Vol. 57 No. 5 Oct. 2025

DOI:10.16356/j.1005-2615.2025.05.007

# 视觉语义增强的联合小样本开集识别分类器

丁相舒1,2, 耿传兴1,2, 陈松灿1,2

(1. 南京航空航天大学计算机科学与技术学院,南京 211106; 2. 模式分析与机器智能工业和信息化部重点实验 室,南京 211106)

摘要:探究了视觉-语言预训练模型对比语言-图像预训练(Contrastive language-image pre-training, CLIP)在小样本开集识别(Few-shot open-set recognition, FSOR)任务中的潜力。实验发现基于CLIP图像编码特征的视觉原型分类器通常不如传统FSOR基线方法;基于CLIP语义编码特征的语义原型分类器虽然在闭集分类上显著优于传统基线,但在开集识别方面表现不佳。本文分析造成这些问题的主要原因可能是CLIP的训练数据与FSOR目标数据之间的分布差异及CLIP语义原型分类器为已知类别划分了过大的决策边界。本文提出了一种简单有效的视觉语义增强的联合小样本开集分类器,其不仅充分利用CLIP语义原型分类器的闭集分类优势,还巧妙挖掘了传统FSOR预训练模型构建的视觉原型分类器的潜力,以更紧密的决策边界进一步提升开集识别的精准度。在4个基准数据集上的实验结果表明,该方法在准确率(Accuracy, ACC)和受试者工作特征曲线下的面积(Area under the receiver operating characteristic, AUROC)指标上相比最优基线平均提升了 2.9% 和 2.6%。

关键词:小样本开集识别;视觉-语言模型;原型分类器;分布差异;决策边界

中图分类号:TP391 文献标志码:A 文章编号:1005-2615(2025)05-0861-09

# Visual-Semantic Enhanced Joint Classifier for Few-Shot Open-Set Recognition

DING Xiangshu<sup>1,2</sup>, GENG Chuanxing<sup>1,2</sup>, CHEN Songcan<sup>1,2</sup>

(1. College of Computer Science and Technology, Nanjing University of Aeronautics & Astronautics, Nanjing 211106, China;2. MIIT Key Laboratory of Pattern Analysis and Machine Intelligence, Nanjing 211106, China)

Abstract: This paper investigates the potential of the vision-language pretrained model contrastive language-image pre-training (CLIP) in few-shot open-set recognition (FSOR). The experiments reveal that the visual-prototype classifier based on CLIP's image encoding features generally performs worse than traditional FSOR baseline methods; although the semantic-prototype classifier based on CLIP's semantic encoding features significantly outperforms traditional baselines in closed-set performance, it underperforms in open-set performance. The primary reasons for these issues may be the gap between CLIP's training data and the FSOR target data, as well as the overly large decision boundaries assigned by the CLIP semantic prototype classifier to known classes. To tackle these problems, a simple yet effective joint few-shot open-set classification advantages of the semantic-prototype classifier based on CLIP but also skillfully exploits the potential of the visual-prototype classifier constructed by traditional FSOR pretrained models, which further enhances the open-set performance by establishing tighter decision boundaries. Experiments on four benchmark datasets

基金项目:江苏省自然科学青年基金(BK20210292)。

收稿日期:2024-08-26;修订日期:2024-11-25

通信作者:陈松灿,男,教授,博士生导师,E-mail: s.chen@nuaa.edu.cn。

引用格式: 丁相舒, 耿传兴, 陈松灿. 视觉语义增强的联合小样本开集识别分类器[J]. 南京航空航天大学学报(自然科学版), 2025, 57(5): 861-869. DING Xiangshu, GENG Chuanxing, CHEN Songcan. Visual-semantic enhanced joint classifier for few-shot open-set recognition[J]. Journal of Nanjing University of Aeronautics & Astronautics (Natural Science Edition), 2025, 57(5): 861-869.

demonstrate that this method achieves average improvements of 2.9% in accuracy (ACC) and 2.6% in area under the receiver operating characteristic (AUROC) compared to the best traditional baselines.

**Key words:** few-shot open-set recognition(FSOR); vision-language model; prototype classifier; distribution discrepancy; decision boundary

小样本学习(Few-shot learning, FSL)旨在模仿人类从少量样本中学习并快速推广到新概念的能力,近年来取得了显著进展<sup>[1-3]</sup>。然而,传统的小样本学习方法主要遵循封闭集假设,即训练样本和测试样本均来自相同的标记空间,因此难以适用于测试样本可能来自未知类别的开放场景。为了解决这一挑战,小样本开集识别(Few-shot open-set recognition, FSOR)作为小样本学习和开集识别(Open-set recognition, OSR)的交集应运而生,引起了广泛的关注<sup>[4-9]</sup>。

小样本开集识别的目标是利用极少量的训练样本,快速适应新的分类任务,在准确分类已知类别样本的同时有效识别并拒绝未知类别的样本,这要求模型具有高度有效的特征表示能力<sup>[4]</sup>。近年来,在大规模的互联网数据集上预训练的视觉-语言模型(Vision-language models, VLMs)<sup>[10-12]</sup>,如对比语言-图像预训练(Contrastive language-image pre-training, CLIP)<sup>[10]</sup>,在目标检测<sup>[13]</sup>和视频理解<sup>[14]</sup>等多种下游视觉任务中表现出色,展示了强大的表示学习能力,这为FSOR的发展提供了新的可能性。

