## 基于 ROC 的三元再编码研究

雷 蕾<sup>①</sup> 王晓丹<sup>\*①</sup> 罗 玺<sup>②</sup> <sup>①</sup>(空军工程大学防空反导学院 西安 710051) <sup>②</sup>(空军工程大学信息与导航学院 西安 710077)

**摘** 要: 针对三元编码矩阵中基分类器不包含被忽略样本类别先验知识的问题,该文提出一种基于接收机工作特性 (ROC)曲线的矩阵再编码方法。首先基于 ROC 曲线寻找构造拒绝域的阈值对,从而获得最优分类器;然后利用最 优分类器对训练样本中被忽略的类别进行分类,将经典的二值输出变为三值输出,从而对初始编码矩阵的码元"0" 进行重新编码。在解码阶段,采用经典的汉明距离解码方法对未知样本进行决策。该方法能够避免基分类器的二次 训练,适用于任意的三元纠错输出编码,具有良好的普适性和实用性。基于人工和 UCI 公共数据集的实验结果表 明该方法简单高效,在不增加训练时间的基础上,能够提高解码的速度和精度,促进分类效果的提升。 关键词: 三元纠错输出编码;二次编码;最优分类器;拒绝域;接收机工作特性

中图分类号: TP391 文献标识码: A 文章编号: 1009-5896(2016)10-2515-08 **DOI**: 10.11999/JEIT151343

# Recoding Error-correcting Output Codes Based on Receiver Operating Characteristics

 ${\rm LEI} \ {\rm Lei}^{\tiny (1)} \qquad {\rm WAGN} \ {\rm Xiaodan}^{\tiny (2)} \qquad {\rm LUO} \ {\rm Xi}^{\tiny (2)}$ 

<sup>®</sup>(Institute of Air and Missile Defense, Aire Force Engineering University, Xi'an 710051, China) <sup>®</sup>(Institute of Information and Navigation, Aire Force Engineering University, Xi'an 710077, China)

Abstract: As to the problem that the base classifiers in ternary Error Correcting Output Codes (ECOC) matrix do not contain the prior information of classes which are ignored in binary splits, a new recoding ECOC based on Receiver Operating Characteristic (ROC) curve is presented. To recode the ternary matrix, the two thresholds of reject region are obtained based on ROC to build the optimal classifiers. Then, the optimal classifiers are used to classify the ignored classes based on bipartition in training phase. In so doing, the classical two-symbol output expands to three-symbol to recode the zeros. Finally, the Hamming decoding strategy is adopted for decision in decoding. This method can avoid a second training and is applied to any kind of ternary matrix. The experiments based on Synthetic and UCI datasets validate the better efficiency and remarkable promotion without increasing training complexity of the proposed approach.

Key words: Ternary error-correcting output codes; Recoding; Optimal classifier; Reject option; Receiver Operating Characteristics (ROC)

## 1 引言

纠错输出编码(Error-Correcting Output Codes, ECOC)<sup>[1]</sup>,作为一种分而治之的多类与二类分类问题的连接桥梁,已成功应用到生物学数据识别<sup>[2]</sup>、疾病诊断<sup>[3]</sup>和计算机视觉识别<sup>[4]</sup>等诸多领域。而三元纠错输出编码(ternary ECOC)<sup>[5]</sup>作为一种更为通用的ECOC 类型,几乎能统一现有的所有多类分类框架。

矩阵中"0"符号的引入增加了二类划分的多样性和 基分类器之间的差异性,使ECOC在多类分类问题 领域的性能得以提升,并已受到众多学者的关注。 文献[6]讨论了一种新的编码设计准则——最大化三 符号编码距离准则(ternary distance),通过该准则 作者提出了一种新的稀疏编码矩阵,实验结果表明 基于该准则获得的ECOC分类效果有显著提高,且 与解码方法的选取相关性较小。2008年,文献[7]针 对样本集线性不可分问题,提出对基类子集再分割 的Subclass ECOC编码方法(SECOC)。文献[8]利 用数据分布的先验知识,基于混淆矩阵和Fisher准 则得到最优的子类划分,对子类采用"一对多"划

