

FINE-TUNE YOUR CLASSIFIER: FINDING CORRELATIONS WITH TEMPERATURE

Benjamin Chamand^{1*} Olivier Risser-Maroux^{2*†}
Camille Kurtz² Philippe Joly¹ Nicolas Loménie²

¹IRIT, Université de Toulouse, CNRS, Toulouse INP, UT3, Toulouse, France

²LIPADE, Université de Paris, France

benjamin.chamand@irit.fr, orissermaroux@gmail.com

ABSTRACT

Temperature is a widely used hyperparameter in various tasks involving neural networks, such as classification or metric learning, whose choice can have a direct impact on the model performance. Most of existing works select its value using hyperparameter optimization methods requiring several runs to find the optimal value. We propose to analyze the impact of temperature on classification tasks by describing a dataset as a set of statistics computed on representations on which we can build a heuristic giving us a default value of temperature. We study the correlation between these extracted statistics and the observed optimal temperatures. This preliminary study on more than a hundred combinations of different datasets and features extractors highlights promising results towards the construction of a general heuristic for temperature.

Index Terms— temperature, hyperparameter, heuristic, softmax, cross-entropy

1. INTRODUCTION

The performance of a machine learning algorithm applied to a computer vision task is highly dependent on the choice of its hyperparameters. Among these, the temperature is a scaling factor often used in a neural network linked to the softmax layers, the latter being usually followed by a cross-entropy (CE) like loss function. Intuitively, the temperature (in allusion to statistical mechanics) is introduced to choose the level of uniformity of the distribution. Since most deep classification models involve both softmax layer and CE like loss functions for their training, determining an optimal temperature for a particular task can then have a broad impact.

For example, this parameter is widely considered in various tasks such as knowledge distillation, classification, text generation, self-supervised and metric learning [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11] Traditionally, in most of these domains and in the underlying applications, the temperature is determined empirically, with a value that can be constant (typically from

a grid search) or evolve dynamically over iterations, in the same vein as the learning rate parameter. Nevertheless, such strategies for determining a *good* temperature may be suboptimal or computationally too cumbersome. Surprisingly, there are very few studies proposing strategies for determining an optimal temperature. In this paper, we focus on the particular problem that, given a classification task, we need to find a correlation between an optimal value for the temperature and statistics describing the dataset such as complexity, dimension, number of classes, etc.

2. RELATED WORKS

The temperature hyperparameter is typically employed in the softmax layer to control the uniformity of the distribution. Although the use of a good temperature has shown its impact in many computer vision tasks, the existing strategies to define such a temperature parameter are quite limited.

The first way to proceed is to consider a constant temperature throughout the training. The choice can be done empirically, as in [1, 8, 9]. It can also be considered as a fixed hyperparameter to be optimized via a grid search in a field of possible values, but this implies significant computational requirements and leads to different hyperparameters for each dataset and architecture. A simple heuristic can also allow to fix the parameter as proposed in the Transformers [12] with \sqrt{d} , d being the dimension of the queries and the key vectors.

Other strategies rely on dynamic temperature adjustment during learning iterations. In this case, the elements of the temperature can evolve at each epoch using a scheduler [4], in the manner of the learning rate to refine the network. In [13], the authors also showed that a batch normalization rescaled by \sqrt{d} , with d the number of dimensions of embeddings, worked slightly better than a simple $L2$ normalization, and can also lead to more embedding vectors. Dynamic adjustment of temperature can also be done by learning it as a standard parameter [14, 15]. This usually requires additional steps like clipping or adding *exp* to avoid negative values. Furthermore, the learned temperature strongly depends on the learning rate hyperparameter.

* Equal contribution.

† Financed by Smiths Detection

An alternative approach is to determine the temperature value analytically. The authors of [16] propose both an evaluation function designed to measure the effectiveness of a temperature parameter and an iterative updating rule to determine the optimal temperature value. However, their work suffers from two drawbacks: (1) Authors introduce a novel hyperparameter λ in the temperature formulation, λ being an improvement factor affecting the number of iterations and the selection of the optimal temperature; (2) It was designed for the *D-armed bandit* problem in reinforcement learning and only tested on synthetic data.

Another work by [17] proposes a theoretical lower bound formulated as a function of the loss value and the number of classes with a loss smaller than ϵ , ϵ supposed to be around $10e - 4$. Interestingly, unlike [12, 13], their solution does not derive any benefit from or rely on any information on the embeddings dimensions. However, since temperature determination was not the main part of their contribution, no benchmark was made to compare the proposed theoretical lower bound with other temperature values. Finally, the assumption on such a low loss value does not correspond to real cases at the beginning of the learning.

