
MITIGATING DISTRIBUTIONAL BIASES IN CONTRASTIVE LEARNING

A PREPRINT

Lidia Troeshestova

Moscow Institute of Physics and Technology
troeshestova.ls@phystech.edu

Roman Isachenko*

ABSTRACT

Recently there has been a renewed interest in using contrastive learning for self-supervised representation learning. A popular method for learning representations without labels is to compare similar (positive) and dissimilar (negative) pairs of samples. However, negative samples are often chosen randomly, which means that they could be of the same class. This paper analyzes various ways to eliminate these biases. Based on the fully-supervised case, we develop debiased contrastive models that account for same-label datapoints without requiring knowledge of true labels, and explore their properties. The experiments are performed on MNIST and CIFAR10 datasets. This will further improve availability of accurate models for classification in tasks where extensive labeling is expensive or inaccessible.

Keywords : Contrastive learning · Representation learning · Self-supervised learning

1 Introduction

Representation learning has become increasingly popular in recent years due to its ability to learn meaningful representations from large amounts of data, without the need for manual feature engineering. A widely used solution for this problem is contrastive learning, a technique that uses the principle of contrasting samples against each other to learn attributes that are common between data classes and attributes that set apart a data class from another. It encourages the representations of similar pairs (x, x^+) to be close, and those of dissimilar pairs (x, x^-) to be more orthogonal.

The basic metric for supervised similarity is *triplet loss* [Schroff et al., 2015], where \mathbf{x} is called *anchor*, \mathbf{x}^+ is a sample with the same label as anchor, and \mathbf{x}^- is a sample with the different label than anchor. An improvement of triplet loss is *Multi-class N-pair Loss* invented by [Sohn, 2016]:

$$\mathcal{L}(\mathbf{x}, \mathbf{x}^+, \{\mathbf{x}_i^-\}_{i=1}^{N-1}; f) = -\log \frac{\exp(f^T f_i^+)}{\exp(f^T f_i^+) + \sum_{i=1}^{N-1} \exp(f^T f_i^-)} \quad (1)$$

Contrastive learning has proven to be effective for eliciting visual representations: SimCLR [Ting Chen, 2020] is a framework that considerably outperforms previous self-supervised and semi-supervised learning methods. Its authors developed NT-Xent loss (normalized temperature-scaled cross entropy loss), that uses cosine similarity instead of scalar products.

Due to the unavailability of actual labels during training, negative instances x^- are often randomly selected from the training data. However, this approach results in sampling bias, where x^- may actually be similar to x . This bias causes a significant drop in performance.

To address this issue, [Prannay Khosla, 2021] extend the self-supervised batch contrastive approach to the fully-supervised setting, allowing to effectively leverage label information. According to the experiments, supervised

*Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

contrastive loss outperforms the base cross entropy, but only by a small amount. It also outperforms the cross entropy on robustness benchmark (ImageNet-C, which applies common naturally occurring perturbations such as noise, blur and contrast changes to the ImageNet dataset) and is less sensitive to hyperparameter changes.

As for the case of self-supervised learning, a debiased contrastive loss was proposed to mitigate the bias of negative sampling by [Ching-Yao Chuang, 2020]. Given N negative and M positive samples, $\mathcal{L}_{\text{Debiased}}^{N,M}(f)$ is the loss function that corrects for False Negative samples. To remain unsupervised in practice, samples are acquired from the data distribution and a positive distribution mimicked by data augmentations, which can contain False Positives.

This work aims to develop a new algorithm for debiased sampling in contrastive learning, which accounts for False Positive as well as False Negative samples. In addition, it is suggested to tackle the issue of choosing a wrong anchor in multimodal applications such as attribute detection.

2 Problem Statement

2.1 Debiased Contrastive Learning

In [Ching-Yao Chuang, 2020] the distribution over classes is denoted by $\rho(c)$, the joint distribution $p_{x,c}(x, c) = p(x|c)\rho(c)$. Let $h : X \rightarrow C$ be the function assigning the latent class labels. Then $p_x^+(x') = p(x'|h(x') = h(x))$ is the probability of observing x' as a positive example for x and $p_x^-(x') = p(x'|h(x') \neq h(x))$ the probability of a negative example. It is assumed that the class probabilities $\rho(x) = \tau^+$ are uniform, and $\tau^- = 1 - \tau^+$ is the probability of observing any different class. Then the “ideal”, or unbiased loss to optimize would be:

