ABOUT CONTRASTIVE UNSUPERVISED REPRESENTA-TION LEARNING FOR CLASSIFICATION AND ITS CON-VERGENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Contrastive representation learning has been recently proved to be very efficient for self-supervised training. These methods have been successfully used to train encoders which perform comparably to supervised training on downstream classification tasks. A few works have started to build a theoretical framework around contrastive learning in which guarantees for its performance can be proven. We provide extensions of these results to training with multiple negative samples and for multiway classification. Furthermore, we provide convergence guarantees for the minimization of the contrastive training error with gradient descent of an overparametrized deep neural encoder, and provide some numerical experiments that complement our theoretical findings.

1 INTRODUCTION

The aim of this work is to provide additional theoretical guarantees for *contrastive learning* (van den Oord et al. 2018), which corresponds to methods allowing to learn useful data representations in an *unsupervised* setting. Unsupervised representation learning was initially approached with a fair amount of success by training through the minimization of losses coming from "pretext" tasks, a technique known as *self-supervision* (Doersch & Zisserman, 2017), where labels can be automatically constructed. Notable examples of pretext tasks in computer vision include colorization (Zhang et al. 2016), transformation prediction (Gidaris et al., 2018; Dosovitskiy et al., 2014) or predicting patch relative positions (Doersch et al., 2015). Some theoretical guarantees (Lee et al., 2020) were recently proposed to support training on pretext tasks.

Contrastive learning is also known to be very effective for pretraining supervised methods (Chen et al. 2020ab) Grill et al. 2020 Caron et al. 2020), where we can observe that, quite surprisingly, the gap between unsupervised and supervised performance has been closed for tasks such as image classification: the use of a pretrained image encoder on top of simple classification layers, that are trained on a fraction of the labels available, allows to achieve an accuracy comparable to that of a fully supervised end-to-end training (Hénaff et al. 2019) Grill et al. 2020). Contrastive methods show also strong success in natural language processing (Logeswaran & Leel 2018) Mikolov et al. 2013; Devlin et al. 2018; van den Oord et al. 2018), video classification (Sun et al. 2019), reinforcement learning (Srinivas et al. 2020) and time-series (Franceschi et al. 2019).

Although the papers cited above introduce methods with considerable variations, they mostly agree on the following basic pretraining approach: provided a dataset, an encoder is trained using a contrastive loss whose minimization allows to learn embeddings that are *similar* for pairs of samples (called the *positives*) that are close to each other (such as pairs of random data augmentations of the same image, see He et al. (2020); Chen et al. (2020a)), while such embeddings are *contrasted* for dissimilar pairs (called the *negatives*).

However, despite growing efforts (Saunshi et al.) [2019] Wang & Isola, [2020], as of today, few theoretical results have been obtained. For instance, there is still no clear theoretical explanation of how a supervised task could benefit from an upstream unsupervised pretraining phase, or of what could be the theoretical guarantees for the convergence of the minimization procedure of the contrastive loss during this pretraining phase. Getting some answers to these questions would undoubtedly be a step towards a better theoretical understanding of contrastive representation learning.

Our contributions in this paper are twofold. In Section 3 we provide new theoretical guarantees for the classification performance of contrastively trained models in the case of multiway classification tasks, using *multiple* negative samples. We extend results from Saunshi et al. (2019) to show that unsupervised training performance reflects on a subsequent classification task in the case of multiple tasks and when a high number of negative samples is used. In Section 4 we prove a convergence result for an *explicit* algorithm (gradient descent), when training overparametrized deep neural network for unsupervised contrastive representation learning. We explain how results from Allen-Zhu et al. (2019) about training convergence of overparametrized deep neural networks can be applied to a contrastive learning objective. The results and major assumptions of both Sections 3 and 4 are illustrated in Section 5 through experiments on a few simple datasets.