收稿日期: 2015-12-01; 改回日期: 2016-06-06; 网络出版: 2016-08-26 \*通信作者: 王晓丹 afeu-wang@163.com

基金项目: 国家自然科学基金(61273275, 61503407)

Foundation Items: The National Natural Science Foundation of China (61273275, 61503407)

分方法; 在单个子类内部每次只选择两种类别作为 正负类,而子类的其他类别将被忽略,即在矩阵中 编码为"0",从而弥补了单个基分类器识别的误差。 文献[9]认为利用编码矩阵中二类划分的先验原始类 结构信息可以提高 ECOC 分类性能,并给出了在流 形假设和聚类假设的情况下将先验结构信息融入基 分类器决策函数的方法。文献[10]针对经典的"一对 一"三符号编码矩阵中码元"0"会引入分类偏差的 问题,提出了避免二次训练的码元"0"再编码方法, 并将该分类结果作为权值融入基于损失函数的解码 过程中。文献[11]利用遗传算法来优化编码矩阵的构 造,将初始 ECOC 编码矩阵看作遗传个体,经过交 叉和变异形成新的编码矩阵。相关的研究还有文献 [12-14]等,这些成果都有力地促进了三元 ECOC 多类分类的发展。

在三元编码矩阵中,码元"0"所对应的类别不 参与训练,因此构造得到的基分类器将不包含此类 样本的先验信息。故当该基分类器作用于测试样本 时,就有可能因为信息缺失而输出片面的结果,导致 分类错误。本文针对此分类偏差问题,提出了一种基 于 ROC 的再编码方法(re-coding error-correcting output codes based on ROC)。该方法能够在不进 行二次训练的基础上,通过基于 ROC 构造的最优 基分类器对码元"0"进行重新编码,使得矩阵尽可 能包含所有类别样本的结构信息;同时拒绝域的构 造实现了三元 ECOC 对样本的选择性分类,矩阵中 的"0"也在解码过程中变得更为具体,从而减小解 码误差。决策时,将输出向量与新的类别编码进行 Hamming 距离解码。该方法采用的初始编码矩阵可 以为任意的三元纠错输出编码矩阵,具有很好的普 适性。

论文结构安排如下:第2节首先简要介绍基于 ROC的三元再编码思想;第3节提出基于ROC再 编码思想的步骤和方法,并对该方法存在的问题进 行一一阐述;第4节给出实验结果和分析;最后给 出总结说明。

## 2 再编码思想

ECOC框架中基分类器不能对样本进行选择性 分类,而对于未参与基分类器训练的样本而言,直 接对其分类将引入误差。

问题1 假设有符合高斯混合分布的5类数据 (ecoli数据集<sup>1)</sup>),如图1所示。当采用"一对一"方 法对其进行编码时,由 $C_1$ 和 $C_5$ 作为训练样本产生的 基分类器的决策面也能在一定程度上对类别 $C_2$ , $C_3$ 

和*C*<sub>4</sub> 进行正确的划分。这样,基分类器 *h*<sub>15</sub> 就能在 测试时对其做出正确的硬判决输出。针对此问题, 文献[10]提出了一种不需要重新训练的再编码方法, 通过已训练的基分类器对码元"0"所对应的类别进 行分类,通过设定分类正确率的阈值,将其重新编 码为{1,-1}。但对于编码矩阵中遗留的码元"0"仍 采用传统的损失函数进行解码。

