

# Contrastive Learning with Hard Negative Samples

Joshua Robinson  
Ching-Yao Chuang  
Suvrit Sra  
Stefanie Jegelka

Massachusetts Institute of Technology

joshrob@mit.edu  
cychuang@mit.edu  
suvrit@mit.edu  
stefje@csail.mit.edu

## Abstract

We consider the question: how can you sample good negative examples for contrastive learning? We argue that, as with metric learning, learning contrastive representations benefits from hard negative samples (i.e., points that are difficult to distinguish from an anchor point). The key challenge toward using hard negatives is that contrastive methods must remain unsupervised, making it infeasible to adopt existing negative sampling strategies that use label information. In response, we develop a new class of unsupervised methods for selecting hard negative samples where the user can control the amount of hardness. A limiting case of this sampling results in a representation that tightly clusters each class, and pushes different classes as far apart as possible. The proposed method improves downstream performance across multiple modalities, requires only few additional lines of code to implement, and introduces no computational overhead.

## 1 Introduction

Owing to their empirical success, contrastive learning methods [7, 16] have become one of the most popular self-supervised approaches for learning representations [34, 46, 4]. For instance, on certain computer vision tasks unsupervised contrastive learning methods have even outperformed supervised pre-training [30, 20].

Contrastive learning relies on two key ingredients: notions of similar (positive)  $(x, x^+)$  and dissimilar (negative)  $(x, x^-)$  pairs of data points. The training objective, typically *noise-contrastive estimation* [15], guides the learned representation  $f$  to map positive pairs to nearby locations, and negative pairs farther apart; other objectives have also been considered [4]. The success of the associated methods depends on the informativeness of the positive and negative pairs, but their selection is particularly challenging without true label supervision.

Much research effort has addressed sampling strategies for positive pairs [2, 53, 47]. For image data, positive sampling strategies often apply transformations that preserve semantic content, e.g., jittering, random cropping, separating color channels, etc. [4, 6, 46]. Such transformations have also been effective in learning control policies from raw pixel data [42]. Positive sampling techniques have also been proposed for sentence, audio, and video data [28, 34, 37, 40].

Surprisingly, the choice of negative pairs has drawn much less attention in contrastive learning. Often, given an “anchor” point  $x$ , a “negative”  $x^-$  is simply sampled uniformly from the training data, independent of how informative it may be for the learned representation. In supervised and metric learning settings, “hard” (negative) examples can help guide a learning method to correct its mistakes more quickly [38, 41]. For representation learning, informative negative examples are intuitively those pairs that are mapped nearby but should be far apart. This idea is successfully applied in metric learning, where true pairs of dissimilar points are available, as opposed to unsupervised contrastive learning.

With this motivation, we address the challenge of selecting informative negatives for contrastive representation learning. In particular, we propose a solution that builds a tunable sampling distribution that prefers negative pairs whose representations are currently very similar. This solution faces two challenges: (1) we do not have access to any true label or dissimilarity information; and (2) we need an efficient sampling strategy for this tunable distribution. We overcome (1) by building on ideas from positive-unlabeled learning [13, 12], and (2) by designing an efficient, easy to implement importance sampling technique that incurs no computational overhead.

Our theoretical analysis shows that, as a function of the tuning parameter, the optimal representations for our new method place similar inputs in tight clusters, whilst spacing the clusters as far apart as possible. Empirically, our hard negative sampling strategy improves the downstream task performance for image, graph and text data, supporting that indeed, our negative examples are more informative.

**Contributions.** In summary, we make the following contributions:

1. We formulate to our knowledge the first sampling distribution for hard negative pairs for contrastive representation learning, and an efficient sampling strategy that takes into account the lack of true dissimilarity information;
2. We theoretically analyze the properties of the objective and optimal representations, showing that they capture desired goals for representations;
3. We empirically observe that the proposed sampling method improves the downstream task performance on image, graph and text data.

Before moving onto the problem formulation and our results, we summarize related work below.

## 1.1 Related Work

**Contrastive Representation Learning.** Various frameworks for contrastive learning of visual representations have been proposed, including SimCLR [4, 5], which uses augmented views of other items in a minibatch as negative samples, and MoCo [20, 6], which uses a momentum updated memory bank of old negative representations to enable the use of very large batches of negative samples. Many works study the role of positive pairs, and, e.g., propose to apply large perturbations for images [4, 6], or argue to minimize the mutual information within positive pairs, apart from relevant information for the ultimate prediction task [47]. Beyond visual data, contrastive methods have been developed for sentence embeddings [28], sequential data [34, 21], graph [44, 18, 27] and node representation learning [48], and learning representations from raw images for off-policy control [42]. The role of negative pairs has been much less studied. Chuang et al. [9] propose a method for “debiasing”, i.e., avoiding the “sampling bias” due to sometimes sampling false negatives that have the same label as the anchor. [24] exploit the specific temporal structure of video to generate negatives for object detection.

**Negative Mining in Deep Metric Learning.** As opposed to the contrastive representation learning literature, selection strategies for negative samples have been successfully applied in (deep) metric learning [38, 41, 17, 52, 14, 43]. Most of these works observe that it is helpful to use negative samples that are difficult for the current embedding to discriminate. [38] qualify this, observing that some examples are simply too hard, and propose selecting “semi-hard” negative samples. The well known importance of negative samples in metric learning, where (partial) true dissimilarity information is available, raises the question of negative samples in contrastive learning, the subject of this paper.

## 2 Contrastive Learning Setup

We begin with the setup and the idea of contrastive representation learning. We wish to learn an embedding  $f : \mathcal{X} \rightarrow \mathbb{S}^{d-1}/t$  that maps an observation  $x$  to a point on a hypersphere  $\mathbb{S}^{d-1}/t$  in  $\mathbb{R}^d$  of radius  $1/t$ , where  $t$  is the “temperature” scaling hyperparameter.

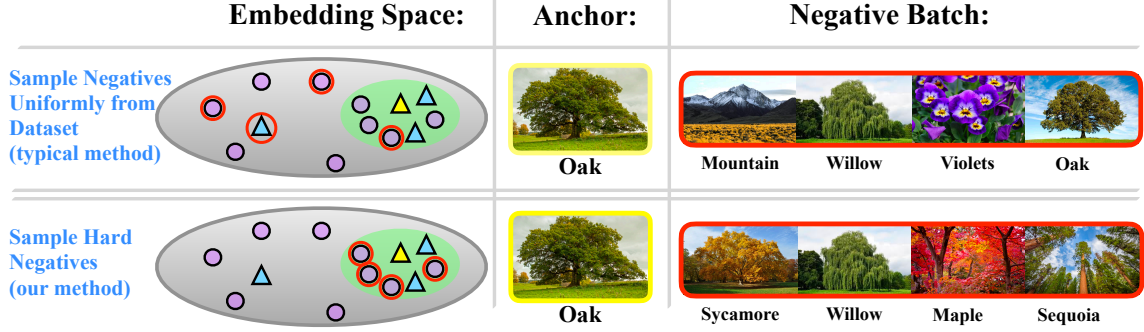


Figure 1: Schematic illustration of negative sampling methods for the example of classifying species of tree. Top row: uniformly samples negative examples (red rings); mostly focuses on very different data points from the anchor (yellow triangle), and may even sample examples from the same class (triangles, vs. circles). Bottom row: Hard negative sampling prefers examples that are (incorrectly) close to the anchor.

Following the setup of [1], we assume an underlying set of discrete latent classes  $\mathcal{C}$  that represent semantic content, so that similar pairs  $(x, x^+)$  have the same latent class. Denoting the distribution over latent classes by  $\rho(c)$  for  $c \in \mathcal{C}$ , we define the joint distribution  $p_{x,c}(x, c) = p(x|c)\rho(c)$  whose marginal  $p(x)$  we refer to simply as  $p$ , and assume  $\text{supp}(p) = \mathcal{X}$ . For simplicity, we assume  $\rho(c) = \tau^+$  is uniform, and let  $\tau^- = 1 - \tau^+$  be the probability of another class. Since the class-prior  $\tau^+$  is unknown in practice, it must either be treated as a hyperparameter, or estimated [8, 23].

Let  $h : \mathcal{X} \rightarrow \mathcal{C}$  be the true underlying hypothesis that assigns class labels to inputs. We write  $x \sim x'$  to denote the label equivalence relation  $h(x) = h(x')$ . We denote by  $p_x^+(x') = p(x'|h(x') = h(x))$ , the distribution over points with same label as  $x$ , and by  $p_x^-(x') = p(x'|h(x') \neq h(x))$ , the distribution over points with labels different from  $x$ . We drop the subscript  $x$  when the context is clear. Following the usual convention, we overload  $\sim$  and also write  $x \sim p$  to denote a point sampled from  $p$ .

For each data point  $x \sim p$ , the noise-contrastive estimation (NCE) objective [15] for learning the representation  $f$  uses a *positive* example  $x^+$  with the same label as  $x$ , and *negative* examples  $\{x_i^-\}_{i=1}^N$  with (supposedly) different labels,  $h(x_i^-) \neq h(x)$ , sampled from  $q$ :

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

Typically, we take  $Q = N$ , the number of negative samples. The negative sample distribution  $q$  is frequently chosen to be the marginal distribution  $p$ , or, in practice, an empirical approximation of it [46, 4, 6, 20, 6, 34, 21]. In this paper we ask: is there a better way to choose  $q$ ?

### 3 Hard Negative Sampling

In this section we describe our approach for hard negative sampling. We begin by asking *what makes a good negative sample?* To answer this question we adopt the following two guiding principles:

**Principle 1.**  $q$  should only sample “negatives”  $x_i^-$  whose labels differ from that of the anchor  $x$ .

**Principle 2.** The most useful negative samples are ones that the embedding currently believes to be similar to the anchor.