$$L_{\text{Unbiased}}^N(f) = \mathbb{E}_{\substack{x \sim p, x^+ \sim p_x^+, \\ x_i^- \sim p_x^-}} \left[-\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + \sum_{i=1}^N e^{f(x)^T f(x_i^-)}} \right] \quad (2)$$

$$L_{\text{Unbiased}}^N(f) = \mathbb{E}_{\substack{x \sim p, x^+ \sim p_x^+, \\ x_i^- \sim p_x^-}} \left[-\log \frac{\exp(\text{sim}(x, x^+))}{\exp(\text{sim}(x, x^+)) + \sum_{i=1}^N \exp(\text{sim}(x, x_i^-))} \right] \quad (3)$$

In practice, $p_x^-(x_i^-)$ is unavailable. Thus, the standard solution is to sample negative examples x_i^- from the unlabeled $p(x)$ instead. When drawn from $p(x)$, the sample x_i^- will be of the same class as x with probability τ^+ .

$$L_{\text{Biased}}^N(f) = \mathbb{E}_{\substack{x \sim p, x^+ \sim p_x^+, \\ x_i^- \sim p_x}} \left[-\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + \sum_{i=1}^N e^{f(x)^T f(x_i^-)}} \right] \quad (4)$$

Problem statement: Derive a loss function that approximates the ideal L_{Unbiased}^N , only having access to M positive samples and samples from the distribution p .

The data distribution can be decomposed:

$$p(x') = \tau^+ p_x^+(x') + \tau^- p_x^-(x') \quad (5)$$

And the missing $p_x^-(x')$ is derived as

$$p_x^-(x') = \frac{p(x') - \tau^+ p_x^+(x')}{\tau^-} \quad (6)$$

However, estimating $p(x')$ and $p_x^+(x')$ empirically is computationally expensive, and the authors propose a cheaper way to estimate $p_x^-(x')$, using N samples $\{u_i\}_{i=1}^N$ from p and M samples $\{v_i\}_{i=1}^M$ from p_x^+ :

$$L_{\text{Debiased}}^{N,M}(f) = \mathbb{E}_{\substack{x \sim p, x^+ \sim p_x^+, \\ \{u_i\}_{i=1}^N \sim p^N, \\ \{v_i\}_{i=1}^M \sim p_x^{+M}}} \left[-\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Ng(x, \{u_i\}_{i=1}^N, \{v_i\}_{i=1}^M)} \right], \quad (7)$$

where

$$g(x, \{u_i\}_{i=1}^N, \{v_i\}_{i=1}^M) = \max \left\{ \frac{1}{\tau^-} \left(\frac{1}{N} \sum_{i=1}^N e^{f(x)^T f(u_i)} - \tau^+ \frac{1}{M} \sum_{i=1}^M e^{f(x)^T f(v_i)} \right), e^{-1/t} \right\} \quad (8)$$

Using this loss on SimCLR framework allowed to achieve higher accuracy compared to the regular biased loss. With a similar theoretical inference, this work is aimed at developing a new algorithm, that would additionally focus on eliminating bias caused by False Positives.

Tasks:

- reproduce Debiased Contrastive Learning results
- develop novel debiased contrastive loss that accounts for False Positives.
- test the algorithm on CIFAR10 and MNIST datasets

2.2 False Positive Samples

False Positive samples can be encountered in contrastive learning due to excessive or inappropriate augmentation. Similar to Debiased Contrastive Learning, we propose to estimate the positive sampling distribution in the loss function.

$$p_x^+(x') = \frac{p(x') - \tau^- p_x^-(x')}{\tau^+} \quad (9)$$

Contrary to Debiased Contrastive Learning, drawing samples from $p(x')$ will also include positive labelled datapoints, since some of them could be actually negative. With $N \rightarrow \infty$:

$$\mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+ \\ \{x_i^-\}_{i=1}^N \sim p_x^-}} \left[-\log \frac{\exp(f^T f^+)}{\exp(f^T f^+) + \sum_{i=1}^N \exp(f^T f_i^-)} \right] \rightarrow \mathbb{E}_{x^- \sim p_x^-} \left[-\log \frac{R}{R + N \mathbb{E}_{x^- \sim p_x^-} \exp(f^T f^-)} \right], \quad (10)$$

where

$$R = \frac{1}{\tau^+} (\mathbb{E}_{x' \sim p} \exp(f^T f') - \tau^- \mathbb{E}_{x^- \sim p_x^-} \exp(f^T f^-)) \quad (11)$$