**问题2** 假设一个三元编码矩阵 **M**<sub>4×7</sub> 和输出向 量 **x**。

从图2中可以看出,根据Hamming距离解码和 欧式距离解码,未知样本x被判定为 $C_1$ 。但事实上x属于第2类。因为只有当C,类在基分类器训练中没 有被忽略时,其类别属性与C,类一致;当C,类被忽 略时基分类器的决策面不能对样本 x 进行正确识 别。出现这样的情况是由于编码矩阵中码元"0"的 干扰, 使训练得到的分类器不包含对应类别的分布 信息,在解码过程中引入了偏差,从而不能做出正 确的判决。本文结合问题1和问题2,将"拒绝域" 引入每一个基分类器,对样本输出值落入拒绝域中 的样本予以拒绝,不对其进行分类识别,使基分类 器的输出由二元值扩展到三元值,从而减少基分类 器对"0"标识的类别样本直接分类带来的误差。具 体来讲: 就是用训练好的基分类器对码元'0'所对 应的样本进行分类,根据一定的拒绝域准则,将"0" 重新编码为"1"或"-1";如果码元"0"所对应样 本的输出落入拒绝域,此时选择性拒判就产生了, 此类别在新的编码矩阵中码元依旧为"0",这样就 从本质上减小了利用Hamming距离解码时带来的 误差,从而克服了问题2。



图 1 5 类数据集分布



<sup>&</sup>lt;sup>1)</sup>对 ecoli 数据集进行了归一化处理,并删除了一些样本数很小的类

要实现选择性拒判,关键是构造拒绝域。作为 分类器性能评估的有效手段,ROC曲线判决直观、 概念清楚,对样本数据的先验分布知识和错分代价 矩阵都不敏感,为解决前面提到的拒绝判决问题提 供了强大的工具。为此本文引入ROC来设计拒绝 域,构造最优分类器进行再编码设计。根据以上分 析,图3给出了基于ROC的三元再编码的多类分类 结构框图。其中*S<sub>T</sub>*和*S<sub>V</sub>*分别代表训练样本的训练子 集和验证子集。值得注意的是,再编码过程仍然在 样本训练阶段进行,这样就避免了二次训练,减少 了训练时间。

## 3 基于 ROC 的三元再编码研究方法

本节利用 ROC 构造阈值对,进而得到带拒绝 域的最优分类器,通过其对训练阶段所忽略的类别 样本进行分类再编码,最终得到新的编码矩阵。

## 3.1 基于 ROC 的最优分类器

为了说明 ROC 的生成及相关特性,首先简要 介绍两类样本混淆矩阵的相关概念。混淆矩阵描绘 了样本数据的真实类别属性与识别结果之间的关 系,是评价分类器性能的一种常用的有效方法。假 定一个二分器的分类结果为一个2×2的混淆矩阵, 如表 1 所示,其中行元素代表样本的真实属性,列 元素代表分类器的分类结果。矩阵中:TP 为正确分 类的正类样本数;FP 为被错误分类的负类样本数; FN 为被错误分类的正类样本数;TN 为正确分类的 负类样本数。定义代价矩阵  $C_{\rm M}$ ,其中 TP 和 TN 的 分类代价为0。cr =  $c_{21}/c_{12}$  为代价矩阵的费效比(cost ratio)。由此定义该分类器的分类损失代价为

$$\cot = \frac{FNc_{12} + FPc_{21}}{TP + FN + FP + TN}$$
(1)



图 3 基于 ROC 的三元再编码多类分类结构框图

表1	混淆矩阵和代价矩阵

	混淆矩阵 <b>C</b>						矩阵 <b>(</b>	<b>у</b> М	
	+	-				+	-	0	
+	TP	$_{\rm FN}$	P		+	0	$c_{12}$	$c_{13}$	
-	$\mathbf{FP}$	TN	N		-	$c_{21}$	0	$c_{23}$	

其中, FP+TN=N, TP+FN=P。ROC 为一个 2 维图形,如图4所示,y轴表示为 tp=TP/(TP+FN), x轴表示为 fp = FP/(FP+TN)。利用式(1)对 FP 求 导,并令其等于 0 可得:  $f'_{ROC}(fp^*) = \operatorname{cr} \frac{N}{P}$ 。

对于带拒绝域的基分类器而言,要找到最优的 阈值对(fp<sup>\*</sup>,tp<sup>\*</sup>),文献[15]探讨了一系列算法,本文 就对其采用的方法进行简单介绍.

