

用于高分辨遥感影像场景分类的迁移学习 混合专家分类模型

龚希¹, 陈占龙^{1,2}, 吴亮^{1,2}, 谢忠^{1,2*}, 徐永洋^{1,2}

¹中国地质大学(武汉)地理与信息工程学院, 湖北 武汉 430074;

²国家地理信息系统工程技术研究中心, 湖北 武汉 430074

摘要 针对小样本遥感影像场景数据集中地物多样性和分布复杂性引起的分类精度低下的问题, 提出一种基于迁移学习的混合专家(TLMoE)分类模型。该模型通过多通道充分利用包含场景全局信息的全连接层特征和包含场景局部细节信息的卷积层特征, 能够实现更精确的场景分类。基于全连接层特征的预判通道, 利用场景全局信息完成对全部类别场景的初判; 通过专家通道为每类场景训练专属专家网络, 针对性地挖掘各类场景卷积层特征中蕴含的关键局部信息, 提取可区分相似场景间细微差异的局部特征, 完成细粒度的识别; 结合预判权重实现顾及场景全局及局部差异的分类。在小样本数据集上的实验表明, 本文方法可有效识别易混淆场景, 能够取得较好的分类效果。

关键词 遥感; 高分辨率遥感影像; 场景分类; 混合专家系统; 迁移学习

中图分类号 P237

文献标志码 A

doi: 10.3788/AOS202141.2301003

Transfer Learning Based Mixture of Experts Classification Model for High-Resolution Remote Sensing Scene Classification

Gong Xi¹, Chen Zhanlong^{1,2}, Wu Liang^{1,2}, Xie Zhong^{1,2*}, Xu Yongyang^{1,2}

¹ School of Geography and Information Engineering, China University of Geosciences, Wuhan, Hubei 430074, China;

² National Engineering Research Center of Geographic Information System, Wuhan, Hubei 430074, China

Abstract To tackle the low classification accuracy caused by the diversity and distribution complexity of surface objects in small-sample datasets of remote sensing image scenes, this paper proposes a transfer learning based mixture of experts (TLMoE) classification model. The model can achieve more accurate scene classification by taking full advantage of the features from the convolution layer containing the local details and the fully-connected layer containing the global information of scenes through multi-channels. First, a pre-judgment channel based on the fully-connected layer features is established to preliminarily judge all kinds of scenes with global scene information; then exclusive expert networks are trained for each kind of scenes via the expert channel, which can mine the key local details contained in the convolution layer features of all categories of scenes targetedly and extract the local features used to distinguish the subtle differences between similar scenes to complete fine-grained identification. Finally, combined with the pre-judged weight, the model realizes the scene classification considering the global and local differences. Experiments on small-sample datasets show that the proposed method can effectively identify confusing scenes and achieve good classification results.

收稿日期: 2021-01-25; 修回日期: 2024-04-20; 录用日期: 2021-06-10

基金项目: 国家自然科学基金(42001340, U1711267, 41871305)、国家重点研发计划(2018YFB0505500, 2018YFB0505504)、地质探测与评估教育部重点实验室开放基金(GLAB2020ZR05)、自然资源部城市国土资源监测与仿真重点实验室开放基金(KF-2020-05-068)

通信作者: *xiezhong@cug.edu.cn

Key words remote sensing; high-resolution remote sensing image; scene classification; mixture of experts system; transfer learning

OCIS codes 010.0280; 100.3008; 100.2960

1 引言

随着遥感传感器技术和制图技术的快速发展,遥感影像的质量与数量得到极大提升,利用遥感影像解译技术挖掘海量遥感数据中的兴趣知识并实现信息的增值具有重要意义。遥感影像场景分类作为一项重要的解译技术,可自动提取并识别遥感场景中丰富的语义信息,在遥感数据知识挖掘中的重要性凸显。然而遥感影像场景具有颜色纹理信息丰富、地物种类多样、空间分布复杂等特点,如何对遥感影像场景进行区分度更强、准确度更高的场景特征提取与分类成为极具挑战的课题,已引起众多研究者的关注^[1-2]。

早期分类方法主要依赖人工设计特征完成分类,这类方法使用的中低层特征缺乏对高层语义信息的表达,难以跨越遥感影像场景的“语义鸿沟”^[3],存在一定的局限性。近年来,深度学习在多个领域取得巨大成功,包括遥感影像场景分类领域^[4-5]。尤其是在图像领域大放异彩的卷积神经网络(CNN),大幅提升了遥感影像场景分类的准确度^[6]。但基于CNN的遥感影像场景分类仍面临两个方面的问题。问题一:CNN的诸多优点建立于大量标注的训练样本之上,而当前标注的遥感影像场景集的样本量往往较小,难以达到百万数量级的要求;问题二:遥感影像地物多样、分布复杂,遥感影像场景数据集中存在诸多外观差异显著的场景类别,同时场景的类内多样性与类间相似性使得部分类别场景间外观差异非常细微,传统CNN模型仅使用单个分类器,难以顾及各类场景间不同程度的差异,无法在考虑全部场景类别的同时实现对相似类别的细粒度分类。

通过迁移学习利用预训练CNN特征可有效解决问题一中样本量不足的问题,如:文献[7]通过多个预训练CNN模型证明全连接(FC)层特征可充分挖掘遥感影像场景的全局语义信息;文献[5]证明预训练CNN的卷积层特征可对遥感影像场景的局部信息进行表达,将BoVW编码后的卷积层特征与全连接层特征融合,能够实现对遥感影像场景的高效表达。由于CNN的卷积层特征和全连接层特征是对图像局部信息和全局信息的抽象表达^[8],基于预训练CNN特征的方法可快速地提升对小样本数

据集中场景特征的表达能力,但也会导致很多预训练网络层特征并未经过针对性训练,因此未能充分挖掘有用信息。而针对问题二中类内多样性和类间相似性引起的分类精度低下的问题,学者们尝试从网络的结构设计^[9]、输入增强^[10]及训练约束^[11]等多个方面对CNN模型进行改进,在一定限度内增强了模型对场景特征的描述能力,获得更优的分类结果,但仍面临计算复杂度过高、子通道训练未充分结合场景语义信息、小样本数据集仍需额外的数据增强来提升样本量等问题。

为充分利用预训练CNN特征中蕴含的场景全局信息和局部信息来克服小样本遥感影像场景数据集类内差异大、类间相似度高造成的分类精度低下的问题,本文提出了一种基于迁移学习的混合专家(TLMoE)分类模型,分别通过预判通道和专家通道挖掘场景全连接层特征中的全局信息和卷积层特征中的局部信息。TLMoE利用场景全局特征来实现对全部遥感场景类别的初判,同时为弥补单一预判分类器对差异细微、易混淆场景识别能力的不足,通过专家通道建立多个专家网络并针对性地学习各类场景的局部细节信息,提升对相似场景间细微差别的分辨能力,最后结合预判权重完成更准确的场景分类。该模型实现了由全局信息挖掘到局部信息增强、由整体类别识别到单一类别判断的特征学习和分类过程,它利用多个通道针对性地学习不同表达层次和抽象程度的场景特征,顾及了不同类别场景间的显著差异和细微差异,实现了精度更高的遥感影像场景分类。