尽管视觉-语言模型在FSOR中具有巨大的应 用潜力,目前针对这一交叉领域的研究仍然较为有 限。为填补这一空白,本文以CLIP作为研究对象, 探讨了其在FSOR中的应用,通过两种无需训练的 原型分类器进行了研究。一种是通过CLIP的图像 编码器结合支撑样本生成的视觉原型分类器,充分 利用CLIP的视觉表示能力,另一种是通过CLIP的 语义编码器结合文本提示生成的语义原型分类器, 利用其文本语义理解能力增强图像识别,如图1所 示。实验结果显示,视觉原型分类器通常不及传统 最优方法,而语义原型分类器虽然在闭集分类方面 显著优于传统最优方法,但在开集识别性能方面表 现不佳。对于前者,主要原因可能是CLIP的训练 数据与FSOR目标任务数据之间的分布差异。对 于后者,语义原型不依赖于数据分布的特性使得分 类器为已知类别形成了更广泛的决策边界,这在闭 集设定下具有优势,但对需要紧凑决策边界以有效 拒绝未知类别的开集场景却是不利的[15]。上述两 种原型分类器决策边界的对比如图2所示,其中蓝 色和绿色的圆分别代表了语义原型分类器和视觉 原型分类器的决策边界,其半径之比由图1中的余 弦相似度获得。可以看出,语义原型分类器较宽的 央策边界虽然能够更好地分类已知类别样本,容纳 更大的类内差异性,但也削弱了模型区分已知和未 知类别的能力,增加了混淆的风险。

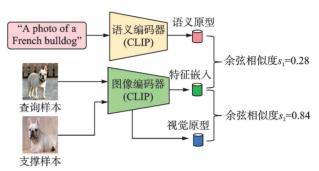


图 1 两种原型分类器示意图

Fig.1 Overview of the two prototype classifiers

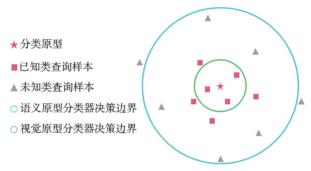


图 2 两种原型分类器决策边界对比

Fig.2 Comparison of decision boundaries of two prototype classifiers

为了解决这些问题,本文提出了一种简单有效且无需训练的视觉语义增强的联合小样本开集分类器(Visual-semantic enhanced joint classifier, VS-JC),巧妙结合了两种不同的原型分类器。一种是基于文本提示增强的CLIP语义原型分类器,具有出色的闭集分类能力;另一种是利用传统FSOR预训练模型构建的先验增强的视觉原型分类器,包含更多对目标任务数据分布的先验信息,能够为已知类别提供更紧凑的决策边界。通过两者的有机结合实现优势互补,获得了性能的双重提升。

# 1 相关工作

## 1.1 小样本开集识别

作为小样本学习与开集识别的交叉领域,小样本开集识别因其内在的挑战性和实际意义,正逐渐成为研究的热点。现有方法大致可以分为两类:基于阈值的方法和基于负原型的方法。基于阈值的

方法通常设置分类置信度阈值,当模型预测的置信 度低于该阈值时,便将该样本视为未知类。Liu等[4] 借鉴并扩展了在小样本学习中已被证明有效的元 学习策略,同时提出了一种基于最大熵的开集识别 损失,显著增强了模型对未知类别的鲁棒性。 Wang 等[5]提出了一种基于全局能量分数的阈值方 法,该方法通过计算类间和像素间的相似性来获得 能量得分,并利用该得分设置阈值,从而实现对未 知类别样本的有效识别。CHE 等[6]引入了一种即 插即用的多关系边缘训练损失,以更好地捕捉已知 类和未知类的分布特征。Jeong 等[9]利用变换一致 性原理,通过比较原始原型与经过变换后将查询样 本特征替换的原型之间的距离差异以有效地检测 未知类别样本。基于负原型的方法为未知类生成 一个或多个分类原型,将原本的N类开集识别任务 转换为(N+1)类闭集分类任务。Huang等[7]创新 性地引入了自注意力机制,通过将已知类原型输入 Transformer模型以生成任务自适应的未知类原型, 显著提升了模型在动态任务场景下的表现。Song 等[8]则提出了一种独特的训练策略,将已知类别样 本的背景区域视为未知类别,通过对背景分布的建 模,进一步丰富了分类器对未知类别的区分能力。