2 RELATED WORK

A growing literature attempts to build a theoretical framework around contrastive learning and to provide justifications for its success beyond intuitive ideas. In Saunshi et al. (2019) a formalism is proposed together with results on classification performance based on unsupervisedly learned representation. However, these results do not explain the performance gain that is observed empirically (Chen et al., 2020a; He et al., 2020) when a high number of negative samples are used, while the results proposed in Section 3 below hold for an arbitrary large number of negatives (and decoupled from the number of classification tasks). A more recent work (Wang & Isola, 2020) emphasizes the two tendencies encouraged by the contrastive loss: the encoder's outputs are incentivized to spread evenly on the unit hypersphere, and encodings of same-class samples are driven close to each other while those of different classes are driven apart. Interestingly, this work also shows how the tradeoff between these two aspects can be controlled, by introducing weight factors in the loss leading to improved performance. Chuang et al. (2020) considers the same setting as Saunshi et al. (2019) and addresses the bias problem that comes from collisions between positive and negative sampling in the unsupervised constrastive loss. They propose to simulate unbiased negative sampling by assuming, among other things, extra access to positive sampling. However, one has to keep in mind that an excessive access to positive sampling gets the setting closer to that of supervised learning.

In a direction that is closer to the result proposed in Section 4 below, Wen (2020) provides a theoretical guarantee on the training convergence of gradient descent for an overparametrized model that is trained with an unsupervised contrastive loss, using earlier works by Allen-Zhu et al. (2019). However, two separate encoders are considered instead of a single one: one for the query, which corresponds to a sample from the dataset, and one for the (positive and negative) samples to compare the query to. In this setting, it is rather unclear how the two resulting encoders are to be used for downstream classification. In Section 4 below, we explain how the results from Allen-Zhu et al. (2019) can be used for the more realistic setting of a single encoder, by introducing a reasonable assumption on the encoder outputs.

3 UNSUPERVISED TRAINING IMPROVES SUPERVISED PERFORMANCE

In this section, we provide new results in the setting previously considered in Saunshi et al. (2019). We assume that data are distributed according to a finite set C of latent classes, and denote $N_C = \operatorname{card}(C)$ its cardinality. Let ρ be a discrete distribution over C that is such that

$$\sum_{c \in \mathcal{C}} \rho(c) = 1 \quad \text{and} \quad \rho(c) > 0$$

for all $c \in C$. We denote D_c a distribution over the feature space \mathcal{X} from a class $c \in C$. In order to perform unsupervised contrastive training, on the one hand we assume that we can sample *positive* pairs (x, x^+) from the distribution

$$\mathcal{D}_{\rm sim}(x, x^+) = \sum_{c \in \mathcal{C}} \rho(c) \mathcal{D}_c(x) \mathcal{D}_c(x^+), \tag{1}$$

namely, (x, x^+) is sampled as a mixture of independent pairs conditionally to a shared latent class, sampled according to ρ . On the other hand, we assume that we can sample *negative* samples x^- from the distribution

$$\mathcal{D}_{\text{neg}}(x^{-}) = \sum_{c \in \mathcal{C}} \rho(c) \mathcal{D}_c(x^{-}).$$
(2)

Given $k \leq N_{\mathcal{C}} - 1$, a (k + 1)-way classification task is a subset $\mathcal{T} \subseteq \mathcal{C}$ of cardinality $|\mathcal{T}| = k + 1$, which induces the conditional distribution

$$\mathcal{D}_{\mathcal{T}}(c) = \rho(c \mid c \in \mathcal{T})$$

for $c \in \mathcal{C}$ and we define

$$\mathcal{D}_{\mathcal{T}}(x,c) = \mathcal{D}_{\mathcal{T}}(c)\mathcal{D}_c(x).$$

In particular, we denote as C, whenever there is no ambiguity, the N_{C} -way classification task where the labels are sampled from ρ , namely $\mathcal{D}_{C}(x, c) = \rho(c)\mathcal{D}_{c}(x)$.

