

## 基于正交约束和最大类内特征判别性的 分层分类特征选择算法

殷金钊<sup>1</sup>, 郑文利<sup>2</sup>, 钱婷<sup>2,3\*</sup>, 折延宏<sup>2,3</sup>

(1. 西安石油大学 计算机学院, 陕西 西安 710065;

2. 西安石油大学 理学院, 陕西 西安 710065;

3. 西北大学 概念、认知与智能研究中心, 陕西 西安 710127)

**摘要:**在大数据时代,数据正呈现出指数级增长趋势。数据间的类别层次结构使得分类学习任务更有效率。现有的分层分类特征选择算法未充分体现类内特征的判别性,因此本文提出了一种基于正交约束和最大化类内特征判别性的分层分类特征选择算法(Hierarchical Classification Feature Selection Algorithm Based on Orthogonal Constraints and Intra-class Maximum Feature Discriminability, HFSOC)。该算法在使用稀疏正则化项去除不相关特征后,利用改进后的正交约束公式来度量类间独立性,并将每个内部节点特征矩阵的各个列向量互相正交,以提高类内特征的判别性。最后,利用递归正则项优化输出特征权重矩阵。实验结果表明,本文所提算法在5个数据集上取得了一定的效果,其分类准确率在DD数据集上相比于HFisher算法提高约17%,在F194数据集和CLEF数据集上相比于基于 $\ell_{2,1}$ 范数最小化的高效鲁棒的特征选择算法(HFSNM)均提高约10%,在ILSVRC数据集上相比于HFSNM算法提高约1%。

**关键词:**特征选择;稀疏学习;层次树结构;正交约束;递归正则化

中图分类号:TP311 文献标志码:A 文章编号:0253-2395(2025)01-0144-14

## Hierarchical Classification Feature Selection Algorithm Based on Orthogonal Constraints and Intra-class Maximum Feature Discriminability

YIN Jinzhao<sup>1</sup>, ZHENG Wenli<sup>2</sup>, QIAN Ting<sup>2,3\*</sup>, SHE Yanhong<sup>2,3</sup>

(1. College of Computer, Xi'an Shiyou University, Xi'an 710065, China;

2. College of Science, Xi'an Shiyou University, Xi'an 710065, China;

3. Institute of Concepts, Cognition and Intelligence, Northwest University, Xi'an 710127, China)

**Abstract:** In the era of big data, data is showing an exponential growth trend. The hierarchical structure of categories among data makes the classification learning task more efficient. However, existing hierarchical classification feature selection algorithms do not fully reflect the discriminative nature of intra-class features. This paper proposes a hierarchical classification feature selection algorithm based on orthogonal constraints and maximizing the discriminative nature of intra-class features(HFSOC). The algorithm utilizes an improved orthogonality constraint formula to measure inter-class independence and orthogonalizes the individual column

收稿日期:2024-06-14;接受日期:2024-10-24

基金项目:国家自然科学基金(12171388; 61976244; 12101478; 12171294);陕西省自然科学基金基础研究计划项目(2023JCYB027);浙江海洋大学海洋大数据挖掘与应用重点实验室(OBDMA202101);陕西数理基础科学研究项目(23JSZ008);研究生创新项目(YCX2413143)

作者简介:殷金钊(2000—),女,山东聊城人,硕士研究生,研究方向为分层分类。E-mail:yin2991802968@163.com

\* 通信作者:钱婷(QIAN Ting),E-mail:qiant2000@126.com

引文格式:殷金钊,郑文利,钱婷,等.基于正交约束和最大类内特征判别性的分层分类特征选择算法[J].山西大学学报(自然科学版),2025,48(1):144-157. DOI:10.13451/j.sxu.ns.2024137.

vectors of each internal node's feature matrix to each other to improve the discriminative property of intra-class features after using sparse regularization terms to remove irrelevant features. Finally, the output feature weight matrix is optimized using recursive regularization terms. The experimental results show that the algorithm proposed in this paper achieves certain results on five datasets, and its classification accuracy is improved by about 17% compared to HFisher's algorithm on the DD dataset, by about 10% compared to efficient and robust feature selection via joint  $\ell_{2,1}$ -norms minimization algorithm (HFSNM) on both the F194 dataset and the CLEF dataset, and by about 1% compared to HFSNM's algorithm on the ILSVRC dataset.

**Key words:** feature selection; sparse learning; hierarchical tree structure; orthogonal constraint; recursive regularization

## 0 引言

在大数据时代,数据正呈现出指数级增长趋势,数据样本急剧增加,数据种类也由最初的数千发展至数十万数百万乃至更多<sup>[1]</sup>。例如,在 ImageNet 中图像数据的类别就有上万个,图像数据样本更是达到了上千万之多<sup>[2]</sup>。在超大规模的数据下,人们面临的分类学习任务也遇到了新的挑战<sup>[3]</sup>。面对这种情况,学者通常采用分而治之策略在类别间建立层次结构进行分类学习。这种具有层次结构的分类任务被称为分层分类<sup>[4]</sup>。此后,许多学者针对分层分类问题展开了研究。Wang 等<sup>[5]</sup>认为在层次结构分类时每个非叶节点存在保守风险和冒进风险,提出了局部 Bayes 风险最小化的分类模型,即通过在每个非叶节点比较保守风险和冒进风险来决定停止或者继续向下进行分类。Zhou 等<sup>[6]</sup>提出了一种深度超类学习模型来解决长尾分布式数据分类问题,该模型将组稀疏与卷积神经网络相结合,类别之间的视觉相似度越高,它们出现在同一组中的可能性越大,进而提高分类精度。分层分类学习已成为一个热点研究问题<sup>[7-12]</sup>。

在分类任务中,特征选择是一项不可或缺的环节<sup>[13-14]</sup>。研究人员提出了各种分层分类特征选择算法。Freeman 等<sup>[15]</sup>提出了一种多类分类模型,该模型一边构建层次树结构,一边对局部分类器进行特征选择,二者同时进行。该模型可以为每个分类子任务选取互相独立的特征子集。Grimaud 等<sup>[16]</sup>提出了基于层次结构的分类器,在每个非叶节点设立分类器,每个分类器相互独立,利用 mRMR 特征选择算法为每个分类器选择有代表性的特征。Song 等<sup>[17]</sup>提出了一种基于信息增益和特征位置信息及其频率分布的分层分类特征选择算法。该算法使用类别的语义结构进行特征降维,并且为每个类别选取适当且互相独立的特征子集。上述算法解决了传统特征选择中只能选择统一特征子集的问题,但均未考虑层次结构间的各种关系。

父子关系及兄弟关系是层次结构中常被考虑的类别关系<sup>[3]</sup>。Tuo 等<sup>[18]</sup>为了探索不同类别之间的双向依赖关系,提出了基于子树的图正则化分层分类特征选择算法。Zhao 等<sup>[19]</sup>提出了类别的层次结构与递归正则化结合的特征选择算法。该算法利用父子节点之间的相关性以及兄弟节点之间的独立性分别构造正则项,进而递归地为每个非叶节点选择不同的特征子集。Lin 等<sup>[20]</sup>认为层次结构中的内部节点之间存在共有特征和固有特征,进而提出利用损失函数以及类别相似度矩阵构建特征选择目标函数来学习层次结构中不同类别的固有特征和共有特征。但上述这些算法均未考虑类内特征之间的关系。Shi 等<sup>[21]</sup>提出了基于类间最大独立性和类内最小冗余性的分层分类特征选择算法,其利用兄弟关系正交约束提高类间独立性,用内积项度量类内相关性进而减小类内特征冗余性。

运用正交约束理论刻画特征关系是一种有效的手段,但我们发现 Shi 等<sup>[21]</sup>提出的正交约束会导致兄弟节点之间选取的特征较为一致,影响兄弟节点之间的独立性。例如对于特征权重矩阵为  $2 \times 2$  的布尔矩阵,为使得  $\|W_j^T W_i - E\|_F^2$  取最小,而极端情况下的最小值就是  $\|W_j^T W_i - E\|_F^2 = 0$ ,我们通过穷举法发现,只有  $W_i = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} W_j = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$  或者  $W_i = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} W_j = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$  这两种组合情况能满足最小约束,而这两种组合情况均说明兄弟节点之间所选特征一致,故不能较好地反映兄弟节点