## 1.2 视觉-语言模型

近年来,将图像和文本映射到共同特征空间进 行多模态表示学习的VLMs在计算机视觉领域引 起了广泛关注。其中,CLIP[10]是一个具有代表性 的模型,它通过自监督对比学习,在包含数百万图 像文本对的互联网数据集上训练,并在各种下游任 务中取得了出色的表现[16-18]。许多后续研究致力 于开发更有效的方法,以便更好地将CLIP应用于 下游的小样本分类任务。例如,Zhou等[19]提出用 可学习的文本标记替代传统手工设计的文本提示 词,通过小样本训练数据进行微调,以自适应生成 更符合任务需求的文本提示。Zhou等[20]在此基础 上加入了一个轻量级网络,为每个输入图像生成条 件标记,从而动态适应不同图像的内容,有效缓解 了对基类的过拟合问题。Gao等[21]引入了一个轻 量级适配模块,该模块通过利用目标任务中的训练 数据进行微调,从而生成更加适配任务需求的特征 表示。Zhang等[22]通过引入键值缓存模型,在微调 过程中加速了模型的收敛速度,以实现更高效的训 练过程。为了充分挖掘更多的先验知识, Zhang 等[23]将 CLIP 与 DALL-E[11]和 DINO[24]等预训练模 型进行级联,通过集成这些强大的基础模型,有效 提升了模型的泛化能力。He等[25]在CLIP的图像 编码器中集成了一个基于 Transformer 的网络模 块,并通过微调以生成判别性强且任务自适应的原 型特征,提升了模型在特定任务中的表示能力。

尽管上述研究在小样本学习领域取得了显著进展,但基于 VLMs的研究中,专门针对 FSOR问题的工作仍然较少。因此本文以 CLIP 为例,深入探讨其在 FSOR问题中的应用潜力。

# 2 问题定义

FSOR旨在通过有限的标注数据,既能精确分 类已知类别样本,又能有效检测和识别未知类别样 本。在一个具体的FSOR任务中,假设总类别数为 N,每个类别的支撑样本(训练样本)数为K,则该 任务被称为 N-way K-shot 任务。具体地,一个 N-way K-shot 的 FSOR 任务可以表示为 G=  $(D^{\mathrm{s}},D^{\mathrm{q}})$ , 其中  $D^{\mathrm{s}} = \{ x_i^{\mathrm{s}} \in X^{\mathrm{s}}, y_i^{\mathrm{s}} \in Y^{\mathrm{s}} \}_{-}^{N \times K}$  表示 支撑集(训练集), $D^{Q} = \{x_{i}^{q} \in X^{Q}, y_{i}^{q} \in Y^{Q}\}_{i=1}^{m}$ 表示 查询集(测试集), m为查询样本的总数。与传统闭 集小样本学习不同,查询集 $D^{q}$ 不仅包括已知类别 的查询集 $D^{K} = \{x_i^{K} \in X^{K}, y_i^{K} \in Y^{K}\}_{i=1}^{N \times n},$ 还包括未 知类别的查询集 $D^{\text{U}} = \{x_i^{\text{u}} \in X^{\text{U}}, y_i^{\text{u}} \in Y^{\text{U}}\}_{i=1}^{N^{\text{U}} \times n},$ 其 中  $Y^{s} \cap Y^{U} = \emptyset$ ,  $N^{U}$  和 n 分别表示未知类别的总 数和每个类别的查询样本数量。FSOR任务的目 标是通过有限的支撑集 D<sup>s</sup>构建分类器,以准确分 类来自DK中的已知类查询样本,同时检测并拒绝 来自D<sup>U</sup>中的未知类查询样本。

# 3 CLIP在小样本开集识别任务中的 研究

本章节主要探讨CLIP在FSOR中的表现,采用了两种无需训练的原型分类器:视觉原型分类器(CLIP-vis)和语义原型分类器(CLIP-sem)。

### 3.1 视觉原型分类器 CLIP-vis

CLIP-vis 是一个利用 CLIP 的图像编码器  $F_{img}$  和支撑样本构建的视觉原型分类器。具体来说,给定一个 N-way K-shot 的 FSOR 任务,首先通过支撑样本计算每个已知类别  $c_i$  的视觉分类原型。具体步骤是将类别  $c_i$  的支撑样本输入 CLIP 的图像编码器,对获得的所有特征嵌入进行平均,得到的特征即为视觉原型  $V_i$ ,表示为

$$V_{i} = \frac{1}{K} \sum_{i=1}^{K} F_{img} (x_{i,j}^{s})$$
 (1)

在获得每个已知类别的视觉原型后,通过比较它们与查询样本特征之间的余弦相似度进行分类。具体而言,查询样本 $x^q$ 属于类别 $c_i$ 的概率计算为

$$P(c_{i}|x^{q}) = \frac{\exp(\cos(F_{img}(x^{q}), V_{i})/\tau)}{\sum_{i=1}^{N} \exp(\cos(F_{img}(x^{q}), V_{j})/\tau)}$$
(2)

式中: $\cos$ 表示余弦相似度; $\tau$ 为一个温度缩放因子,设置为0.01,与CLIP的设定一致。

## 3.2 语义原型分类器 CLIP-sem

CLIP-sem 是由 CLIP 的语义编码器  $F_{\text{text}}$  结合人工编写的文本提示构建的语义原型分类器。具体来说,对于每个已知类别  $c_i$ ,首先生成一个形式为"a photo of a [CLASS]"的文本提示  $T_i$ ,其中[CLASS]对应该类别的具体名称,如"cat""plane"或"apple"。然后将该文本提示输入 CLIP 的语义编码器,以获得该类别的语义分类原型  $S_i$ ,即

$$S_i = F_{\text{text}}(T_i) \tag{3}$$

通过比较语义原型与查询样本特征之间的余弦相似度,查询样本*x*°属于类别*c*;的概率计算为

$$P(c_i|x^{q}) = \frac{\exp(\cos(F_{\text{text}}(x^{q}), S_i)/\tau)}{\sum_{j=1}^{N} \exp(\cos(F_{\text{text}}(x^{q}), S_j)/\tau)}$$
(4)