2 基于迁移学习的混合专家分类模型

TLMoE模型的整体网络结构如图1所示,它主要由预训练CNN特征提取器、预判通道及专家通道三个部分组成,其中预判通道和专家通道分别利用预训练CNN提取的场景全局特征和局部特征进行整体类别的识别和单一类别的判断。具体而言:预判通道针对遥感影像场景种类丰富、外观多样的特点,利用场景全局特征实现场景类别的初分类,获得场景属于各类的预判概率,并基于此为各专家网络分配训练样本及权重;专家通道则进一步针对场景的类内多样性和类间相似性的特点进行细粒度

分类,该通道构建了多个网络结构相同的二分类专家网络,并通过相应类别样本及其相似类别样本进行专家化训练,针对性地学习对应类别场景与易混

淆场景的局部细节信息,利用细节差异判断场景是否属于该类;最后通过预判概率值对专家网络的结果加权,得到综合全局信息和局部信息的分类结果。

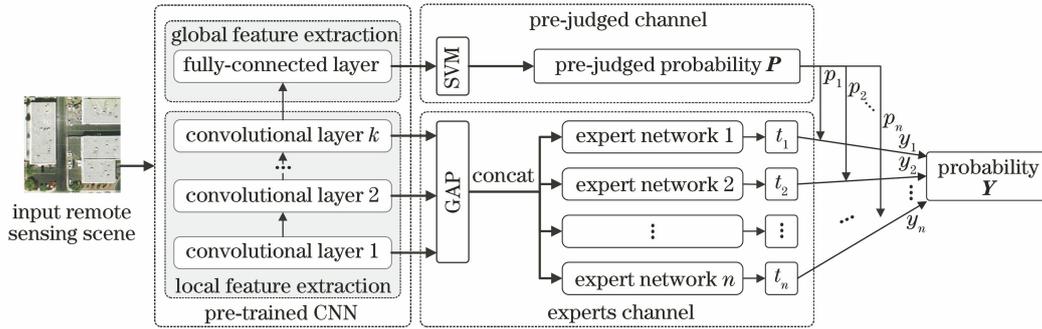


图 1 TLMoE 方法流程图

Fig. 1 Flow chart of TLMoE

2.1 场景全局与局部特征的提取

相较浅层网络,深度卷积神经网络(DCNN)对地物类型多样、空间分布复杂的遥感影像具有更强的抽象与表达能力,因此本研究通过 DCNN 进行场景特征的提取,分别以 VGG19^[12] 和 Resnet50^[13] 两个 DCNN 为例,对 TLMoE 方法进行阐述,并将采用 VGG19 和 Resnet50 的 TLMoE 模型分别记为 TLMoE-VGG19 和 TLMoE-Resnet50。VGG19 和 Resnet50 的网络层数分别为 19 和 50,它们的组成可大致归纳为表 1。两个模型结构均包含了 5 组卷

积层,且各组卷积层数目及输出的特征大小略有差异。对于任意卷积层 l ,可将其输出特征尺寸表示为 $d^{(l)} \times n^{(l)} \times n^{(l)}$, $n^{(l)} \times n^{(l)}$ 为单个特征图的大小, $d^{(l)}$ 为特征图数目。表 1 中列出了 VGG19 各组第一个卷积层及 Resnet50 各组 shortcut 层的输出特征大小,其中 shortcut 是 Resnet50 解决网络的退化问题、完成恒等映射时的重要步骤。在全连接层方面,VGG19 包含两个中间层和一个分类层,而 Resnet50 通过全局平均池化(GAP)^[14]层替代中间层,仅利用了一层分类全连接层。

表 1 VGG19 和 Resnet50 组成对比

Table 1 Structure comparison between VGG19 and Resnet50

No.	Layer group	VGG19		Resnet50	
		Layer number	Feature size	Layer number	Feature size
1	conv1	2	$64 \times 224 \times 224$	1	$64 \times 112 \times 112$
2	conv2	2	$128 \times 112 \times 112$	9	$256 \times 56 \times 56$
3	conv3	4	$256 \times 56 \times 56$	12	$512 \times 28 \times 28$
4	conv4	4	$512 \times 28 \times 28$	18	$1024 \times 14 \times 14$
5	conv5	4	$512 \times 14 \times 14$	9	$2048 \times 7 \times 7$
6	FC/GAP	2	4096	1	2048
7	output FC	1	1000	1	1000

由于预训练的 VGG19 和 Resnet50 是针对语义分类任务训练的 DCNN,其卷积层特征对图像局部细节信息的表达有较强的不变性,利用不同的卷积层特征可提取不同抽象程度下的局部细节信息^[5]。本研究将提取场景的多个卷积层特征并使其联合为多级局部特征。两个网络中,相较底层的 conv1~conv2,中高层的 conv3~conv5 的卷积层对场景的抽象程度更高,TLMoE 将选用各组第一个卷积层或 shortcut 层作为局部特征提取层,提取 3 组大小不同的特征图。对任意遥感影像场景 s ,其

多级局部特征表示为

$$F_L(s) = \{f_{conv3}, f_{conv4}, f_{conv5}\}, \quad (1)$$

式中: f_{convl} 是 conv l 组局部特征提取层输出的特征图。

全局特征表达方面,FC 层特征作为整图高度抽象的顶层特征,被广泛应用于图像的全局描述^[8]。同时在 VGG19 中,第一个 FC 层特征被证明在遥感场景分类中具有更好的表达效果^[7,15],因此 TLMoE 采用第一个 FC 层为全局特征提取层。而 Resnet50 中除分类 FC 层外,仅有 GAP 层可产生高维特征向

量,本研究采用该层作为全局特征提取层。对于任意遥感影像场景 s ,其全局特征可表示为

$$F_G(s) = \{f_{FC}\}. \quad (2)$$

采用 VGG19 和 Resnet50 时,它的维度分别为 4096 和 2048 维。

2.2 TLMoE 的预测

TLMoE 模型构建了预判通道和专家通道两类通道。在预判通道中利用支持向量机对任意场景 s 的全局特征 $F_G(s)$ 进行分类,得到预判概率,表达式为

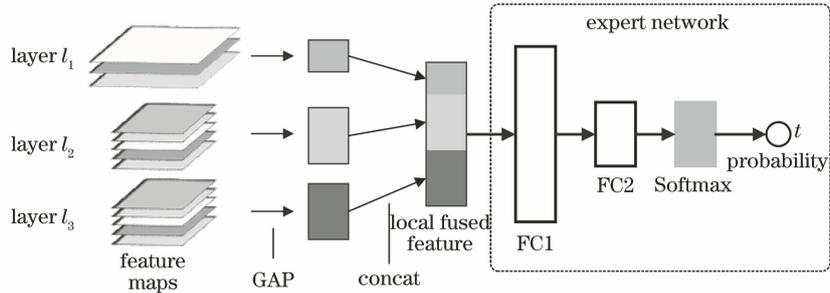


图 2 专家网络迁移学习的过程

Fig. 2 Transfer learning process of expert network

由于原始卷积层的特征维度较高,计算代价和存储成本较大,专家通道采用 GAP 对其降维。卷积层输出的每个特征图是对场景不同局部信息的特征表达,GAP 通过求取各特征图的平均值,在保持特征图对局部信息描述的相对侧重点的同时,直观地建立特征图间的联系,汇聚蕴含在局部区域间的空间信息。conv l 组局部特征提取层的输出特征 f_{convl} 可表示为特征集合 $C^{(l)} = \{c_1^{(l)}, c_2^{(l)}, \dots, c_{d^{(l)}}^{(l)}\}$,其包含了 $d^{(l)}$ 个 $n^{(l)} \times n^{(l)}$ 大小的特征图,通过 GAP 可将该卷积层特征表示为 $d^{(l)}$ 维的特征向量 \mathbf{X}^{convl} ,表达式为

$$\mathbf{X}^{convl} = (x_1^{convl}, x_2^{convl}, \dots, x_{d^{(l)}}^{convl}) = \text{GAP}(C^{(l)}) = [\text{ave}(c_1^{(l)}), \text{ave}(c_2^{(l)}), \dots, \text{ave}(c_{d^{(l)}}^{(l)})], \quad (4)$$