In short, negative samples that have different label from the anchor, but that are embedded nearby are likely to be most useful and provide significant gradient information during training. In metric learning there is access to true negative pairs, automatically fulfilling the first principle.

In contrastive learning there is no supervision, so upholding Principle 1 is impossible to do exactly. In this paper we propose a method that upholds Principle 1 approximately, and simultaneously combines this idea with the key additional conceptual ingredient of “hardness” (encapsulated in Principle 2). The level of “hardness” in our method can be smoothly adjusted, allowing the user to select the hardness that best trades-off between an improved learning signal from hard negatives, and the harm due to the correction of false negatives being only approximate. This is important since the hardest points are those closest to the anchor, and are expected to have a high propensity to have the same label. Therefore the damage from the approximation not removing all false negatives becomes larger for harder samples, creating the trade-off. As a special case of our method, when the hardness level is tuned fully down, we obtain the method proposed in [9] that only upholds Principle 1 (approximately) but not Principle 2. Finally, beyond Principles 1 and 2, we wish to design an efficient sampling method that does not add additional computational overhead during training.

### 3.1 Proposed Hard Sampling Method

Our first goal is to design a distribution  $q$  on  $\mathcal{X}$  that is allowed to depend on the embedding  $f$  and the anchor  $x$ . From  $q$  we sample a batch of negatives  $\{x_i^-\}_{i=1}^N$  according to the principles noted above. We propose sampling negatives from the distribution  $q_\beta^-$  defined as

$$q_\beta^-(x^-) := q_\beta(x^- | h(x) \neq h(x^-)), \quad \text{where} \quad q_\beta(x^-) \propto e^{\beta f(x)^\top f(x^-)} \cdot p(x^-),$$

for  $\beta \geq 0$ . The exponential term in  $q_\beta$  is an unnormalized von Mises–Fisher distribution with mean direction  $f(x)$  and “concentration parameter”  $\beta$  [29]. There are two key components to  $q_\beta^-$ , corresponding to each principle: 1) conditioning on the event  $\{h(x) \neq h(x^-)\}$  which guarantees that  $(x, x^-)$  correspond to different latent classes (Principle 1); 2) the concentration parameter  $\beta$  term controls the degree by which  $q_\beta$  up-weights points  $x^-$  that have large inner product (similarity) to the anchor  $x$  (Principle 2). Since  $f$  lies on the surface of a hypersphere of radius  $1/t$ , we have  $\|f(x) - f(x')\|^2 = 2/t^2 - 2f(x)^\top f(x')$  so preferring points with large inner product is equivalent to preferring points with small squared Euclidean distance.

Although we have designed  $q_\beta^-$  to have all of the desired components, it is not clear how to sample efficiently from it. To work towards a practical method, note that we can rewrite this distribution by adopting a PU-learning viewpoint [13, 12, 9]. That is, by conditioning on the event  $\{h(x) = h(x^-)\}$  we can split  $q_\beta(x^-)$  as

$$q_\beta(x^-) = \tau^- q_\beta^-(x^-) + \tau^+ q_\beta^+(x^-), \quad (2)$$

where  $q_\beta^+(x^-) = q_\beta(x^- | h(x) = h(x^-)) \propto e^{\beta f(x)^\top f(x^-)} \cdot p^+(x^-)$ . Rearranging equation 2 yields a formula  $q_\beta^-(x^-) = (q_\beta(x^-) - \tau^+ q_\beta^+(x^-)) / \tau^-$  for the negative sampling distribution  $q_\beta^-$  in terms of two distributions that are tractable since we have samples from  $p$  and can approximate samples from  $p^+$  using a set of semantics-preserving transformations, as is typical in contrastive learning methods.

It is possible to generate samples from  $q_\beta$  and (approximately from)  $q_\beta^+$  using rejection sampling. However, rejection sampling involves adding an additional step to sampling batches, which could be slow. Instead note that fixing the number  $Q$  and taking the limit  $N \rightarrow \infty$  in the objective (1) yields,

$$\mathcal{L}(f, q) = \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^\top f(x^+)}}{e^{f(x)^\top f(x^+)} + Q \mathbb{E}_{x^- \sim q} [e^{f(x)^\top f(x^-)}]} \right]. \quad (3)$$

The original objective (1) can be viewed as a finite negative sample approximation to  $\mathcal{L}(f, q)$

(note implicitly  $\mathcal{L}(f, q)$  depends on  $Q$ ). Inserting  $q = q_\beta^-$  and using the rearrangement of equation (2) we obtain the following hardness-biased objective:

$$\mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + \frac{Q}{\tau^-} (\mathbb{E}_{x^- \sim q_\beta} [e^{f(x)^T f(x^-)}] - \tau^+ \mathbb{E}_{v \sim q_\beta^+} [e^{f(x)^T f(v)}])} \right]. \quad (4)$$

This objective suggests that we need only to approximate *expectations*  $\mathbb{E}_{x^- \sim q_\beta} [e^{f(x)^T f(x^-)}]$  and  $\mathbb{E}_{v \sim q_\beta^+} [e^{f(x)^T f(v)}]$  over  $q_\beta$  and  $q_\beta^+$  (rather than explicitly sampling). This can be achieved using classical Monte-Carlo importance sampling techniques using samples from  $p$  and  $p^+$  as follows:

$$\begin{aligned} \mathbb{E}_{x^- \sim q_\beta} [e^{f(x)^T f(x^-)}] &= \mathbb{E}_{x^- \sim p} [e^{f(x)^T f(x^-)} q_\beta / p] = \mathbb{E}_{x^- \sim p} [e^{(\beta+1)f(x)^T f(x^-)} / Z_\beta], \\ \mathbb{E}_{v \sim q_\beta^+} [e^{f(x)^T f(v)}] &= \mathbb{E}_{v \sim p^+} [e^{f(x)^T f(v)} q_\beta^+ / p^+] = \mathbb{E}_{v \sim p^+} [e^{(\beta+1)f(x)^T f(v)} / Z_\beta^+], \end{aligned}$$

where  $Z_\beta, Z_\beta^+$  are the partition functions of  $q_\beta$  and  $q_\beta^+$  respectively. The right hand terms readily admit empirical approximations by replacing  $p$  and  $p^+$  with  $\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i^-}(x)$  and  $\hat{p}^+(x) = \frac{1}{M} \sum_{i=1}^M \delta_{x_i^+}(x)$  respectively ( $\delta_w$  denotes the Dirac delta function centered at  $w$ ). The only unknowns left are the partition functions,  $Z_\beta = \mathbb{E}_{x^- \sim p} [e^{\beta f(x)^T f(x^-)}]$  and  $Z_\beta^+ = \mathbb{E}_{x^+ \sim p^+} [e^{\beta f(x)^T f(x^+)}]$  which themselves are expectations over  $p$  and  $p^+$  and therefore admit empirical estimates,

$$\hat{Z}_\beta = \frac{1}{N} \sum_{i=1}^N e^{\beta f(x_i^-)^T f(x_i^-)}, \quad \hat{Z}_\beta^+ = \frac{1}{M} \sum_{i=1}^M e^{\beta f(x_i^+)^T f(x_i^+)}.$$

It is important to emphasize the simplicity of the implementation of our proposed approach. Since we propose to reweight the objective instead of modifying the sampling procedure, only two extra lines of code are needed to implement our approach, with no additional computational overhead. PyTorch-style pseudocode for the objective is given in Figure 7 in Appendix D.

## 4 Analysis of Hard Negative Sampling

### 4.1 Hard Sampling Interpolates Between Marginal and Worst-Case Negatives

Intuitively, the concentration parameter  $\beta$  controls the level of “hardness” of the negative samples. It is clear that taking  $\beta = 0$  (the “least hard” negatives) recovers the population distribution  $q_0 = p$  (and  $q_0^-$  recovers the debiasing method of [9]). But what interpretation does large  $\beta$  admit? Specifically, what does the distribution  $q_\beta^-$  converge to in the limit  $\beta \rightarrow \infty$ , if anything? We show that in the limit  $q_\beta^-$  approximates an inner solution to the following zero-sum two player game.

$$\min_f \sup_{q \in \Pi} \left\{ \mathcal{L}(f, q) = \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Q \mathbb{E}_{x^- \sim q} [e^{f(x)^T f(x^-)}]} \right] \right\}. \quad (5)$$

where  $\Pi = \{q = q(\cdot; x, f) : \text{supp}(q(\cdot; x, f)) \subseteq \{x' \in \mathcal{X} : x' \approx x\}, \forall x \in \mathcal{X}\}$  is the set of distributions with support that is disjoint from points with the same class as  $x$  (without loss of generality we assume  $\{x' \in \mathcal{X} : x' \approx x\}$  is non-empty). Since  $q = q(\cdot; x, f)$  depends on  $x$  and  $f$  it can be thought of as a family of distributions. The formal statement is as follows.