### 3.3 案例研究

以两个FSOR基准数据集 TieredImageNet<sup>[26]</sup>和 CIFAR-FS<sup>[27]</sup>为例对 CLIP-vis 和 CLIP-sem 的应用潜力进行了探究,结果如图 3 所示。与传统基线方法 GEL<sup>[5]</sup>相比,CLIP-vis 在 TieredImageNet 数据集 5-way 1-shot 任务上的闭集分类 ACC 和开集识别 AUROC 指标均明显落后,在 CIFAR-FS 数据集 5-way 5-shot 任务中同样表现不佳,这可能是由于 CLIP 的训练数据与 FSOR 目标任务数据之间存在分布差异。尽管 CLIP-sem 在两种任务设定下的 ACC 指标均明显优于 GEL,但在 AUROC 指标上却明显不足,其主要原因可能是未能为已知类别建立紧凑的决策边界,这表明可能需要获得更多的数据分布信息以进一步提升性能。

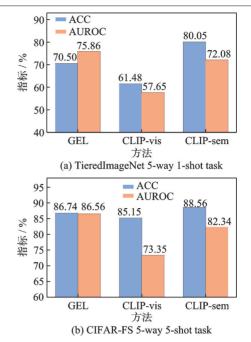


图 3 两种原型分类器与传统基线 GEL 在不同任务设定下 的性能对比

Fig.3 Performance comparison of two prototype classifiers with GEL across different settings

# 4 视觉语义增强的联合小样本开集 分类器

为了克服上述的局限性,本文提出了一种简单有效且无需训练的 VSJC,其核心思想是将基于CLIP的语义原型分类器和传统 FSOR 预训练模型的优势充分结合。前者在闭集分类中表现出色,而后者对目标任务的数据分布有更丰富的先验知识,能够为已知类别形成更紧凑的决策边界,从而提升开集识别性能。为实现这一目标,该方法结合了经过文本提示增强的 CLIP 语义原型分类器与由传统 FSOR 预训练模型生成的先验增强的视觉原型分类器,进行联合推理。图4展示了方法的整体框架。

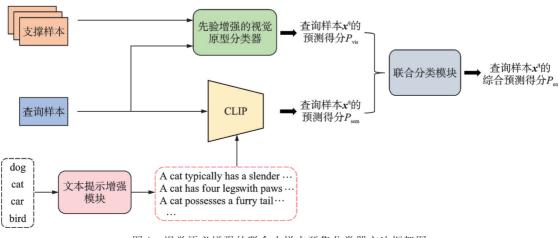


图 4 视觉语义增强的联合小样本开集分类器方法框架图

Fig.4 Framework of VSJC method

### 4.1 先验增强的视觉原型分类器

鉴于 CLIP-vis 由于分布偏移问题而表现不佳,本文未将其作为 VSJC 的视觉原型分类器。相反,本文旨在挖掘传统 FSOR 预训练模型的潜力来缓解数据分布差异。具体而言,本文使用在基类数据集上预训练的图像编码器  $F_{base}$ 来构建视觉原型分类器,充分利用其对目标任务数据分布的先验知识以增强性能。本文提出了一种先验增强的视觉原型分类器 FSOR-vis。给定  $1 \land N$ -way K-shot 的FSOR 任务,首先使用支撑样本计算每个类别  $c_i$  的的视觉分类原型  $V_i$ ,即

$$V_{i} = \frac{1}{K} \sum_{i=1}^{K} F_{\text{base}} (x_{i,j}^{s})$$
 (5)

然后,通过计算查询样本特征与各类别视觉原型之间的余弦相似度,得到其预测得分。具体而言,对查询样本 $x^q$ 的预测得分 $P_{vis}$ 可以表示为

$$P_{\text{vis}} = \left\{ P_{\text{vis}}^{i} \right\}_{i=1}^{N} \tag{6}$$

式中: $P_{vis}^{i}$ 表示查询样本 $x^{q}$ 属于第i个类别的预测得分,具体计算为

$$P_{\text{vis}}^{i} = \cos(F_{\text{base}}(x^{q}), V_{i}) \tag{7}$$

## 4.2 文本提示增强的语义原型分类器

为了更好地将 CLIP-sem 应用于 FSOR 问题,本文在其基础上引入了一个文本提示增强模块,利用语言模型 GPT- $3^{[28]}$ 来生成更细致的文本提示,这有利于丰富上下文语义信息,获得更具代表性的语义原型,从而形成更紧凑的决策边界。本文提出了一种文本提示增强的语义原型分类器。具体而言,给定一个 N-way K-shot 的 FSOR 任务,使用统一的模板作为输入 Input (例如"what a [CLASS] looks like")来查询 GPT-3,为每个已知类别获得  $n_i$ 种更加丰富具体的文本提示。将为类别  $c_i$ 生成的文本提示记作  $T_i^{aug}$ ,有

$$T_i^{\text{aug}} = \text{GPT-3}(\text{input}_c)$$
 (8)

