# A Guide to InfoNCE and Its Use in SimCLR

**Abstract**

Self-supervised learning has revolutionized representation learning in deep neural networks by leveraging unlabeled data to learn meaningful features. At the heart of many such methods lies the InfoNCE (Information Noise-Contrastive Estimation) loss, a tractable approximation to mutual information maximization. This article provides a comprehensive, mathematically rigorous treatment of InfoNCE, tracing its theoretical foundations from information theory to its practical embodiment in SimCLR (A Simple Framework for Contrastive Learning of Visual Representations). We derive key equations step-by-step, analyze the roles of components like the critic function, partition function approximations, and projection heads, and discuss design choices that enhance empirical performance. Spanning theoretical derivations, practical implementations, and insights into why these methods work, this exposition aims to bridge the gap between abstract information-theoretic principles and state-of-the-art computer vision algorithms. The article is structured to build progressively from fundamentals to advanced topics, with a focus on mathematical precision.

# 1 Introduction: The Theoretical Triangle of Mutual Information, InfoNCE, and SimCLR

Self-supervised learning (SSL) addresses the challenge of learning representations without explicit labels by creating supervisory signals from the data itself. A cornerstone of modern SSL is contrastive learning, which pulls together representations of similar (positive) samples while pushing apart those of dissimilar (negative) ones. This paradigm is theoretically grounded in maximizing mutual information (MI) between different views of the same data point.

The "theoretical triangle" encapsulates this progression:

- **Mutual Information (MI)**: An intractable but ideal objective for capturing dependencies between views.

- **InfoNCE**: A tractable lower bound on MI, transforming the problem into a noise-contrastive estimation task.

- **SimCLR**: A practical implementation of InfoNCE tailored for visual representation learning, incorporating architectural and augmentation strategies.

This article expands on these connections with rigorous derivations. We begin with MI, derive InfoNCE as its approximation, and detail SimCLR's specifics. Throughout, we emphasize mathematical formulations, including proofs where applicable, to ensure rigor.

Consider the core goal: Given an image $\mathbf{x}$, we generate two augmented views $\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{x}}_j$ (positive pair) via stochastic transformations $t \sim \mathcal{T}$, where $\mathcal{T}$ is a distribution over augmentations like random crops, color jitter, and flips. The objective is to learn an encoder $f(\cdot)$ such that representations $\mathbf{h}_i = f(\tilde{\mathbf{x}}_i)$ and $\mathbf{h}_j = f(\tilde{\mathbf{x}}_j)$ are close in embedding space, while being distant from representations of other images (negatives).

Mathematically, this aligns with maximizing the MI $I(\tilde{\mathbf{x}}_i; \tilde{\mathbf{x}}_j)$, but direct optimization is challenging. InfoNCE provides a solution.

# 2    Mutual Information – The Intractable Foundation

Mutual information quantifies the reduction in uncertainty about one random variable given knowledge of another. For two random variables $X$ and $Y$, it is defined as:

$$I(X;Y) = H(X) - H(X \mid Y) = \mathbb{E}_{p(x,y)} \left[ \log \frac{p(x,y)}{p(x)p(y)} \right], \tag{1}$$

where $H(\cdot)$ is entropy, $p(x,y)$ is the joint distribution, and $p(x)$, $p(y)$ are marginals.

In SSL, we aim to maximize $I(f(\tilde{\mathbf{x}}_i); f(\tilde{\mathbf{x}}_j))$, encouraging the encoder $f$ to preserve information shared between views while being invariant to augmentations. However, computing MI is intractable because:

1. It requires knowledge of the high-dimensional joint and marginal distributions.

2. The expectation involves integration over potentially infinite spaces.

3. In deep learning, representations are continuous and high-dimensional, exacerbating computational issues.

To illustrate, suppose we parameterize the encoder with $\theta$. Optimizing $\max_\theta I(f_\theta(\tilde{\mathbf{x}}_i); f_\theta(\tilde{\mathbf{x}}_j))$ demands estimating density, often via variational approximations. InfoNCE emerges as a variational lower bound, drawing from Noise-Contrastive Estimation (NCE).