**Proposition 3.** *Let  $\mathcal{L}^*(f) = \sup_{q \in \Pi} \mathcal{L}(f, q)$ . Then for any  $t > 0$  and  $f : \mathcal{X} \rightarrow \mathbb{S}^{d-1}/t$  we observe the convergence  $\mathcal{L}(f, q_\beta^-) \rightarrow \mathcal{L}^*(f)$  as  $\beta \rightarrow \infty$ .*

*Proof.* See Appendix A.1.  $\square$

To develop a better intuitive understanding of the worst case negative distribution objective  $\mathcal{L}^*(f) = \sup_{q \in \Pi} \mathcal{L}(f, q)$ , we note that the supremum can be characterized analytically. Indeed,

$$\begin{aligned} \sup_{q \in \Pi} \mathcal{L}(f, q) &= -\mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} f(x)^T f(x^+) + \sup_{q \in \Pi} \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \log \left\{ e^{f(x)^T f(x^+)} + Q \mathbb{E}_{x^- \sim q} [e^{f(x)^T f(x^-)}] \right\} \\ &= -\mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} f(x)^T f(x^+) + \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \log \left\{ e^{f(x)^T f(x^+)} + Q \cdot \sup_{q \in \Pi} \mathbb{E}_{x^- \sim q} [e^{f(x)^T f(x^-)}] \right\}. \end{aligned}$$

The supremum over  $q$  can be pushed inside the expectation since  $q$  is a family of distribution indexed by  $x$ , reducing the problem to maximizing  $\mathbb{E}_{x^- \sim q} [e^{f(x)^T f(x^-)}]$ , which is solved by any  $q^*$  whose support is a subset of  $\arg \sup_{x^-: x^- \sim x} e^{f(x)^T f(x^-)}$  if the supremum is attained. If it is not attained, distributions supported on approximations of the supremum yield approximate solutions to  $\sup_{q \in \Pi} \mathcal{L}(f, q)$ . However, computing such points involves maximizing a neural network. Instead of taking this challenging route,  $q_\beta$  represents a tractable approximation for large  $\beta$  (Proposition 3).

## 4.2 Embeddings on the Hypersphere for Worst-Case Negative Samples

What desirable properties does an optimal contrastive embedding (global minimizer of  $\mathcal{L}$ ) possess that make the representation generalizable? To study this question, we analyze the distribution of an optimal embedding  $f^*$  on the hypersphere when negatives are sampled from the adversarial worst-case distribution. We consider a different limiting viewpoint of objective (1) as the number of negative samples  $N \rightarrow \infty$ . Following the formulation of [50] we take  $Q = N$  in (1), and subtract  $\log N$ . This changes neither the set of minimizers, nor the geometry of the loss surface. Taking the number of negative samples  $N \rightarrow \infty$  yields the limiting objective,

$$\mathcal{L}_\infty(f, q) = \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^T f(x^+)}}{\mathbb{E}_{x^- \sim q} [e^{f(x)^T f(x^-)}]} \right]. \quad (6)$$

**Theorem 4.** Suppose the downstream task is classification (i.e.  $\mathcal{C}$  is finite), and let  $\mathcal{L}_\infty^*(f) = \sup_{q \in \Pi} \mathcal{L}_\infty(f, q)$ . The infimum  $\inf_{f: \text{measurable}} \mathcal{L}_\infty^*(f)$  is attained, and any  $f^*$  achieving the global minimum is such that  $f^*(x) = f^*(x^+)$  almost surely. Furthermore, letting  $\mathbf{v}_c = f^*(x)$  for any  $x$  such that  $h(x) = c$  (so  $\mathbf{v}_c$  is well defined up to a set of  $x$  of measure zero),  $f^*$  is characterized as being any solution to the following ball-packing problem,

$$\max_{\{\mathbf{v}_c \in \mathbb{S}^{d-1}/t\}_{c \in \mathcal{C}}} \sum_{c \in \mathcal{C}} \rho(c) \cdot \min_{c' \neq c} \|\mathbf{v}_c - \mathbf{v}_{c'}\|^2. \quad (7)$$

*Proof.* See Appendix A.2.  $\square$

**Interpretation:** The first component of the result is that  $f^*(x) = f^*(x^+)$  almost surely for an optimal  $f^*$ . That is, an optimal embedding  $f^*$  must be invariant across pairs of similar inputs  $x, x^+$ . The second component is characterizing solutions via the classical Tammes Ball-Packing Problem from geometry [45] (Eq. 7) that has only been solved exactly for uniform  $\rho$ , for specific of  $|\mathcal{C}|$  and typically for  $\mathbb{S}^2$ , [39, 32, 45]. When the distribution  $\rho$  over classes is uniform this problem is solved by a set of  $|\mathcal{C}|$  points on the hypersphere such that the average squared- $\ell_2$  distance from a point to the nearest other point is as large as possible. In other words, suppose we wish to place  $|\mathcal{C}|$  number of balls<sup>1</sup> on  $\mathbb{S}^{d-1}$  so that they do not intersect. Then solutions to Tammes' Problem (7) expresses (twice) the largest possible average squared radius that the balls can have. So, we have a ball-packing problem where instead of trying to pack as many balls as possible of

<sup>1</sup>For a manifold  $\mathcal{M} \subseteq \mathbb{R}^d$ , we say  $C \subset \mathcal{M}$  is a ball if it is connected, and there exists a Euclidean ball  $B = \{x \in \mathbb{R}^d : \|x\|_2 \leq R\}$  for which  $C = \mathcal{M} \cap B$ .

a fixed size, we aim to pack a fixed number of balls (one for each class) to have as big radii as possible. Non-uniform  $\rho$  adds importance weights to each fixed ball. This interpretation shows that solutions of the problem  $\min_f \mathcal{L}^*(f)$  are a type of maximum margin clustering.

## 5 Empirical Results

Next, we evaluate our hard negative sampling method empirically, and apply it as a modification to state-of-the-art contrastive methods on image, graph, and text data. In all cases, we use one positive example  $x^+$  per data point and treat the class-prior  $\tau^+$  as a hyperparameter.

### 5.1 Image Representations

We begin by testing the hard sampling method on vision tasks using the STL10, CIFAR100 and CIFAR10 data. We use SimCLR [4] as the baseline method. The results in Figure 2 show consistent improvement over SimCLR ( $q = p$ ) and the particular case of our method with  $\beta = 0$  proposed in [9] (called debiasing) on STL10 and CIFAR100. For  $N = 510$  negative examples per data point we observe absolute improvements over the best debiased baseline of 1.9% and 3.2%, respectively, on CIFAR100 and STL10, and absolute improvements of 3% and 7.3% over SimCLR. On CIFAR10 there is a slight improvement for smaller  $N$ , which disappears at larger  $N$ . See Section 6.1 for further discussion, and Appendix D.1 for setup details.

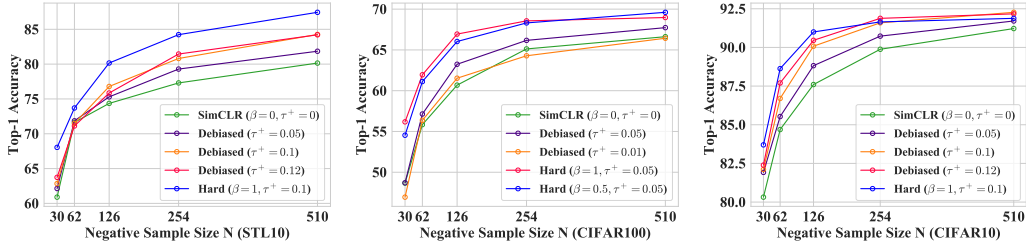


Figure 2: **Classification accuracy on downstream tasks.** Embeddings trained using hard ( $\beta > 0$ ), debiased ( $\beta = 0$ ), and standard ( $\beta = 0, \tau^+ = 0$ ) versions of SimCLR, and evaluated using linear readout accuracy.

### 5.2 Graph Representations

Second, we consider hard negative sampling in the context of learning graph representations. We use the state-of-the-art InfoGraph method introduced by [44] as the baseline, which is suitable for downstream graph-level classification. The objective is of a slightly different form from the NCE loss: it has two terms, the first encourages alignment of embeddings of the whole graph and nodes from the same graph, while the second contrasts embeddings for a graph with node embeddings from different graphs. We only modify the second term, using a generalization of the formulation presented in Section 3 (See Appendix B for details). In doing so, we illustrate that it is easy to adapt our hard sampling method to other contrastive objectives beyond the NCE loss.

Figure 3 shows the results of fine-tuning an SVM [3, 10] on the fixed, learned embedding for a range of different values of  $\beta$ . Hard sampling does as well as InfoGraph in all cases, and better in 6 out of 8 cases. For ENZYMES and REDDIT, hard negative samples improve the accuracy by 3.2% and 2.4%, respectively, for DD and PTC by 1 – 2%, and for IMDB-B and MUTAG by at least 0.5%. Usually, multiple different choices of  $\beta > 0$  were competitive with the InfoGraph baseline: 17 out of the 24 values of  $\beta > 0$  tried (across all 8 datasets) achieve accuracy as high or better than InfoGraph (i.e. the special case  $\beta = 0$ ).

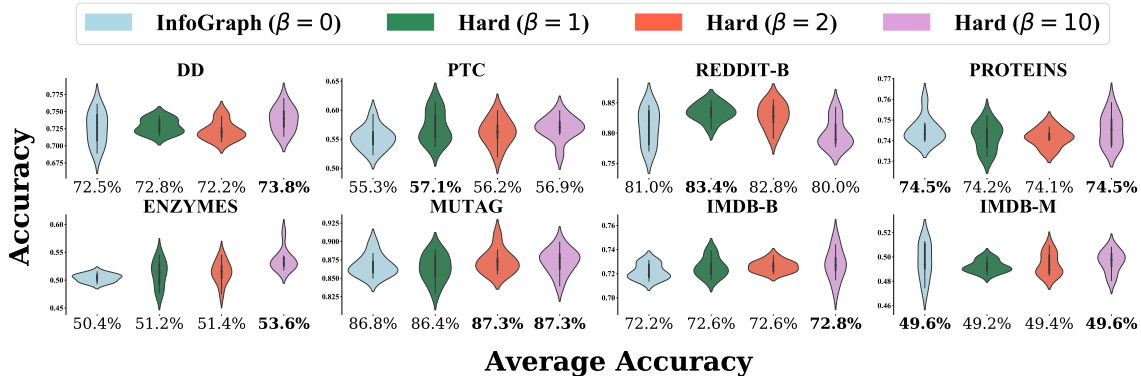


Figure 3: **Classification accuracy on downstream tasks.** We compare graph representations on four classification tasks. Accuracies are obtained by fine-tuning an SVM readout function, and are the average of 10 runs, each using 10-fold cross validation. Results in **bold** indicate best performer.