然后通过对所生成文本提示编码后的特征取 均值,计算每个类别的精细化语义原型,即

$$S_{i} = \frac{1}{n_{t}} \sum_{j=1}^{n_{t}} F_{\text{text}} (T_{i,j}^{\text{aug}})$$
 (9)

和 FSOR-vis 类似, 查询样本  $x^{q}$ 的预测得分  $P_{sem}$ 可以表示为

$$P_{\text{sem}} = \left\{ P_{\text{sem}}^{i} \right\}_{i=1}^{N} \tag{10}$$

$$P_{\text{sem}}^{i} = \cos(F_{\text{img}}(x^{q}), S_{i})$$
 (11)

#### 4.3 视觉-语义联合推断

为充分发挥上述两种原型分类器的优势,本文将其结合进行联合推断。具体来说,通过对两个分类器的预测得分 $P_{vis}$ 和 $P_{sem}$ 进行加权求和,然后通过 softmax 归一化,得到查询样本 $x^{q}$ 的综合预测得

$$\mathcal{P}_{en} = \{P_{en}^i\}_{i=1}^N$$
,具体计算为

 $P_{\text{en}} = \operatorname{softmax}(\alpha \cdot P_{\text{vis}} + (1 - \alpha) \cdot P_{\text{sem}})$  (12) 式中: $\alpha$ 为一个取值范围在 $0 \sim 1$ 的权重超参数。

最后,通过对综合预测得分 $P_{en}$ 进行阈值判断,确定查询样本 $x^{q}$ 的最终分类结果 $y^{q}$ 。将阈值表示为n,有

$$y^{q} = \begin{cases} \hat{n} & P_{\text{en}}^{\hat{n}} > \eta \\ \text{unknown} & \text{Otherwise} \end{cases}$$
 (13)

$$\hat{n} = \operatorname{argmax}_{i} \{ P_{\text{en}}^{i} \}_{i=1}^{N}$$
 (14)

## 5 实验与分析

## 5.1 实验设置

### 5.1.1 数据集与评估指标

为验证所提方法的有效性,本文在FSOR领域 广泛认可的 4 个基准数据集: MiniImageNet<sup>[29]</sup>、TieredImageNet<sup>[26]</sup>、CIFAR-FS<sup>[27]</sup>和 CUB<sup>[30]</sup>上进行了实验。本文采用准确率(Accuracy, ACC)和受试者工作特征曲线下的面积(Area under the receiver operating characteristic, AUROC)作为主要评估指标,其中ACC用于衡量模型对已知类查询样本的闭集分类准确度,AUROC用于评估模型检测未知类别的能力。根据现有研究的设置<sup>[4-9]</sup>,本文在600个5-way 1-shot或5-way 5-shot的FSOR任务上对模型进行测试,并计算平均ACC和AUROC指标。

## 5.1.2 实现细节

本文采用传统方法 ATT- $G^{[7]}$ 中基于残差网络(Residual network12, ResNet-12) $^{[31]}$ 的预训练模型作为 FSOR-vis 中的图像编码器  $F_{base}$ 。对于CLIP $^{[10]}$ ,使用预训练的 ResNet- $50^{[31]}$ 作为图像编码器  $F_{img}$ ,并使用与之对应的 Transformer 作为文本编码器  $F_{text}$ 。对于文本提示增强策略,将由GPT- $3^{[28]}$ 为每个类别生成的文本提示数量  $n_i$ 设置为 20。关于权重超参数  $\alpha$ ,其选择通过在验证集上进行网格搜索确定。对于 5-way 1-shot 任务,在MiniImageNet、TieredImageNet、CIFAR-FS 和CUB数据集上的设置分别为 0.04、0.15、0.06 和 0.07; 对于 5-way 5-shot 任务,分别设置为 0.07、0.28、0.1 和 0.05。

### 5.2 方法对比与分析

### 5.2.1 与无需训练的方法对比

为了验证本文方法的优势,将 VSJC 与两种典型的无需训练的方法 CLIP-vis 和 CLIP-sem,以及 6种经典的传统 FSOR 方法 (PEELER<sup>[4]</sup>, SnaTCH-er<sup>[9]</sup>, ATT-G<sup>[7]</sup>, MRM<sup>[6]</sup>和 GEL<sup>[5]</sup>, FSOR-vis)进行

%

了比较,结果如表1和表2所示。通过结合由传统FSOR预训练模型构建的先验增强的视觉原型分类器FSOR-vis进行联合推理,VSJC在ACC和AUROC方面相比CLIP-sem取得了显著提升。例如,在TieredImageNet 5-way 1-shot任务中,VSJC在

ACC上提升了5.5%,在AUROC上提升了6.9%;在CIFAR-FS5-way1-shot任务中,VSJC在ACC和AUROC上分别达到92.14%和86.51%,相较于CLIP-sem分别提高了3.6%和4.2%,充分体现了联合推理策略在小样本开集识别任务中的有效性。

表 1 在 MiniImageNet 和 TieredImageNet 上与无需训练方法的对比

Table 1 Comparison results with training-free methods on MiniImageNet and TieredImageNet