式中:ave($c_i^{(l)}$)为对第 i 个特征图元素的平均运算,其结果为 \mathbf{X}^{convl} 的第 i 个元素 x_i^{convl} 。

对于任意场景 s ,TLMoE 融合了 $F_L(s)$ 中 3 个卷积层的输出特征,定义该特征为场景的局部融合特征,即

$$\mathbf{X}(s) = (X^{conv3}, X^{conv4}, X^{conv5}) = (x_1^{conv3}, \dots, x_{d^{conv3}}^{conv3}, x_1^{conv4}, \dots, x_{d^{conv4}}^{conv4}, x_1^{conv5}, \dots, x_{d^{conv5}}^{conv5}), \quad (5)$$

式中: $\mathbf{X}(s)$ 为一个 $d^{conv3} + d^{conv4} + d^{conv5}$ 维的特征向量,在 VGG19 和 Resnet50 中分别为 1280 和 3584 维,该特征将作为各专家网络的输入特征。

$$\mathbf{P}(s) = (p_1, p_2, \dots, p_n), \quad (3)$$

式中: n 为场景集类别数目, p_i ($1 \leq i \leq n$) 为场景 s 属于第 i 类的概率。

在专家通道中建立了 n 个具有相同网络结构的专家网络。基于该网络分别对各类场景的局部特征进行迁移学习,将不同层次的局部特征融合,针对性地挖掘相应类别场景的关键局部信息并进一步地抽象化和降维,获取样本属于本类的概率。图 2 展示了单个专家网络对多级局部特征的迁移学习的过程。

简单的网络结构可降低模型对数据量的需求。本研究设计的专家网络仅由 2 个 FC 层及 1 个 Softmax 函数组成(图 2)。其中第 1 个 FC 层 FC1 将输入特征 $\mathbf{X}(s)$ 进一步优化并降维,再通过修正线性单元(ReLU)^[16] 激活获取 512 维的局部专家特征,即

$$\mathbf{O}_{FC1}(s) = \sigma[\mathbf{W}_1 \mathbf{X}(s) + \mathbf{b}_1], \quad (6)$$

式中: $\sigma(x) = \max(0, x)$ 为 ReLU 激活函数, \mathbf{W}_1 和 \mathbf{b}_1 为 FC1 的权重与偏置项组成的矩阵。由于专家网络为二分类网络,第 2 个 FC 层 FC2 包含 2 个神经元,Softmax 函数可获取当前场景属于本类的概率,即

$$\mathbf{O}_{FC2}(s) = \varphi[\mathbf{W}_2 \mathbf{O}_{FC1}(s) + \mathbf{b}_2], \quad (7)$$

式中: $\varphi(x) = e^x / \sum e^x$ 为 Softmax 函数, \mathbf{W}_2 和 \mathbf{b}_2 为 FC2 的权重与偏置项组成的矩阵。 $\mathbf{O}_{FC2}(s) = (1 - t, t)$ 为一个二维特征向量,其中 t 为样本属于当前类的概率。通过 n 个专家网络则可获取当前场景 s 属于各类的概率,即

$$\mathbf{T} = (t_1, t_2, \dots, t_n), \quad (8)$$

式中: t_i ($1 \leq i \leq n$) 为第 i 个专家网络判断当前场景属于第 i 类的概率。在专家通道中利用场景的局部信息完成对场景的精细分类,结合预判通道的初分类结果,可得最终的分分类概率为

$$\mathbf{Y} = (y_1, y_2, \dots, y_n) = (p_1 * t_1, p_2 * t_2, \dots, p_n * t_n). \quad (9)$$

2.3 专家通道的样本筛选

针对外观差异显著、种类繁多的遥感影像场景, 预判通道利用场景的全局特征完成初判。但对类内多样性和类间相似性产生的相似场景, 需进一步通过专家网络学习的场景关键局部信息进行分辨。因此训练对局部关键细节信息敏感的专家网络对提升 TLMoE 模型整体的分类识别能力至关重要。为对各类场景实现更准确的二分类, 各专家网络利用对应类别样本及其相似类别样本的局部融合特征进行针对性地特征优化与信息挖掘。在这一过程中, TLMoE 模型将根据预判概率动态筛选各类样本及其易混淆的相似类别样本, 以供特定专家网络进行训练。

设训练集中任意场景 s 的真实标签为 i ($1 \leq i \leq n$), 预判通道对它的预判概率为 $\mathbf{P} = (p_1, p_2, \dots, p_n)$ 。 \mathbf{P} 中较大的 k 个元素对应的类别, 为 s 在预判通道最易被误判的 k 个类, 它们对应的 k 个专家网络都需学习场景 s 的局部细节特征。显然, 当这

k 类包含第 i 类时, s 将作为第 i 个专家网络的正样本以及剩余 $k-1$ 个专家网络的负样本; 而当 k 类不包含第 i 类时, s 将作为这 k 个专家网络的负样本。

k 值过大或过小都会造成正负样本的极度失衡, 导致专家网络对正样本或负样本中的信息学习不足。仅当预判结果的 top- k 合理且 k 值适中时, 网络可充分学习正负样本间的局部特征, 挖掘两者间的细节差异, 此时正负样本的比例接近 $1:k-1$ 。在深度学习中, $1:3$ 是广泛应用于目标检测如 SSD 网络^[17]训练中的正负样本比, 此时模型可以更快地优化和更稳定地训练^[14]。基于此, 本研究采用 $k=4$ 进行样本动态筛选, 从而使专家网络获取的正负样本比接近 $1:3$ 。图 3 展示了这一样本的筛选过程, 可发现对于真实标签为 i 的任意样本, 其第 i 类的预判概率 p_i 属于预判概率 top-4 的元素时, 它可作为第 i 个专家网络的正样本(标签 P)以及其他 3 个专家网络的负样本(标签 N), 从而将总样本集的正负样本比例维持在 $1:3$ 左右。