With finite N:

$$L_{\text{Debiased}}^N(f) = \mathbb{E}_{\substack{x \sim p \\ \{u_i\}_{i=1}^N \sim p^- \\ v \sim p_x^+}} \left[-\log \frac{R' - \tau^- R^-}{R' + (N\tau^+ - \tau^-)R^-} \right], \quad (12)$$

where

$$R^- = \frac{1}{N} \sum_{i=1}^N \exp(f(x)^T f(u_i)); \quad (13)$$

$$R' = \exp(f(x)^T f(v)) + R^- \quad (14)$$

When the number of positive samples $M > 1$, similar to [Prannay Khosla, 2021] we will aggregate using the summation over positives located outside of the log (which yields better accuracy for supervised contrastive losses rather than summation inside of the log):

$$L_{\text{Debiased}}^{N,M}(f) = \mathbb{E}_{\substack{x \sim p \\ \{u_i\}_{i=1}^N \sim p^- \\ \{v_i\}_{i=1}^M \sim p^+}} \frac{1}{M} \sum \left[-\log \frac{R' - \tau^- R^-}{R' + (N\tau^+ - \tau^-)R^-} \right] \quad (15)$$

2.3 Wrong Anchor

This issue arises in downstream applications of contrastive learning such as image attribute detection, e.g. discrepancy between the image (image representation f) and its description (description representation t) in online shopping. To solve this problem, when calculating the losses for each picture, we propose to weigh the losses by the cosine similarity of the anchor image vector and the description text vector.

$$L_{\text{Unbiased}}(x, x^+, \{x_i^-\}_{i=1}^{N-1}; f) = -\text{sim}(f^T, t^T) \log \frac{\exp(f^T f_i^+)}{\exp(f^T f_i^+ + \sum_{i=1}^N \exp(f^T f_i^-))} \quad (16)$$

3 Experiment

3.1 Baseline Reproduction

The goal of the experiment is to compare Multi-class N-pair Loss 1 with Negative-Debiased Loss 7. We use STL10 [Coates et al., 2011] dataset, an image recognition dataset for developing unsupervised feature learning algorithms. It contains 100000 unlabeled images of 96x96 pixels for unsupervised learning. The baseline negative-debiased model is SimCLR [Ting Chen, 2020] with ResNet-18 [He et al., 2016] as the encoder architecture, with the Adam optimizer [Diederik P. Kingma, 2014] with learning rate 0.001 and batch size 512. All models are trained for 50 epochs and evaluated by training a linear classifier after fixing the learned embedding.

The results are shown in Figure 1. While the difference between Top-5 accuracy for the debiased model and the biased model is insignificant, Top-1 accuracy of the debiased model is higher than that of the biased model.

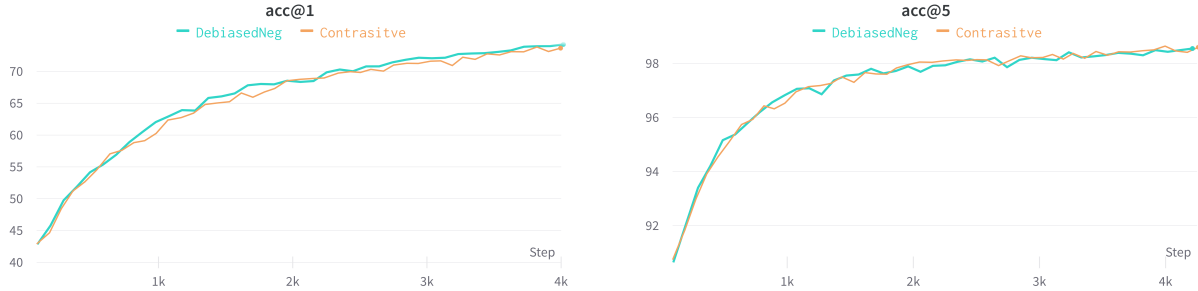


Figure 1: Classification accuracy on CIFAR10. Negative-Debiased loss is denoted as debiasedNeg.

3.2 Testing New Loss

We run the same experiment for Positive-Debiased loss 15. Figure 2 shows a tangible difference between the performances of the novel loss and the previous losses. The accuracy advantage is prevalent from the very first steps of training both for Top-1 and Top-5 accuracy, meaning Positive-Debiased loss is more efficient in the early stages of training.