While some of the previously mentioned heuristics are based on feature dimensions, others use the number of classes or class separability measures. None of them have been designed specifically for use in a classification task or have been evaluated on this particular hyperparameter to demonstrate the effectiveness of the proposed heuristic. Other heuristics could be derived from other criteria reflecting information such as the difficulty of the dataset to be classified. For example, [18] proposes to estimate the difficulty of classifying datasets from six classes of measures based on information such as feature-based, neighborhood or dimensionality measures. However, most of proposed measures have a complexity at least equal to $O(n^2)$, with n the number of points in the dataset, making these measures difficult to scale up to larger datasets. Similarly, we seek to describe each dataset by a set of statistics computable in a reasonable amount of time. We then propose to determine which variables / statistics are really correlated with the best empirical temperature, in order to propose a simple heuristic based on these dataset statistics.

3. METHODOLOGY

3.1. Rescaling Cross-Entropy with temperature scaling

Inspired by the formulation of [13], we start from the same basic observation as they do. We define a set of N samples labeled $\{(x_1, y_1), \dots, (x_N, y_N)\}$, where $x_i \in \mathbb{R}^d$ is the representation (embedding) of the i -th sample, d being the dimensionality of x_i , and $y_i \in \{1, \dots, C\}$ is the category label of the sample x_i , C being the total number of categories. Let us consider $W = [w_1, \dots, w_C]$ where $w_j \in \mathbb{R}^d$ is the weights associated with the class C_j , we define $z_i = x_i W$

with $i \in \{1, \dots, N\}$. In our case, we focus on the learning of the weights W . In the same vein as [2, 7, 8, 17], we removed the bias term, and we consider the inputs x and weights W L_2 normalized. We optimize the cosine similarity since this choice is both popular in classification and metric learning.

The probability that a sample x belongs to the category $c \in \{1, \dots, C\}$ can be predicted by the softmax function as:

$$p(c|x, \alpha) = \frac{\exp(\alpha z_c)}{\sum_{j=1}^C \exp(\alpha z_j)} \quad (1)$$

To simplify the notation, we note $\alpha = 1/T$ as the inverse of the temperature T to choose the level of uniformity of the softmax output distribution.

Assuming that the ground truth distribution of the training sample is $q(c|x)$, generally encoded in a one-hot vector (which equals 1 if $c = y$ and 0 otherwise), the cross-entropy loss with respect to x is defined as:

$$\mathcal{L}(x, \alpha) = - \sum_{c=1}^C \log(p(c|x, \alpha)) q(c|x) \quad (2)$$

and the gradient with respect to the weight w_c is:

$$\frac{\partial \mathcal{L}}{\partial w_c} = \alpha(p(c|x, \alpha) - q(c|x)) x^\top \quad (3)$$

From Eq. 3, we can observe that the temperature has two effects in the gradients. The first is, as mentioned earlier, to control the probability distribution $p(c|x, \alpha) \in [0, 1]$. The second effect is simply to multiply gradients by the value of the temperature α . As α often lies (empirically) in $[1, 250]$, this last effect is harmful during learning since it rescales the learning rate with this value; this can lead to divergence. To cancel this effect and study only the impact of the choice of distribution during training, we propose to normalize the cross-entropy loss function by a constant of value equal to α .

3.2. Finding correlations

As previously mentioned and illustrated, the temperature has a strong impact on the final accuracy but poses different difficulties in finding the optimal value. We then look for a heuristic h , a universal rule, to select a temperature α close to its optimal value, denoted α^* , that is general across datasets and representations. We need to represent each dataset of embeddings $e \in \mathcal{E}$ in a common space \mathcal{S} by using some statistical features $s \in \mathbb{R}^m$ computed over \mathcal{E} with m the number of statistical features. Our goal is to find the best correlation with the observed optimal temperature in order to design a heuristic for temperature.

To ensure our heuristic will be sufficiently general and not just specialized for particular cases such as small number of classes or large embedding sizes, we need to cover as many as possible different cases. To this end, we construct

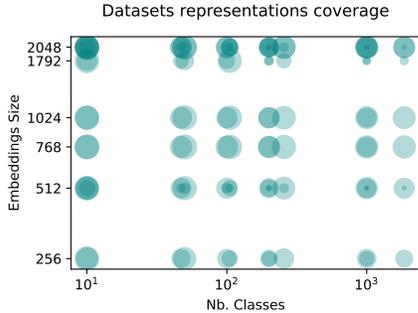


Fig. 1. Coverage between all combinations of feature extractors and datasets. The size of the circles represents the best accuracy achieved for each dataset.