假设 ROC 曲线上两点  $(fp_{\alpha}, tp_{\alpha})$  和  $(fp_{\beta}, tp_{\beta})$  对 应的基分类器为  $C_{\alpha}$  和  $C_{\beta}$ , 设  $fp_{\alpha} \leq fa_{\beta}$ , 如图 4。 $C_{\alpha}$ 和  $C_{\beta}$  是分类器根据不同阈值设定而得到的,属于一 个家族,其数据分布具有一致性,即使存在分类不 一致的情况,对分类结果也不会产生太大的影响。 故通过此 ROC 曲线得到的最优分类器  $C_{opt}$  能满足 下列条件:

$$C_{\rm opt}(x) = \begin{cases} +, \ C_{\alpha} = + \land C_{\beta} = + \\ 0, \ C_{\alpha} = - \land C_{\beta} = + \\ -, \ C_{\alpha} = - \land C_{\beta} = - \end{cases}$$
(2)

同时假设 $C_{\alpha}$ 和 $C_{\beta}$ 的混淆矩阵为(TP\_{\alpha}, FN\_{\alpha}, FP\_{\alpha}, TN\_{\alpha})和(TP<sub> $\beta$ </sub>, FN<sub> $\beta$ </sub>, FP<sub> $\beta$ </sub>, TN<sub> $\beta$ </sub>), 分类代价矩阵为  $C_{M}$ ,式(1)的分类损失代价函数可改写为

$$\begin{aligned} \operatorname{rc}(N+P) &= \underbrace{\left(\operatorname{FP}_{\beta} - \operatorname{FP}_{\alpha}\right)c_{23}}_{C_{\alpha},C_{\beta}\operatorname{disagree},-} + \underbrace{\left(\operatorname{FN}_{\alpha} - \operatorname{FN}_{\beta}\right)c_{13}}_{C_{\alpha},C_{\beta}\operatorname{disagree},+} \\ &+ \underbrace{\operatorname{FP}_{\alpha}c_{21}}_{\operatorname{FP}\operatorname{for}\operatorname{both}} + \underbrace{\operatorname{FN}_{\beta}c_{12}}_{\operatorname{FN}\operatorname{for}\operatorname{both}} \\ &= \left(\operatorname{FN}_{\alpha}c_{13} + \operatorname{FP}_{\alpha}\left(c_{21} - c_{23}\right)\right) \\ &+ \operatorname{FN}_{\beta}\left(c_{12} - c_{13}\right) + \operatorname{FN}_{\beta}c_{23}\right) \\ &= P\left(1 - f_{\operatorname{ROC}}\left(\frac{\operatorname{FP}_{\alpha}}{N}\right)\right)c_{13} + \operatorname{FP}_{\alpha}\left(c_{21} - c_{23}\right) \\ &+ P\left(1 - f_{\operatorname{ROC}}\left(\frac{\operatorname{FP}_{\beta}}{N}\right)\right)\left(c_{12} - c_{13}\right) + \operatorname{FP}_{\beta}c_{23} \end{aligned}$$

$$(3)$$

可以看出式(3)是关于  $FP_{\alpha}$ 和  $FP_{\beta}$ 两个变量的函数, 为了得到局部最小值利用式(3)分别对该两个变量 求偏导可得:



图 4 ROC 曲线特性

$$\frac{\partial \operatorname{rc}}{\partial \operatorname{FP}_{\alpha}}(N+P) = -\frac{P}{N} f_{\operatorname{ROC}}' \left(\frac{\operatorname{FP}_{\alpha}}{N}\right) c_{13} + c_{21} - c_{23}$$

$$\frac{\partial \operatorname{rc}}{\partial \operatorname{FP}_{\alpha}} \left(\operatorname{FP}_{\alpha}\right) \left(\operatorname$$

$$\frac{\partial \mathcal{FC}}{\partial \mathcal{FP}_{\beta}}(N+P) = -\frac{P}{N} f_{\text{ROC}}' \left(\frac{\mathcal{F1}_{\beta}}{N}\right) (c_{12} - c_{13}) + c_{23}$$