之间特征的独立性。

其次 Shi 等<sup>[21]</sup>通过利用行与行之间的正交关系考虑类内特征冗余性,降低了特征维数,但其忽略了分类的需求。我们仍以  $2 \times 2$  的布尔矩阵为例,假设非叶节点  $i$  的类内特征权重矩阵  $W_i$  为  $\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ ,从矩阵可以看出两个特征对两个类别都重要且重要程度相同,并且二者之间的相关性强,冗余性大。经过内积项对特征(即行与行)之间的相关性进行惩罚之后,非叶节点  $i$  的特征权重矩阵  $W_i$  被约束为  $\begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix}$ 。从该矩阵中我们可以看出特征的冗余性变小。但是所选特征仍然不能进行有效分类,这就是所选类内特征判别性小。基于选择特征的最终目的是分类,所以类内特征选择时首要考虑筛选后的特征具有最大判别性。

因此本文提出基于正交约束和最大化类内特征判别性的分层分类特征选择算法(Hierarchical Classification Feature Selection Algorithm Based on Orthogonal Constraints and Intra-class Maximum Feature Discriminability, HFSOC)。并做了相关对比实验,实验结果表明,本文所提出的 HFSOC 算法具有一定的有效性。

本文相关工作如下:第1节介绍相关理论知识;第2节介绍基于正交约束和最大化类内特征判别性的分层分类特征选择算法;第3节将所提算法与其他算法进行比较,并使用多种评价指标对结果进行分析;第4节进行总结与展望。

## 1 相关理论

### 1.1 类别的层次树结构

类别的层次结构主要分为两种类型,树型和有向无环图型<sup>[22]</sup>。为了便于理解,本论文重点讨论树型层次结构。层次树结构将类别标签组织成树状结构,以表示类别之间的一种“从属”关系<sup>[4]</sup>。用  $(Y, <)$  表示层次结构,  $Y$  表示标签集合,“ $<$ ”表示“从属”关系。类别间“从属”关系的三种性质描述如下:

- (1)不可逆性:若  $i < j$ ,则对于  $\forall i, j \in Y$ ,有  $j \not< i$ ;
- (2)反自反性:对  $\forall i \in Y$ ,有  $i \not< i$ ;
- (3)传递性:若  $i < k$  且  $k < j$ ,则对  $\forall i, j, k \in Y$ ,有  $i < j$ 。

### 1.2 基于稀疏学习的特征选择

在有监督的特征学习中,由于稀疏学习具有良好的可解释性,所以其常为特征选择首选的技术手段<sup>[23]</sup>。分层特征选择的目标是确保在一定的精度损失下挑选出少数有代表性的特征,因此将分层特征选择与稀疏学习相结合可以达到更好的效果<sup>[24]</sup>。设  $X \in R^{n \times m}$  是样本矩阵,其中  $n$  表示样本数量, $m$  表示特征数量。设  $T = \{0, 1, \dots, N\}$  表示类别层次树结构的非叶子节点集,当  $N=0$  时表示根节点。令  $X_0, X_1, \dots, X_N$  表示每个非叶节点的样本矩阵,其中  $X_i = [x_i^1, x_i^2, \dots, x_i^{n_i}] \in R^{n_i \times m}$  ( $0 \leq i \leq N, n_i \leq n$ );  $Y_0, Y_1, \dots, Y_N$  表示每个非叶节点对应的标签矩阵,其中  $Y_i = [y_i^1, y_i^2, \dots, y_i^{n_i}] \in R^{n_i \times d}$ , 并且  $y_i^j = \{0, 1\}^d$  ( $1 \leq j \leq n_i$ ),  $d$  是非叶节点中类别数目的最大值。

基于稀疏学习的特征选择目标函数<sup>[25]</sup>的一般形式为

$$\min_w L(XW, Y) + \lambda R(W), \quad (1)$$

其中  $L(\cdot)$  是一个损失项,  $R(W)$  是一个稀疏正则项,  $\lambda > 0$ 。在常见的损失函数<sup>[26]</sup>中,最小二乘损失具有计算简单、容易实现、最优解唯一等优点<sup>[19]</sup>,因此本文所提算法的损失函数使用最小二乘损失进行计算,如公式(2):

$$L(XW, Y) = \|XW - Y\|_F^2. \quad (2)$$

在稀疏学习方法中,常见的正则项包括  $l_0$  范数、 $l_1$  范数、 $l_2$  范数、 $l_{2,0}$  范数和  $l_{2,1}$  范数等<sup>[20]</sup>。其中

$\ell_0$  范数可以实现特征稀疏化,但由于其非凸、不光滑且不连续的特性,以至于不方便求解; $\ell_1$  范数不能进行类的结构稀疏化; $\ell_2$  范数不能实现特征稀疏化; $\ell_{2,0}$  范数预先定义了所选特征的数量,容易陷入局部最优;相比之下  $\ell_{2,1}$  范数可以实现类的结构稀疏化,并且是凸的,便于优化<sup>[27]</sup>。基于此,本文采用  $\ell_{2,1}$  范数作为稀疏学习的正则项。因此稀疏正则项可表示为

$$R(\mathbf{W}) = \|\mathbf{W}\|_{2,1} \quad (3)$$

综上,结合公式(2)和公式(3)可将基于稀疏学习的特征选择目标函数公式(1)重新表示为

$$\min_{\mathbf{W}} \sum_{i=0}^N \left( \|\mathbf{X}_i \mathbf{W}_i - \mathbf{Y}_i\|_F^2 + \|\mathbf{W}_i\|_{2,1} \right), \quad (4)$$

其中  $\mathbf{W}_i$  表示非叶节点  $i$  的权重矩阵,并且  $\mathbf{W}_i = [\mathbf{w}_i^1, \mathbf{w}_i^2, \dots, \mathbf{w}_i^m] \in R^{m \times d}$ 。

## 2 算法模型

### 2.1 基于正交约束的兄弟关系正则项

类别层次结构中兄弟关系是节点间的重要依赖关系<sup>[3]</sup>。一般来说,兄弟节点来自同一个父节点,相比于来自不同父节点的类别共享更多的特征,但是每个内部节点都有自己的子树,在每个内部节点挑选的特征都应具备判别子类别的特性。因此兄弟节点间应减少共享特征,增加判别性特征,进而加大类间的独立性。

正交约束是指在求解优化问题时,强制让变量在某个规范化的空间中正交。在机器学习中,正交约束优化方法将权重矩阵分解为两个正交矩阵的乘积,这样做有助于减少模型中的冗余信息,提高模型的泛化能力和稳定性。本文将修改 Shi 等<sup>[21]</sup>所提出的结构关系正则化项,使得经过正交约束后兄弟节点之间的特征更具有独立性,以及避免出现上述所分析的极端情况。设  $S_i$  为第  $i$  个非叶节点的兄弟节点,  $j \in S_i$  表示第  $i$  个非叶节点的第  $j$  个兄弟节点。因此兄弟关系正则化定义为

$$\sum_{j \in S_i} \|\mathbf{W}_j^T \mathbf{W}_i\|_F^2 \quad (5)$$

基于以上讨论,将兄弟关系正则项加入稀疏学习特征选择公式(4)中,以便提高特征的独立性。因此目标函数定义为

$$T(\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_N) = \min_{\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_N} \sum_{i=0}^N \left( \|\mathbf{X}_i \mathbf{W}_i - \mathbf{Y}_i\|_F^2 + \lambda \|\mathbf{W}_i\|_{2,1} \right) + \alpha \sum_{i=1}^N \sum_{j \in S_i} \|\mathbf{W}_j^T \mathbf{W}_i\|_F^2, \quad (6)$$

其中第一项是损失函数,用来度量预测标签和真实标签之间的误差程度;第二项是稀疏正则项,用于最小化特征权重值,增强模型的稀疏性,避免过拟合;第三项是兄弟关系正则项,用于加大类间的独立性,  $\alpha$  是惩罚兄弟节点间依赖度的非负参数。