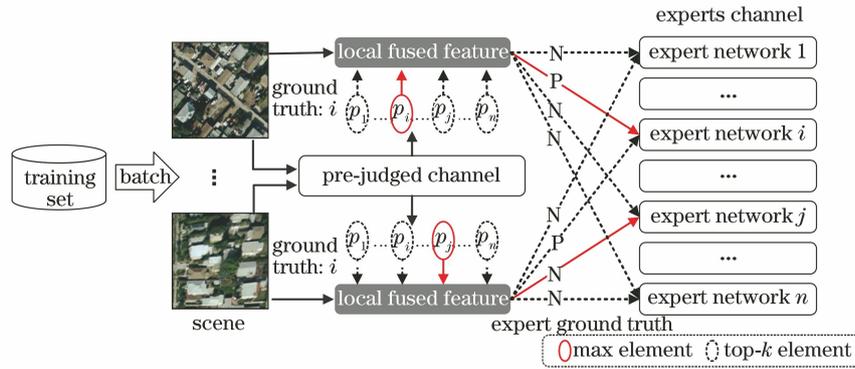


图 3 专家网络的训练样本筛选过程

Fig. 3 Training sample filter for expert networks

专家网络经过针对性训练后, 对本类场景及其在预判通道中类别易被误判的场景有较高的区分能力, 综合两类通道的结果可有效地改善单一分类器对相似场景区分能力不足的问题。如当预判通道可正确预测真实标签为 i 的场景 s 时, 经过训练的第 i 类专家网络将进一步确认场景 s 属于本类并产生较高的 t_i 值, 该值与预判概率 p_i 结合后产生的第 i 类概率值仍保持最高, 最终 TLMoE 维持正确判断; 而当预判通道将场景 s 误判为第 j ($1 \leq j \leq n$) 类时, 第 j 个专家网络在训练阶段已充分学习第 i 类与第 j 类别样本的特征知识, 据此可判断出场景 s 不属于本类并产生更低的 t_j 值, 该值与预判概率 p_j 相乘后产生的第 j 类概率值大幅降低, 以此避免 TLMoE 将场景 s 误分为第 j 类。

2.4 模型训练

在通过预判通道获取样本数据后, 专家通道将利用联合损失函数进行训练, 联合损失函数可表示为

$$L = \frac{1}{n} \sum_{i=1}^n L_{\text{expert}}^{(i)} + L_{\text{TLMoE}}, \quad (10)$$

式中, 第一项为专家网络分类损失, 第二项为融合了预判通道和专家通道结果的分损失, 两者分别顾及各专家通道的分类精度及模型整体的分类表现。 n 个专家网络均为二分类网络, 其中任意专家网络的分类损失表达式为

$$L_{\text{expert}}^{(i)} = - \sum_{x \in D_i} t' \log[\mathbf{O}_{\text{FC2}}(s)] + \frac{\lambda}{2} \|\mathbf{w}_i\|^2, \quad (11)$$

式中:第一项为交叉熵损失,其中 D_i 为第 i 个专家网络的训练集,由预判通道动态筛选的场景局部特征组成, $t' \in \{0, 1\}$ 为当前样本在第 i 个专家网络中的真实标签;第二项为权值的 L2 正则项,用以避免网络过拟合,其中 W_i 为第 i 个专家网络 FC 层的权重矩阵, λ 为正则项系数,用于平衡两项损失,它的值由权重的衰减系数的乘积决定。

由于最终分类结果融合了专家通道结果与预判结果,完成单个专家网络的优化后, TLMoE 模型整体分类结果还需进一步优化,即

$$L_{\text{TLMoE}} = - \sum_{i=1}^n Y'_i \log Y_i, \quad (12)$$

式中: Y_i 为 TLMoE 中输出分类概率 Y 中的第 i 个元素, Y'_i 为该场景的真实类别。(12)式通过(10)式与(11)式共享正则项损失。

3 实验与分析

首先在三个小样本数据集上对 TLMoE 模型的表现进行分析,并对比了其他多种模型方法,结果证明所提 TLMoE 模型在深度学习方法中仍不失优越

性;然后对 TLMoE 模型框架运行机制进行解析,通过消融实验证明 TLMoE 模型采用的混合专家模型框架可高效整合多通道信息,提升网络整体表现;最后对比多种特征表现,证明专家通道对场景局部信息的针对性学习可有效提升局部特征的表达能力,且可增强 TLMoE 模型的整体表现。

3.1 实验设置及数据集

本研究实验均在载有 NVIDIA GeForce GTX1060 的显卡、Inter \times core i5-3479 CPU \times 3.20 GHz, RAM 为 16.0 GB 的工作站上进行。所有实验均在 TensorFlow^[18] 框架中完成,使用的预训练模型来自 TF-slim^[19], TF-slim 是一个轻量级的 TensorFlow 库,它提供了丰富的、可方便调用的预训练 CNN 模型,因此本研究通过 TF-slim 中封装的 VGG19 和 Resnet50 模型提取场景的全局特征和多级局部特征。

实验数据采用了 UCM、SIRI、RSSCN7 三个小样本遥感影像场景集,各数据集的场景类别如图 4 所示。UCM 数据集包含 21 类遥感影像场景,每类有 100 幅尺寸为 256 pixel \times 256 pixel 的图像,空间



图 4 遥感影像场景示例。(a) UCM 数据集;(b) SIRI 数据集;(c) RSSCN7 数据集

Fig. 4 Image examples of remote sensing scenes. (a) UCM dataset; (b) SIRI dataset; (c) RSSCN7 dataset

分辨率为 0.3 m;SIRI 数据集覆盖了 12 类,每类包含 200 幅尺寸为 200 pixel×200 pixel 的遥感影像场景,空间分辨率为 2 m;RSSCN7 数据集由 7 类场景组成,每类场景包含 400 幅尺寸为 700 pixel×700 pixel 的场景图像,覆盖了 4 种不同空间分辨率下的场景影像。

为公平地与其他方法的效果作横向对比,本文实验数据集的划分策略与已有研究保持一致,在 UCM 和 SIRI 数据集中随机选取 80% 的样本作为训练集,在 RSSCN7 中随机选取 50% 的样本作为训练集,剩余样本为测试集。本文方法不对样本进行旋转、平移、随机采样等操作,避免了复杂的数据预处理,只在原有小样本量下完成训练与分类。实验部分采用分类精度和混淆矩阵作为评价指标,分类精度和混淆矩阵各元素的值域范围均为 $[0, 1]$,其中分类精度值及混淆矩阵主对角元素值越高,方法分类表现越好,对不同类别遥感影像场景的识别能力越强。

3.2 TLMoE 分类表现对比与分析

1) UCM 数据集。表 2 列出 TLMoE 模型及多种前沿方法在 UCM 数据集上的分类结果。随机森林(RF)^[20]直接用于原图分类时的分类精度仅达 44.77%,结合 SIFT 特征的 BoVW^[21]和 SPMK^[22]方法,可将分类精度提升至 75% 以上,但这与深度学习方法的表现仍有一定差距。在 UCM 上经过重新训练的 VGG19 和 Resnet50 的分类精度可达 80% 以上,其相较传统方法有很大提升,但这在深度学习方法中并不突出,这是由于这两个网络层次较深、参数量较多,在小样本数据集上难以充分训练。因此降低网络结构的复杂度(第 6 行),或借助预训练网络特征(第 7~11 行)都更利于改善网络在小样本数据集上的表现。特别地,第 7~8 行的结果是直接利用预训练 CNN 提取场景特征并通过 SVM 分类所得,它们均采用与 TLMoE 内两类通道输入相同的迁移特征,通过简单连接操作得到融合的特征向量,以便于 SVM 直接分类。对比第 7 行与第 10 行的分类结果时,发现尽管两类方法的输入特征均来自预训练网络 VGG19 中相同的网络层,但两者的表现仍存在一定的差异,直接利用预训练 VGG19 特征迁移的分类方法可获得 94.29% 的分类精度,但 TLMoE-VGG19 在预训练 VGG19 特征的基础上建立了混合专家结构网络,并充分利用多层卷积层与全连接层特征进行了不同层次的分类,最终将分类结果提升至 98.10%。类似的情况也出