令式(4)等于 0 得到最终结果为  

$$f'_{\text{ROC}}(\text{fp}^*_{\alpha}) = \frac{c_{21} - c_{23}}{c_{13}} \frac{N}{P}$$
  
 $f'_{\text{ROC}}(\text{fp}^*_{\beta}) = \frac{c_{23}}{c_{12} - c_{13}} \frac{N}{P}$ 
(5)

由式(5)就可获得 ROC 曲线上满足要求的两点,以此 作为该基分类器产生拒绝域的阈值对( $\mathbf{fp}_{\alpha}^{*}, \mathbf{fp}_{\beta}^{*}$ )。由前 面的假设可知  $\mathbf{fp}_{\alpha}^{*} \leq \mathbf{fp}_{\beta}^{*}, f_{ROC}(\mathbf{fp}_{\alpha}^{*}) \leq f_{ROC}(\mathbf{fp}_{\beta}^{*})$ 。由 于 ROC 的凸包是递增和凸的,所以它的一阶导数 是非负和非增的,也就是  $f_{ROC}(\mathbf{fp}_{\alpha}^{*}) \geq f_{ROC}(\mathbf{fp}_{\beta}^{*}) \geq 0$ 。 所以代价矩阵  $C_{M}$ 需要满足以下条件:

$$(c_{21} \ge c_{23}) \land (c_{12} > c_{13}) \land (c_{21}c_{12} \ge c_{21}c_{13} + c_{23}c_{12})$$
(6)

#### 3.2 新编码矩阵的生成

上一节通过 ROC 曲线构造出带拒绝域的基分 类器并获得阈值对 $(fp_{\alpha}^{*}, fp_{\beta}^{*})$ 。其再编码过程为

$$M(i,j) = \begin{cases} 1, & c_i^j \ge f \mathbf{p}_{\beta}^* \\ 0, & f \mathbf{p}_{\alpha}^* < c_i^j < f \mathbf{p}_{\beta}^* \\ -1, & c_i^j \le f \mathbf{p}_{\alpha}^* \end{cases}$$
(7)

其中, *ci* 为利用编码矩阵第 *j* 列对应的最优分类器 对"0"标识的第 *i* 类的分类正确率。表 2 给出了基 于再编码的多类分类方法的具体实现步骤。初始编 码矩阵经过带拒绝域的基分类器重新编码以后,矩 阵中码元"0"的个数会有所下降,而基分类器对未 知样本的输出不再只有两种输出,即正类输出和负 类输出,而是出现第 3 种输出结果,即拒绝做出判 决(用"0"标示)。对此我们可以直接利用经典的汉 明距离进行解码。

下面将重点通过实验验证该方法所得的编码矩阵在分类中的应用效果。

#### 4 实验

本节采用人工数据集和 UCI 数据集来验证本文 方法的分类效果。

#### 4.1 实验数据

实验中所用的第1类数据集为5类2维正态分布人工数据集<sup>[16]</sup>,各类别数据的先验概率相同且样本个数为400,其类别分布参数如表3所示。

第2类数据集为UCI数据集,其各类数据描述如表4所示,针对部分UCI数据集,本文对其进行 了归一化处理,并删除了一些样本数很小的类,同 时为了提高分类速度,实验中对高维数据使用主成 分分析法(PCA)进行了降维处理。

#### 表 2 基于 ROC 三元再编码的多类分类方法

输入:初始编码矩阵 $M_{init} \in \{-1,0,+1\}^{K \times L}$ ,训练样本集  $\Im$  步骤 1 从训练样本集  $\Im$  中随机抽取部分样本作为验证集  $\Im_{validate}$ ,更新训练样本集  $\Im' = \Im - \Im_{validate}$ ,令编码矩阵