### 2.2 基于正交约束的类内判别正则项

类别层次结构中每个非叶节点都有自己的子类别,因此在非叶节点  $i$  的特征子集应该具有最大程度判别子类别的作用,即具备最大类内特征判别性。已知特征权重矩阵  $\mathbf{W}_i$  的每个列向量  $\mathbf{w}_i^j$  对应着非叶节点  $i$  的第  $j$  个子类别的特征的权重,因此  $\mathbf{W}_i$  中列与列之间的相关性即表示非叶节点  $i$  的子类别之间的相关性。相关性越大,选取的特征之间越相似,对子类别的判别性越弱;相关性越小,选取的特征之间越独立,对子类别的判别性越强。受正交约束方法的启发,将特征权重矩阵  $\mathbf{W}_i$  的每个列向量互相正交,可以惩罚非叶节点  $i$  上子类别之间的相关性,进而提高子类别特征之间的独立性,使得在非叶节点  $i$  挑选出的特征可以最大化地判别子类别。因此类内特征关系正则化定义为

$$\|\mathbf{W}_i^T \mathbf{W}_i - \mathbf{E}\|_F^2, \quad (7)$$

其中  $\mathbf{E}$  是单位矩阵,减  $\mathbf{E}$  是为了忽略各个列向量与自身的正交约束。

综上讨论,我们在稀疏学习特征选择公式(4)中添加类内判别正则项来加强类内特征的判别性,因此目标函数定义为

$$T(W_0, W_1, \dots, W_N) = \min_{W_0, W_1, \dots, W_N} \sum_{i=0}^N \left( \|X_i W_i - Y_i\|_F^2 + \lambda \|W_i\|_{2,1} + \beta \|W_i^T W_i - E\|_F^2 \right), \quad (8)$$

其中前两项分别是损失函数和稀疏正则项;第三项是类内判别正则项,用于提高类内特征的判别性, $\beta$ 是惩罚类内特征相关性的非负参数。

### 2.3 基于正交约束和最大化类内特征判别性的分层分类特征选择算法

在本节中,结合考虑兄弟关系公式(5)、类内特征关系公式(7)和稀疏学习特征选择公式(4),可得最终目标函数公式

$$T(W_0, W_1, \dots, W_N) = \min_{W_0, W_1, \dots, W_N} \sum_{i=0}^N \left( \|X_i W_i - Y_i\|_F^2 + \lambda \|W_i\|_{2,1} + \beta \|W_i^T W_i - E\|_F^2 \right) + \alpha \sum_{i=1}^N \sum_{j \in S_i} \|W_j^T W_i\|_F^2. \quad (9)$$

通过优化公式(9),得到第*i*个非叶节点的特征权重矩阵 $W_i$ ,我们对相应任务中的特征权重按降序排序,选择排名靠前的特征作为最优特征。

### 2.4 算法优化

由于 $\ell_{2,1}$ 范数的非光滑性导致优化时很难推导出封闭解。针对这个问题,Argyriou等<sup>[27]</sup>提出了解决方法,当 $w_i \neq 0, i = 1, \dots, d$ 时, $\|W\|_{2,1}$ 对 $W$ 的导数为

$$\frac{\partial \|W\|_{2,1}}{\partial W} = \frac{\partial \text{Tr}(W^T D W)}{\partial W} = 2D W, \quad (10)$$

其中 $D \in R^{d \times d}$ 是对角矩阵,第*j*个对角元素是 $D_{jj} = \frac{1}{2\|w_j\|_2}$ ,若 $w_j = 0$ ,则 $D_{jj} = \epsilon$ 。

本文对 $\ell_{2,1}$ 范数使用上述优化方法,由于在层次结构中根节点没有兄弟节点,因此根节点和非根节点的内部节点需要分开优化计算,因此最终目标函数公式(9)可变换为

$$T(W_0, W_1, \dots, W_N) = \min_{W_0, W_1, \dots, W_N} \left( \|X_0 W_0 - Y_0\|_F^2 + \lambda \|W_0\|_{2,1} + \beta \|W_0^T W_0 - E\|_F^2 \right) + \sum_{i=1}^N \left( \|X_i W_i - Y_i\|_F^2 + \lambda \|W_i\|_{2,1} + \alpha \sum_{j \in S_i} \|W_j^T W_i\|_F^2 + \beta \|W_i^T W_i - E\|_F^2 \right). \quad (11)$$

对于根节点,目标函数为

$$T(W_0) = \min_{W_0} \|X_0 W_0 - Y_0\|_F^2 + \lambda \|W_0\|_{2,1} + \beta \|H_0 W_0 - E\|_F^2. \quad (12)$$

为了方便维度计算,令 $H_0 = W_0^T$ 。然后将(12)式中 $W_0$ 的导数设为0,即

$$\begin{aligned} \frac{\partial T}{\partial W_0} &= 2X_0^T(X_0 W_0 - Y_0) + 2\lambda D_0 W_0 + 2\beta H_0^T(H_0 W_0 - E) = \\ &2(X_0^T X_0 + \lambda D_0 + \beta H_0^T H_0)W_0 - 2(X_0^T Y_0 + \beta H_0^T E) = 0. \end{aligned} \quad (13)$$

因此,求得 $W_0$ 为

$$W_0 = (X_0^T X_0 + \lambda D_0 + \beta H_0^T H_0)^{-1} (X_0^T Y_0 + \beta H_0^T E). \quad (14)$$

对于中间节点,目标函数为

$$T(W_1, \dots, W_N) = \min_{W_1, \dots, W_N} \sum_{i=1}^N \left( \|X_i W_i - Y_i\|_F^2 + \lambda \|W_i\|_{2,1} + \alpha \sum_{j \in S_i} \|W_j^T W_i\|_F^2 + \beta \|H_i W_i - E\|_F^2 \right). \quad (15)$$

同样为了便于优化,令 $H_i = W_i^T$ 。然后将(15)式中 $W_i$ 的导数设为0,即

$$\begin{aligned} \frac{\partial T}{\partial W_i} &= 2X_i^T(X_i W_i - Y_i) + 2\lambda D_i W_i + 2\alpha \sum_{j \in S_i} W_j (W_j^T W_i) + 2\beta H_i^T(H_i W_i - E) = \\ &2 \left( X_i^T X_i + \lambda D_i + \alpha \sum_{j \in S_i} W_j W_j^T + \beta H_i^T H_i \right) W_i - 2(X_i^T Y_i + \beta H_i^T E) = 0, \end{aligned} \quad (16)$$

因此,求得 $W_i$ 为

$$\mathbf{W}_i = \left( \mathbf{X}_i^T \mathbf{X}_i + \lambda \mathbf{D}_i + \alpha \sum_{j \in \mathcal{S}_i} \mathbf{W}_j \mathbf{W}_j^T + \beta \mathbf{H}_i^T \mathbf{H}_i \right)^{-1} (\mathbf{X}_i^T \mathbf{Y}_i + \beta \mathbf{H}_i^T \mathbf{E}). \quad (17)$$

算法伪代码如算法1 HFSOC 所示,非叶节点的特征权重矩阵更新计算在第2行到第12行,对角矩阵在第3行到第5行更新计算,根节点的特征权重矩阵  $\mathbf{W}_0$  在第6行更新计算,中间节点的特征权重矩阵  $\mathbf{W}_i$  在第7行到第9行更新计算。算法1的时间复杂度主要取决于特征权重矩阵  $\mathbf{W}$  的更新时间,每个中间节点的  $\mathbf{X}_i^T \mathbf{X}_i$  与  $\mathbf{X}_i^T \mathbf{Y}_i$  各需要计算一次,时间复杂度分别为  $O(m^2 n_i)$ ,  $O(m d n_i)$ ; 其中  $m$  表示特征数,  $n_i$  表示第  $i$  个内部节点的样本数,  $d$  表示内部节点的最大类别数。因此,所有中间节点更新的时间复杂度为  $O(m^2 n + m d n)$ ; 根节点更新的时间复杂度为  $O(m^3 + m^2 d)$ 。假设总迭代次数为  $K$ , 那么 HFSOC 算法的总时间复杂度为  $O(K(m^3 + m^2 d) + m^2 n + m d n)$ 。

**算法1 HFSOC**

输入: 样本矩阵  $\mathbf{X}_i \in R^{n_i \times m}$ , 标签  $\mathbf{Y} \in \{0, 1\}^{n_i \times d}$ , 其中  $i=0, 1, \dots, N$ ,  $d$  是内部节点的最大类别数, 正则化参数  $\lambda, \alpha, \beta$ , 迭代次数  $F$

输出: 特征权重矩阵  $\mathbf{W} \in R^{m \times d}$

1: 设置迭代次数  $f$ , 随机初始化矩阵  $\mathbf{W} \in R^{m \times d}$ ,  $\mathbf{W} = [\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_N]$