	MiniImageNet				TieredImageNet			
方法	1-shot		5-shot		1-shot		5-shot	
	ACC	AUROC	ACC	AUROC	ACC	AUROC	ACC	AUROC
PEELER <sup>[4]</sup>	65.86	60.57	80.61	67.35	69.51	65.20	84.10	73.27
SnaTCHer <sup>[9]</sup>	67.02	68.27	82.02	72.42	70.52	74.28	84.74	82.02
$ATT-G^{[7]}$	68.11	72.41	83.12	79.85	70.58	73.43	<u>85.38</u>	81.64
$MRM^{[6]}$	67.03	71.20	82.00	80.39	71.13	75.59	85.27	83.03
$\operatorname{GEL}^{\scriptscriptstyle{[5]}}$	68.26	73.70	83.05	82.29	70.50	<u>75.86</u>	84.60	81.95
FSOR-vis	65.38	63.64	83.65	73.37	68.99	67.71	84.67	75.71
CLIP-vis	80.30	67.47	93.34	75.80	61.48	57.65	79.61	64.63
CLIP-sem	96.82	92.23	96.82	92.23	80.05	72.08	80.05	72.08
VSJC	97.56	94.74	97.80	94.94	85.57	78.95	90.21	82.43

表 2 在 CIFAR-FS 和 CUB 上与无需训练方法的对比

Table 2 Comparison results with training-free methods on CIFAR-FS and CUB

CIFAR-FS **CUB** 方 法 1-shot 5-shot 1-shot 5-shot ACC AUROC ACC AUROC ACC AUROC ACC AUROC PEELER<sup>[4]</sup> 71.47 71.28 85.46 75.97 62.62 57.26 82.44 65.01 SnaTCHer<sup>[9]</sup> 75.09 78.15 87.18 85.81 69.03 74.34 83.77 84.57 ATT-G<sup>[7]</sup> 90.02 72.43 76.72 86.52 84.64 79.27 81.73 90.01  $MRM^{[6]}$ \_ 70.00 76.30 84.5185.91  $GEL^{\scriptscriptstyle [5]}$ 76.77 86.56 79.00 89.75 89.95 78.67 86.74 81.78 FSOR-vis 70.45 70.46 86.47 80.35 67.5468.10 85.38 82.72 CLIP-vis 65.25 61.67 85.1566.69 94.61 79.24 73.35 78.67 CLIP-sem 81.98 88.56 82.34 <u>88.56</u> 82.34 95.63 81.98 95.63 VSJC 92.14 86.51 94.30 88.65 96.20 88.39 97.16 90.41

此外,与传统FSOR方法相比,VSJC在大多数数据集上均表现出显著的性能优势。例如,与传统最优方法GEL和ATT-G相比,VSJC在4个基准数据集上的ACC和AUROC指标分别平均提升了14%、6.8%和14.4%、8.1%。这些结果充分证明了VSJC在小样本开集识别任务中的卓越性能与广泛适用性。

# 5.2.2 与基于微调的方法对比

为了进一步证明本文方法的有效性,将两种流行 的 CLIP 微 调 策 略 (CLIP-Adapter<sup>[21]</sup> 和 Tip-Adapter<sup>[22]</sup>)适配到 CLIP-sem上,并进行了对比分析,如表3和表4所示。结果显示,尽管 CLIP-Adapter和 Tip-Adapter通过在有限的支撑集上进行微调可以显著提升性能,但它们在4个基准

数据集上的表现仍然落后于本文提出的 VSJC 方法,并且后者无需额外的训练开销。例如,在 TieredImageNet 数据集上,VSJC 的 ACC 相较于CLIP-Adapter 和 Tip-Adapter 分别平均提升了4.5%和5.9%, AUROC 分别提升了7.0%和7.7%;在 CIFAR-FS 数据集上,VSJC 的 ACC 和AUROC 分别相较于 CLIP-Adapter提高了3.2%和3.8%,相较于 Tip-Adapter提高了3.4%和4.1%。这些结果进一步验证了 VSJC 在分类精度和开集识别性能方面的显著优势。

## 5.3 消融实验

为了评估所提出的 VSJC 模型各个组件的有效性,本文在 TieredImageNet 的 5-way 1-shot 任务和 CIFAR-FS 的 5-way 1-shot 任务上进行了消融

%

%

### 表 3 MiniImageNet 和 TieredImageNet 上与基于微调方法的对比

Table 3 Comparison results with fintuning methods on MiniImageNet and TieredImageNet

方法	MiniImageNet			TieredImageNet				
	1-shot		5-shot		1-shot		5-shot	
	ACC	AUROC	ACC	AUROC	ACC	AUROC	ACC	AUROC
CLIP-sem	96.82	92.23	96.82	92.23	80.05	72.08	80.05	72.08
CLIP-Adapter	96.99	92.37	97.08	92.45	81.85	<u>72.82</u>	<u>85.03</u>	<u>74.63</u>
Tip-Adapter	97.19	92.60	<u>97.25</u>	92.70	81.39	72.50	82.56	73.45
VSJC	97.56	94.74	97.80	94.94	85.57	78.95	90.21	82.43