 M = M<sub>init</sub>。

 步骤2
 基于编码矩阵 M 和训练样本集 3' 训练得到 L 个

 基分类器 l<sub>i</sub> (i = 1,2,...,L), 利用每个基分类器对验证集 3<sub>validate</sub>

 进行分类,保留分类结果。

步骤3 针对上一步每个基分类器画出相应的ROC曲线,同时利用3.1节所示的方法确定构造每个基分类器拒绝域的阈值对 $(fp^*_{\alpha}, fp^*_{\beta})$ ,得到编码矩阵各列所对应的最优基分类器 $f_i(i = 1, 2, \cdots, L)$ 。

步骤4 利用  $f_i$  对编码矩阵 M 中"0" 对应的类别进行分类, 根据式(7) 对其进行重新编码,得到新的编码矩阵  $M_{\text{new}}$ 。

步骤 5 利用  $f_i$ 对未知样本进行分类,将输出向量与  $M_{new}$ 进行 Hamming 距离解码。

输出:分类正确率

表3 5 类人工数据集各类别分布参数

class	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$
$\mu_i$	(0,0)	(3,0)	(0,5)	(7,0)	(0,9)
$\sigma_i^2$	1	4	9	25	64

表 4 UCI 数据集及数据描述

数据集	样本书	类别数	维数
Balance	625	3	4
Ecoli	336	8	8
Glass	214	6	9
Iris	150	3	4
Satimag	6435	6	36
Segment	2310	7	19
Thyroid	215	3	5
vehicle	846	4	18
Vowel	990	11	10
Wine	178	3	13
Yeast	1484	10	8
Zoo	101	7	18

#### 4.2 实验设计

实验比较了原始"一对一"编码、稀疏随机编 码和基于两者的再编码矩阵的分类效果。在选择随 机编码方法时,我们将分别从已产生的 2000 个稀疏 随机编码阵集(对应各码元 p(-1)=1/3, p(0)=1/3, p(+1)=1/3)中随机选择所需要的编码阵。基分类 器采用了线性逻辑分类器(LOGLC)和多项式核函数支持向量机(C=2)。采用 Hamming 距离进行解码,输出最后分类决策。

在估计分类错误率时为保证估计的准确性,样本数据个数大于 500 时我们采用 10 重交叉验证,小于 500 时采用 5 重交叉验证,双边估计 t 检验法见 文献[17]。

ROC 曲线的绘制采用算法 2 实现<sup>[18]</sup>,如图 5 所示,其中 *p*, *n* 分别为正负样本个数。根据 3.1 节的方法求出该基分类器对应的阈值对和最优分类器,用于下一步测试样本的分类。

#### 4.3 实验结果及分析

**4.3.1 人工数据集** 图 6 展示了基于一对一编码矩阵的线性逻辑基分类器所绘制的 ROC 曲线。从图 6 中可以看出,将所有的点根据不同的决策阈值进行分类,得到一条从点(0,0)向点(1,1)方向延伸的 ROC 曲线,并最终趋于稳定。而对于有限数据集合而言,生成的 ROC 曲线是阶梯形状的,只有当样本数据趋近无穷时,才能得到理论上的真实 ROC 曲线。

对于人工数据集而言,正负类样本比例大致为 1:1,所以得到的最优分类器的阈值为(0.2754, 0.8751),表5为验证集的分类正确率。其中 f<sub>i</sub>为基 分类器, class,为类别。



图 6 ROC 实验曲线

假阳性率fp

表 5 基于验证集的分类正确率

	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$	$f_7$	$f_8$	$f_9$	$f_{10}$
${\rm class}_1$	0	0	0	0	0.8802	1	1	0.9323	1	0.7188
${\rm class}_2$	0	0.9100	0.6050	0.9800	0	0	0	0.3200	1	0.9100
${\rm class}_3$	0.6963	0	0.9159	0.3785	0	0.9907	0.3411	0	0	0.1355
${\rm class}_4$	0.1256	0.7236	0	0.8543	0.8794	0	0.9095	0	0.9548	0
$class_5$	0.6443	0.2577	0.7577	0	0.2268	0.8247	0	0.8299	0	0