2: WHILE  $f < F$  DO

3: FOR  $i=0:N$  DO

4: 通过  $D_{jj} = 1 / (2 \|w_j\|_2)$  求得矩阵  $\mathbf{D}_i^{(f)}$

5: END FOR

6: 更新  $\mathbf{W}_0$ :  $\mathbf{W}_0^{(f+1)} = (\mathbf{X}_0^T \mathbf{X}_0 + \lambda \mathbf{D}_0^{(f)} + \beta \mathbf{H}_0^T \mathbf{H}_0)^{-1} (\mathbf{X}_0^T \mathbf{Y}_0 + \beta \mathbf{H}_0^T \mathbf{E})$

7: FOR  $i=1:N$  DO

8: 更新  $\mathbf{W}_i$ :  $\mathbf{W}_i^{(f+1)} = (\mathbf{X}_i^T \mathbf{X}_i + \lambda \mathbf{D}_i^{(f)} + \alpha \sum_{j \in \mathcal{S}_i} \mathbf{W}_j \mathbf{W}_j^T + \beta \mathbf{H}_i^T \mathbf{H}_i)^{-1} (\mathbf{X}_i^T \mathbf{Y}_i + \beta \mathbf{H}_i^T \mathbf{E})$

9: END FOR

10: 更新  $\mathbf{W}^{(f+1)} = [\mathbf{W}_0^{(f+1)}, \mathbf{W}_1^{(f+1)}, \dots, \mathbf{W}_N^{(f+1)}]$

11:  $f = f + 1$

12: END WHILE

13: 返回  $\mathbf{W}$

### 3 实验分析

#### 3.1 数据集

本实验使用了5个数据集,所有数据集均具有层次结构,其中有两个蛋白质数据集: DD<sup>[28]</sup>和 F194<sup>[29]</sup>, 三个图像数据集: CLEF<sup>[30]</sup>、VOC (The PASCAL Visual Object Classes dataset)<sup>[31]</sup>和 ILSVRC<sup>[32]</sup>。表1给出了数据集的详细信息。

#### 3.2 评价指标

为了能更准确地衡量算法,除了传统的度量精度的方法以外,实验还引用了以下两种分层分类特有的评价指标:

表1 数据集描述

Table 1 Description of the datasets

数据集	训练集	测试集	特征数	节点数	叶子节点数	层数
DD	3 020	605	473	32	27	3
F194	7 105	1 420	473	202	194	3
CLEF	8 368	939	80	88	63	3
VOC	7 178	5 105	1 000	30	20	5
ILSVRC	12 346	11 845	4 096	65	57	4

树诱导误差(Tree Induced Error, TIE)<sup>[33]</sup>能够反映测试样本在层次结构上的误差程度,用真实标签节点和预测标签节点在层次结构中的总边数表示,即节点*i*的TIE表示为:

$$i_{\text{TIE}}(y, \hat{y}) = \sum_E (y, \hat{y}). \quad (18)$$

Hierarchical-F1 measure<sup>[25]</sup>在传统的F1 measure的基础上充分考虑了真实类别和预测类别的祖先和后代,它代表了算法在准确率和召回率上的总体分类表现。用准确率和召回率的调和平均数表示,即

$$F_H = \frac{2 \cdot P_H \cdot R_H}{P_H + R_H}. \quad (19)$$

准确率 $P_H$ 和召回率 $R_H$ 分别为

$$P_H = \frac{|Y_{\text{aug}} \cap \hat{Y}_{\text{aug}}|}{|Y_{\text{aug}}|}, R_H = \frac{|Y_{\text{aug}} \cap \hat{Y}_{\text{aug}}|}{|\hat{Y}_{\text{aug}}|},$$

其中 $Y_{\text{aug}}$ 表示真实标签扩展集, $\hat{Y}_{\text{aug}}$ 表示预测标签扩展集,且有 $Y_{\text{aug}} = y \cup \text{Anc}(y)$ , $\hat{Y}_{\text{aug}} = \hat{y} \cup \text{Anc}(\hat{y})$ , $\text{Anc}(y)$ 和 $\text{Anc}(\hat{y})$ 分别表示真实标签 $y$ 和预测标签 $\hat{y}$ 的祖先节点集合。

### 3.3 对比算法

在本实验中选取的六种对比算法如下所示:

(1)HFisher<sup>[10]</sup>:Fisher分数是一种用于多元分类和数据降维的有效方法。HFisher算法是在Fisher分数的基础上进行改进后与层次结构结合的分层分类特征选择算法。

(2)HFSNM(Efficient and Robust Feature Selection via Joint  $\ell_{2,1}$ -norms Minimization)<sup>[11]</sup>:基于 $\ell_{2,1}$ 范数最小化的高效鲁棒的特征选择算法,该算法对损失函数和正则化联合使用 $\ell_{2,1}$ 范数最小化,即对所有数据进行 $\ell_{2,1}$ 范数正则化,选择具有联合稀疏性的特征,实现了鲁棒的特征选择目标。

(3)HFSDK(Robust Hierarchical Feature Selection Driven by Data and Knowledge)<sup>[12]</sup>:基于数据和知识驱动的鲁棒的分层分类特征选择算法,该算法采用自顶向下的方式从粗粒度到细粒度为每个非叶节点选择具有鲁棒性和判别性的局部特征子集。

(4)Hier-FS(Hierarchical Feature Selection)<sup>[19]</sup>:基于层次结构的特征选择算法,该算法只考虑节点间的层次结构不考虑层次结构间的依赖关系。

(5)HiRRfam-FS(Family Relationship Based Hierarchical Feature Selection with Recursive Regularization)<sup>[19]</sup>:基于家庭关系的递归正则化层次特征选择算法,该算法提出父子关系正则项和兄弟关系正则项,进而为每个节点选取不同的特征子集。

(6)HFS-MIMR(Feature Selection via Maximizing Inter-class Independence and Minimizing Intra-Class Redundancy for Hierarchical Classification)<sup>[21]</sup>:基于最大化类间独立性和最小化类内冗余性的分层分类特征选择算法,该算法不仅充分考虑类间结构关系同时也考虑类内特征关系来进行特征选择。

### 3.4 实验设置

本实验采用线性支持向量机作为分类器。实验参数 $\lambda$ 设置为10,参数 $\alpha$ 和 $\beta$ 的设置:蛋白质数据集的 $\alpha$ 和 $\beta$ 分别是0.1和1,图像数据集的 $\alpha$ 和 $\beta$ 分别是1和10。据文献[10-12,19,21],本实验所有对比算法均使用其调整好的参数,为了与对比算法保持一致,本实验保持分别采用10%的蛋白质数据集和20%的图像数据集的特征。另外,为了保证算法的精度,本实验采用十折交叉验证法作为模型评估方法。

实验环境:主机环境为16 GB内存、3.30 GHz的AMD Ryzen 5 5600H CPU和Windows 11系统,编程环境为Matlab 2021a软件。

### 3.5 实验结果及分析

#### 3.5.1 性能比较

本节分别从三个评价指标列表分析,表2、表3和表4分别给出了7个算法在各个数据集上的准确率、Hierarchical-F1 measure和TIE。“↑”表示评价指标结果数值越大越好,“↓”表示评价指标结果数据越小越好。表中黑色粗体数值表示在各个数据集上最好的结果。

从三个表中可以得出,本论文所提算法在大部分数据集上都能表现出良好效果。例如在DD数据集上,本论文所提算法比HFisher的准确率高出0.1784, Hierarchical-F1 measure高出0.0909。与其他算法相比,本论文所提算法在VOC数据集上表现略差于HFSDK和HFS-MIMR,这可能是由于所提算法实验设置的参数不适合VOC数据集,并且VOC数据集中的噪声数据较多,算法具有鲁棒性对VOC数据集分类学习更有优势。总体来说,实验结果证明本论文所提算法在分类任务的特征选择方面具有良好的性能,并且能够在分类过程中带来更好的预测效果。

为了评价算法性能,引入统计学检验Friedman检验和Bonferroni-Dunn后验检验<sup>[34]</sup>。非参数检验Friedman检验的统计量表示为

$$F_F = \frac{(N-1)\chi_F^2}{N(k-1) - \chi_F^2}, \text{其中 } \chi_F^2 = \frac{12N}{k(k+1)} \left( \sum_{i=1}^k R_i^2 - \frac{k(k+1)^2}{4} \right),$$

其中 $k$ 和 $N$ 分别是算法和数据集的数量; $R_i$ 表示给定算法在所有数据集中的平均排名; $F_F$ 服从Fisher分布,自由度为 $F(k-1, k-1(N-1))$ 。