现在 Resnet50 迁移特征分类结果和 TLMoE-Resnet50 分类结果上(第 8、11 行),由此可表明 TLMoE 有效地结合了迁移学习和混合专家系统结构的优点,可获得更好的分类表现。

表 2 UCM 数据集上的分类精度对比

Table 2 Classification accuracy comparison on the

UCM dataset		
No.	Method	Accuracy /%
1	RF ^[20]	44.77
2	SIFT+BoVW ^[21]	76.81
3	SIFT+SPMK ^[22]	75.29
4	VGG19 (training from scratch)	83.48
5	Resnet50 (training from scratch)	85.71
6	DCT-CNN ^[1]	95.76
7	Pre-trained VGG19 features+SVM	94.29
8	Pre-trained Resnet50 features+SVM	97.14
9	GLDFB ^[5]	97.62
10	TLMoE-VGG19	98.10
11	TLMoE-Resnet50	98.33

图 5 的混淆矩阵展示了 TLMoE-VGG19 对各类场景的识别情况。21 类场景中 20 类场景的识别准确率达到 95% 及以上,其中 14 类达到 100%,唯一未达 95% 以上的 buildings 类的识别准确率亦达到 90%,它主要被误判为 tennis court 类。这是由于 tennis court 类场景中部分网球场分布在建筑物周围且尺度较小,整体外观与建筑物相关类别场景相似,因而产生少数混淆情况,但它们中大部分仍能被准确区分,tennis court 类及 dense/medium residential 等建筑物相关类的识别准确度均达到 95% 及以上。UCM 数据集中主要存在两个易混淆大类,即包含了 dense/medium residential、buildings 的建筑物相关类和包含了 overpass、intersection、freeway 的道路相关类。表 3 对比了 TLMoE-VGG19 与同样使用预训练 VGG19 特征优化分类的 GLDFB 方法对这类易混淆场景的分类精度。由表 2 可知,GLDFB 模型的总体分类精度达 97.62%,且对这几类场景已具备较好的识别能力,但对 overpass 的分类精度仍低于 TLMoE-VGG19, TLMoE-VGG19 对道路相关类中三类场景的识别准确率达到 100%;同时 TLMoE-VGG19 对建筑物相关类中 medium residential 的识别准确率相较 GLDFB 模型也提升了 5 个百分点;此外, TLMoE-VGG19 对 tennis court 类的识别准确率也有所提升。整体而言, TLMoE 模型对不同复杂程度与相似程度的场景都有较好的识别区分能力。

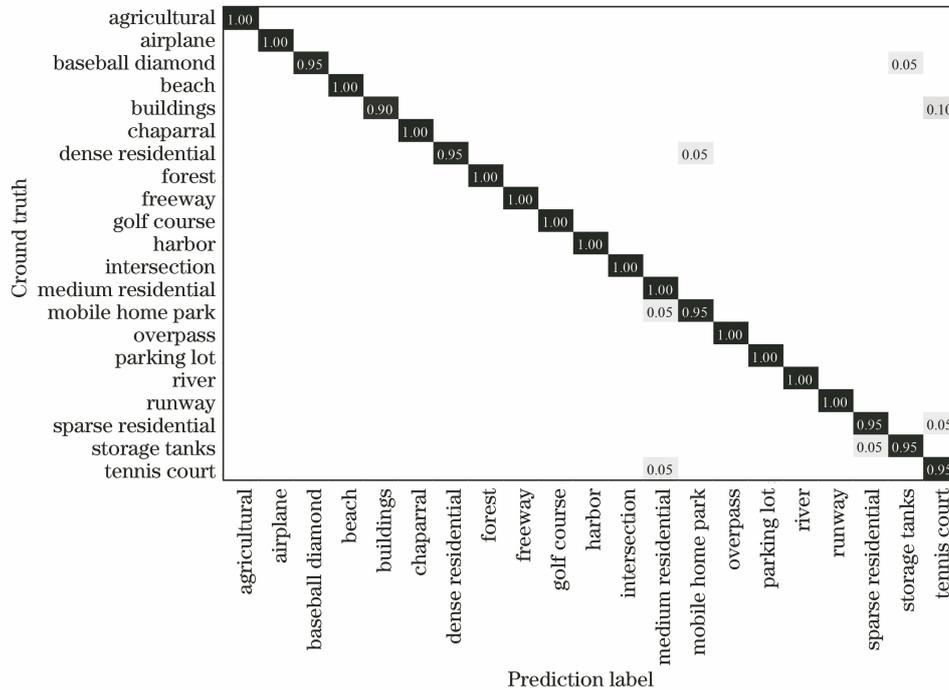


图 5 TLMoE-VGG19 模型在 UCM 数据集上的分类混淆矩阵

Fig. 5 Classification confusion matrix of TLMoE-VGG19 on UCM dataset

表 3 UCM 数据集易混淆类别的分类精度对比

Table 3 Classification accuracy comparison on the confusing classes of UCM dataset

unit: %

Type	Road type			Building type			Other
	Freeway	Intersection	Overpass	Buildings	Lense residential	Medium residential	Tennis court
GLDFB(VGG19)	100	100	95	90	95	95	90
TLMoE-VGG19	100	100	100	90	95	100	95

2) SIRI 数据集。表 4 对比了多种方法在 SIRI 数据集上的分类表现,其中:RF 直接用于原图分类时的分类精度仅达 49.90%;基于中低层特征的场景分类方法(第 2~3 行)可将分类精度提升至 60%~80%;而深度学习类方法(第 4~9 行)利用抽象程度更高的特征获取更优的表现,在数据量不足的情况下,对 VGG19 和 Resnet50 重新训练的分类精度仍达 85%以上;MCNN^[23]对图像多尺度采样后送入 CNN,有效改善了分类效果,但需要额外的数据增强处理;而迁移学习方法(第 7~11 行)汲取了大样本数据集训练的 CNN 特征对局部和全局信息高效表达的优点,可降低对迁移任务的数据需求。在同等数据量下,利用预训练的 VGG19 和 Resnet50 提取场景特征并通过 SVM 分类分别获得了 94.79%和 96.25%的分类结果,TLMoE 模型则更进一步通过不同网络通道充分利用预训练 CNN 中不同网络层次的特征,故 TLMoE-VGG19 和 TLMoE-Resnet50 将分类精度提高到 97.29%和 97.50%。

表 4 SIRI 数据集上的分类精度比较

Table 4 Classification accuracy comparison on the SIRI dataset

No.	Method	Accuracy /%
1	RF ^[20]	49.90
2	SIFT+BoVW ^[21]	75.63
3	SIFT+SPMK ^[22]	77.69±1.01
4	VGG19 (training from scratch)	86.13
5	Resnet50 (training from scratch)	89.26
6	MCNN ^[23]	93.75±1.13
7	Pre-trained VGG19 features+SVM	94.79
8	Pre-trained Resnet50 features+SVM	96.25
9	GLDFB ^[5]	96.67
10	TLMoE-VGG19	97.29
11	TLMoE-Resnet50	97.50