得到新的编码矩阵如表6所示。

表 7 所示结果为各编码类型在人工数据集下的 分类正确率及其置信区间。ROC+编码方法为对应

And on Watching a strength and a	表	6	新编码矩阵行列
----------------------------------	---	---	---------

4					3	利				
仃	1	2	3	4	5	6	7	8	9	10
1	1	1	1	1	1	1	1	1	1	0
2	-1	1	0	1	1	1	1	0	1	1
3	0	-1	1	0	-1	1	0	1	1	-1
4	-1	0	-1	0	1	-1	1	-1	1	1
5	0	-1	0	-1	-1	0	-1	0	-1	-1

的基于 ROC 的再编码多类分类方法。

从表 7 可以看出,基于 ROC 的再编码方法是 有效的,其分类得到的结果要优于原始编码方法, 而基于一对一编码方法分类精度的提高要高于随机 编码,这是因为在一对一编码中,码元"0"所占的 比例高,对其进行再编码,能明显地减少"0"的数 量,缩小测试样本输出向量与真实类码字的距离, 从而提高分类正确率。同时再编码的复杂度也没有 明显的增加,所需时间消耗在可接受范围内。 4.3.2 UCI 数据集 表 8 和表 9 列出了分别以 SVM 和 LOGLC 作为基分类器的不同编码方法的分类结 果比较。其中 Re-为文献[10]所采用的再编码方法。 在每张表中加粗的数据为最大分类正确率。

	基分类器								
编码类型	线性逻辑分类器	P.C.	支持向量机						
	分类正确率&置信区间(%)	时间复杂度(s)	分类正确率&置信区间(%)	时间复杂度(s)					
一对一	$61.45 \pm 2.19$	57	$69.80{\pm}1.68$	115.00					
ROC 一对一	$63.25 \pm 1.08$	102	$72.85{\pm}4.79$	136.73					
一对多	$59.70 \pm 1.30$	6	$67.75 \pm 3.59$	59.00					
密集编码	$55.20 {\pm} 0.12$	18	$60.35{\pm}0.12$	156.00					
稀疏编码	$48.90{\pm}3.10$	26	$63.10{\pm}7.06$	84.00					
ROC 稀疏编码	$63.90{\pm}1.27$	50	$64.45 \pm 14.57$	117.54					

表7 人工数据集分类正确率及置信水平为 0.95 的置信区间(%)

表 8 基于 SVM 和 Hamming 距离解码的各数据集分类正确率及置信水平为 0.95 的置信区间(%)