Supervised loss and mean classifier. For an encoder function $f : \mathcal{X} \to \mathbb{R}^d$, we define a supervised loss (cross-entropy with the best possible linear classifier on top of the representation) over task \mathcal{T} as

$$L_{\sup}(f,\mathcal{T}) = \inf_{W \in \mathbb{R}^{|\mathcal{T}| \times d}} \mathbb{E}_{(x,c) \sim \mathcal{D}_{\mathcal{T}}} \left[-\log\left(\frac{\exp\left(Wf(x)\right)_{c}}{\sum_{c' \in \mathcal{T}} \exp\left(Wf(x)\right)_{c'}}\right) \right].$$
 (3)

Then, it is natural to consider the *mean* or *discriminant* classifier with weights W^{μ} which stacks, for $c \in \mathcal{T}$, the vectors

$$W_{c,:}^{\mu} = \mathbb{E}_{x \sim \mathcal{D}_c} \left[f(x) \right]$$
(4)
and whose corresponding (supervised) loss is given by

$$L^{\mu}_{\sup}(f,\mathcal{T}) = \mathbb{E}_{(x,c)\sim\mathcal{D}_{\mathcal{T}}}\left[-\log\left(\frac{\exp\left(W^{\mu}f(x)\right)_{c}}{\sum_{c'\in\mathcal{T}}\exp\left(W^{\mu}f(x)\right)_{c'}}\right)\right].$$
(5)

Note that, obviously, one has $L_{\sup}(f, \mathcal{T}) \leq L_{\sup}^{\mu}(f, \mathcal{T})$.

Unsupervised contrastive loss. We consider the unsupervised contrastive loss with N negative samples given by

$$L_{\rm un}^{N}(f) = \mathbb{E}_{\substack{(x,x^{+}) \sim \mathcal{D}_{\rm sim} \\ X^{-} \sim \mathcal{D}_{\rm neg}^{\otimes N}}} \left[-\log\left(\frac{\exp\left(f(x)^{T}f(x^{+})\right)}{\exp\left(f(x)^{T}f(x^{+})\right) + \sum_{x^{-} \in X^{-}}\exp\left(f(x)^{T}f(x^{-})\right)}\right) \right], \quad (6)$$

where \mathcal{D}_{sim} is given by Equation [1] and where $\mathcal{D}_{neg}^{\otimes N}$ stands for the N tensor product of the \mathcal{D}_{neg} distribution given by Equation [2]. When a single negative sample is used (N = 1), we will use the notation $L_{un}(f) = L_{un}^1(f)$. In the rest of the paper, N will stand for the number of negatives used in the unsupervised loss [6].

3.1 INEQUALITIES FOR UNSUPERVISED TRAINING WITH MULTIPLE CLASSES

The following Lemma states that the unsupervised objective with a single negative sample can be related to the supervised loss for which the target task is classification over the whole set of latent classes C.

Lemma 3.1. For any encoder $f : \mathcal{X} \to \mathbb{R}^d$, one has

$$L_{\sup}(f,\mathcal{C}) \le L_{\sup}^{\mu}(f,\mathcal{C}) \le \frac{1}{p_{\min}^{\rho}} L_{\mathrm{un}}(f) + \log N_{\mathcal{C}},\tag{7}$$

where $p_{\min}^{\rho} = \min_{c} \rho(c)$.

The proof of Lemma 3.1 is given in Supplementary Material, and uses a trick from Lemma 4.3 in Saunshi et al. (2019) relying on Jensen's inequality. This Lemma relates the unsupervised and the supervised losses, a shortcoming being the introduction of p_{\min}^{ρ} , which is small for a large $N_{\mathcal{C}}$ since obviously $p_{\min}^{\rho} \leq 1/N_{\mathcal{C}}$.

The analysis becomes more difficult with a larger number of negative samples. Indeed, in this case, one needs to carefully keep track of how many distinct classes will be represented by each draw. This is handled by Theorem B.1 of Saunshi et al. (2019), but the bound given therein only estimates an expectation of the supervised loss w.r.t. the random subset of classes considered (so called tasks). For multiple negative samples, the approach adopted in the proof of Lemma 3.1 above further degrades, since p_{\min}^{ρ} would be replaced by the minimum probability among tuple draws, an even much smaller quantity.

We propose the following Lemma, which assumes that the number of negative samples is large enough compared to the number of latent classes. **Lemma 3.2.** Consider the unsupervised objective with N negative samples as defined in Equation (6) and assume that N satisfies $N = \Omega(N_C \log N_C)$. Then, we have

$$L_{\sup}(f,\mathcal{C}) \le L_{\sup}^{\mu}(f,\mathcal{C}) \le \frac{1}{p_{cc}^{\rho}(N)} L_{\mathrm{un}}^{N}(f), \tag{8}$$

where $p_{cc}^{\rho}(N)$ is the probability to have all coupons after N draws in an $N_{\mathcal{C}}$ -coupon collector problem with draws from ρ .