## 2.1    Derivation of MI Lower Bounds

A common approach to bounding MI is the Barber-Agakov bound or its variants. For InfoNCE, we consider the following: Let $\mathbf{c} = f(\tilde{\mathbf{x}}_i)$ be a "context" representation, and $\mathbf{x} = f(\tilde{\mathbf{x}}_j)$ the target. We seek $I(\mathbf{c}; \mathbf{x})$.

Using the chain rule and non-negativity of KL-divergence:

$$I(\mathbf{c}; \mathbf{x}) = \mathbb{E}_{p(\mathbf{x}|\mathbf{c})} \left[ \log p(\mathbf{x} \mid \mathbf{c}) \right] + H(\mathbf{x}). \tag{2}$$

Since $H(\mathbf{x})$ is constant, we focus on the conditional log-likelihood. But estimating $p(\mathbf{x} \mid \mathbf{c})$ is hard. InfoNCE approximates this via classification among positives and negatives.

# 3 InfoNCE as a Tractable Lower Bound on Mutual Information

InfoNCE, introduced in Contrastive Predictive Coding (CPC), provides a lower bound on MI by framing the problem as distinguishing a positive sample from noise (negatives). Consider a positive pair $(\mathbf{c}, \mathbf{x}^+)$ drawn from $p(\mathbf{x} \mid \mathbf{c})$, and $K - 1$ negatives $\{\mathbf{x}_k^-\}_{k=1}^{K-1}$ from the marginal $p(\mathbf{x})$.

The InfoNCE loss for a single instance is:

$$\mathcal{L}_{\text{InfoNCE}} = -\log \left( \frac{\exp(s(\mathbf{c}, \mathbf{x}^+))}{\exp(s(\mathbf{c}, \mathbf{x}^+)) + \sum_{k=1}^{K-1} \exp(s(\mathbf{c}, \mathbf{x}_k^-))} \right), \tag{3}$$

where $s(\cdot, \cdot)$ is a similarity function (the "critic"), often bilinear or dot product.

## 3.1 Rigorous Derivation of InfoNCE as MI Lower Bound

To derive this as a bound on MI, start with the expectation over positives and negatives:

$$\mathbb{E}[\mathcal{L}_{\text{InfoNCE}}] = -\mathbb{E}_{\mathbf{c}, \mathbf{x}^+, \{\mathbf{x}_k^-\}} \left[ \log \frac{\exp(s(\mathbf{c}, \mathbf{x}^+))}{\sum_{j=1}^{K} \exp(s(\mathbf{c}, \mathbf{x}_j))} \right], \tag{4}$$

where $\mathbf{x}_1 = \mathbf{x}^+$, $\mathbf{x}_{2:K} = \{\mathbf{x}_k^-\}$.

This is the negative log-probability of correctly classifying the positive in a $K$-way softmax classifier. Under the optimal critic $s^*(\mathbf{c}, \mathbf{x}) = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} + C$ (where $C$ is constant), it can be shown that:

$$I(\mathbf{c}; \mathbf{x}) \geq \log K - \mathbb{E}[\mathcal{L}_{\text{InfoNCE}}]. \tag{5}$$

**Proof Sketch**:

The probability of correct classification is:

$$p(\text{correct}) = \mathbb{E} \left[ \frac{\exp(s(\mathbf{c}, \mathbf{x}^+))}{\sum_j \exp(s(\mathbf{c}, \mathbf{x}_j))} \right]. \tag{6}$$

By Jensen's inequality (since $-\log$ is convex):

$$-\log p(\text{correct}) \geq \mathbb{E}[\mathcal{L}_{\text{InfoNCE}}]. \tag{7}$$

For the optimal critic, the expected log-ratio approximates the density ratio, and as $K \to \infty$, the bound tightens to MI. Precisely, Poole et al. (2019) show the bound is $I(\mathbf{c}; \mathbf{x}) \geq \log K + \mathbb{E}[-\mathcal{L}_{\text{InfoNCE}}]$, with equality in the limit.

This derivation highlights InfoNCE's elegance: It turns MI maximization into multiclass classification, where negatives approximate the marginal.