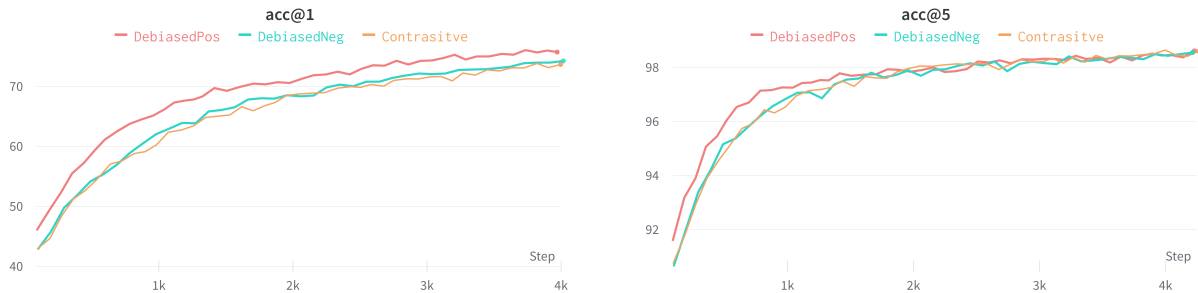


Figure 2: Classification accuracy on CIFAR10, Positive-Debiased loss is denoted as debiasedPos.

To test the performance of the loss under the assumptions made in 2.2, we artificially exclude False Negative errors via accessing ground truth labels during training, and discover that Positive-Debiased loss still outperforms Negative-Debiased loss [Figure 3].

To understand if Positive-Debiased loss is robust to noise, we apply Gaussian blur on train images with probability 0.3, and thus increase the False Positive rate. Figure 4 demonstrates that Positive-Debiased loss is able to mitigate False Positive errors significantly better than Negative-Debiased and Contrastive losses.

Additionally we test all three losses with increased positive sample size M [Figure 5]. Larger M leads to better estimate of the loss both for Positive- and Negative-Debiased losses, yielding better accuracy as a result. It is worth noting that in all of the experiments Positive-Debiased loss achieves better accuracy in the earlier stages of training both for Top-1 and Top-5 scores.