表 4 CIFAR-FS和CUB上与基于微调方法的对比

Table 4 Comparison results with fintuning methods on CIFAR-FS and CUB

	CIFAR-FS				CUB			
方 法	1-shot		5-shot		1-shot		5-shot	
	ACC	AUROC	ACC	AUROC	ACC	AUROC	ACC	AUROC
CLIP-sem	88.56	82.34	88.56	82.34	95.63	81.98	95.63	81.98
CLIP-Adapter	<u>89.62</u>	83.41	90.38	84.20	95.79	82.13	<u>96.16</u>	<u>87.25</u>
Tip-Adapter	89.50	83.23	90.18	83.78	<u>95.96</u>	82.25	96.09	82.30
VSJC	92.14	86.51	94.30	88.65	96.20	88.39	97.16	90.41

实验,结果如表5和表6所示。

表 5 TieredImageNet 5-way 1-shot任务上的消融实验
Table 5 Ablation experiment on TieredImageNet 5-way
1-shot task

	组成部分	性能			
CLIP-sem	文本提示增强	联合推断	ACC/%	AUROC/%	
$\overline{}$	_	_	80.05	72.08	
$\checkmark$	$\checkmark$	_	81.09	72.71	
$\checkmark$	_	$\checkmark$	85.08	78.67	
	$\checkmark$	$\checkmark$	85.57	78.95	

表 6 CIFAR-FS 5-way 1-shot任务上的消融实验
Table 6 Ablation experiment on CIFAR-FS 5-way
1-shot task

	组成部分	性能		
CLIP-sem	文本提示增强	联合推断	ACC/%	AUROC/%
$\overline{\hspace{1cm}}$	_	_	88.56	82.34
$\checkmark$	$\checkmark$		90.10	83.88
$\checkmark$	_	$\checkmark$	91.22	86.20
	$\checkmark$	$\checkmark$	92.14	86.51

实验表明,文本提示增强与联合推断策略均能有效提升闭集分类准确率 ACC 和开集识别性能 AUROC。具体而言,文本提示增强策略通过借助 GPT-3丰富语义上下文信息,在不降低 ACC 的前提下,使 AUROC 在两项任务中平均提升了 1.1%。同时,通过联合推断策略,充分结合两种原型分类器的优势,ACC 和 AUROC 指标分别在两项任务中平均提升了 3.8% 和 5.2%。此外,结合这两种策略可以实现最佳性能表现。

## 5.4 超参数敏感性分析

为了评估 VSJC 模型对于权重超参数  $\alpha$  的敏感性,本文在 MiniImageNet 的 5-way 5-shot 任务和 CIFAR-FS 的 5-way 1-shot 任务上进行了实验,其中 $\alpha$  的值分别被设置为 $\{0.05,0.06,0.07,0.08\}$ ,结果如图 5所示。实验结果表明,在给定范围内,权重超参数 $\alpha$ 的选择对性能的影响总体上并不显著。然而,精细调整 $\alpha$ 的取值仍然能够带来性能的提升。

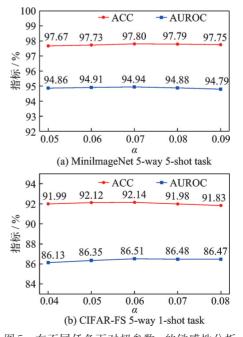


图 5 在不同任务下对超参数 α 的敏感性分析 Fig. 5 Sensitivity analysis of hyperparameter α across different tasks

# 6 结 论

本文探究了视觉-语言预训练模型 CLIP 在FSOR中的应用潜力,实验发现:

- (1)基于CLIP图像编码特征的视觉原型分类器通常不如传统FSOR基线方法,这可能是由于CLIP的预训练数据集与FSOR目标任务数据集之间存在数据分布差异。
- (2)基于CLIP语义编码特征的语义原型分类器虽然在闭集分类上显著优于传统基线,但开集识别性能不佳,这可能是因为语义原型分类器为已知类别划分了过大的决策边界。

针对这些问题,本文提出了一种简单有效且无需训练的视觉语义增强的联合小样本开集分类器。该方法充分利用了基于文本提示增强的CLIP语义原型分类器在闭集分类中的优势,同时发挥了由传统FSOR预训练模型构建的先验增强的视觉原型分类器的潜力,通过其对目标任务数据的先验知识,缓解数据分布差异,并为已知类别建立更紧密的决策边界,从而更好地检测未知类别,实现性能的双重提升。实验结果表明,本文所提出的方法在ACC和AUROC指标上取得了显著的提升,并在大多数数据集上取得了较先进的性能。

### 参考文献:

- [1] 张玉,尚志华,郭晓楠,等.小样本图像分类中的类别信息融合网络[J].南京航空航天大学学报,2022,54(4):715-722.
  - ZHANG Yu, SHANG Zhihua, GUO Xiaonan, et al. A category information fusion network for few-shot image classification[J]. Journal of Nanjing University of Aeronautics & Astronautics, 2022, 54(4): 715-722.
- [2] 徐惠灵,尚政国,董胜波,等.面向深度神经网络应用的小样本学习技术研究[J].南京航空航天大学学报,2022,54(S1):80-86.
  - XU Huiling, SHANG Zhengguo, DONG Shengbo, et al. A study on few-shot learning techniques for deep neural network applications[J]. Journal of Nanjing University of Aeronautics & Astronautics, 2022, 54 (S1): 80-86.
- [3] ZHOU Y, HAO J, HUO S, et al. Automatic metric search for few-shot learning [J]. IEEE Transactions on Neural Networks and Learning Systems, 2023, 35(7): 10098-10109.
- [4] LIU B, KANG H, LI H, et al. Few-shot open-set recognition using meta-learning[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2020.
- [5] WANG H, PANG G, WANG P, et al. Glocal energy-based learning for few-shot open-set recognition