# 4 SimCLR – InfoNCE in Practice for Visual Representations

SimCLR operationalizes InfoNCE for image data. For a batch of $N$ images, apply two augmentations each, yielding $2N$ views. Let $\mathbf{z}_i = g(f(\tilde{\mathbf{x}}_i))$ be projected representations, where $f$ is the encoder (e.g., ResNet) and $g$ is the projection head (MLP).

The similarity is cosine:

$$\text{sim}(\mathbf{z}_i, \mathbf{z}_j) = \frac{\mathbf{z}_i^\top \mathbf{z}_j}{\|\mathbf{z}_i\| \|\mathbf{z}_j\|}. \tag{8}$$

The loss for pair $(i, j)$ (positives from same image) is:

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbf{1}_{k \neq i} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}, \tag{9}$$

where $\tau$ is temperature. The total loss is symmetrized:

$$\mathcal{L}_{\text{SimCLR}} = \frac{1}{2N} \sum_{i=1}^{2N} (\ell_{i,j(i)} + \ell_{j(i),i}), \tag{10}$$

with $j(i)$ the mate of $i$.

This is InfoNCE with $K = 2N - 1$, batch negatives, and cosine critic. Augmentations define positives, ensuring invariance to transformations.

## 4.1 Why Cosine Similarity?

Cosine normalizes for magnitude, focusing on direction, which is robust to scale variations in representations. Temperature $\tau$ controls distribution sharpness: Low $\tau$ emphasizes hard negatives.

# 5 The Critic Function – Judging Similarity

The critic $s(\mathbf{c}, \mathbf{x})$ evaluates pair compatibility. In SimCLR, it's $s(\mathbf{z}_i, \mathbf{z}_j) = \text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau$.

## 5.1 Adversarial Interpretation

Borrowed from GANs, the critic "judges" encoder quality. The encoder maximizes positive scores while the critic (implicitly learned via joint optimization) discriminates positives from negatives.

In the loss, the critic appears in both numerator (positive score) and denominator (all scores):

$$\mathcal{L} = -\log \frac{\exp(s^+)}{\exp(s^+) + \sum \exp(s_k^-)}. \tag{11}$$

This creates contrast: To minimize loss, positives must dominate negatives.

Derivation of optimal critic: At optimum, $s^*(\mathbf{c}, \mathbf{x}) \approx \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}$, aligning with density ratio estimation.

# 6 The Partition Function Problem and Its Approximation

In probabilistic models, densities involve partition functions:

$$p(\mathbf{x} \mid \mathbf{c}) = \frac{\exp(s(\mathbf{c}, \mathbf{x}))}{Z(\mathbf{c})}, \quad Z(\mathbf{c}) = \int \exp(s(\mathbf{c}, \mathbf{y})) p(\mathbf{y}) d\mathbf{y}. \tag{12}$$

$Z(\mathbf{c})$ is intractable due to integration over representation space and unknown $p(\mathbf{y})$.

InfoNCE sidesteps this by approximating $Z$ with batch samples:

$$Z_{\text{approx}}(\mathbf{c}) = \sum_{k=1}^{K-1} \exp(s(\mathbf{c}, \mathbf{x}_k^-)) + \exp(s(\mathbf{c}, \mathbf{x}^+)), \tag{13}$$

a Monte Carlo estimate. Larger batches (higher $K$) improve approximation, tightening the MI bound.

## 6.1 Bias and Variance in Approximation

The estimate is unbiased in expectation for fixed critic, but in practice, joint optimization introduces bias. However, empirical results show large batches (e.g., 4096 in Sim-CLR) yield near state-of-the-art performance.

# 7 Deep Dive into the InfoNCE Loss Formula

Expanding the loss:

$$\mathcal{L} = -\log\left(\frac{\exp(s^+)}{\exp(s^+) + \sum \exp(s_k^-)}\right) = -s^+ + \log\left(\exp(s^+) + \sum \exp(s_k^-)\right). \tag{14}$$

This is cross-entropy between a one-hot label (positive) and softmax probabilities.

Interpretation: It computes the negative log-probability that the positive is selected among candidates, encouraging $p(\mathbf{y}^+ \mid \mathbf{x}, \text{candidates}) \approx 1$.