The proof of Lemma 3.2 is given in Supplementary Material. In this result, $p_{cc}^{\rho}(N)$ is related to the following coupon collector problem. Assume that ρ is the uniform distribution over C and let T be the random number of necessary draws until each $c \in C$ is drawn at least once. It is known (see for instance Motwani & Raghavan (1995)) that the expectation and variance of T are respectively given by $N_{C}H_{N_{C}}$ and $(N_{C}\pi)^{2}/6$, where H_{n} is the *n*-th harmonic number $H_{n} = \sum_{i=1}^{n} 1/i$. This entails using Chebyshev's inequality that

$$\mathbb{P}\left(|T - N_{\mathcal{C}}H_{N_{\mathcal{C}}}| \ge \beta N_{\mathcal{C}}\right) \le \frac{\pi^2}{6\beta^2}$$

for any $\beta > 0$, so that whenever ρ is sufficiently close to a uniform distribution and $N = \Omega(N_C \log N_C)$, the probability p_{cc}^{ρ} is reasonably high. Due to the randomness of the classes sampled during training, it is difficult to obtain a better inequality than Lemma 3.2 if we want to upper bound $L_{un}^N(f)$ by the supervised $L_{sup}(f, C)$ on all classes. However, the result can be improved by considering the average loss over tasks $L_{sup,k}(f)$, as explained in the next Section.

3.2 GUARANTEES ON THE AVERAGE SUPERVISED LOSS

In this Section, we bound the average of the supervised classification loss on tasks that are subsets of C. Towards this end, we need to assume (only in this Section) that ρ is uniform. We consider supervised tasks consisting in distinguishing one latent class from k other classes, given that they are distinct and uniformly sampled from C. We define the average supervised loss of f for (k + 1)-way classification as

$$L_{\sup,k}(f) = \mathbb{E}_{\mathcal{T}\sim\mathcal{D}^{k+1}}\left[L_{\sup}\left(f,\mathcal{T}\right)\right],\tag{9}$$

where \mathcal{D}^{k+1} is the uniform distribution over (k + 1)-way tasks, which means uniform sampling of $\{c_1, \dots, c_{k+1}\}$ distinct classes in \mathcal{C} . We define also the average supervised loss of the mean classifier

$$^{\mu}_{\sup,k}(f) = \mathbb{E}_{\mathcal{T}\sim\mathcal{D}^{k+1}}\left[L^{\mu}_{\sup}\left(f,\mathcal{T}\right)\right],\tag{10}$$

where we recall that $L_{\sup}^{\mu}(f, \mathcal{T})$ is given by (5). The next Proposition is a generalization to arbitrary values of k and N of Lemma 4.3 from Saunshi et al. (2019), where it is assumed k = 1 and N = 1. **Proposition 3.3.** Consider the unsupervised loss $L_{un}^{N}(f)$ from Equation (6) with N negative samples. Assume that ρ is uniform over C and that $2 \leq k + 1 \leq N_{c}$. Then, any encoder function $f : \mathcal{X} \to \mathbb{R}^{d}$ satisfies

$$L_{\sup,k}(f) \le L_{\sup,k}^{\mu}(f) \le \frac{\kappa}{1 - \tau_N^+} \left(L_{\operatorname{un}}^N(f) - \tau_N^+ \log(N+1) \right)$$

with $\tau_N^+ = \mathbb{P}\left[c_i = c, \forall i \mid (c, c_1, \cdots, c_N) \sim \rho^{\otimes N+1} \right].$

The proof of Proposition 3.3 is given in Supplementary Material. This Proposition states that, in a setting similar to that of Saunshi et al. (2019), on average, the (k + 1)-way supervised classification loss is upper-bounded by the unsupervised loss (both with N = 1 negative or N > 1 negatives), that contrastive learning algorithms actually minimize. Therefore, these results give hints for the performances of the learned representation for downstream tasks.