利用表3的Hierarchical-F1 measure值进行排序,在每个数据集上性能最好的算法排名为1,第

表2 算法HFSOC与6个对比算法在各个数据集上的分类准确率(↑)

Table 2 Classification accuracy of algorithm HFSOC versus six comparison algorithms on each dataset

数据集	HFisher	HFSNM	HFSDK	Hier-FS	HiRRfam-FS	HFS-MIMR	HFSOC
F194	0.257 7	0.245 1	0.335 2	0.331 0	0.338 0	0.345 1	<b>0.345 8</b>
DD	0.520 7	0.682 8	0.687 7	0.682 7	0.686 1	0.687 8	<b>0.699 1</b>
CLEF	0.580 4	0.524 0	0.625 1	0.621 9	0.616 6	0.618 8	<b>0.629 4</b>
VOC	0.393 4	0.403 3	<b>0.430 2</b>	0.419 0	0.424 3	0.428 8	0.425 1
ILSVRC	0.850 1	0.843 1	0.850 9	0.850 9	0.849 7	0.851 0	<b>0.852 1</b>

注:↑表示评价指标结果数值越大越好。

表3 算法HFSOC与6个对比算法在各个数据集上的Hierarchical-F1 measure(↑)

Table 3 Hierarchical-F1 measure of algorithm HFSOC versus six comparison algorithms on each dataset

数据集	HFisher	HFSNM	HFSDK	Hier-FS	HiRRfam-FS	HFS-MIMR	HFSOC
F194	0.675 8	0.646 2	0.709 2	0.708 9	0.711 7	0.714 6	<b>0.716 2</b>
DD	0.774 1	0.852 4	0.859 0	0.858 4	0.860 6	0.861 2	<b>0.865 0</b>
CLEF	0.739 5	0.710 0	0.772 5	0.765 7	0.762 8	0.766 9	<b>0.773 9</b>
VOC	0.657 6	0.673 9	0.677 7	0.675 4	0.675 8	<b>0.679 4</b>	0.678 3
ILSVRC	0.958 0	0.956 3	0.959 3	0.959 1	0.958 8	0.958 9	<b>0.959 5</b>

注:↑表示评价指标结果数值越大越好。

表4 本文所提算法HFSOC与6个对比算法在各个数据集上的TIE(↓)

Table 4 TIE of algorithm HFSOC versus six comparison algorithms on each dataset

数据集	HFisher	HFSNM	HFSDK	Hier-FS	HiRRfam-FS	HFS-MIMR	HFSOC
F194	1.945	2.123	1.750	1.746	1.730	1.728	<b>1.721</b>
DD	1.230	0.830	0.850	0.850	0.836	0.826	<b>0.818</b>
CLEF	2.010	2.240	1.760	1.796	1.758	1.794	<b>1.746</b>
VOC	2.279	2.220	2.120	2.143	2.138	<b>2.115</b>	2.126
ILSVRC	0.336	0.350	0.334	0.328	0.329	0.327	<b>0.323</b>

注:↓表示评价指标结果数值越小越好。

二好的算法排名为2,以此类推。在相等数据值的情况下,分配平均排名。按照以上规则排好序后计算每个算法在所有数据集上的平均排名。原假设为所有特征选择算法性能都相同,通过计算得出  $F_F=16.270$ 。对于七个算法和五个数据集的临界值  $F(7-1, (7-1) \times (5-1))=F(6, 24)$ , 大于  $\alpha=0.05$  时的  $F$  检验临界值 2.508, 因此拒绝原假设。即被比较的七种算法并不相同并且算法之间存在显著差异。

本文采用 Bonferroni-Dunn 后验检验来进一步比较算法之间的差异。Bonferroni-Dunn 后验检验的临界值计算方式为  $CD_\alpha = q_\alpha \sqrt{\frac{k(k+1)}{6N}}$ 。通过计算,在显著性水平  $\alpha = 0.1$  时,有  $q_\alpha = 2.394$ , 因此可计算出  $CD=3.2708 (k=7, N=5)$ 。为了直观地展示这些差异,使用带有临界值的特殊图形来连接这些算法。图 1 为使用 Bonferroni-Dunn 后续检验比较所提算法 HFSOC 与其他六种对比算法性能的检验结果。

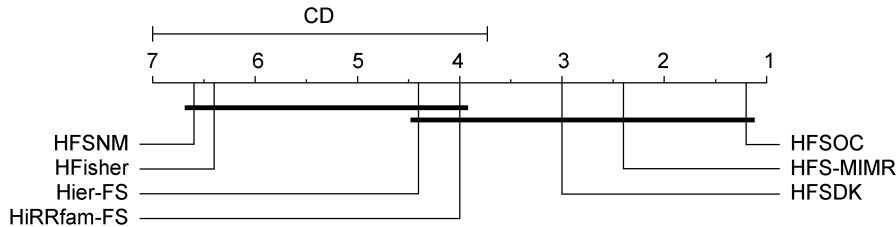


图 1 通过 Bonferroni-Dunn 检验验证 HFSOC 算法性能

(图中数字横轴表示各个算法的平均性能排名,CD线表示两个算法间在统计上可被认为有显著差异的最小差异值。图中如果两个算法之间没有连线,表示这两个算法之间存在显著差异;反之,如果存在连线,则差异不显著。本图说明 HFSOC 算法与 HFisher 算法和 HFSNM 算法之间存在显著差异。)

Fig. 1 Verification of HFSOC algorithm performance by Bonferroni-Dunn test

(The numerical horizontal axis of the figure indicates the average performance ranking of each algorithm, and the CD line indicates the value of the smallest difference between two algorithms that can be considered statistically significant. If there is no line between two algorithms in the figure, it means that there is a significant difference between these two algorithms; conversely, if there is a line, the difference is not significant. This figure illustrates that there is a significant difference between the HFSOC algorithm and the HFisher and HFSNM algorithms.)

### 3.5.2 消融实验

本节通过消融实验验证算法中兄弟关系正则化和类内特征关系正则化的有效性,公式(11)各部分可以重新组合表示如下:

(1)Hspar: 该公式为公式(11)减去兄弟关系正则化和特征关系正则化,只剩下稀疏学习公式,即

$$\min_W \sum_{i=0}^N (\|X_i W_i - Y_i\|_F^2 + \|W_i\|_{2,1}) \tag{20}$$

(2)HFSOC- $\alpha$ : 该公式为公式(11)减去兄弟关系正则化,由稀疏学习公式和类内特征关系正则化组合表示,即

$$\min_{W_0, W_1, \dots, W_N} \sum_{i=0}^N (\|X_i W_i - Y_i\|_F^2 + \lambda \|W_i\|_{2,1} + \beta \|W_i^T W_i - E\|_F^2) \tag{21}$$

(3)HFSOC- $\beta$ : 该公式为公式(11)减去类内特征关系正则化,由稀疏学习公式和兄弟关系正则化组合表示,即

$$\min_{W_0, W_1, \dots, W_N} \sum_{i=0}^N (\|X_i W_i - Y_i\|_F^2 + \lambda \|W_i\|_{2,1}) + \alpha \sum_{i=1}^N \sum_{j \in S_i} \|W_j^T W_i\|_F^2 \tag{22}$$

图 2 中(a-e)分别展示了在各个数据集中 HFSOC 算法不同组合的 Hierarchical-F1 measure 结果。从图中我们可以看出,类间独立性正则项对 F194 数据集起到重要作用,而 CLEF 数据集和