Objective	MR	CR	SUBJ	MPQA	TREC	MSRP	
						(Acc)	(F1)
QT ( $\beta = 0, \tau^+ = 0$ )	76.8	81.3	86.6	93.4	<b>89.8</b>	73.6	81.8
Debiased ( $\tau^+ = 0.01$ )	76.2	82.9	86.9	<b>93.7</b>	89.1	<b>74.7</b>	<b>82.7</b>
Hard ( $\beta = 1, \tau^+ = 0$ )	77.1	82.5	<b>87.0</b>	92.9	89.2	73.9	82.2
Hard ( $\beta = 2, \tau^+ = 0$ )	<b>77.4</b>	<b>83.6</b>	86.8	93.4	88.7	73.5	82.0

Table 1: **Classification accuracy on downstream tasks.** Sentence representations are learned using quick-thoughts (QT) vectors on the BookCorpus dataset and evaluated on six classification tasks. Evaluation of binary classification tasks (MR, CR, SUBJ, MPQA) uses 10-fold cross validation.

### 5.3 Sentence Representations

Third, we test hard negative sampling on learning representations of sentences using the *quick-thoughts* (QT) vectors framework introduced by [28], which uses adjacent sentences (before/after) as positive samples. Embeddings are trained using the unlabeled BookCorpus dataset [26], and evaluated following the protocol of [28] on six downstream tasks. The results are reported in Table 1. Hard sampling outperforms or equals the QT baseline in 5 out of 6 cases, the  $\beta = 0$  (“debiased” [9]) baseline in 4 out of 6, and both in 3 out of 6 cases. Setting  $\tau^+ > 0$  led to numerical issues in optimization for hard sampling.

## 6 A Closer Look at Hard Sampling

### 6.1 Are Harder Samples Necessarily Better?

By setting  $\beta$  to large values, one can focus on only the hardest samples in a training batch. But is this desirable? Figure 4 (left, middle) shows that for vision problems, taking larger  $\beta$  does not necessarily lead to better representations. In contrast, when one uses true positive pairs during training (green curve, uses label information for positive but not negative pairs), the downstream performance monotonically increases with  $\beta$  until convergence (Figure 4, middle). Interestingly, this is achieved without using label information for the negative pairs. This observation suggests an explanation for why bigger  $\beta$  hurts performance in practice. Conditioning on the event  $\{h(x) \neq h(x^-)\}$  using the true  $p^+$  corrects for sampling  $x^-$  with the same label as  $x$ . However, since in practice we approximate  $p^+$  using a set of data transformations, we can only partially correct. This is harmful for large  $\beta$  since this regime strongly prefers  $x^-$  for which  $f(x^-)$  is close to  $f(x)$ , many of whom will have the same label as  $x$  if not corrected for. By annealing  $\beta$  (gradually decreasing  $\beta$  to 0 throughout training) it is possible to be more robust to the choice of initial  $\beta$ , although downstream accuracy is lower than using the best value of  $\beta$ .

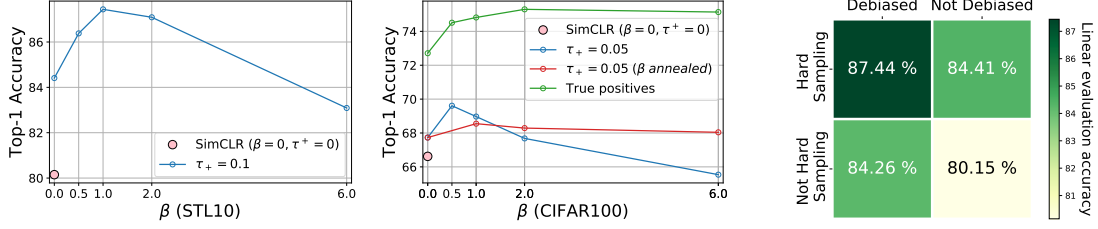


Figure 4: Left: the effect of varying concentration parameter  $\beta$  on linear readout accuracy. Middle: linear readout accuracy as concentration parameter  $\beta$  varies, in the case of contrastive learning (fully unsupervised), using true positive samples (uses label information), and an annealing method that improves robustness to the choice of  $\beta$  (see Appendix D.1 for details). Right: STL10 linear readout accuracy for hard sampling with and without debiasing, and non-hard sampling ( $\beta = 0$ ) with and without debiasing. Best results come from using both simultaneously.

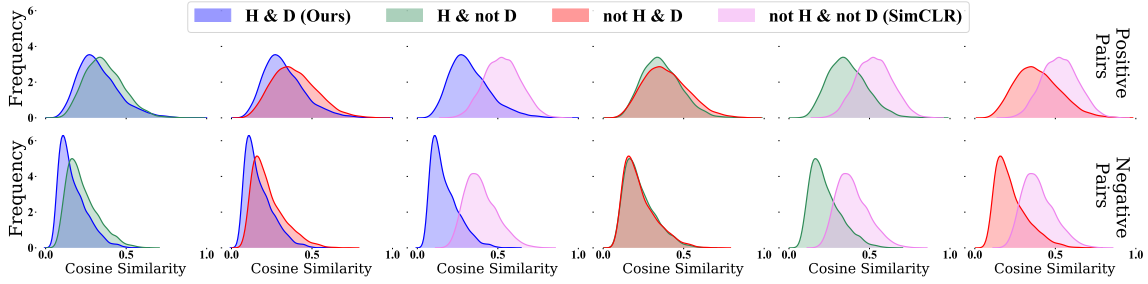


Figure 5: Histograms of cosine similarity of pairs of points with the same label (top) and different labels (bottom) for embeddings trained on STL10 with four different objectives. H=Hard Sampling ( $\beta > 0$ ), D=Debiasing (conditioning on  $\{h(x) \neq h(x^-)\}$ ). Histograms overlaid pairwise to allow for convenient comparison.

## 6.2 Does Avoiding False Negatives Improve Hard Sampling?

Our proposed hard negative sampling method conditions on the event  $\{h(x) \neq h(x^-)\}$  in order to avoid false negatives (termed “debiasing” [9]). But does this help? To test this, we train four embeddings: hard sampling with and without debiasing, and uniform sampling ( $\beta = 0$ ) with and without debiasing. The results in Figure 4 (right) show that hard sampling with debiasing obtains the highest linear readout accuracy on STL10, only using hard sampling or only debiasing yields (in this case) similar accuracy. All improve over the SimCLR baseline.

Figure 5 compares the histograms of cosine similarities of positive and negative pairs for the four learned representations. The representation trained with hard negatives and debiasing assigns much lower similarity score to a pair of negative samples than other methods. On the other hand, the SimCLR baseline assigns higher cosine similarity scores to pairs of positive samples. However, to discriminate positive and negative pairs, a key property is the amount of *overlap* of positive and negative histograms. Our hard sampling method appears to obtain less overlap than SimCLR, by better trading off higher dissimilarity of negative pairs with less similarity of positive pairs. Similar tradeoffs are observed for the debiased objective, and hard sampling without debiasing.

## 7 Conclusion

We have initiated a discussion on the value of hard negative sampling for unsupervised contrastive representation learning. Our work connects two major lines of work: contrastive learning, and negative mining in metric learning. Doing so required overcoming an apparent roadblock: negative mining in metric learning uses pairwise similarity information as a core component, while contrastive learning is unsupervised. By showing this roadblock can be passed, we open the door to the possibility of designing other negative sampling strategies.

## References

1. Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. In *Int. Conference on Machine Learning (ICML)*, pages 5628–5637, 2019.
2. Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100, 1998.
3. Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152, 1992.
4. Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Int. Conference on Machine Learning (ICML)*, pages 10709–10719, 2020.
5. Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv:2006.10029*, 2020.
6. Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv:2003.04297*, 2020.
7. Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 539–546, 2005.
8. Marthinus Christoffel, Gang Niu, and Masashi Sugiyama. Class-prior estimation for learning from positive and unlabeled data. In *Asian Conference on Machine Learning*, pages 221–236, 2016.
9. Ching-Yao Chuang, Joshua Robinson, Lin Yen-Chen, Antonio Torralba, and Stefanie Jegelka. Debaised Contrastive Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
10. Corinna Cortes and Vladimir Vapnik. Support-Vector Networks. *Machine learning*, 20(3): 273–297, 1995.
11. Bill Dolan, Chris Quirk, and Chris Brockett. Unsupervised construction of large paraphrase corpora: Exploiting massively parallel news sources. In *Proceedings of the 20th international conference on Computational Linguistics*, page 350. Association for Computational Linguistics, 2004.
12. Marthinus C Du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 703–711, 2014.
13. Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In *ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 213–220, 2008.
14. Weifeng Ge. Deep metric learning with hierarchical triplet loss. In *Europ. Conference on Computer Vision (ECCV)*, pages 269–285, 2018.
15. Michael Gutmann and Aapo Hyvärinen. Noise-Contrastive Estimation: A new estimation principle for unnormalized statistical models. In *Proc. Int. Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 297–304, 2010.

16. Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1735–1742, 2006.
17. Ben Harwood, Vijay Kumar BG, Gustavo Carneiro, Ian Reid, and Tom Drummond. Smart mining for deep metric learning. In *Int. Conference on Computer Vision (ICCV)*, pages 2821–2829, 2017.
18. Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive multi-view representation learning on graphs. In *Int. Conference on Machine Learning (ICML)*, pages 3451–3461, 2020.
19. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
20. Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9729–9738, 2020.
21. Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. In *Int. Conference on Machine Learning (ICML)*, pages 6661–6671, 2020.
22. Minqing Hu and Bing Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177, 2004.
23. Shantanu Jain, Martha White, and Predrag Radivojac. Estimating the class prior and posterior from noisy positives and unlabeled data. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2693–2701, 2016.
24. SouYoung Jin, Aruni RoyChowdhury, Huaizu Jiang, Ashish Singh, Aditya Prasad, Deep Chakraborty, and Erik Learned-Miller. Unsupervised hard example mining from videos for improved object detection. In *Europ. Conference on Computer Vision (ECCV)*, pages 307–324, 2018.
25. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Int. Conf. on Learning Representations (ICLR)*, 2015.
26. Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-Thought Vectors. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3294–3302, 2015.
27. Yujia Li, Chenjie Gu, Thomas Dullien, Oriol Vinyals, and Pushmeet Kohli. Graph matching networks for learning the similarity of graph structured objects. In *Int. Conference on Machine Learning (ICML)*, pages 3835–3845, 2019.
28. Lajanugen Logeswaran and Honglak Lee. An efficient framework for learning sentence representations. In *Int. Conf. on Learning Representations (ICLR)*, 2018.
29. K. V. Mardia and P. Jupp. *Directional Statistics*. John Wiley and Sons Ltd., second edition, 2000.
30. Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6707–6717, 2020.
31. Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *Graph Representation Learning and Beyond, ICML Workshop*, 2020.

32. Oleg R Musin and Alexey S Tarasov. The tammes problem for  $n=14$ . *Experimental Mathematics*, 24(4):460–468, 2015.
33. Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-GAN: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 271–279, 2016.
34. Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv:1807.03748*, 2018.
35. Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd annual meeting on Association for Computational Linguistics*, page 271. Association for Computational Linguistics, 2004.
36. Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd annual meeting on association for computational linguistics*, pages 115–124. Association for Computational Linguistics, 2005.
37. Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases. *arXiv:2007.13916*, 2020.
38. Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 815–823, 2015.
39. K Schütte and BL Van der Waerden. Auf welcher kugel haben 5, 6, 7, 8 oder 9 punkte mit mindestabstand eins platz? *Mathematische Annalen*, 123(1):96–124, 1951.
40. Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1134–1141, 2018.
41. Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4004–4012, 2016.
42. Aravind Srinivas, Michael Laskin, and Pieter Abbeel. CURL: Contrastive unsupervised representations for reinforcement learning. In *Int. Conference on Machine Learning (ICML)*, pages 10360–10371, 2020.
43. Yumin Suh, Bohyung Han, Wonsik Kim, and Kyoung Mu Lee. Stochastic class-based hard example mining for deep metric learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7251–7259, 2019.
44. Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. InfoGraph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *Int. Conf. on Learning Representations (ICLR)*, 2020.
45. Pieter Merkus Lambertus Tammes. On the origin of number and arrangement of the places of exit on the surface of pollen-grains. *Recueil des travaux botaniques néerlandais*, 27(1):1–84, 1930.
46. Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive Multiview Coding. In *Europ. Conference on Computer Vision (ECCV)*, pages 770–786, 2019.
47. Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? *arXiv:2005.10243*, 2020.

- 48. Petar Velickovic, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. Deep Graph Infomax. In *Int. Conf. on Learning Representations (ICLR)*, 2019.
- 49. Ellen M Voorhees and Donna Harman. Overview of trec 2002. 2002.
- 50. Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Int. Conference on Machine Learning (ICML)*, pages 9574–9584, 2020.
- 51. Janyce Wiebe, Theresa Wilson, and Claire Cardie. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2-3):165–210, 2005.
- 52. Chao-Yuan Wu, R Manmatha, Alexander J Smola, and Philipp Krahenbuhl. Sampling matters in deep embedding learning. In *Int. Conference on Computer Vision (ICCV)*, pages 2840–2848, 2017.
- 53. Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning. *arXiv:1304.5634*, 2013.
- 54. Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *Int. Conf. on Learning Representations (ICLR)*, 2019.

## A Analysis of Hard Sampling

### A.1 Hard Sampling Interpolates Between Marginal and Worst-Case Negatives

We begin by proving Proposition 3. Recall that the proposition stated the following.

**Proposition 5.** *Let  $\mathcal{L}^*(f) = \sup_{q \in \Pi} \mathcal{L}(f, q)$ . Then for any  $t > 0$  and measurable  $f : \mathcal{X} \rightarrow \mathbb{S}^{d-1}/t$  we observe the convergence  $\mathcal{L}(f, q_\beta^-) \rightarrow \mathcal{L}^*(f)$  as  $\beta \rightarrow \infty$ .*

*Proof.* Consider the following essential supremum,

$$M(x) = \operatorname{ess\,sup}_{x^- \in \mathcal{X}: x^- \sim x} f(x)^T f(x^-) = \sup\{m > 0 : m \geq f(x)^T f(x^-) \text{ a.s. for } x^- \sim p^-\}.$$

The second inequality holds since  $\operatorname{supp}(p) = \mathcal{X}$ . We may rewrite

$$\begin{aligned} \mathcal{L}^*(f) &= \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Qe^{M(x)}} \right], \\ \mathcal{L}(f, q_\beta^-) &= \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Q\mathbb{E}_{x^- \sim q_\beta^-}[e^{f(x)^T f(x^-)}]} \right]. \end{aligned}$$

The difference between these two terms can be bounded as follows,

$$\begin{aligned} \left| \mathcal{L}^*(f) - \mathcal{L}(f, q_\beta^-) \right| &\leq \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left| -\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Qe^{M(x)}} + \log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Q\mathbb{E}_{x^- \sim q_\beta^-}[e^{f(x)^T f(x^-)}]} \right| \\ &= \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left| \log \left( e^{f(x)^T f(x^+)} + Q\mathbb{E}_{x^- \sim q_\beta^-}[e^{f(x)^T f(x^-)}] \right) - \log \left( e^{f(x)^T f(x^+)} + Qe^{M(x)} \right) \right| \\ &\leq \frac{e^{1/t}}{Q+1} \cdot \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left| e^{f(x)^T f(x^+)} + Q\mathbb{E}_{x^- \sim q_\beta^-}[e^{f(x)^T f(x^-)}] - e^{f(x)^T f(x^+)} - Qe^{M(x)} \right| \\ &= \frac{e^{1/t}Q}{Q+1} \cdot \mathbb{E}_{x \sim p} \left| \mathbb{E}_{x^- \sim q_\beta^-}[e^{f(x)^T f(x^-)}] - e^{M(x)} \right| \\ &\leq e^{1/t} \cdot \mathbb{E}_{x \sim p} \mathbb{E}_{x^- \sim q_\beta^-} \left| e^{M(x)} - e^{f(x)^T f(x^-)} \right| \end{aligned}$$

where for the second inequality we have used the fact that  $f$  lies on the hypersphere of radius  $1/t$  to restrict the domain of the logarithm to values greater than  $(Q+1)e^{-1/t}$ . Because of this the logarithm is Lipschitz with parameter  $e^{1/t}/(Q+1)$ . Using again the fact that  $f$  lies on the hypersphere we know that  $|f(x)^T f(x^-)| \leq 1/t^2$  and hence have the following inequality,

$$\mathbb{E}_{x \sim p} \mathbb{E}_{q_\beta^-} \left| e^{M(x)} - e^{f(x)^T f(x^-)} \right| \leq e^{1/t^2} \mathbb{E}_{x \sim p} \mathbb{E}_{q_\beta^-} \left| M(x) - f(x)^T f(x^-) \right|$$

Let us consider the inner expectation  $E_\beta(x) = \mathbb{E}_{q_\beta^-} \left| M(x) - f(x)^T f(x^-) \right|$ . Note that since  $f$  is bounded,  $E_\beta(x)$  is uniformly bounded in  $x$ . Therefore, in order to show the convergence

$\mathcal{L}(f, q_\beta^-) \rightarrow \mathcal{L}^*(f)$  as  $\beta \rightarrow \infty$ , it suffices by the dominated convergence theorem to show that  $E_\beta(x) \rightarrow 0$  pointwise as  $\beta \rightarrow \infty$  for arbitrary fixed  $x \in \mathcal{X}$ .

From now on we denote  $M = M(x)$  for brevity, and consider a fixed  $x \in \mathcal{X}$ . From the definition of  $q_\beta^-$  it is clear that  $q_\beta^- \ll p^-$ . That is, since  $q_\beta^- = c \cdot p^-$  for some (non-constant)  $c$ , it is absolutely continuous with respect to  $p^-$ . So  $M(x) \geq f(x)^T f(x^-)$  almost surely for  $x^- \sim q_\beta^-$ , and we may therefore drop the absolute value signs from our expectation. Define the following event  $\mathcal{G}_\varepsilon = \{x^- : f(x)^T f(x^-) \geq M - \varepsilon\}$  where  $\mathcal{G}$  refers to a “good” event. Define its complement  $\mathcal{B}_\varepsilon = \mathcal{G}_\varepsilon^c$  where  $\mathcal{B}$  is for “bad”. For a fixed  $x \in \mathcal{X}$  and  $\varepsilon > 0$  consider,

$$\begin{aligned} E_\beta(x) &= \mathbb{E}_{x^- \sim q_\beta^-} \left| M(x) - f(x)^T f(x^-) \right| \\ &= \mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{G}_\varepsilon) \cdot \mathbb{E}_{x^- \sim q_\beta^-} \left[ \left| M(x) - f(x)^T f(x^-) \right| \mid \mathcal{G}_\varepsilon \right] + \mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon) \cdot \mathbb{E}_{x^- \sim q_\beta^-} \left[ \left| M(x) - f(x)^T f(x^-) \right| \mid \mathcal{B}_\varepsilon \right] \\ &\leq \mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{G}_\varepsilon) \cdot \varepsilon + 2\mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon) \\ &\leq \varepsilon + 2\mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon). \end{aligned}$$