Also, while Saunshi et al. (2019) only considers an unsupervised loss with N = k negatives along with (k + 1)-way tasks for evaluation, the quantities N and k are decoupled in Proposition 3.3 Furthermore, whenever ρ is uniform, one has $\tau_N^+ = \sum_{c \in \mathcal{C}} \rho(c)^{N+1} = N_c^{-N}$, which decreases to 0 as $N \to +\infty$, so that a larger number of negatives N makes $k/(1 - \tau_N^+)$ smaller and closer to k. This provides a step towards a better understanding of what is actually done in practice with unsupervised contrastive learning. For instance, N = 65536 negatives are used in He et al. (2020).

While we considered a generic encoder f and a generic setting in this Section, the next Section deconsiders a more realistic setting of an unsupervised objective with a fixed available dataset, and the study of an *explicit* algorithm for the training of f.

4 CONVERGENCE OF GRADIENT DESCENT FOR CONTRASTIVE UNSUPERVISED LEARNING

This section leverages results from Allen-Zhu et al. (2019) to provide convergence guarantees for gradient-descent based minimization of the contrastive training error, where the unsupervisedly trained encoder is an overparametrized deep neural network.

Deep neural network encoder. We consider a family of encoders f defined as a deep feed-forward neural network following Allen-Zhu et al. (2019). We quickly restate its structure here for the sake of completeness. A deep neural encoder f is parametrized by matrices $A \in \mathbb{R}^{m \times d_x}$, $B \in \mathbb{R}^{d \times m}$ and $W_1, \ldots, W_L \in \mathbb{R}^{m \times m}$ for some depth L. For an input $x \in \mathbb{R}^{d_x}$, the feed-forward output $y \in \mathbb{R}^d$ is given by

$$g_0 = Ax, \quad h_0 = \phi(g_0), \quad g_l = W_l h_{l-1}, \quad h_l = \phi(g_l) \text{ for } l = 1, \dots, L,$$

 $y = Bh_L,$

where ϕ is the ReLU activation function. Note that the architecture can also include residual connections and convolutions, as explained in Allen-Zhu et al. (2019).

We know from Allen-Zhu et al. (2019) that, provided a δ -separation condition on the dataset (x_i, y_i) for $i = 1, \ldots, n$ with $\delta > 0$ and sufficient overparametrization of the model $(m = \Omega (\operatorname{poly}(n, L, \delta^{-1}) \cdot d))$, the optimisation of the least-squares error $\frac{1}{2} \sum_{i=1}^{n} \|\widehat{y}_i - y_i\|_2^2$ using gradient descent provably converges to an arbitrarily low value $\epsilon > 0$, where $\widehat{y}_i = f(x_i)$ are the network outputs. Moreover, the convergence is linear i.e. the number of required epochs is $T = O(\log(1/\epsilon))$, although involving a constant of order $\operatorname{poly}(n, L, \delta^{-1})$. Although this result does not directly apply to contrastive unsupervised learning, we explain below how it can be adapted provided a few additional assumptions.

Ideally, we would like to prove a convergence result on the unsupervised objective defined in Equation 6. However, we need to define an objective through an explicitly given dataset so that it falls within the scope of Allen-Zhu et al. (2019). Regarding this issue, we assume in what follows that we dispose of a set of fixed triplets $(x, x^+, x^-) \in (\mathbb{R}^{d_x})^3$ we train on.

Objective function. Let us denote this fixed training set $\{(x_i, x_i^+, x_i^-)\}_{i=1}^n$. Each element leads to an output $z_i = (f(x_i), f(x_i^+), f(x_i^-))$ by the encoder and we optimize the empirical objective

$$\widehat{L}_{un}(f) = \sum_{i=1}^{n} \zeta(f(x_i)^T (f(x_i^-) - f(x_i^+))) = \sum_{i=1}^{n} \ell(z_i),$$
(11)

where we introduced the loss function $\ell(z_i) = \ell(z_{i,1}, z_{i,2}, z_{i,3}) = \zeta(z_{i,1}^T(z_{i,3} - z_{i,2}))$ with $\zeta(x) = \log(1 + e^x)$. Note that $\widehat{L}_{un}(f)/n$ is the empirical counterpart of the unsupervised loss (6). Our management of the set of training triplets can be compared to that of Wen (2020) who similarly fixes them in advance but uses multiple negatives and the same x_i as a positive. However, two distinct encoders are trained therein, one for the reference sample x_i and another for the rest. We consider here the more realistic case where a single encoder is trained. Our approach also applies to multiple negatives, but we only use a single one here for simplicity. We need the following data separation assumption from Allen-Zhu et al. (2019).