图 6 的混淆矩阵展示了 TLMoE-VGG19 模型对 SIRI 数据集各类的识别情况。12 类场景中 10 类的识别准确率达到 95%及以上,其中 commercial、industrial、residential 等组成复杂、差

别较小的建筑相关类场景,其识别准确率达到 97.5%或 100%。易混淆的水系相关类中,river 和 water 的识别准确率达到 100%,但 pond 类中产生少量误分,这是由于该类场景中的部分池塘形状走向与河流类似,且周围常伴有绿地,小部分 pond 类场景被误判为 river、meadow 类,但该类大部分场景

可正确识别,整体识别准确率达到 90%。此外 TLMoE-VGG19 对 overpass、industrial 等人工建造场景类别的识别准确率相较 GLDFB 模型提升 2.5~5 个百分点。整体而言,TLMoE-VGG19 对各类场景的识别准确率均高于 90%,整体分类精度达到 97.29%。

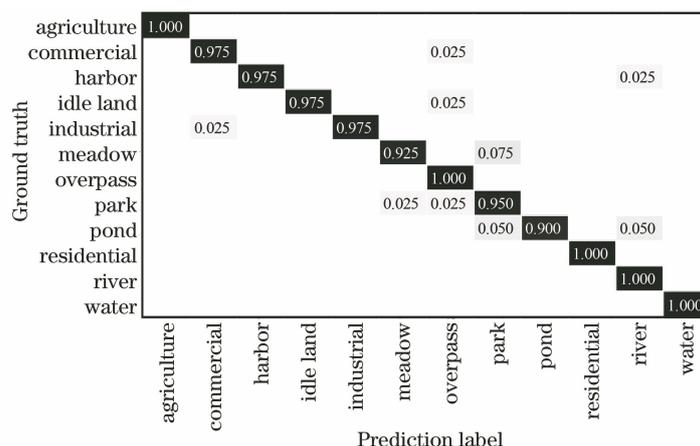


图 6 TLMoE-VGG19 在 SIRI 数据集上的分类混淆矩阵

Fig. 6 Classification confusion matrix of TLMoE-VGG19 on SIRI dataset

3) RSSCN7 数据集。表 5 列出了 RSSCN7 数据集上的分类结果。与在 UCM 和 SIRI 数据集上的结果类似,RF 直接用于原图分类时的分类精度仍较低,仅达 55.43%,同样受限于训练样本的不足,重新训练的 VGG19 和 Resnet50 模型的分精度略高于 80%,而采用预训练 VGG19 及 Resnet50 提取场景特征的分精度提升至 91.93%和 89.92%。更进一步地,采用预训练 VGG19 和 Resnet50 的 TLMoE 模型的分精度达到 93.21%和 93.29%,在深度学习类方法(第 2~8 行)中仍具有一定优势。尽管训练集样本占比相较前两个数据集降为 50%,但在图 7 的混淆矩阵中,7 类场景中仍有 6 类场景的识别准确率达到 90%以上,仅有 Industry 的混淆比率相对较高。

表 5 RSSCN7 数据集上的分类精度比较

Table 5 Classification accuracy comparison on the

RSSCN7 dataset

No.	Method	Accuracy / %
1	RF ^[20]	55.43
2	VGG19 (training from scratch)	82.50
3	Resnet50 (training from scratch)	81.70
4	Deep filter bank ^[24]	90.04±0.6
5	Pre-trained VGG19 features+SVM	91.93
6	Pre-trained Resnet50 features+SVM	89.92
7	TLMoE-VGG19	93.21
8	TLMoE-Resnet50	93.29

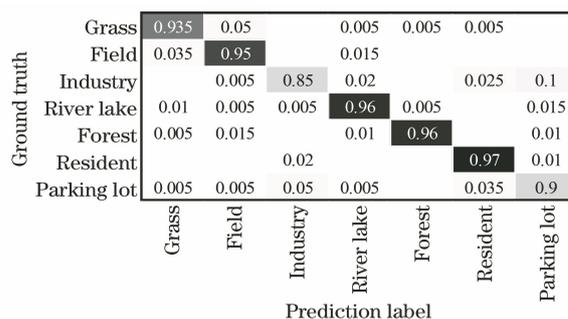


图 7 TLMoE-VGG19 在 RSSCN7 数据集上的分类混淆矩阵

Fig. 7 Classification confusion matrix of TLMoE-VGG19

on RSSCN7 dataset

3.3 TLMoE 模型结构性能分析

为更深入地了解 TLMoE 模型结构及各通道的性能,表 6 对比了 TLMoE 中预判通道、专家通道与整体模型的分结果(见第 1, 2, 4 行),可发现在 3 个数据集上,TLMoE 模型相较单独的预判通道和专家通道,均能取得最优结果;同时表 6 列出了两类通道的输入特征 $F_L(s)$ 、 $X(s)$ 直接连接融合后通过 SVM 分类的结果(第 3 行),可发现直接融合方式下获得的分类精度全部低于 TLMoE 的结果,部分结果甚至低于单通道的结果。由此可见,简单的特征连接融合并不能充分利用两类特征中的信息。相较之下,TLMoE 模型采用的混合专家系统结构能够更有效地结合不同通道挖掘的场景信息,利用不同通道的识别优势提升了分类精度。

表 6 TLMoE 通道分类精度对比

Table 6 Classification accuracy comparison between TLMoE channels

No.	Channel	Accuracy (pre-trained VGG19) /%			Accuracy (pre-trained Resnet50) /%		
		UCM	SIRI	RSSCN7	UCM	SIRI	RSSCN7
1	Pre-judged channel	94.60	93.13	87.50	97.62	96.25	89.21
2	Expert channel	93.81	96.04	92.14	96.67	97.08	92.93
3	$(F_L(s), X(s))$ -SVM	94.29	94.79	91.93	97.14	96.25	89.92
4	TLMoE	98.10	97.29	93.21	98.33	97.50	93.29

TLMoE 模型由预训练 CNN 特征提取器、预判通道、专家通道三个部分组成。图 8 对比了这三个部分及 TLMoE 模型整体的训练耗时,所有实验 batch size 均采用 512,最大 epoch 采用 150。从图 8 中可发现,不论是 VGG19 还是 Resnet50 作为预训练 CNN 特征提取器,预训练 CNN 对场景特征的提取速度均较快。而将预训练 CNN 特征送入对应通道后,预判通道和专家通道的训练耗时逐渐递增。三个部分组合后的训练耗时均略小于实际 TLMoE 模型整体的训练耗时,这是由于 TLMoE 模型训练

时,预判通道产生预测值并为专家通道进行动态样本筛选的过程仍需时间。但总体而言,它们之间的时间差较小,采用 VGG19 的时间差不超过 30 s,采用 Resnet50 的时间差不超过 11 s,其中 VGG19 由于产生的全局特征维度更高,计算耗时更长。同步对比表 5 中各通道及 TLMoE 模型的整体训练时长,可发现尽管 TLMoE 整体的训练时长有小幅增加,但对数据集的分类精度提升效果较为明显,且 TLMoE 方法的整体训练耗时仍处在较低范围内,整体而言 TLMoE 在时间性能与分类精度方面均有较好表现。