a list of datasets with different numbers of classes and a list of features extractors with different feature sizes / discrimination powers. From each pair (dataset, feature extractor) we build a dataset of embeddings e_i divided in training and validation sets e_i^{train}, e_i^{val} . For each set of embeddings e_i , we compute a description by extracting the statistical features s and empirically find the corresponding optimal temperature $\hat{\alpha}_i^* = \arg \max_{\alpha, W_\alpha} Accuracy(e_i^{val})$ for α selected from a given set of possible temperatures and W_α the weights learned on e_i^{train} with a given temperature α . Thus, to each pair (dataset, feature extractor) is associated the pair (embedding dataset statistical features, optimal empiric temperature) noted: $(S_i, \hat{\alpha}_i^*)$.

In order to find our heuristic $h(\cdot)$ several options are possible. The simplest one would be to consider an affine function on the form $a \cdot s_j + b$, with s_j the most correlated variable in s to our optimal temperature. If a strong enough correlation exists, this would be the simplest heuristic possible. However, more than one variable may be needed to find this correlation. We therefore investigate the strength of the correlation between each statistic independently and a linear combination of our statistical features to the optimal temperature.

4. EXPERIMENTAL STUDY

4.1. Datasets, feature extractors selection

In order to find a generalizable heuristic covering a wide range of cases for a classification task, we selected 12 datasets and 9 feature extractors (Fig. 1). The number of classes ranges from 10 to 1854 while the dimensionality of the features ranges from 256 to 2048. The selected datasets are MNIST, CIFAR10, DTD, PhotoArt, CIFAR100, 105-PinterestFaces, CUB200, ImageNet-R, Caltech256, FSS1000, ImageNetMini, THINGS, containing respectively 10, 10, 47, 50, 100, 105, 200, 200, 256, 1000, 1000, 1854 classes. Regarding the feature extractors, different architectures have been selected with different pre-

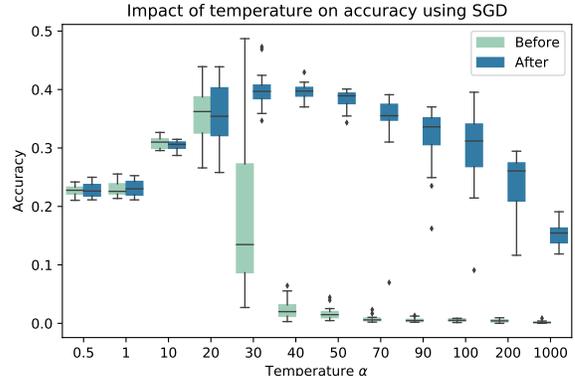


Fig. 2. Impact of temperature on the accuracy of the dataset THINGS [22]. For each temperature, we learned the classification layer and performed 20 learning for 100 epochs.

trainings, in order to cover a large number of dimensions while decorrelating this aspect from the network performance. For example, FaceNet is expected to perform poorly on CIFAR datasets since it is learned on a face recognition task while a ResNet18 pretrained on ImageNet is expected to perform better on natural images while being weaker on 105-PinterestFaces. The feature extractors used are: AlexNet [19] and ResNet- $\{18, 50, 101\}$ [20] pre-trained on ImageNet, ResNet- $\{34, 152\}$ [20] randomly initialized, FaceNet [21] pre-trained on VGGFaces2 and CLIP- $\{RN50, ViT32.b\}$ [14] pre-trained on millions of image-text pairs. The embedding dimensions are respectively: 256, 512, 2048, 2048, 512, 2048, 1792, 1024, 768.

4.2. Statistical features selection

In order to find the hidden relationship between a given dataset and the associated optimal temperature we need to describe each dataset by a feature vector s in a common space \mathcal{S} . Since, as we have seen previously, very different heuristics are proposed to set up the temperature, we selected various features s_i : the dimensionality of embeddings e (dim), the number of output classes ($n_classes$), the mean value of all embeddings values ($mean$), the variance of all embedding values (var), the trace of the average matrix of all intra-class covariance matrices (sb_trace), the trace of the average of all inter-classes covariances matrices (sw_trace), the mean squared error (MSE) between the features correlation matrix and the identity ($feats_corr$), the mean cosine similarity between each dimensions pair ($feats_cos_sim$), the number of samples in the training set ($n_samples$), the average number of samples per class (avg_samp_class) and the percentage of dimensions to be retained for a given explained variance (as in PCA) of 50, 75, 80, 90, 95, 99 ($pca.\%$), the average of all embedding values ($train_mean$) and the standard deviation ($train_std$), the average kurtosis computed on each dimension