We need to control  $\mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon)$ . Expanding,

$$\mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon) = \int_{\mathcal{X}} \mathbf{1} \left\{ f(x)^T f(x^-) < M(x) - \varepsilon \right\} \frac{e^{\beta f(x)^T f(x^-)} \cdot p^-(x^-)}{Z_\beta} dx^-$$

where  $Z_\beta = \int_{\mathcal{X}} e^{\beta f(x)^T f(x^-)} p^-(x^-) dx^-$  is the partition function of  $q_\beta^-$ . We may bound this expression by,

$$\begin{aligned} \int_{\mathcal{X}} \mathbf{1} \left\{ f(x)^T f(x^-) < M - \varepsilon \right\} \frac{e^{\beta(M-\varepsilon)} \cdot p^-(x^-)}{Z_\beta} dx^- &\leq \frac{e^{\beta(M-\varepsilon)}}{Z_\beta} \int_{\mathcal{X}} \mathbf{1} \left\{ f(x)^T f(x^-) < M - \varepsilon \right\} p^-(x^-) dx^- \\ &= \frac{e^{\beta(M-\varepsilon)}}{Z_\beta} \mathbb{P}_{x^- \sim p^-}(\mathcal{B}_\varepsilon) \\ &\leq \frac{e^{\beta(M-\varepsilon)}}{Z_\beta} \end{aligned}$$

Note that

$$Z_\beta = \int_{\mathcal{X}} e^{\beta f(x)^T f(x^-)} p^-(x^-) dx^- \geq e^{\beta(M-\varepsilon/2)} \mathbb{P}_{x^- \sim p^-}(f(x)^T f(x^-) \geq M - \varepsilon/2).$$

By the definition of  $M = M(x)$  the probability  $\rho_\varepsilon = \mathbb{P}_{x^- \sim p^-}(f(x)^T f(x^-) \geq M - \varepsilon/2) > 0$ , and we may therefore bound,

$$\begin{aligned} \mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon) &= \frac{e^{\beta(M-\varepsilon)}}{e^{\beta(M-\varepsilon/2)} \rho_\varepsilon} \\ &= e^{-\beta\varepsilon/2} / \rho_\varepsilon \\ &\longrightarrow 0 \text{ as } \beta \rightarrow \infty. \end{aligned}$$

We may therefore take  $\beta$  to be sufficiently big so as to make  $\mathbb{P}_{x^- \sim q_\beta^-}(\mathcal{B}_\varepsilon) \leq \varepsilon$  and therefore  $E_\beta(x) \leq 3\varepsilon$ . In other words,  $E_\beta(x) \longrightarrow 0$  as  $\beta \rightarrow \infty$ .  $\square$

## A.2 Optimal Embeddings on the Hypersphere for Worst-Case Negative Samples

In order to study properties of global optima of the contrastive objective using the adversarial worst case hard sampling distribution recall that we have the following limiting objective,

$$\mathcal{L}_\infty(f, q) = \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \left[ -\log \frac{e^{f(x)^\top f(x^+)}}{\mathbb{E}_{x^- \sim q_\beta} [e^{f(x)^\top f(x^-)}]} \right]. \quad (8)$$

We may separate the logarithm of a quotient into the sum of two terms plus a constant,

$$\mathcal{L}_\infty(f, q) = \mathcal{L}_{\text{align}}(f) + \mathcal{L}_{\text{unif}}(f, q) - 1/t^2$$

where  $\mathcal{L}_{\text{align}}(f) = \mathbb{E}_{x, x^+} \|f(x) - f(x^+)\|^2/2$  and  $\mathcal{L}_{\text{unif}}(f, q) = \mathbb{E}_{x \sim p} \log \mathbb{E}_{x^- \sim q} e^{f(x)^\top f(x^-)}$ . Here we have used the fact that  $f$  lies on the boundary of the hypersphere of radius  $1/t$ , which gives us the following equivalence between inner products and squared Euclidean norm,

$$2/t^2 - 2f(x)^\top f(x^+) = \|f(x)\|^2 + \|f(x^+)\|^2 - 2f(x)^\top f(x^+) = \|f(x) - f(x^+)\|^2. \quad (9)$$

Taking supremum to obtain  $\mathcal{L}_\infty^*(f) = \sup_{q \in \Pi} \mathcal{L}_\infty(f, q)$  we find that the second expression simplifies to,

$$\mathcal{L}_{\text{unif}}^*(f) = \sup_{q \in \Pi} \mathcal{L}_{\text{unif}}(f, q) = \mathbb{E}_{x \sim p} \log \sup_{x^- \approx x} e^{f(x)^\top f(x^-)} = \mathbb{E}_{x \sim p} \sup_{x^- \approx x} f(x)^\top f(x^-).$$

Using Eqn. (9), this can be re-expressed as,

$$\mathbb{E}_{x \sim p} \sup_{x^- \approx x} f(x)^\top f(x^-) = -\mathbb{E}_{x \sim p} \inf_{x^- \approx x} \|f(x) - f(x^+)\|^2/2 + 1/t^2. \quad (10)$$

The forthcoming theorem exactly characterizes the global optima of  $\min_f \mathcal{L}_\infty^*(f)$

**Theorem 6.** Suppose the downstream task is classification (i.e.  $\mathcal{C}$  is finite), and let  $\mathcal{L}_\infty^*(f) = \sup_{q \in \Pi} \mathcal{L}_\infty(f, q)$ . The infimum  $\inf_{f: \text{measurable}} \mathcal{L}_\infty^*(f)$  is attained, and any  $f^*$  achieving the global minimum is such that  $f^*(x) = f^*(x^+)$  almost surely. Furthermore, letting  $\mathbf{v}_c = f^*(x)$  for any  $x$  such that  $h(x) = c$  (so  $\mathbf{v}_c$  is well defined up to a set of  $x$  of measure zero),  $f^*$  is characterized as being any solution to the following ball-packing problem,

$$\max_{\{\mathbf{v}_c \in \mathbb{S}^{d-1}/t\}_{c \in \mathcal{C}}} \sum_{c \in \mathcal{C}} \rho(c) \cdot \min_{c' \neq c} \|\mathbf{v}_c - \mathbf{v}_{c'}\|^2. \quad (11)$$

*Proof.* Any minimizer of  $\mathcal{L}_{\text{align}}(f)$  has the property that  $f(x) = f(x^+)$  almost surely. So, in order to prove the first claim, it suffices to show that there exist functions  $f \in \arg \inf_f \mathcal{L}_{\text{unif}}^*(f)$  for which  $f(x) = f(x^+)$  almost surely. This is because, at that point, we have shown that  $\arg \min_f \mathcal{L}_{\text{align}}(f)$  and  $\arg \min_f \mathcal{L}_{\text{unif}}^*(f)$  intersect, and therefore any solution of  $\mathcal{L}_\infty^*(f) = \mathcal{L}_{\text{align}}(f) + \mathcal{L}_{\text{unif}}^*(f)$  must lie in this intersection.

To this end, suppose that  $f \in \arg \min_f \mathcal{L}_{\text{unif}}^*(f)$  but that  $f(x) \neq f(x^+)$  with non-zero probability. We shall show that we can construct a new embedding  $\hat{f}$  such that  $f(x) = f(x^+)$  almost surely, and  $\mathcal{L}_{\text{unif}}^*(\hat{f}) \leq \mathcal{L}_{\text{unif}}^*(f)$ . Due to Eqn. (10) this last condition is equivalent to showing,

$$\mathbb{E}_{x \sim p} \inf_{x^- \approx x} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \geq \mathbb{E}_{x \sim p} \inf_{x^- \approx x} \|f(x) - f(x^-)\|^2. \quad (12)$$

Fix a  $c \in \mathcal{C}$ , and let  $x_c \in \arg \max_{x: h(x)=c} \inf_{x^- \sim x} \|f(x) - f(x^-)\|^2$ . Then, define  $\hat{f}(x) = f(x_c)$  for any  $x$  such that  $h(x) = c$  and  $\hat{f}(x) = f(x)$  otherwise. Let us first aim to show that Eqn. (12) holds for this  $\hat{f}$ . Let us begin to expand the left hand side of Eqn. (12),

$$\begin{aligned}
& \mathbb{E}_{x \sim p} \inf_{x^- \sim x} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \\
&= \mathbb{E}_{\hat{c} \sim \rho} \mathbb{E}_{x \sim p(\cdot|\hat{c})} \inf_{x^- \sim x} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \\
&= \rho(c) \mathbb{E}_{x \sim p(\cdot|c)} \inf_{x^- \sim x} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \\
&\quad + (1 - \rho(c)) \mathbb{E}_{\hat{c} \sim \rho(\cdot|\hat{c} \neq c)} \mathbb{E}_{x \sim p(\cdot|\hat{c})} \inf_{x^- \sim x} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \\
&= \rho(c) \mathbb{E}_{x \sim p(\cdot|c)} \inf_{x^- \sim x} \|f(x_c) - f(x^-)\|^2 \\
&\quad + (1 - \rho(c)) \mathbb{E}_{\hat{c} \sim \rho(\cdot|\hat{c} \neq c)} \mathbb{E}_{x \sim p(\cdot|\hat{c})} \inf_{x^- \sim x} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \\
&= \rho(c) \inf_{x^- \sim x_c} \|f(x_c) - f(x^-)\|^2 \\
&\quad + (1 - \rho(c)) \mathbb{E}_{\hat{c} \sim \rho(\cdot|\hat{c} \neq c)} \mathbb{E}_{x \sim p(\cdot|\hat{c})} \inf_{h(x^-) \neq \hat{c}} \|\hat{f}(x) - \hat{f}(x^-)\|^2 \tag{13}
\end{aligned}$$

