

# 关于《自然语言处理：大模型理论与实践》第七章 7.3.2 节 中对比损失函数表述错误的详细分析

## 1. 教材中的错误表述

教材在第 144 页公式 (7.9) 中，将 CLIP 模型的优化目标描述为最大化正样本相似度与最小化负样本相似度的简单线性差值。书中给出的公式形式如下：

$$\mathcal{L}_{textbook} = \min \left[ \sum_{i=1}^N \sum_{j=1, j \neq i}^N (I_i \cdot T_j) - \sum_{i=1}^N (I_i \cdot T_i) \right]$$

该公式仅表示了正负样本点积数值的直接相减。在深度学习的对比表示学习中，这种线性的损失函数不仅难以收敛，更缺失了概率归一化和温度系数调节这两个关键机制，无法体现 CLIP 模型构建高质量语义空间的核心原理。

## 2. 正确的数学表述

CLIP 模型实际采用的是基于 InfoNCE (Information Noise Contrastive Estimation) 的对称交叉熵损失函数。其核心思想是将相似度分值通过 Softmax 转换为概率分布，并最大化正样本对的对数似然概率。

对于一个包含 N 个图像-文本对的训练批次，图像  $I_i$  到文本  $T$  的损失函数  $\mathcal{L}_{I \rightarrow T}$  的标准定义应为：

$$\mathcal{L}_{I \rightarrow T}^{(i)} = -\log \frac{\exp(\text{sim}(I_i, T_i)/\tau)}{\sum_{k=1}^N \exp(\text{sim}(I_i, T_k)/\tau)}$$

最终的总损失函数  $\mathcal{L}_{total}$  通常为图像到文本损失与文本到图像损失的平均值：

$$\mathcal{L}_{total} = \frac{1}{N} \sum_{i=1}^N \left( \frac{\mathcal{L}_{I \rightarrow T}^{(i)} + \mathcal{L}_{T \rightarrow I}^{(i)}}{2} \right)$$

参数含义：

**N (Batch Size):**训练批次的大小，即一个 Batch 中包含的“图像-文本”对的数量。

作用：N 决定了对比学习中负样本的数量。对于第  $i$  张图片，它有 1 个正样本（对应的文本  $T_i$ ）和  $N - 1$  个负样本（同一批次中的其他文本）。

**$I_i, T_k$ (Normalized Embeddings):** $I_i$  表示第  $i$  张图像经过图像编码器输出并归一化后的特征向量； $T_k$  表示第  $k$  个文本经过文本编码器输出并归一化后的特征向量。

**$sim(I_i, T_k)$  (Cosine Similarity):**图像特征与文本特征的余弦相似度。由于向量已归一化，余弦相似度等价于点积，即  $sim(I_i, T_k) = I_i \cdot T_k$ 。

**$\tau$  (Temperature Parameter):**温度系数，一个可学习的标量参数。这是教材公式中完全遗漏的关键参数。当  $\tau \rightarrow 0$  时，Softmax 分布变得非常尖锐，模型会极度关注那些很难区分的负样本，产生巨大的梯度信号。当  $\tau \rightarrow \infty$  时，分布变得平滑，模型对所有负样本一视同仁。 $\tau$  用于调节模型对“难样本”的挖掘力度，是 CLIP 训练成功的核心超参数。