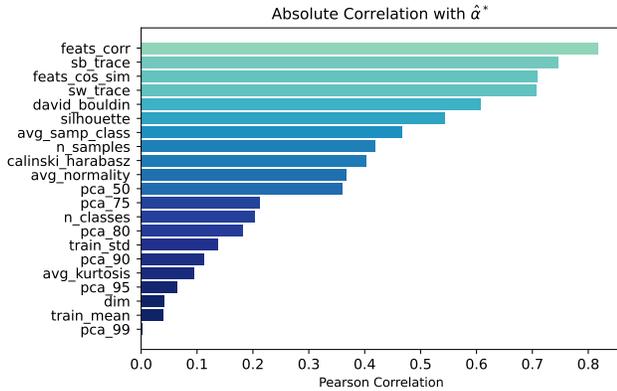


Fig. 3. Absolute value of Pearson correlation between each dataset statistic and the optimal empirical temperature.

(*avg_kurtosis*), and the average Shapiro-Wilk value testing the normality of each dimension (*avg_normality*). Three other popular metrics used in clustering are used such as the Silhouette (*silhouette*), Calinski Harabasz (*calinski_harabasz*) and the David Bouldin (*david_bouldin*) score using the true labels as cluster prediction to obtain measures of the quality of the representation.

4.3. Empiric study of correlations

Once we have extracted the embeddings from various datasets and feature extractors, we need to find the best temperature usable for each case. To do this, we split each dataset of embeddings into training and testing sets and trained the model during 1000 epochs with a batch of size 2048 for a given temperature. We used the default temperature of 1 and temperatures ranging from 5 to 250 with a step of 5: $\alpha \in \{1, 5, \dots, 245, 250\}$. By tracking the accuracy on the test set, we are able to observe the best achievable accuracy for each temperature. We used the rescaled CE loss presented in Sec. 3.1 which allows strong improvements in accuracy over high temperatures using the SGD optimizer as shown in Fig. 2. The latter allows us to observe experimentally the advantage of isolating the peaking distribution effect. However, we found that this had no impact during training when using a smarter optimizer like Adam [23].

In order to find a heuristic for setting a default temperature, we need to find strong correlations from the pairs of optimal empiric temperatures and datasets statistics. Fig. 3 shows the absolute correlation between each statistical value and the temperature. We found that the most interesting variable was the measure of correlation between embedding features. To increase the correlation, we propose to learn a linear regression from our statistics and the optimal temperature, using a cross-validation strategy. The latter omits all sets of embeddings of a dataset (e.g. MNIST) during the learning phase in order to use them in the validation of the found linear combi-

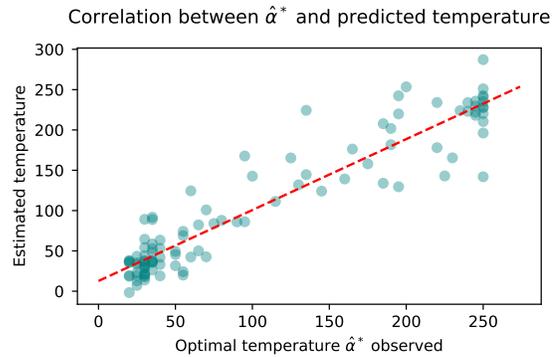


Fig. 4. Correlation between observed optimal temperature versus predicted temperature.

Method	Med. Corr.	Avg. Corr. (\pm std)	<i>p</i> -value
Best Stat	0.9262	0.851 (\pm 0.159)	0.031
4-Best Stats	0.9265	0.865 (\pm 0.131)	0.018
All Stats	0.9563	0.884 (\pm 0.124)	0.014

Table 1. Observed correlations between the most correlated variable, a linear combination of 4 variables and a linear combination of all variables with the optimal empiric temperature.

nation. After that, we repeated this procedure on a subset of the most correlated statistics. The scores are shown in Tab. 1. Finally, we fitted a linear regression on all points whose correlation between our predicted temperature and the empirical optimal temperature is shown in Fig. 4. Obtained results are promising and statistically significant with a Pearson’s correlation of 0.9563 and a *p*-value of $0.014 < 0.05$.