By construction, the first term can be lower bounded by  $\inf_{x^- \sim x_c} \|f(x_c) - f(x^-)\|^2 \geq \inf_{x^- \sim x} \|f(x) - f(x^-)\|^2$  for any  $x$  such that  $h(x) = c$ . To lower bound the second term, consider any fixed  $\hat{c} \neq c$  and  $x \sim p(\cdot|\hat{c})$  (so  $h(x) = \hat{c}$ ). Define the following two subsets of the input space  $\mathcal{X}$

$$\mathcal{A} = \{f(x^-) : f(x^-) \neq \hat{c} \text{ for } x^- \in \mathcal{X}\} \quad \hat{\mathcal{A}} = \{f(x^-) \in \mathcal{X} : \hat{f}(x^-) \neq \hat{c} \text{ for } x^- \in \mathcal{X}\}.$$

Since by construction the range of  $\hat{f}$  is a subset of the range of  $f$ , we know that  $\hat{\mathcal{A}} \subseteq \mathcal{A}$ . Combining this with the fact that  $\hat{f}(x) = f(x)$  whenever  $h(x) = \hat{c} \neq c$  we see,

$$\begin{aligned}
\inf_{h(x^-) \neq \hat{c}} \|\hat{f}(x) - \hat{f}(x^-)\|^2 &= \inf_{h(x^-) \neq \hat{c}} \|f(x) - \hat{f}(x^-)\|^2 \\
&= \inf_{u \in \hat{\mathcal{A}}} \|f(x) - u\|^2 \\
&\geq \inf_{u \in \mathcal{A}} \|f(x) - u\|^2 \\
&= \inf_{h(x^-) \neq \hat{c}} \|f(x) - f(x^-)\|^2
\end{aligned}$$

Using these two lower bounds we may conclude that Eqn. (13) can be lower bounded by,

$$\rho(c) \inf_{x^- \sim x_c} \|f(x_c) - f(x^-)\|^2 + (1 - \rho(c)) \mathbb{E}_{\hat{c} \sim \rho(\cdot|\hat{c} \neq c)} \mathbb{E}_{x \sim p(\cdot|\hat{c})} \inf_{h(x^-) \neq \hat{c}} \|f(x) - f(x^-)\|^2$$

which equals  $\mathbb{E}_{x \sim p} \inf_{x^- \sim x} \|f(x) - f(x^-)\|^2$ . We have therefore proved Eqn. (12). To summarize the current progress; given an embedding  $f$  we have constructed a new embedding  $\hat{f}$  that attains lower  $\mathcal{L}_{\text{unif}}$  loss and which is constant on  $x$  such that  $\hat{f}$  is constant on  $\{x : h(x) = c\}$ . Enumerating  $\mathcal{C} = \{c_1, c_2, \dots, c_{|\mathcal{C}|}\}$ , we may repeatedly apply the same argument to construct a sequence of embeddings  $f_1, f_2, \dots, f_{|\mathcal{C}|}$  such that  $f_i$  is constant on each of the following sets  $\{x : h(x) = c_j\}$  for  $j \leq i$ . The final embedding in the sequence  $f^* = f_{|\mathcal{C}|}$  is such that  $\mathcal{L}_{\text{unif}}^*(f^*) \leq \mathcal{L}_{\text{unif}}^*(f)$  and therefore  $f^*$  is a minimizer. This embedding is constant on each of  $\{x : h(x) = c_j\}$  for  $j = 1, 2, \dots, |\mathcal{C}|$ . In other words,  $f^*(x) = f^*(x^+)$  almost surely. We have proved the first claim.

Obtaining the second claim is a matter of manipulating  $\mathcal{L}_\infty^*(f^*)$ . Indeed, we know that  $\mathcal{L}_\infty^*(f^*) = \mathcal{L}_{\text{unif}}^*(f^*) - 1/t^2$  and defining  $\mathbf{v}_c = f^*(x)$  for any  $x$  such that  $h(x) = c$ , this expression is minimized if and only if  $f^*$  attains,

$$\begin{aligned} \max_f \mathbb{E}_{x \sim p} \inf_{x^- \sim x} \|f(x) - f(x^-)\|^2 &= \max_f \mathbb{E}_{c \sim \rho} \mathbb{E}_{x \sim p(\cdot|c)} \inf_{h(x^-) \neq c} \|f(x) - f(x^-)\|^2 \\ &= \max_f \sum_{c \in \mathcal{C}} \rho(c) \cdot \inf_{h(x^-) \neq c} \|f(x) - f(x^-)\|^2 \\ &= \max_{\{\mathbf{v}_c \in \mathbb{S}^{d-1}/t\}} \sum_{c \in \mathcal{C}} \rho(c) \cdot \min_{c' \neq c} \|\mathbf{v}_c - \mathbf{v}_{c'}\|^2 \end{aligned}$$

□

## B Hard Sampling for Graph Representation Learning

We describe in detail the hard sampling method for graphs whose results are reported in Section 5.2. Before getting that point, in the interests of completeness we cover some required background details on the InfoGraph method of [44]. For further information see the original paper [44].

### B.1 Background on Graph Representations

We observe a set of graphs  $\mathbf{G} = \{G_j \in \mathcal{G}\}_{j=1}^n$  sampled according to a distribution  $p$  over an ambient graph space  $\mathcal{G}$ . Each node  $u$  in a graph  $G$  is assumed to have features  $h_u^{(0)}$  living in some Euclidean space. We consider a  $K$ -layer graph neural network, whose  $k$ -th layer iteratively computes updated embeddings for each node  $v \in G$  in the following way,

$$h_v^{(k)} = \text{COMBINE}^{(k)} \left( h_v^{(k-1)}, \text{AGGREGATE}^{(k)} \left( \left\{ \left( h_v^{(k-1)}, h_u^{(k-1)}, e_{uv} \right) : u \in \mathcal{N}(v) \right\} \right) \right)$$

where  $\text{COMBINE}^{(k)}$  and  $\text{AGGREGATE}^{(k)}$  are parameterized learnable functions and  $\mathcal{N}(v)$  denotes the set of neighboring nodes of  $v$ . The  $K$  embeddings for a node  $u$  are collected together to obtain a single final summary embedding for  $u$ . As recommended by [54] we use concatenation,  $h^u = h^u(G) = \text{CONCAT}(\{h_u^{(k)}\}_{k=1}^K)$  to obtain an embedding in  $\mathbb{R}^d$ . Finally, the node representations are combined together into a length  $d$  graph level embedding using a readout function,

$$H(G) = \text{READOUT}(\{h^u\}_{u \in G})$$

which is typically taken to be a simple permutation invariant function such as the sum or mean. The InfoGraph method aims to maximize the mutual information between the graph level embedding  $H(G)$  and patch-level embeddings  $h^u(G)$  using the following objective,

$$\max_h \mathbb{E}_{G \sim p} \frac{1}{|G|} \sum_{u \in G} I(h^u(G); H(G))$$

In practice the population distribution  $p$  is replaced by its empirical counterpart, and the mutual information  $I$  is replaced by a variational approximation  $I_T$ . In line with [44] we use the

Jensen-Shannon mutual information estimator as formulated by [33]. It is defined using a neural network discriminator  $T : \mathbb{R}^{2d} \rightarrow \mathbb{R}$  as,

$$I_T(h^u(G); H(G)) = \mathbb{E}_{G \sim p} \left[ -\text{sp}(-T(h^u(G), H(G))) \right] - \mathbb{E}_{(G, G') \sim p \times p} \left[ \text{sp}(T(h^u(G), H(G'))) \right]$$

where  $\text{sp}(z) = \log(1 + e^z)$  denotes the softplus function. The final objective is the joint maximization over  $h$  and  $T$ ,

$$\max_{\theta, \psi} \mathbb{E}_{G \sim p} \frac{1}{|G|} \sum_{u \in G} I_T(h^u(G); H(G))$$

## B.2 Hard Sampling for Graph Representations

In order to derive a simple modification of the NCE hard sampling technique that is appropriate for use with InfoGraph, we first provide a mildly generalized view of hard sampling. Recall that the NCE contrastive objective can be decomposed into two constituent pieces,

$$\mathcal{L}(f, q) = \mathcal{L}_{\text{align}}(f) + \mathcal{L}_{\text{unif}}(f, q)$$

where  $q$  is in fact a family of distributions  $q(x^-; x)$  over  $x^-$  that is indexed by the possible values of the anchor  $x$ .  $\mathcal{L}_{\text{align}}$  performs the role of “aligning” positive pairs (embedding near to one-another), while  $\mathcal{L}_{\text{unif}}$  repels negative pairs. The hard sampling framework aims to solve,

$$\inf_f \sup_q \mathcal{L}(f, q).$$

In the case of NCE loss we take,

$$\begin{aligned} \mathcal{L}_{\text{align}}(f) &= -\mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} f(x)^T f(x^+), \\ \mathcal{L}_{\text{unif}}(f, q) &= \mathbb{E}_{\substack{x \sim p \\ x^+ \sim p_x^+}} \log \left\{ e^{f(x)^T f(x^+)} + Q \mathbb{E}_{x^- \sim q} [e^{f(x)^T f(x^-)}] \right\}. \end{aligned}$$

View this view, we can easily adapt to the InfoGraph framework, taking

$$\begin{aligned} \mathcal{L}_{\text{align}}(h, T) &= -\mathbb{E}_{G \sim p} \frac{1}{|G|} \sum_{u \in G} \text{sp}(-T(h^u(G), H(G))), \\ \mathcal{L}_{\text{unif}}(h, T, q) &= -\mathbb{E}_{G \sim p} \frac{1}{|G|} \sum_{u \in G} \mathbb{E}_{G' \sim q} \text{sp}(T(h^u(G), H(G'))) \end{aligned}$$