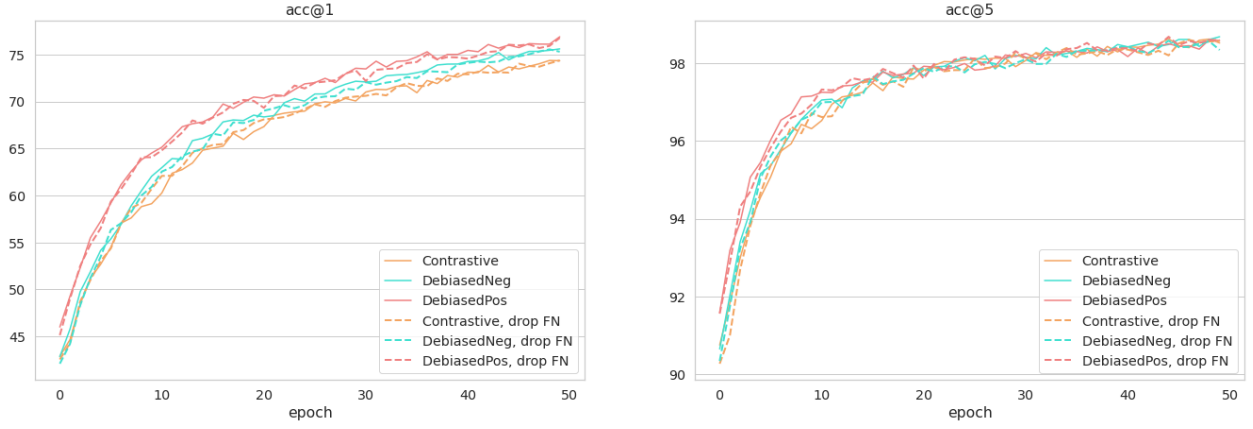


Figure 3: Classification accuracy on CIFAR10 with dropping False Negatives

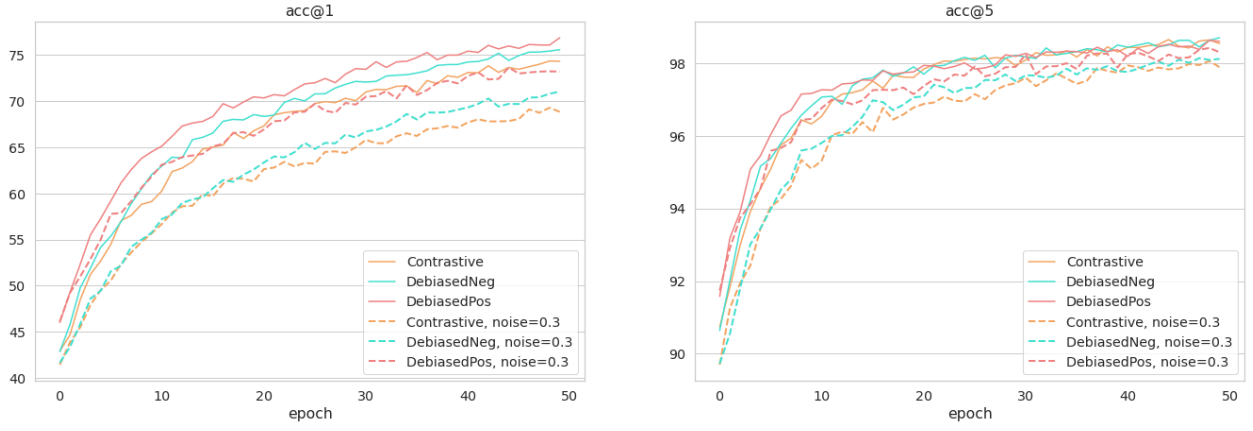


Figure 4: Classification accuracy on CIFAR10 with increased False Positive rate

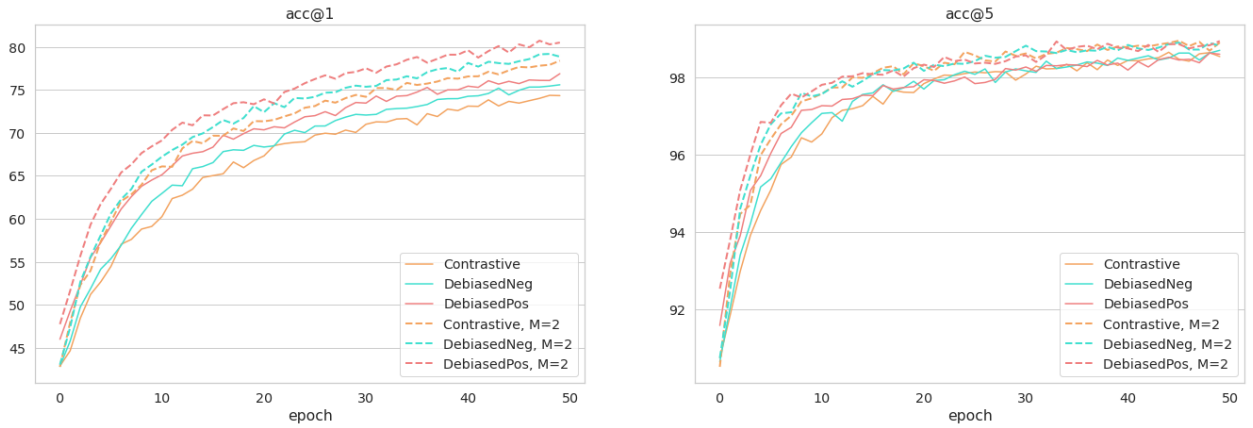


Figure 5: Classification accuracy on CIFAR10 with increased positive sample size M

References

- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Face recognition using triplet loss. *arXiv preprint arXiv:1503.03832*, 2015.
- Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/6b180037abbebea991d8b1232f8a8ca9-Paper.pdf>.
- Mohammad Norouzi Geoffrey Hinton Ting Chen, Simon Kornblith. A simple framework for contrastive learning of visual representations. *arXiv:2002.05709*, 2020.
- Chen Wang Aaron Sarna Yonglong Tian Phillip Isola Aaron Maschinot Ce Liu Dilip Krishnan Prannay Khosla, Piotr Teterwak. Supervised contrastive learning. *arXiv:2004.11362*, 2021.
- Lin Yen-Chen Antonio Torralba Stefanie Jegelka Ching-Yao Chuang, Joshua Robinson. Debiased contrastive learning. *arXiv:2007.00224*, 2020.
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 215–223, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. doi:10.1109/CVPR.2016.90.
- Jimmy Ba Diederik P. Kingma. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.