5. CONCLUSION

In this paper, we have shown the importance of the temperature hyperparameter for finetuning a linear classifier on learned representation. We showed that cross-entropy loss can suffer from high temperature if not properly re-scaled. After re-scaling the cross-entropy, we proposed to study the correlations between the optimal empirical temperature observed on many datasets, over a wide range of classes and dimensions, and the statistics computed on the representations of the dataset. In this way, we revealed that some heuristics (such as the dimensionality of embeddings) had little correlation with the optimal temperature while a measure of correlation between features showed strong correlations. We found that appropriate selection and combination of statistics could improve the correlation with the best temperature. We suggest enhancing this pipeline in subsequent work [24] and applying it to other issues, such as identifying the elements of representation learning that will result in high accuracy by predicting it with symbolic regression.

6. REFERENCES

- [1] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean, “Distilling the knowledge in a neural network,” *CoRR*, vol. abs/1503.02531, 2015.
- [2] Andrew Zhai and Hao-Yu Wu, “Classification is a strong baseline for deep metric learning,” in *BMVC, Procs.*, 2019, p. 91.
- [3] Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin, “Language gans falling short,” in *ICLR, Procs.*, 2020.
- [4] Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing, “Toward controlled generation of text,” in *ICML, Procs.*, 2017, p. 1587–1596.
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, “A simple framework for contrastive learning of visual representations,” in *ICML, Procs.*, 2020, pp. 1597–1607.
- [6] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger, “On calibration of modern neural networks,” in *ICML, Procs.*, 2017, pp. 1321–1330.
- [7] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan, “Supervised contrastive learning,” in *NIPS, Procs.*, 2020.
- [8] Zhirong Wu, Alexei A. Efros, and Stella X. Yu, “Improving generalization via scalable neighborhood component analysis,” in *ECCV, Procs.*, 2018, pp. 712–728.
- [9] Olivier Risser-Maroux, Camille Kurtz, and Nicolas Loménie, “Learning an adaptation function to assess image visual similarities,” in *ICIP, Procs.*, 2021, pp. 2498–2502.
- [10] Pei-Hsin Wang, Sheng-Iou Hsieh, Shih-Chieh Chang, Yu-Ting Chen, Jia-Yu Pan, Wei Wei, and Da-Chang Juan, “Contextual temperature for language modeling,” *CoRR*, vol. abs/2012.13575, 2020.
- [11] Xuezhe Ma, Pengcheng Yin, Jingzhou Liu, Graham Neubig, and Eduard H. Hovy, “Softmax q-distribution estimation for structured prediction: A theoretical interpretation for RAML,” *CoRR*, vol. abs/1705.07136, 2017.
- [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *NIPS, Procs.*, 2017, pp. 5998–6008.
- [13] Xu Zhang, Felix X. Yu, Svebor Karaman, Wei Zhang, and Shih-Fu Chang, “Heated-up softmax embedding,” *CoRR*, vol. abs/1809.04157, 2018.
- [14] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, et al., “Learning transferable visual models from natural language supervision,” in *ICML, Procs.*, 2021, pp. 8748–8763.
- [15] Rajeev Ranjan, Carlos Domingo Castillo, and Rama Chellappa, “L2-constrained softmax loss for discriminative face verification,” *CoRR*, vol. abs/1703.09507, 2017.
- [16] Yu-Lin He, Xiao-Liang Zhang, Wei Ao, and Joshua Zhexue Huang, “Determining the optimal temperature parameter for softmax function in reinforcement learning,” *Applied Soft Computing*, vol. 70, pp. 80–85, 2018.
- [17] Yu Liu, Hongyang Li, and Xiaogang Wang, “Rethinking feature discrimination and polymerization for large-scale recognition,” *CoRR*, vol. abs/1710.00870, 2017.
- [18] Ana C Lorena, Luís PF Garcia, Jens Lehmann, Marcilio CP Souto, and Tin Kam Ho, “How complex is your classification problem? a survey on measuring classification complexity,” *ACM Computing Surveys*, vol. 52, no. 5, pp. 1–34, 2019.
- [19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, “Imagenet classification with deep convolutional neural networks,” in *NIPS, Procs.*, 2012.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *CVPR, Procs.*, 2016, pp. 770–778.
- [21] Florian Schroff, Dmitry Kalenichenko, and James Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *CVPR, Procs.*, 2015, pp. 815–823.
- [22] Martin N. Hebart, Adam H. Dickter, Alexis Kidder, Wan Y. Kwok, Anna Corriveau, Caitlin Van Wicklin, and Chris I. Baker, “THINGS: A database of 1, 854 object concepts and more than 26, 000 naturalistic object images,” *PLOS One*, vol. 14, no. 10, pp. e0223792, 2019.
- [23] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *ICLR, Procs.*, 2015.
- [24] Olivier Risser-Maroux and Benjamin Chamand, “What can we learn by predicting accuracy?,” in *IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2023, pp. XXX–XXX.