Denote by  $\hat{p}$  the distribution over nodes  $u \in \mathbb{R}^s$  defined by first sampling  $G \sim p$ , then sampling  $u \in G$  uniformly over all nodes of  $G$ . Then these two terms can be simplified to

$$\begin{aligned} \mathcal{L}_{\text{align}}(h, T) &= -\mathbb{E}_{u \sim \hat{p}} \text{sp}(-T(h^u(G), H(G))), \\ \mathcal{L}_{\text{unif}}(h, T, q) &= -\mathbb{E}_{(u, G') \sim \hat{p} \times q} \text{sp}(T(h^u(G), H(G'))) \end{aligned}$$

At this point it becomes clear that, just as with NCE, a distribution  $q^* \in \arg \max_q \mathcal{L}(f, q)$  in the InfoGraph framework if it is supported on  $\arg \max_{G' \in G} \text{sp}(T(h^u(G), H(G')))$ . Although this is still hard to compute exactly, it can be approximated by,

$$q_u^\beta(G') \propto \exp(\beta T(h^u(G), H(G))) \cdot p(G').$$

## C Additional Ablations

To study the affect of varying the concentration parameter  $\beta$  on the learned embeddings Figure 6 plots cosine similarity histograms of pairs of similar and dissimilar points. The results show that for  $\beta$  moving from 0 through 0.5 to 2 causes both the positive and negative similarities to gradually skew left. In terms of downstream classification, an important property is the *relative* difference in similarity between positive and negative pairs. In this case  $\beta = 0.5$  find the best balance (since it achieves the highest downstream accuracy). When  $\beta$  is taken very large ( $\beta = 6$ ), we see a change in conditions. Both positive and negative pairs are assigned higher similarities in general. Visually it seems that the positive and negative histograms for  $\beta = 6$  overlap a lot more than for smaller values, which helps explain why the linear readout accuracy is lower for  $\beta = 6$ .

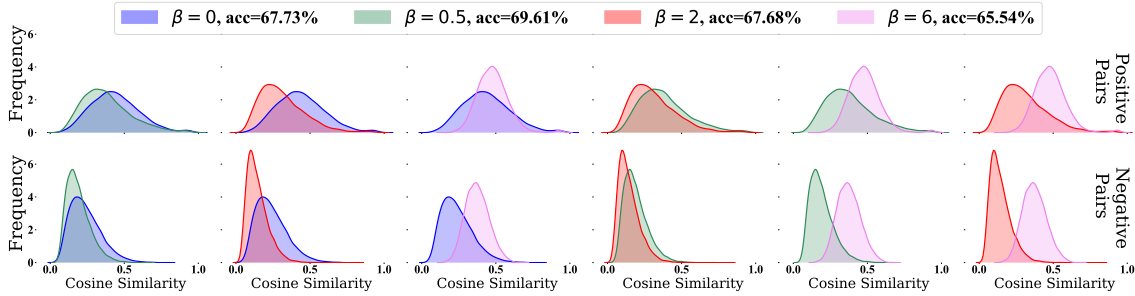


Figure 6: Histograms of cosine similarity of pairs of points with different label (bottom) and same label (top) for embeddings trained on CIFAR100 with different values of  $\beta$ . Histograms overlaid pairwise to allow for easy comparison.

## D Experimental Details

Figure 7 shows PyTorch-style pseudocode for the standard objective, the debiased objective [9], and the hard sampling objective. The proposed hard-sample loss is very simple to implement, requiring only two extra lines of code compared to the standard objective.

### D.1 Visual Representations

We implement SimCLR in PyTorch. We use a ResNet-50 [19] as the backbone with embedding dimension 128. We use the Adam optimizer [25] with learning rate 0.001 and weight decay  $10^{-6}$ . Official code will be released. Since we adopt the SimCLR framework, the number of negative samples  $N = 2(\text{batch size} - 1)$ . Since we always take the batch size to be a power of 2 (16, 32, 64, 128, 256) the negative batch sizes are 30, 62, 126, 254, 510 respectively.

**Annealing  $\beta$  Method:** We detail the annealing method whose results are given in Figure 4. The idea is to reduce the concentration parameter down to zero as training progresses. Specifically, suppose we have  $e$  number of total training epochs. We also specify a number  $\ell$  of “changes” to the concentration parameter we shall make. We initialize the concentration parameter  $\beta_1 = \beta$  (where this  $\beta$  is the number reported in Figure 4), then once every  $e/\ell$  epochs we reduce  $\beta_i$  by  $\beta/\ell$ . In other words, if we are currently on  $\beta_i$ , then  $\beta_{i+1} = \beta_i - \beta/\ell$ , and we switch from  $\beta_i$  to  $\beta_{i+1}$ .

```

1 # pos      : exp of inner products for positive examples
2 # neg      : exp of inner products for negative examples
3 # N        : number of negative examples
4 # t        : temperature scaling
5 # tau_plus : class probability
6 # beta     : concentration parameter
7
8 #Original objective
9 standard_loss = -log(pos.sum() / (pos.sum() + neg.sum()))
10
11 #Debiased objective
12 Neg = max((-N*tau_plus*pos + neg).sum() / (1-tau_plus), e**(-1/t))
13 debiased_loss = -log(pos.sum() / (pos.sum() + Neg))
14
15 #Hard sampling objective (Ours)
16 reweight = (beta*neg) / neg.mean()
17 Neg = max((-N*tau_plus*pos + reweight*neg).sum() / (1-tau_plus), e**(-1/t))
18 hard_loss = -log( pos.sum() / (pos.sum() + Neg))

```

Figure 7: Pseudocode for our proposed new hard sample objective, as well as the original NCE contrastive objective, and debiased contrastive objective [9]. In each case we take the number of positive samples to be  $M = 1$ . The implementation of our hard sampling method only requires two additional lines of code compared to the standard objective.

in epoch number  $i \cdot e / \ell$ . The idea of this method is to select particularly difficult negative samples early on order to obtain useful gradient information early on, but later (once the embedding is already quite good) we reduce the “hardness” level so as to reduce the harmful effect of only approximately correcting for false negatives (negatives with the same labels as the anchor).

```

1 # pos      : exp of inner products for positive examples
2 # neg      : exp of inner products for negative examples
3 # N        : number of negative examples
4 # t        : temperature scaling
5 # tau_plus : class probability
6 # beta     : concentration parameter
7
8 #Clipping negatives trick before computing reweighting
9 reweight = 2*neg / max( neg.max().abs(), neg.min().abs() )
10 reweight = (beta*reweight) / reweight.mean()
11 Neg = max((-N*tau_plus*pos + reweight*neg).sum() / (1-tau_plus), e**(-1/t))
12 hard_loss = -log( pos.sum() / (pos.sum() + Neg))

```

Figure 8: In cases where the learned embedding is not normalized to lie on a hypersphere we found that clipping the negatives to live in a fixed range (in this case  $[-2, 2]$ ) stabilizes optimization.

## D.2 Graph Representations

All datasets we benchmark on can be downloaded at [www.graphlearning.io](http://www.graphlearning.io) from the TUDataset repository of graph classification problems [31]. Information on basic statistics of the datasets is included in Tables 2 and 3. For fair comparison to the original InfoGraph method, we adopt the official code, which can be found at <https://github.com/fanyun-sun/InfoGraph>. We modify only the `gan_losses.py` script, adding in our proposed hard sampling via reweighting. For simplicity we trained all models using the same set of hyperparameters: we used the GIN architecture [54] with  $K = 3$  layers and embedding dimension  $d = 32$ . Each model is trained for 200 epochs with batch size 128 using the Adam optimizer [25]. with learning rate 0.001, and weight decay of  $10^{-6}$ . Each embedding is evaluated using the average accuracy 10-fold cross-validation using an SVM as the classifier (in line with the approach taken by [31]). Each experiment is repeated from scratch 10 times, and the distribution of results from these 10 runs is plotted in

Figure 3.

Since the graph embeddings are not constrained to lie on a hypersphere, for a batch we clip all the inner products to live in the interval  $[-2, 2]$  while computing the reweighting, as illustrated in Figure 8. We found this to be important for stabilizing optimization.

Dataset	DD	PTC	REDDIT-B	PROTEINS
No. graphs	1178	344	2000	1113
No. classes	2	2	2	2
Avg. nodes	284.32	14.29	429.63	39.06
Avg. Edges	715.66	14.69	497.75	72.82

Table 2: Basic statistics for graph datasets.

Dataset	ENZYMES	MUTAG	IMDB-B	IMDB-M
No. graphs	600	188	1000	1500
No. classes	6	2	2	3
Avg. nodes	32.63	17.93	19.77	13.00
Avg. Edges	62.14	19.79	96.53	65.94

Table 3: Basic statistics for graph datasets.

### D.3 Sentence Representations

We adopt the official quick-thoughts vectors experimental settings, which can be found at <https://github.com/lajanugen/S2V>. We keep all hyperparameters at the default values and change only the `s2v-model.py` script. Since the official BookCorpus dataset [26] is not available, we use an unofficial version obtained using the following repository: <https://github.com/soskek/bookcorpus>. Since the sentence embeddings are also not constrained to lie on a hypersphere, we use the same clipping trick as for the graph embeddings, illustrated in Figure 8.

After training on the BookCorpus dataset, we evaluate the embeddings on six different classification tasks: paraphrase identification (MSRP) [11], question type classification (TREC) [49], opinion polarity (MPQA) [51], subjectivity classification (SUBJ) [35], product reviews (CR) [22], and sentiment of movie reviews (MR) [36].