is loss weighting. This approach has proven itself to be a good way of achieving a lower variance in the per-class accuracies of unbalanced datasets and will be adapted to this research. In general, the weights are chosen to be inversely proportional to the amount of samples in the classes, in order to weight smaller classes more than bigger ones. An example for a weight calculation function, proposed by Yin Cui et al. [4], is shown in formula 1 and 2:

$$w_c = \frac{1 - \beta^{N_c}}{1 - \beta}, \quad (1)$$

further normalization:

$$w_c = \frac{C * w_c}{\sum_{i=1}^C w_i}. \quad (2)$$

This method is described as the *effective number of samples* weighting, where  $N_c$  is the amount of samples of class  $c$  and  $\beta$  is a tuneable hyper-parameter with  $\beta \in [0; 1[$ .

Formula 3 shows how to use such a weight along with the cross entropy loss function for one sample:

$$CE(x, c) = w_c \cdot (-\log(\hat{y}_c)) , \quad (3)$$

where  $\hat{y}$  is the output of a model  $F(x, \theta)$  with input sample  $x$  and learnable parameter  $\theta$ . It can be described as a probability distribution for which  $\hat{y}_i$  is the confidence of the model, that the input  $x$  corresponds to the  $i$ -th class. Therefore,  $\hat{y}_c$  represents the probability that the model correctly classified  $x$  as the wanted class  $c$ . As shown, the weight  $w_c$ , that corresponds to the correct class, is multiplied by the normal error value (in the case of cross entropy: the negative log of the probability).

Another loss weighting method is the so-called *focal loss* and was introduced by Tsung-Yi Lin et al. [6] in order to handle class imbalance in object detection. This approach is especially interesting for this research because it already uses the true positive rates instead of the sizes of the classes to generate the weights. Thus, it can be applied unchanged to the learnability imbalance problem. Since it is also mainly tested and practiced with class imbalanced datasets however, this research will investigate its effectiveness on class balanced datasets that have a high learnability variance. The weight for the focal loss is defined as the following ( $\gamma$  is a positive hyperparameter for scaling):

$$w_c = (1 - \hat{y}_c)^\gamma . \quad (4)$$

When multiplying this weight by the loss, it will be bigger the smaller the confidence of the correct class was. The advantage of the focal loss over the proposed weight generation of this research is that it does not require a pre-training in order to figure out which classes are harder to learn, since it does this during the training process. Though, it can be hypothesized that this weighting could have less impact on learnability imbalanced datasets because it figures out the imbalance while training. In contrast to this, when training unweighted first, calculating the weights afterwards and then train again with those weights, the loss function can start penalizing the classes differently right from the beginning.

Another method is *augmentation*, in which the input samples get randomly modified in order to synthetically generate more training samples. This can help the model to generalize better. For example, in image classification it is common to flip, crop or add white noise to the images. The idea, originally proposed as *SMOTE* (Synthetic Minority Over-sampling Technique) by Nitesh V. Chawla et al. [7], is that such transformations can be used on smaller classes to provide the model with more of their samples. The same technique can also be used on learnability imbalanced datasets but by applying it on classes with a lower per-class accuracy.

Another way to try to compensate for class imbalance is by doing *oversampling* [7], where the model is fed samples of the minority classes more often. This method should be used if the data is not augmentable and therefore *SMOTE* can't be applied.

From the approaches mentioned above, the focal loss is the only one that can be tested on balanced datasets as it is. The others require a pre-training in order to determine the learnability differences and therefore the true positive rates, which then can be used to calculate how much these methods should affect each class.

### III. METHODOLOGY

In order to produce more interpretable results, the proposed approaches have to be applied on different models with different tasks. There are three main disciplines: image, sequential data and tabular data classification. When testing an approach on each task, the results can be interpreted more generalized. But it can also reveal if any methods work better for one specific discipline than for others. If there is such a suspicion, the methods will be tested on even more datasets.

As already mentioned in the introduction, to test the approaches on image classification, the CIFAR-10 dataset will be used. It contains 60,000 32×32 pixel RGB (i.e. channels for red, green and blue color) images, that are equally distributed among ten classes. This dataset is freely available and commonly used in computer vision research. The architecture of the model that will be trained on this dataset will be a convolutional neural network [8] with residual connections [9]. For testing the approach on sequential data, a private name-ethnicity dataset will be used. This dataset was gathered by the government of Great Britain and provided to the author for a prior project based on name-ethnicity-classification. Within this project, the maximum accuracy for this classification

problem was reached by using a *ConvLSTM* architecture, i.e. the combination of a convolutional neural network and a LSTM (long-short term memory) [10], which will also be used in this research. Finally, the tabular dataset, which will be used, has yet to be chosen.