Assumption 1. We assume that all the samples $x \in \mathcal{X}_{data} = \bigcup_{i=1}^{n} \{x_i, x_i^+, x_i^-\}$ are normalized ||x|| = 1 and that there exists $\delta > 0$ such that $||x - x'||_2 \ge \delta$ for any $x, x' \in \mathcal{X}_{data}$.

Note that sampling the positives and negatives x_i^+, x_i^- need not to be made through simple draws from the dataset. A common practice in contrastive learning (Chen et al., 2020a) is to use data augmentations, where we replace x_i^{\pm} by $\psi(x_i^{\pm})$ for an augmentation function ψ also drawn at random. Such an augmentation can include, whenever inputs are color images, Gaussian noise, cropping, resizing, color distortion, rotation or a combination thereof, with parameters sampled at random in prescribed intervals. The setting considered here allows the case where x_i^{\pm} are actually augmentations (we won't write $\psi(x_i^{\pm})$ but simply x_i^{\pm} to simplify notations), provided that Assumption 1 is satisfied and that such augmentations are performed and fixed before training. Note that, in practice, the augmentations are themselves randomly sampled at each training iteration (Chen et al.) 2020a). Unfortunately, this would make the objective intractable and the convergence result we are about to derive does not apply in that case.

In order to apply the convergence result from <u>Allen-Zhu et al.</u> (2019), we need to prove that the following gradient-Lipschitz condition

$$\ell(z+z') \le \ell(z) + \langle \nabla \ell(z), z' \rangle + \frac{L_{\text{smooth}}}{2} \left\| z' \right\|^2$$
(12)

holds for any $z, z' \in \mathbb{R}^{3d}$, for some constant $L_{\text{smooth}} > 0$, where ℓ is the loss given by (11). However, as defined previously, ℓ does not satisfy (12) without extra assumptions. We propose to bypass this problem by making the following additional assumption on the norms of the outputs of the encoder.

Assumption 2. For each element $x \in \mathcal{X}_{data}$, the output $z = f(x) \in \mathbb{R}^d$ satisfies

$$\eta < \|z\| < C$$

during and at the end of the training of the encoder f, for some constants $0 < \eta < C < +\infty$.

In Section 5 we check experimentally on three datasets (see Figure 3 herein) that this assumption is rather realistic. The lower bound $\eta > 0$ is necessary and used in Lemma 4.2 below, while the upper bound C is used in the next Lemma 4.1 which establishes the gradient-Lipschitz smoothness of the unsupervised loss ℓ and provides an estimation of L_{smooth} .

Lemma 4.1. Consider the unsupervised loss ℓ given by (11), grant Assumption 2 and define the set

$$B^{3} = \left\{ z = (z_{1}, z_{2}, z_{3}) \in (\mathbb{R}^{d})^{3} : \max_{j=1,2,3} \|z_{j}\|_{2}^{2} \leq C^{2} \right\}$$

where C > 0 is defined in Assumption 2 Then, the restriction of ℓ to B^3 satisfies (12) with a constant $L_{\text{smooth}} \leq 2 + 8C^2$.

The proof of Lemma 4.1 is given in Supplementary Material. Now, we can state the main result of this Section.