For batch average in SimCLR, this scales to multi-view settings.

# 8 Positive and Negative Samples – The Anchor Perspective

In contrastive learning, $p(\mathbf{y} \mid \mathbf{x})$ is the probability that $\mathbf{y}$ matches anchor $\mathbf{x}$. Positives are augmentation-defined: $\mathbf{y}^+ \sim p(\mathbf{y} \mid \mathbf{x}) \approx \delta(\text{same image})$.

Negatives provide contrast, forcing discrimination. Without them, collapse to constant representations occurs.

The anchor $\mathbf{x}$ fixes the reference; normalization is per-anchor.

# 9 The Partition Function in the Denominator – Competition Dynamics

The denominator includes the positive term, creating self-competition:

$$p^+ = \frac{\exp(s^+)}{\exp(s^+) + \sum \exp(s^-)}. \tag{15}$$

Including positive scales $p^+ < 1$, but encourages stronger separation. Variants exclude it for easier optimization, but SimCLR includes it.

Softmax induces zero-sum competition: Probability mass is contested, driving learning.

# 10 The Projection Head in SimCLR – A Key Architectural Choice

SimCLR uses a projection head $g$: A 2-layer MLP mapping $\mathbf{h}$ (2048-dim) to $\mathbf{z}$ (128-dim).

## 10.1 Why Required?

Empirically, it boosts accuracy by $\sim$10%. Theoretically:

- Separates contrastive space from representation space, preventing information loss.
- Acts as regularization, reducing dimensionality and noise.
- Nonlinearity allows complex transformations for better contrast.

Ablations show nonlinear > linear > none. During inference, discard $g$, using $\mathbf{h}$ for tasks.

# 11 Connecting Theory to Practice – The Full Pipeline

The pipeline:

1. Augment data to create views.
2. Encode and project.
3. Compute InfoNCE to maximize MI bound.
4. Transfer $f$ to downstream (e.g., linear classification).

Success stems from inductive biases in augmentations, scalability with data, and transferability.

# 12 Key Insights and Takeaways

- **Mathematical Elegance**: InfoNCE approximates intractable MI via batch Monte Carlo.

- **Design Matters**: Batch size tightens bound; temperature tunes hardness; augmentations define invariances.

- **Contrastive Principle**: Discrimination task yields robust features.

Impacts: Rivals supervised learning, enables data-efficient training, inspires methods like BYOL, CLIP.

# 13 Sign Conventions in Energy vs. Similarity Functions

Papers use similarity (positive for close pairs) or energy (negative):

Similarity: $\exp(s/\tau)$, $s > 0$ for positives.

Energy: $\exp(-E/\tau)$, $E < 0$ for positives.

Equivalent via $s = -E$, but affects interpretation.

# 14 Advanced Topics and Extensions

To expand rigorously, consider variances in negatives: Hard negatives (similar but not positive) improve learning, as in MoCo's queue.

Temperature scaling: Derives from logit adjustment; optimal $\tau$ balances exploration.

Collapse prevention: Normalization (cosine) and large batches help; alternatives like stop-gradients in BYOL avoid negatives.

Related bounds: Other MI estimators like MINE, NWJ exist, but InfoNCE excels in scalability.

Empirical analysis: SimCLR achieves 76.5% top-1 on ImageNet with ResNet-50, nearing supervised baselines.

## 14.1 Derivation of Temperature's Role

Temperature arises in scaled softmax:

$$p_i = \frac{\exp(l_i/\tau)}{\sum \exp(l_j/\tau)}, \tag{16}$$

Low $\tau$: Sharp, focuses on max; high $\tau$: Uniform. In InfoNCE, $\tau = 0.07$ empirically optimal for emphasizing distinctions.

# 15  Conclusion: The Beauty of Theory Meeting Practice

InfoNCE elegantly solves MI's intractability through contrastive approximations, enabling SimCLR's breakthroughs. This framework's rigor—rooted in information theory—and practicality have transformed SSL. Future directions include multi-modal extensions and theoretical tightenings of bounds.

The interplay of mathematics and empirics exemplifies deep learning's progress.