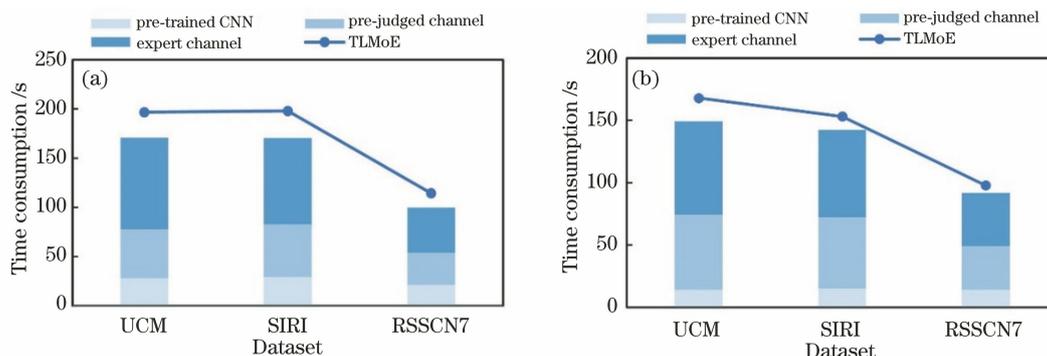


图 8 TLMoE 模型中各通道及预训练 CNN 在组合前后的训练耗时对比。(a) VGG19; (b) Resnet50

Fig. 8 Time consumption comparison before and after the combination of channels and pre-trained CNN in TLMoE.

(a) VGG19; (b) Resnet50

3.4 专家通道特征优化能力分析

分类模型对场景特征的表达能极大地影响着最终的分类表现。TLMoE 通过专家通道构建多个专家网络,针对性地学习各类场景与相似场景的局部信息,进一步提升了场景特征的表达能力和区分能力。表 7 列出专家通道的输入特征即局部融合特征 $X(s)$ 及其他多种视觉图像特征在 SVM 下的分类精度,并对比了专家通道优化后的分类结果。由表 7 可发现:传统中低层图像特征如方向梯度直方图特征(HOG)、尺度不变特征变换特征(SIFT)和局部二值模式特征(LBP)等对遥感影像场景的图像的表达能有限(第 1~3 行),在三个小样本数据集上的分类精度均不超过 60%;相比之下,采用迁移学习方法提取的局部融

合特征 $X(s)$ 的表达能已有较大提升,其在三个数据集上的分类结果均达到 90% 以上;而专家通道则进一步优化局部融合特征,将三个数据集的分类精度提高了 0.48%~1.43%。

图 9 通过 t-SNE^[25] 方法对三个数据集上各类别场景的多种特征进行降维展示,更加直观地反映了不同特征间表达能与区分能的差异,图中局部融合特征和局部专家特征均以 TLMoE-VGG19 为例进行采集。图中每种颜色代表一种场景类别,从中可发现底层特征点的分布较为混乱,尤其是 HOG 特征难以形成明显的聚类中心,SIFT 特征和 LBP 特征中同色特征点可形成少量团簇,但整体的分布仍较为分散。通过预训练 VGG19 提取的局部融合特征的区分能显著增强,同类特征点间的聚

集性更强。而专家网络进一步将局部融合特征优化为局部专家特征。由于各类专家网络是通过不同场景样本来完成针对性训练,是在不同的数据空间完成相应样本知识的专家化学习,从而可将输入数据映射到不同的特征空间,产生区分能力更强的特征

表达。因此局部专家特征中各类特征间分界更明显,同类特征点的聚集性更突出,这也证明 TLMoE 中专家通道对包括相似场景在内的多类场景均具有较强的表达能力,对场景间不同程度的差异有较好的区分能力。

表 7 多种特征的分类精度对比

Table 7 Classification accuracy comparison of several kinds of features

No.	Feature	Accuracy / %		
		UCM	SIRI	RSSCN7
1	HOG	52.14	44.79	35.79
2	SIFT	58.33	53.96	54.14
3	LBP	31.43	46.25	56.14
4	$X(s)$ -VGG19	93.33	94.38	90.71
5	Expert channel-VGG19	93.81	96.04	92.14
6	$X(s)$ -Resnet50	95.48	96.46	92.14
7	Expert channel-Resnet50	96.67	97.08	92.93

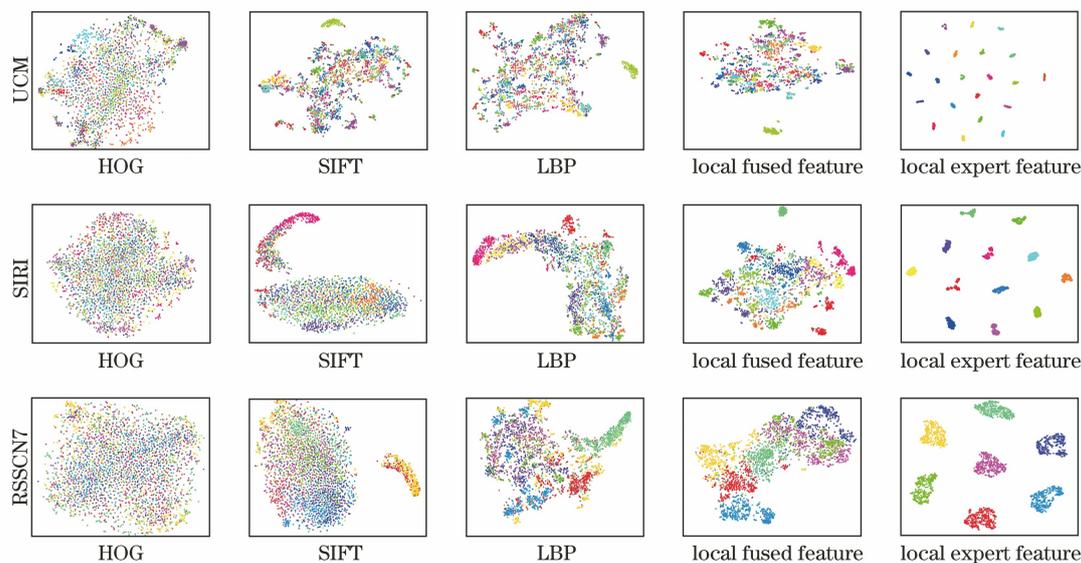


图 9 在三个小样本数据集上的特征二维映射后分布对比

Fig. 9 Comparison of different features on the 3 datasets by 2-dimensional feature visualization

4 结 论

为提升对小样本遥感影像场景数据集的分类精度,提出了一种基于迁移学习的混合专家分类模型 TLMoE。该模型利用预训练卷积神经网络特征对遥感影像场景的信息进行表达,通过建立预判通道和专家通道,分别挖掘场景全连接层特征与卷积层特征中蕴含的场景全局信息与场景局部信息,从而依据不同层次的信息完成对全部类遥感影像场景类别的初判与单一类遥感影像场景类别的细判。通过综合两类通道的判断,TLMoE 能够实现对不同类遥感影像场景间的显著差异和细微差异的区分,可以完成精度更高的遥感影像场景分类。