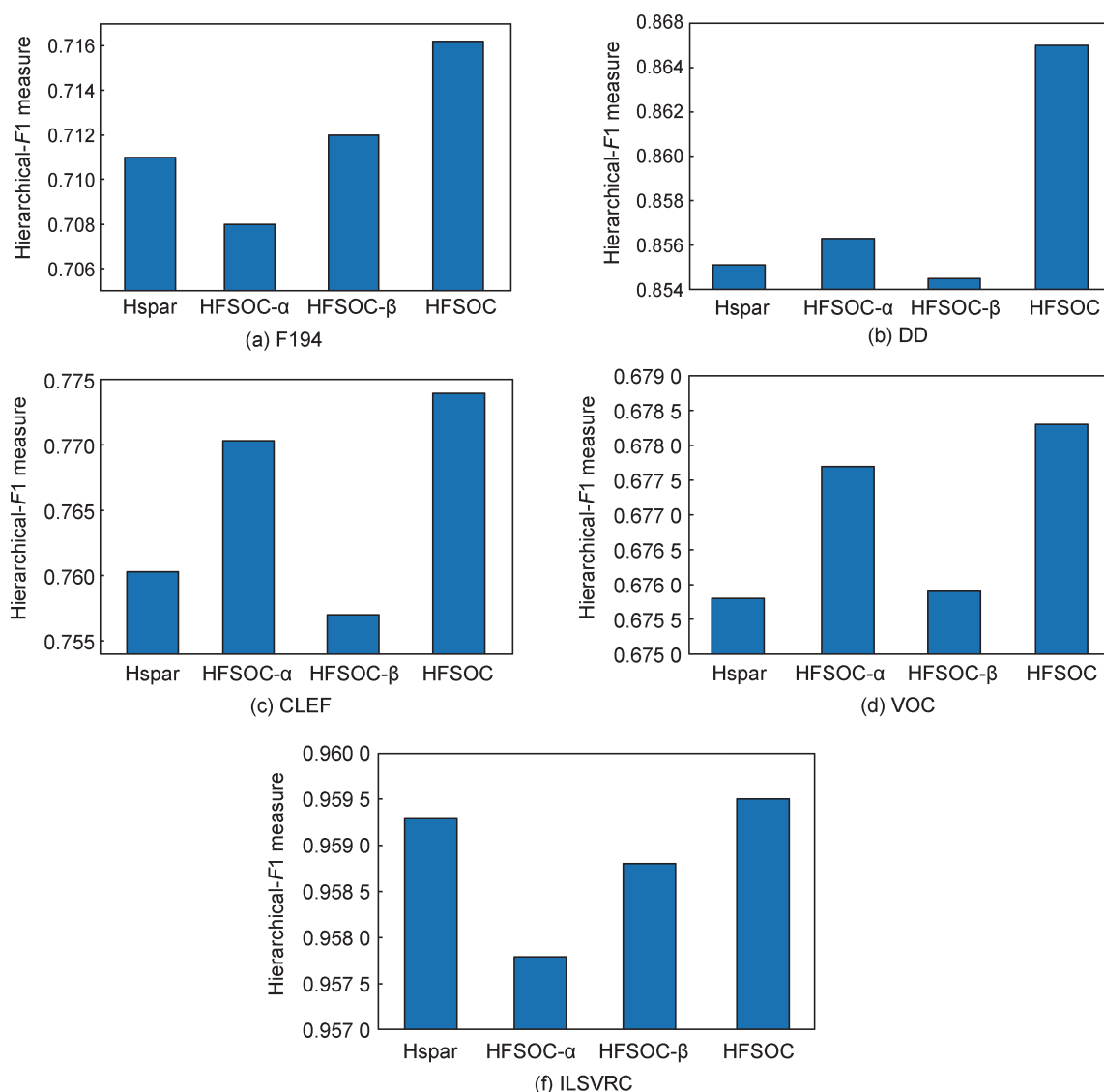


图2 各个数据集上的消融实验结果

((a—d)和(e)分别表示在F194数据集、DD数据集、CLEF数据集、VOC数据集和ILSVRC数据集上对不同模型组合的Hierarchical-F1 measure值的统计对比图,进而验证HFSOC算法的有效性。(a—d)和(e)中横轴均表示四个不同模型组合的公式,纵轴均表示Hierarchical-F1 measure值。)

Fig. 2 Results of ablation experiments on various datasets

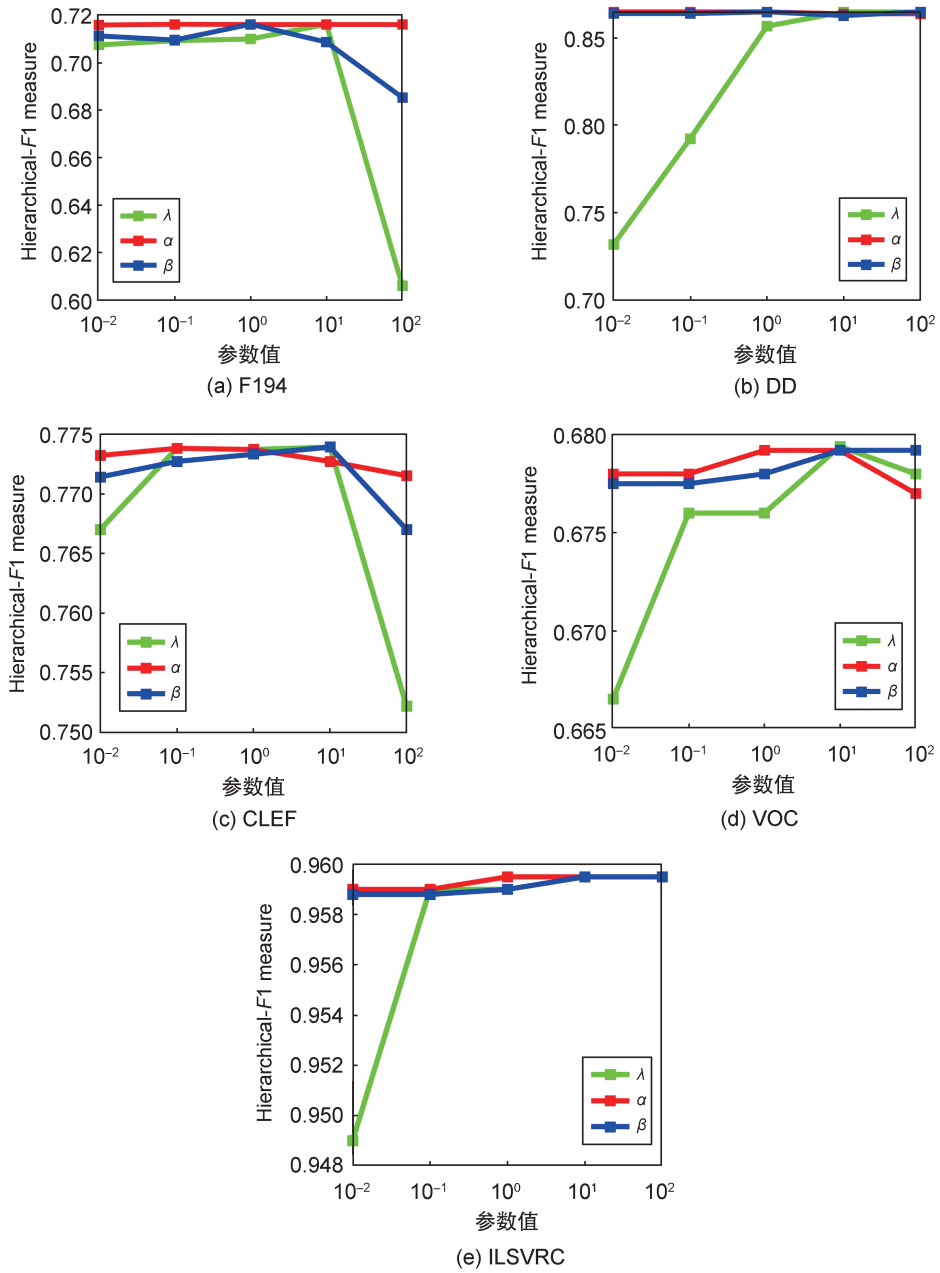
((a—d) and (e) represent the statistical comparison plots of the Hierarchical-F1 measure values for different model combinations on the F194 dataset, the DD dataset, the CLEF dataset, the VOC dataset and the ILSVRC dataset, respectively, which in turn validate the effectiveness of the HFSOC algorithm. The horizontal axes in (a—d), and (e) all denote the formulas for the four different model combinations, and the vertical axes all denote the Hierarchical-F1 measure values.)