Theorem 1. Grant both Assumptions l and 2 let $\epsilon > 0$ and let $\hat{L}_{un}(f)$ be the loss given by (11). Then, assuming that

$$m \ge \Omega\Big(\frac{\operatorname{poly}(n,L,\delta^{-1}) \cdot d}{\epsilon}\Big)$$

the gradient descent algorithm with a learning rate ν and a number of steps T such that

$$\nu = \Theta\Big(\frac{d\delta}{\operatorname{poly}(n,L) \cdot m}\Big) \quad and \quad T = O\Big(\frac{\operatorname{poly}(n,L)}{\delta^2\epsilon^2}\Big),$$

finds a parametrization of the encoder f satisfying

$$\widehat{L}_{\mathrm{un}}(f) \le \epsilon.$$

The proof of Theorem 1 is given in Supplementary Material. Although it uses Theorem 6 from Allen-Zhu et al. (2019), it is actually *not* an immediate consequence of it. Indeed, in our case, the Theorem 6 therein only allows us to conclude that $\|\nabla \hat{L}_{un}(f)\| \le \epsilon$, where the gradient is taken w.r.t. the outputs of f. The convergence of the objective itself is obtained thanks to the following Lemma whose proof is given in Supplementary Material.

Lemma 4.2. Grant Assumption 2 and assume that the parameters of the encoder f are optimized so that $\|\nabla \widehat{L}_{un}(f)\| \le \epsilon$ with $\epsilon < \eta/2$, where η is defined in Assumption 2 Then, for any $i = 1, \ldots, n$, we have $\ell(z_i) \le 2\epsilon/\eta$ where $z_i = (f(x_i), f(x_i^+), f(x_i^-))$.

This Lemma is crucial for proving Theorem 1 as it allows to show, in this setting, that the reached critical point is in fact a global minimum.

A natural idea would be then to combine Theorem 1 with Proposition 3.3 in order to prove that gradient descent training of the encoder using the unsupervised contrastive loss helps to minimize the supervised loss. This paper makes a step towards such a result, but let us stress that it requires

much more work, to be considered in future papers, the technical problems to be addressed being as follows. Firstly, the result of Theorem 1 applies to $\hat{L}_{un}(f)$ and cannot be directly extrapolated on $L_{un}(f)$. Doing so would require a sharp control of the generalization error, while Theorem 1 is about the training error only. Secondly, Assumption 1 requires that all samples are separated and, in particular, distinct. This cannot hold when the objective is defined through an expectation as we did in Section 3 Indeed, it would be invalidated simply by reusing a sample in two different triples.

5 EXPERIMENTS

In this section, we report experiments that illustrate our theoretical findings.

Datasets and Experiments. We use a small convolutional network as encoder on MNIST (LeCun & Cortes 2010) and FashionMNIST (Xiao et al.) 2017), and VGG-16 (Simonyan & Zisserman, 2015) on CIFAR-10 (Krizhevsky et al., 2009). Experiments are performed with PyTorch (Paszke et al., 2019).

Results. Figure 1 provides an illustration of Lemma 3.1 where we display the values of $L_{\rm un}$ (i.e., $L_{\rm un}^N$ with N = 1) and $L_{\rm sup}^{\mu}(f, C)$ along training iterations over 5 separate runs (and their average). We observe that Inequality (7) is satisfied on these experiments, even when the $\log N_C$ term is discarded. Moreover, both losses follow a similar trend. Figure 2 illustrates Lemma 3.2 for several values of N. Once again, we observe that both losses behave similarly, and that Inequality (8) seems to hold even without the $1/p_{cc}^{\rho}$ term (removed for these displays).



Figure 1: Illustration of Lemma 3.1 we observe that Inequality (7) is satisfied on these examples, even without the log $N_{\mathcal{C}}$ term, and that both losses behave similarly (5 runs are displayed together with their average).



Figure 2: Illustration of Lemma 3.2 with N = 15, 25, 35 on MNIST. We observe again that both the unsupervised and supervised losses behave similarly and that Inequality (8) is satisfied in these experiments, even without the $1/p_{cc}^{\rho}$ factor (5 runs are displayed together with their average).

Finally, Figure 3 displays the minimum and maximum Euclidean norms of the outputs of the encoder along training. On these examples, we observe that one can indeed assume these norms to be lower and upper bounded by constants, as stated in Assumption 2



Figure 3: Minimum and maximum Euclidean norms of the outputs of the encoder along contrastive unsupervised training. We observe that Assumption 2 is satisfied on these examples (5 runs are displayed together with their average), the dashed line shows that the minimum norms are away from 0 even in the early iterations.