不同于使用单分类器的一般迁移学习方法,TLMoE 模型在网络结构上进行了新的探索,它将混合专家系统结构与迁移特征有效结合,建立了不同网络通道,并使不同网络通道负责不同相似程度下场景的区分。其中专家通道训练了多个专家网络,通过针对性地学习各类场景的关键局部细节信息,提升了对相似场景的区分能力。由此产生的多分类器架构弥补了单一分类器区分能力的不足,并扩展了网络宽度。同时各专家网络的简单结构也缩减了网络深度,极大降低了训练复杂度,因此无须对训练样本集进行额外的数据增强即可完成充分训练,在三个小样本数据集上的实验亦表明,该方法可有效提高对遥感影像场景的分类精度。未来的研究

工作将重点从局部特征优化、专家网络训练样本的筛选以及网络的损失函数调优等方面展开,以期获取区分能力更强的分类模型。

参 考 文 献

- [1] Liu F, Lu L X, Huang G W, et al. Landform image classification based on discrete cosine transformation and deep network[J]. *Acta Optica Sinica*, 2018, 38(6): 0620001.
刘芳, 路丽霞, 黄光伟, 等. 基于离散余弦变换和深度网络的地貌图像分类[J]. *光学学报*, 2018, 38(6): 0620001.
- [2] Wu C, Wang H W, Yuan Y W, et al. Image feature fusion based remote sensing scene zero-shot classification algorithm [J]. *Acta Optica Sinica*, 2019, 39(6): 0610002.
吴晨, 王宏伟, 袁昱纬, 等. 基于图像特征融合的遥感场景零样本分类算法[J]. *光学学报*, 2019, 39(6): 0610002.
- [3] Hauptmann A, Yan R, Lin W H, et al. Can high-level concepts fill the semantic gap in video retrieval? A case study with broadcast news [J]. *IEEE Transactions on Multimedia*, 2007, 9(5): 958-966.
- [4] He Q, Li Y, Song W, et al. Multimodal remote sensing image classification with small sample size based on high-level feature fusion [J]. *Laser & Optoelectronics Progress*, 2019, 56(11): 111001.
贺琪, 李瑶, 宋巍, 等. 小样本的多模态遥感影像高层特征融合分类[J]. *激光与光电子学进展*, 2019, 56(11): 111001.
- [5] Gong X, Wu L, Xie Z, et al. Classification method of high-resolution remote sensing scenes based on fusion of global and local deep features [J]. *Acta Optica Sinica*, 2019, 39(3): 0301002.
龚希, 吴亮, 谢忠, 等. 融合全局和局部深度特征的高分辨率遥感影像场景分类方法[J]. *光学学报*, 2019, 39(3): 0301002.
- [6] Castelluccio M, Poggi G, Sansone C, et al. Land use classification in remote sensing images by convolutional neural networks [EB/OL]. (2015-08-01) [2021-01-10]. <https://arxiv.org/abs/1508.00092>.
- [7] Hu F, Xia G S, Hu J W, et al. Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery[J]. *Remote Sensing*, 2015, 7(11): 14680-14707.
- [8] Yandex A B, Lempitsky V. Aggregating local deep features for image retrieval [C] // 2015 IEEE International Conference on Computer Vision (ICCV), December 7-13, 2015, Santiago, Chile. New York: IEEE Press, 2015: 1269-1277.
- [9] Zhang X N, Zhong X, Zhu R F, et al. Scene classification of remote sensing images based on integrated convolutional neural networks [J]. *Acta Optica Sinica*, 2018, 38(11): 1128001.
张晓男, 钟兴, 朱瑞飞, 等. 基于集成卷积神经网络的遥感影像场景分类[J]. *光学学报*, 2018, 38(11): 1128001.
- [10] Zheng Z, Fang F, Liu Y Y, et al. Joint multi-scale convolution neural network for scene classification of high resolution remote sensing imagery [J]. *Acta Geodaetica et Cartographica Sinica*, 2018, 47(5): 620-630.
郑卓, 方芳, 刘袁缘, 等. 高分辨率遥感影像场景的多尺度神经网络分类法[J]. *测绘学报*, 2018, 47(5): 620-630.
- [11] Cheng G, Yang C Y, Yao X W, et al. When deep learning meets metric learning: remote sensing image scene classification via learning discriminative CNNs [J]. *IEEE Transactions on Geoscience and Remote Sensing*, 2018, 56(5): 2811-2821.
- [12] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition[EB/OL]. (2015-04-10) [2021-01-10]. <https://arxiv.org/abs/1409.1556>.
- [13] He K M, Zhang X Y, Ren S Q, et al. Deep residual learning for image recognition [C] // 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 27-30, 2016, Las Vegas, NV, USA. New York: IEEE Press, 2016: 770-778.
- [14] Lin M, Chen Q, Yan S. Network in network[EB/OL]. (2013-12-16) [2021-01-10]. <https://arxiv.org/abs/1312.4400>.
- [15] Gong X, Xie Z, Liu Y, et al. Deep salient feature based anti-noise transfer network for scene classification of remote sensing imagery[J]. *Remote Sensing*, 2018, 10(3): 410.
- [16] Glorot X, Bordes A, Bengio Y. Deep sparse rectifier neural networks [C] // Proceedings of the 14th International Conference on Artificial Intelligence and Statistics, April 11-13, 2011, Lauderdale, Florida, USA. Cambridge: JMLR, 2011: 315-323.
- [17] Liu W, Anguelov D, Erhan D, et al. SSD: single shot MultiBox detector[M]//Leibe B, Matas J, Sebe N, et al. Computer vision-ECCV 2016. Lecture notes in computer science. Cham: Springer, 2016, 9905: 21-37.
- [18] Abadi M, Barham P, Chen J M, et al. TensorFlow: a system for large-scale machine learning [C] // Proceedings of the 12th Usenix Conference on

- Operating Systems Design and Implementation, November 2-4, 2016, Savannah, Georgia, USA. Berkeley: USENIX ASSociation, 2016: 265-283.
- [19] Silberman N, Guadarrama S. TensorFlow-Slim image classification model library [EB/OL]. (2020-05-27) [2021-01-10]. <https://github.com/tensorflow/models/blob/master/research/slim/README.md#pre-trained-models>.
- [20] Breiman L. Random forests[J]. Machine Learning, 2001, 45(1): 5-32.
- [21] Csurka G. Visual categorization with bags of keypoints[C]//Agapito, Bronstein M M, Rother C. Workshop on Statistical Learning in Computer Vision, September 6-7, 2014. Zurich, Switzerland: Springer, 2014: 1-22.
- [22] Lazebnik S, Schmid C, Ponce J. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories [C] //2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), June 17-22, 2006, New York, NY, USA. New York: IEEE Press, 2006: 2169-2178.
- [23] Liu Y F, Zhong Y F, Qin Q Q. Scene classification based on multiscale convolutional neural network[J]. IEEE Transactions on Geoscience and Remote Sensing, 2018, 56(12): 7109-7121.
- [24] Wu H, Liu B Z, Su W H, et al. Deep filter banks for land-use scene classification[J]. IEEE Geoscience and Remote Sensing Letters, 2016, 13(12): 1895-1899.
- [25] Laurens V D M, Hinton G. Visualizing Data using t-SNE [J]. Journal of Machine Learning Research, 2008, 9: 2579-2605.