The functions, which will produce loss weights or oversampling chances, are created similarly to the *effective number of samples* loss weighting technique. For a class with a lower true positive score, the functions should generate a larger value. The following formula (5) shows a possible function that this research will use to generate loss weights ( $\gamma$  is a positive hyperparameter for scaling):

$$w_c = \left(1 - \left(\frac{s_c - \bar{s}}{s_c}\right)\right)^\gamma. \quad (5)$$

If the true positive score  $s_c$  of class  $c$  is below the average score  $\bar{s}$ , the corresponding weight is greater than one, otherwise it is less than one.

The metric used for measuring how much a potential learnability imbalance impacted the model is the variance of the per-class accuracies. Therefore, the effectiveness of the proposed approaches can be determined by the difference between the variances generated with and without the approaches:

$$\text{Var}(s) = \frac{1}{C} \sum_{i=1}^C (s_i - \bar{s})^2, \quad (6)$$

$$D = \text{Var}(s) - \text{Var}(s'), \quad (7)$$

where  $s'$  are the true positive scores produced by a model using one of the proposed methods and  $s$  of a model that was trained normally. If  $\text{Var}(s')$  is smaller than  $\text{Var}(s)$  and  $D$  is therefore positive, an approach showed a potential effectiveness. But if the difference is not significant, a statistical test must be performed in order to investigate if the reason for this difference was indeed the used method.

The experiments will be implemented in the *Python* programming language [11] and the models will be built and trained using the *Pytorch* deep learning framework [12]. In addition, an experiment manager, such as *Weights&Biases* [13], will keep the results in an orderly and accessible manner.

### IV. RESEARCH PLAN

This research is expected to be completed in a little over a year and will not consume many resources. The reason for that is that the capacities, i.e. number of parameters, of the models used for the experiments are not high, since they do not need to reach state-of-the-art results. Additionally, the datasets (which are already provided) are not particularly high dimensional and consist of a manageable number of samples. Therefore, the experiments do not require high-end hardware and can be run in a short amount of time.

TABLE I: schedule of research

stage	description	duration
1	investigation	4 months
2	execution	5 months
3	revision	2 months
4	writing	3 months

The estimated duration of this research is 14 months and will be performed in four stages, as shown in the table above.

In the first stage, the investigation, theoretical approaches are conceptualized by doing a more in depth literature review and adapting the proposed methods to counteract class imbalance to the learnability imbalance problem. The aim is to find several candidates for functions that produce loss weights, augmentation and oversampling chances based on the true positive scores. This investigation also gives the opportunity to encounter previously unconsidered methods with which the research can be extended.

The execution stage will begin by setting up the project code repository (link [14]). It will contain the model architectures, train scripts and implementations of the conceived methods. Then the experiments will be run, the results collected and finally analyzed.

In the next stage, the revision, the used approaches and the results will be reviewed and compared. After this process is finished, the last stage, the writing of the research paper can begin. Presumably after three months, the paper can be submitted.

#### REFERENCES

- [1] N. Japkowicz and S. Stephen, "The class imbalance problem: A systematic study," *Intelligent data analysis*, vol. 6, no. 5, pp. 429–449, 2002.
- [2] P. Mateos, "A review of name-based ethnicity classification methods and their potential in population studies," *Population, Space and Place*, vol. 13, no. 4, pp. 243–263, 2007.
- [3] A. Krizhevsky and G. Hinton. "Learning multiple layers of features from tiny images." (2009).
- [4] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9268–9277.
- [5] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," California Univ San Diego La Jolla Inst for Cognitive Science, Tech. Rep., 1985.
- [6] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [7] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: Synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [8] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [10] F. A. Gers, J. Schmidhuber, and F. Cummins, "Learning to forget: Continual prediction with lstm," *Neural computation*, vol. 12, no. 10, pp. 2451–2471, 2000.
- [11] G. Van Rossum and F. L. Drake Jr, "Python reference manual." Centrum voor Wiskunde en Informatica Amsterdam, 1995.
- [12] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *arXiv preprint arXiv:1912.01703*, 2019.
- [13] W. bibinitperiod Biases. "Weights & biases: Tools for ml - build better ml models faster." (2018), [Online]. Available: <https://wandb.ai/site> (visited on 09/23/2021).
- [14] T. Peifer. "Thi-research-proposal." (2021), [Online]. Available: <https://github.com/hollowcodes/thi-research-proposal> (visited on 09/23/2021).