VOC数据集对类内特征判别性正则项更为敏感。DD数据集在重新组合的三个公式上表现均不佳,ILSVRC数据集的类别独立性高,易于分类,因此ILSVRC数据集在稀疏学习公式上表现与HFSOC不相上下。总体来说,HFSOC算法明显优于其他组合。

### 3.5.3 参数敏感性分析

本节进行参数敏感性分析,在本文所提算法HFSOC中共有三个参数,分别是 $\lambda$ 、 $\alpha$ 和 $\beta$ 。其中 $\lambda$ 控制着稀疏正则项, $\alpha$ 和 $\beta$ 分别控制着兄弟关系正则项和类内特征关系正则项。我们采用网格搜索法在 $\{0.01, 0.1, 1, 10, 100\}$ 中调整三个参数。实验通过固定两个参数改变一个参数来得出Hierarchical-F1 measure的结果,并根据此结果验证本文所提算法对变化参数的敏感性。

图3给出了各个数据集上的参数敏感性分析结果,可以得出当 $\lambda$ 过小时,特征的稀疏性约束小,选择的特征较多,影响分类效果;当 $\lambda$ 太大时,特征的稀疏性约束较大,选择的特征较少并且特征重要性降低,无法达到想要的效果。 $\alpha$ 和 $\beta$ 与 $\lambda$ 相比较为稳定,但当 $\alpha$ 和 $\beta$ 过小时,兄弟节点之间选择的特征和类内选择的特征的独立性都较小,冗余性较大;当 $\alpha$ 和 $\beta$ 太大时,无法选出紧凑的特征子集,影响分类效果。总体而言,在各个数据集上进行参数调整时,所提算法HFSOC的分类效果不受太大影响。



注:图中绿线表示固定参数 $\alpha$ 和 $\beta$ ,调整参数 $\lambda$ ;红线表示固定参数 $\lambda$ 和 $\beta$ ,调整参数 $\alpha$ ;蓝线表示固定参数 $\lambda$ 和 $\alpha$ ,调整参数 $\beta$ 。

图3 参数敏感性分析

(图(a—d)和(e)分别表示在F194数据集、DD数据集、CLEF数据集、VOC数据集和ILSVRC数据集上的参数敏感性分析,各个图中横轴均表示参数值,纵轴均表示Hierarchical-F1 measure值。)

Fig. 3 Parameter sensitivity analysis

(Figures (a—d), and (e) represent the parameter sensitivity analysis on the F194 dataset, the DD dataset, the CLEF dataset, the VOC dataset, and the ILSVRC dataset, respectively, and the horizontal axis of each figure represents the parameter values, and the vertical axis represents the Hierarchical-F1 measure values.)

### 3.5.4 收敛性分析

本节我们重点研究了HFSOC算法的收敛性,实验结果如图4所示。本实验设置每个数据集的最大迭代次数为10。从图4中可以看出,所有数据集的图像都是递减的,并且都在10次迭代以内收敛。

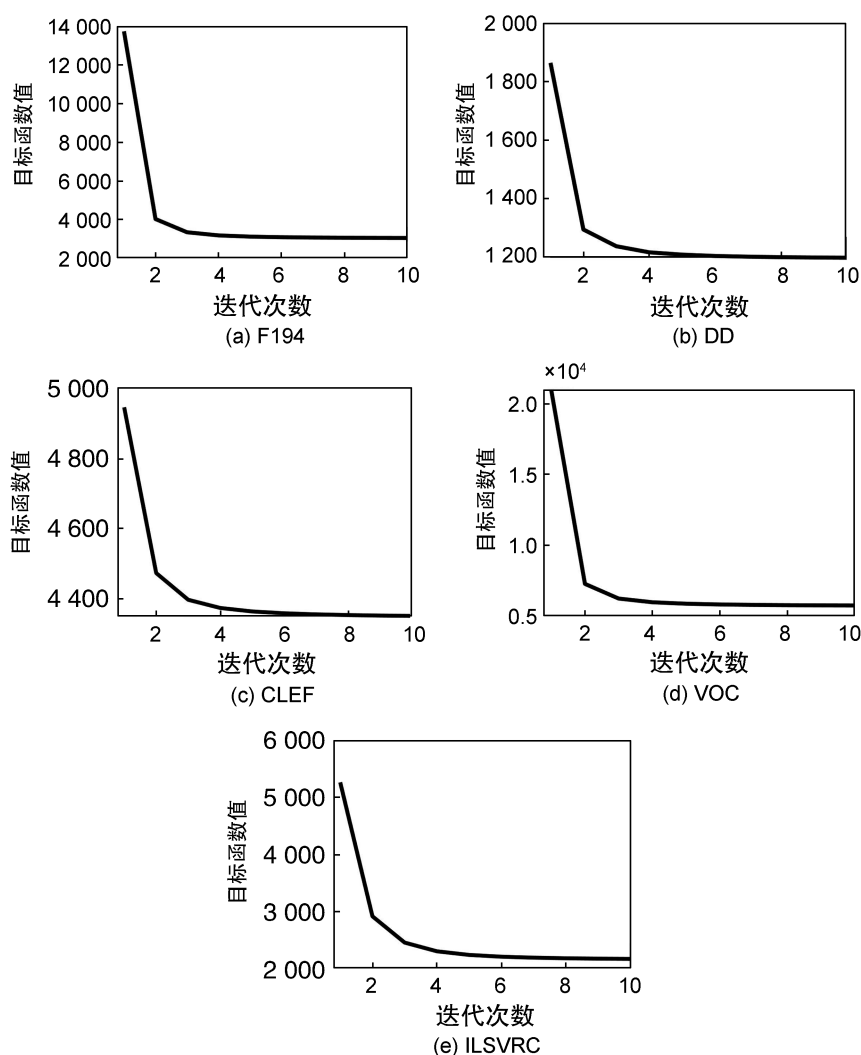


图4 目标函数值的迭代次数

(图(a—d)和(e)分别表示在F194数据集、DD数据集、CLEF数据集、VOC数据集和ILSVRC数据集上目标函数的收敛性,各个图中横轴均表示迭代次数,纵轴均表示目标函数值。)

Fig. 4 Number of iterations for the value of the objective function

(Figures (a—d) and (e) represent the convergence of the objective function on the F194 dataset, the DD dataset, the CLEF dataset, the VOC dataset, and the ILSVRC dataset, respectively, and the horizontal axes of the individual plots denote the number of iterations, and the vertical axes denote the objective function values.)

## 4 总结与展望

本文提出了一种基于正交约束和最大化类内特征判别性的分层分类特征选择算法。本文使用稀疏正则化项去除不相关特征,同时考虑类间特征独立性和类内特征判别性,利用递归正则化优化输出特征权重矩阵。通过7种算法在5个数据集上的对比实验显示,本文所提出的HFSOC算法能够做出良好表现。在本文的基础上,后续将在类内同时联合构造最小化冗余性和最大化判别性的正则项,并考虑去除噪声数据,设计鲁棒性特征选择算法。在所提算法中,手动调整参数较为