- [C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2023: 7507-7516.
- [6] CHE Y, AN Y, XUE H. Boosting few-shot open-set recognition with multi-relation margin loss[C]//Proceedings of the IJCAI. China: IJCAI, 2023: 3505-3513.
- [7] HUANG S, MA J, HAN G, et al. Task-adaptive negative envision for few-shot open-set recognition [C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA; IEEE, 2022; 7171-7180.
- [8] SONG N, ZHANG C, LIN G. Few-shot open-set recognition using background as unknowns[C]//Proceedings of the 30th ACM International Conference on Multimedia. New York: ACM, 2022: 5970-5979.
- [9] JEONG M, CHOIS, KIM C. Few-shot open-set recognition by transformation consistency[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2021; 12566-12575.
- [10] RADFORD A, KIM J W, HALLACY C, et al. Learning transferable visual models from natural language supervision[C]//Proceedings of the International Conference on Machine Learning. New York: PMLR, 2021: 8748-8763.
- [11] REDDY M D M, BASHA M S M, HARI M M C, et al. DALL-E: Creating images from text[J]. UGC Care Group I Journal, 2021, 8(14): 71-75.
- [12] WUX, ZHUF, ZHAOR, et al. Cora: Adapting clip for open-vocabulary detection with region prompting and anchor pre-matching[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2023: 7031-7040.
- [13] ZHAO S, ZHANG Z, SCHULTER S, et al. Exploiting unlabeled data with vision and language models for object detection [C]//Proceedings of the European Conference on Computer Vision. Heidelberg: Springer-Verlag GmbH, 2022: 159-175.
- [14] WU W, SUN Z, OUYANG W. Revisiting classifier: Transferring vision-language models for video recognition [C]//Proceedings of the AAAI Conference on Artificial Intelligence. Washington, DC: AAAI Press, 2023.
- [15] ROADY R, HAYES T L, KEMKER R, et al. Are open set classification methods effective on large-scale datasets[J]. PLoS ONE, 2020, 15(9): e0238302.
- [16] WANG J, CHAN K C, LOY C C. Exploring clip for assessing the look and feel of images [C]//Proceedings of the AAAI Conference on Artificial Intelligence. Washington, DC: AAAI, 2023.

- [17] WASIM S T, NASEER M, KHAN S, et al. Vita-clip: Video and text adaptive clip via multimodal prompting[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2023: 23034-23044.
- [18] LUO H, JI L, ZHONG M, et al. CLIP4Clip: An empirical study of CLIP for end to end video clip retrieval and captioning [J]. Neurocomputing, 2022, 508(C): 293-304
- [19] ZHOU K, YANG J, LOY C C, et al. Learning to prompt for vision-language models[J]. International Journal of Computer Vision, 2022, 130(9): 2337-2348.
- [20] ZHOU K, YANG J, LOY C C, et al. Conditional prompt learning for vision-language models[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2022.
- [21] GAO P, GENG S, ZHANG R, et al. Clip-adapter: Better vision-language models with feature adapters [J]. International Journal of Computer Vision, 2024, 132(2): 581-595.
- [22] ZHANG R, FANG R, ZHANG W, et al. Tip-adapter: Training-free clip-adapter for better vision-language modeling [EB/OL]. (2021-11-15).https://arxiv.org/abs/2111.03930.
- [23] ZHANG R, HU X, LI B, et al. Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2023: 15211-15222.
- [24] CARON M, TOUVRON H, MISRA I, et al.

- Emerging properties in self-supervised vision transformers[C]//Proceedings of the IEEE/CVF International Conference on Computer Vision. Los Alamitos, CA: IEEE, 2021: 9650-9660.
- [25] HEF, LIG, SIL, et al. Prototype Former: Learning to explore prototype relationships for few-shot image classification [EB/OL]. (2023-10-07). https://arxiv.org/abs/2310.03517.
- [26] REN M, TRIANTAFILLOU E, RAVI S, et al. Meta-learning for semi-supervised few-shot classification [EB/OL]. (2018-03-02). https://arxiv.org/abs/1803.00676.
- [27] KRIZHEVSKY A, NAIR V, HINTON G. Cifar-10 (canadian institute for advanced research)[EB/OL]. (2025-10-03). http://www.cs.toronto.edu/kriz/cifar.html
- [28] BROWN T B. Language models are few-shot learners [EB/OL]. (2020-05-28). https://arxiv. org/abs/2005.14165.
- [29] VINYALS O, BLUNDELL C, LILLICRAP T, et al. Matching networks for one shot learning [C]//Proceedings of the 30th International Conference on Neural Information Processing Systems. Barcelona: Curran Associates Inc., 2016: 3637-3645.
- [30] WAH C, BRANSON S, WELINDER P, et al. The caltech-ucsd birds-200-2011 dataset: CNS-TR-2011-001[R].[S.l.]:[s.n.], 2011.
- [31] HE K, ZHANG X, REN S, et al. Deep residual learning for image recognition[C]//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Los Alamitos, CA: IEEE, 2016: 770-778.

(编辑:刘彦东)