6 CONCLUSION

This work provides extensions to previous results on contrastive unsupervised learning, in order to somewhat improve the theoretical understanding of the performance that is empirically observed with pre-trained encoders used for subsequent supervised task. The main hindrance to tighter bounds in Section 3 is the blind randomness of negative sampling, which is unavoidable in the unsupervised setting. Section 4 explains how recent theoretical results about gradient descent training of overparametrized deep neural networks can be used for unsupervised contrastive learning, and concludes with an explanation of why combining the results from Sections 3 and 4 requires many extra technicalities to be considered in future works. Let us conclude by stressing, once again, our motivations for doing this: unsupervised learning theory is much less developed than supervised learning theory, and recent empirical results (see Section 1) indicate that some forms of contrastive learning enable the learning of powerful representations without supervision. In many fields of application, labels are too difficult, too expensive or too invasive to obtain (in medical applications, see for instance Ching et al. (2018)). We believe that a better understanding of unsupervised learning is therefore of utmost importance.

REFERENCES

- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A Convergence Theory for Deep Learning via Over-Parameterization. In *International Conference on Machine Learning*, pp. 242–252. PMLR, 2019.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments, 2020.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020a.
- Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big selfsupervised models are strong semi-supervised learners, 2020b.
- Travers Ching, Daniel S Himmelstein, Brett K Beaulieu-Jones, Alexandr A Kalinin, Brian T Do, Gregory P Way, Enrico Ferrero, Paul-Michael Agapow, Michael Zietz, Michael M Hoffman, et al. Opportunities and obstacles for deep learning in biology and medicine. *Journal of The Royal Society Interface*, 15(141):20170387, 2018.
- Ching-Yao Chuang, Joshua Robinson, Lin Yen-Chen, Antonio Torralba, and Stefanie Jegelka. Debiased contrastive learning, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL http://arxiv.org/abs/1810.04805.

- Carl Doersch and Andrew Zisserman. Multi-task self-supervised visual learning. CoRR, abs/1708.07860, 2017. URL http://arxiv.org/abs/1708.07860,
- Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. *CoRR*, abs/1505.05192, 2015. URL http://arxiv.org/abs/1505.05192.
- Alexey Dosovitskiy, Jost Tobias Springenberg, Martin A. Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. *CoRR*, abs/1406.6909, 2014. URL http://arxiv.org/abs/1406.6909.
- Jean-Yves Franceschi, Aymeric Dieuleveut, and Martin Jaggi. Unsupervised scalable representation learning for multivariate time series. In *Advances in Neural Information Processing Systems*, pp. 4650–4661, 2019.
- Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *CoRR*, abs/1803.07728, 2018. URL http://arxiv.org/abs/ 1803.07728.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning, 2020.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aäron van den Oord. Data-efficient image recognition with contrastive predictive coding. *CoRR*, abs/1905.09272, 2019. URL http://arxiv.org/abs/1905.09272.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL http://yann. lecun.com/exdb/mnist/
- Jason D. Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo. Predicting what you already know helps: Provable self-supervised learning, 2020.
- Lajanugen Logeswaran and Honglak Lee. An efficient framework for learning sentence representations. In International Conference on Learning Representations, 2018.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- Rajeev Motwani and Prabhakar Raghavan. Randomized Algorithms. Cambridge University Press, 1995. doi: 10.1017/CBO9780511814075.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems 32, pp. 8024–8035. Curran Associates, Inc., 2019. URL http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf
- Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. A theoretical analysis of contrastive unsupervised representation learning. In *International Conference on Machine Learning*, pp. 5628–5637, 2019.

- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- Aravind Srinivas, Michael Laskin, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning, 2020.
- Chen Sun, Fabien Baradel, Kevin Murphy, and Cordelia Schmid. Learning video representations using contrastive bidirectional transformer. *arXiv preprint arXiv:1906.05743*, 2019.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. CoRR, abs/1807.03748, 2018. URL http://arxiv.org/abs/1807.03748
- Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere, 2020.
- Zixin Wen. Convergence of end-to-end training in deep unsupervised contrasitive learning, 2020.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful image colorization. *CoRR*, abs/1603.08511, 2016. URL http://arxiv.org/abs/1603.08511.