耗时,后续将建立自适应调参模型选择最优参数。并且在损失函数和优化方法上进行改进,使该模型能适用于有向无环图结构。

### 参考文献:

- [1] BENGIO S, WESTON J, GRANGIER D. Label Embedding Trees for Large Multi-class Tasks[J]. *Adv Neural Inf Process Syst*, 2010, **23**: 163-171. DOI: US20120082371 A1.
- [2] DENG J, DONG W, SOCHER R, *et al.* ImageNet: a Large-scale Hierarchical Image Database[C]//2009 IEEE Conference on Computer Vision and Pattern Recognition. New York: IEEE, 2009: 248-255. DOI: 10.1109/CVPR.2009.5206848.
- [3] 胡清华, 王煜, 周玉灿, 等. 大规模分类任务的分层学习方法综述[J]. *中国科学(信息科学)*, 2018, **48**(5): 487-500. DOI: 10.1360/N112017-00246.  
HU Q H, WANG Y, ZHOU Y C, *et al.* Review on Hierarchical Learning Methods for Large-scale Classification Task[J]. *Sci Sin Inform*, 2018, **48**(5): 487-500. DOI: 10.1360/N112017-00246.
- [4] SILLA C N, FREITAS A A. A Survey of Hierarchical Classification across Different Application Domains[J]. *Data Min Knowl Discov*, 2011, **22**(1): 31-72. DOI: 10.1007/s10618-010-0175-9.
- [5] WANG Y, HU Q H, ZHOU Y C, *et al.* Local Bayes Risk Minimization Based Stopping Strategy for Hierarchical Classification[C]//2017 IEEE International Conference on Data Mining (ICDM). New York: IEEE, 2017: 515-524. DOI: 10.1109/ICDM.2017.61.
- [6] ZHOU Y C, HU Q H, WANG Y. Deep Super-class Learning for Long-tail Distributed Image Classification[J]. *Pattern Recognit*, 2018, **80**: 118-128. DOI: 10.1016/j.patcog.2018.03.003.
- [7] LIN Y J, LIU H Y, ZHAO H, *et al.* Hierarchical Feature Selection Based on Label Distribution Learning[J]. *IEEE Trans Knowl Data Eng*, 2023, **35**(6): 5964-5976. DOI: 10.1109/TKDE.2022.3177246.
- [8] WANG Y, HU Q H, ZHU P F, *et al.* Deep Fuzzy Tree for Large-scale Hierarchical Visual Classification[J]. *IEEE Trans Fuzzy Syst*, 2020, **28**(7): 1395-1406. DOI: 10.1109/TFUZZ.2019.2936801.
- [9] WANG Y, WANG Z, HU Q H, *et al.* Hierarchical Semantic Risk Minimization for Large-scale Classification[J]. *IEEE Trans Cybern*, 2022, **52**(9): 9546-9558. DOI: 10.1109/TCYB.2021.3059631. [PubMed]
- [10] DUDA R O, HART P E. *Pattern Classification*[M]. Hoboken: John Wiley & Sons, 2006. DOI: 10.1007/978-1-4471-0409-4\_3.
- [11] NIE F P, HUANG H, CAI X, *et al.* Efficient and Robust Feature Selection via Joint  $\ell_{2,1}$ -Norms Minimization[C]//Proceedings of the 23rd International Conference on Neural Information Processing Systems. Vancouver: Curran Associates Inc. 2010: 1813-1821. DOI: 10.1007/978-3-319-10690-8\_12.
- [12] LIU X X, ZHOU Y C, ZHAO H. Robust Hierarchical Feature Selection Driven by Data and Knowledge[J]. *Inf Sci*, 2021, **551**: 341-357. DOI: 10.1016/j.ins.2020.11.003.
- [13] 张子宁, 单甘霖, 段修生, 等. 基于改进遗传算法的支持向量机特征选择[J]. *电子产品世界*, 2010, **17** (Z1): 45-47+51. DOI: 10.3969/j.issn.1005-5517.2010.1.008.  
ZHANG Z N, SHAN G L, DUAN X S. Feature Selection for SVM Based on Improved Genetic Algorithm. [J]. *Electron Eng Prod World*, 2010, **17** (Z1): 45-47+51. DOI: 10.3969/j.issn.1005-5517.2010.1.008.
- [14] PENG H C, LONG F H, DING C. Feature Selection Based on Mutual Information Criteria of Max-dependency, Max-relevance, and Min-redundancy[J]. *IEEE Trans Pattern Anal Mach Intell*, 2005, **27**(8): 1226-1238. DOI: 10.1109/TPAMI.2005.159.
- [15] FREEMAN C, KULIĆ D, BASIR O. Joint Feature Selection and Hierarchical Classifier Design[C]//2011 IEEE International Conference on Systems, Man, and Cybernetics. New York: IEEE, 2011: 1728-1734. DOI: 10.1109/ICSMC.2011.6083921.
- [16] GRIMAUDO L, MELLIA M, BARALIS E. Hierarchical Learning for Fine Grained Internet Traffic Classification[C]//2012 8th International Wireless Communications and Mobile Computing Conference (IWCMC). New York: IEEE, 2012: 463-468. DOI: 10.1109/IWCMC.2012.6314248.
- [17] SONG J, ZHANG P Z, QIN S J, *et al.* A Method of the Feature Selection in Hierarchical Text Classification Based on the Category Discrimination and Position Information [C]//2015 International Conference on Industrial Informatics-Computing Technology, Intelligent Technology, Industrial Information Integration. New York: IEEE, 2015: 132-135. DOI: 10.1109/ICIICII.2015.116.
- [18] TUO Q J, ZHAO H, HU Q H. Hierarchical Feature Selection with Subtree Based Graph Regularization[J]. *Knowl Based Syst*, 2019, **163**: 996-1008. DOI: 10.1016/j.knosys.2018.10.023.
- [19] ZHAO H, HU Q H, ZHU P F, *et al.* A Recursive Regularization Based Feature Selection Framework for

- Hierarchical Classification[J]. *IEEE Trans Knowl Data Eng*, 2021, **33**(7): 2833–2846. DOI: 10.1109/TKDE.2019.2960251.
- [20] 林耀进, 白盛兴, 赵红, 等. 基于标签关联性的分层分类共有与固有特征选择[J]. *软件学报*, 2022, **33**(7): 2667–2682. DOI: 10.13328/j.cnki.jos.006335.
- LIN Y J, BAI S X, ZHAO H, *et al.* Label-correlation-based Common and Specific Feature Selection for Hierarchical Classification[J]. *J Softw*, 2022, **33**(7): 2667–2682. DOI: 10.13328/j.cnki.jos.006335.
- [21] SHI J, LI Z Y, ZHAO H. Feature Selection via Maximizing Inter-class Independence and Minimizing Intra-class Redundancy for Hierarchical Classification[J]. *Inf Sci*, 2023, **626**: 1–18. DOI: 10.1016/j.ins.2023.01.048.
- [22] WU F H, ZHANG J, HONAVAR V. Learning Classifiers Using Hierarchically Structured Class Taxonomies [M]//ZUCKER J D, SAIITA L, eds. *Lecture Notes in Computer Science*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005: 313–320. DOI: 10.1007/11527862\_24.
- [23] LI J D, CHENG K W, WANG S H, *et al.* Feature Selection[J]. *ACM Comput Surv*, 2018, **50**(6): 1–45. DOI: 10.1145/3136625.
- [24] 刘浩阳, 林耀进, 刘景华, 等. 由粗到细的分层特征选择[J]. *电子学报*, 2022, **50**(11): 2778–2789. DOI: 10.12263/DZXB.20211263.
- LIU H Y, LIN Y J, LIU J H, *et al.* Hierarchical Feature Selection from Coarse to Fine[J]. *Acta Electron Sin*, 2022, **50**(11): 2778–2789. DOI: 10.12263/DZXB.20211263.
- [25] LI X P, WANG Y D, RUIZ R. A Survey on Sparse Learning Models for Feature Selection[J]. *IEEE Trans Cybern*, 2022, **52**(3): 1642–1660. DOI: 10.1109/TCYB.2020.2982445.
- [26] 史春雨, 毛煜, 刘浩阳, 等. 基于样本相关性的层次特征选择算法[J]. *山东大学学报(理学版)*, 2024, **59**(3): 61–70. DOI: 10.6040/j.issn.1671-9352.7.2023.1073.
- SHI C Y, MAO Y, LIU H Y, *et al.* Hierarchical Feature Selection Algorithm Based on Instance Correlations[J]. *J Shandong Univ Nat Sci*, 2024, **59**(3): 61–70. DOI: 10.6040/j.issn.1671-9352.7.2023.1073.
- [27] ARGYRIOU A, EVGENIOU T, PONTIL M. Convex Multi-task Feature Learning[J]. *Mach Learn*, 2008: 41–48. DOI: 10.1007/s10994-007-5040-8.
- [28] DING C H Q, DUBCHAK I. Multi-class Protein Fold Recognition Using Support Vector Machines and Neural Networks[J]. *Bioinformatics*, 2001, **17**(4): 349–358. DOI: 10.1093/bioinformatics/17.4.349.
- [29] LI D P, JU Y, ZOU Q. Protein Folds Prediction with Hierarchical Structured SVM[J]. *Curr Proteom*, 2016, **13**(2): 79–85. DOI: 10.2174/157016461302160514000940.
- [30] DIMITROVSKI I, KOCEV D, LOSKOVSKA S, *et al.* Hierarchical Annotation of Medical Images[J]. *Pattern Recognit*, 2011, **44**(10/11): 2436–2449. DOI: 10.1016/j.patcog.2011.03.026.
- [31] EVERINGHAM M, VAN GOOL L, WILLIAMS C K I, *et al.* The Pascal Visual Object Classes (VOC) Challenge[J]. *Int J Comput Vis*, 2010, **88**(2): 303–338. DOI: 10.1007/s11263-009-0275-4.
- [32] JIA D, KRAUSE J, BERG A C, *et al.* Hedging your Bets: Optimizing Accuracy-specificity Trade-offs in Large Scale Visual Recognition[C]//2012 IEEE Conference on Computer Vision and Pattern Recognition. New York: IEEE, 2012: 3450–3457. DOI: 10.1109/CVPR.2012.6248086.
- [33] DEKEL O, KESHET J, SINGER Y. Large Margin Hierarchical Classification[C]//Twenty-first international conference on Machine learning-ICML '04. New York: ACM, 2004. DOI: 10.1145/1015330.1015374.
- [34] DEMSAR J. Statistical Comparisons of Classifiers over Multiple Data Sets[J]. *J Mach Learn Res*, 2006, **7**: 1–30.