	一对一	Re-一对一	ROC 一对一	一对多	密集编码	稀疏编码	Re-稀疏	ROC-稀疏
Balance	$95.52 \pm 8.22$	$95.68 \pm 6.01$	$96.00 \pm 2.11$	$96.96 {\pm} 10.23$	-	$96.28{\pm}6.04$	$95.36 \pm 7.14$	$96.32 {\pm} 6.02$
Ecoli	$74.93{\pm}14.60$	$75.34{\pm}2.68$	$75.88{\pm}5.86$	$68.75 \pm 6.15$	$42.56{\pm}0.56$	$69.64{\pm}9.67$	$70.79{\pm}8.54$	$70.67 \pm 11.81$
Glass	$96.70 \pm 18.60$	$96.23{\pm}24.65$	$97.69 {\pm} 17.26$	$96.26{\pm}0.89$	$32.72 \pm 7.77$	$80.83 \pm 10.49$	$87.95{\pm}8.39$	$88.30{\pm}8.72$
Iris	$92.68{\pm}5.38$	$96.67{\pm}8.47$	$97.33{\pm}16.94$	$94.00{\pm}25.41$	-	$64.67{\pm}5.41$	$64.67{\pm}8.47$	$65.33{\pm}3.85$
Satimag	$77.54 \pm 3.03$	$80.40{\pm}4.59$	$82.59 {\pm} 12.39$	$62.00{\pm}8.19$	$67.80{\pm}6.54$	$72.80 \pm 3.10$	$71.09 \pm 13.59$	$71.63 \pm 1.98$
Segment	$93.46{\pm}3.85$	$93.70{\pm}9.90$	$93.68{\pm}9.35$	$92.81{\pm}14.30$	$54.29{\pm}8.17$	$82.34{\pm}2.86$	$83.85 \pm 14.85$	$87.14 \pm 11.55$
Vehicle	$74.94{\pm}9.76$	$75.53 \pm 12.78$	$75.65{\pm}17.29$	$73.52 \pm 3.80$	$53.52 \pm 0.80$	$58.87 \pm 5.28$	$59.57 {\pm} 10.23$	$60.73 \pm 3.54$
Vowel	$96.26{\pm}6.42$	$98.58 \pm 1.36$	$98.71{\pm}2.89$	$87.68 \pm 28.24$	$82.12 \pm 7.97$	$79.80{\pm}2.48$	$80.03 \pm 1.03$	$79.09{\pm}0.00$
Wine	$73.12{\pm}6.19$	$79.81{\pm}3.95$	$80.35 {\pm} 2.19$	$74.71{\pm}10.75$	-	$78.09 \pm 4.01$	$78.89{\pm}14.12$	$79.18{\pm}3.86$
Yeast	$56.33 \pm 5.67$	$65.73{\pm}2.01$	$66.61 \pm 14.90$	$67.06 {\pm} 2.56$	$46.44{\pm}0.84$	$51.21 \pm 10.47$	$52.97{\pm}8.41$	$53.16 \pm 22.10$
Thyroid	$87.93 {\pm} 5.89$	$92.09{\pm}6.38$	$92.57{\pm}5.32$	$91.62{\pm}12.31$	$83.72 {\pm} 6.87$	$87.44 \pm 6.65$	$88.75 \pm 4.78$	$89.10 \pm 17.26$

表 9 基于 LOGLC 和 Hamming 距离解码的各数据集分类正确率及置信水平为 0.95 的置信区间(%)

	一对一	Re-一对一	ROC 一对一	一对多	密集编码	稀疏编码	Re-稀疏	ROC-稀疏
Balance	$88.48 \pm 4.30$	$90.04 \pm 3.88$	$91.72 \pm 12.01$	$96.96 {\pm} 10.23$	-	$46.08{\pm}0.94$	$47.04{\pm}13.15$	$49.76 \pm 17.29$
Ecoli	$80.07 \pm 17.40$	$82.76 \pm 14.08$	$83.62{\pm}12.01$	$77.71 {\pm} 6.26$	$42.56{\pm}0.56$	$72.31{\pm}2.10$	$71.12 \pm 12.09$	$75.32 \pm 14.72$
Glass	$84.90{\pm}8.04$	$87.79{\pm}3.85$	$89.70 {\pm} 14.32$	$86.91 \pm 3.11$	$32.72 \pm 7.77$	$68.27 \pm 2.89$	$69.62 \pm 1.28$	$70.90 \pm 3.70$
Iris	$94.00{\pm}8.47$	$96.00 \pm 16.94$	$97.33{\pm}16.94$	$94.67{\pm}0.05$	-	$74.67{\pm}5.08$	$75.67{\pm}4.23$	$76.33 \pm 16.94$
Satimag	$86.53 \pm 4.07$	$86.15{\pm}0.51$	$86.76{\pm}6.79$	$79.39{\pm}3.43$	$63.82{\pm}0.06$	$79.78{\pm}0.71$	$80.20{\pm}2.07$	$80.26{\pm}2.25$
Segment	$93.11 \pm 17.05$	$94.03 {\pm} 13.35$	$93.03 \pm 11.00$	$88.01{\pm}2.69$	$74.29 \pm 14.29$	$71.77{\pm}3.30$	$74.16{\pm}2.80$	$74.50{\pm}3.52$
Vehicle	$79.31{\pm}3.21$	$77.42 \pm 3.83$	$78.01{\pm}3.66$	$75.65 \pm 12.75$	$23.52{\pm}0.80$	$71.76{\pm}4.27$	$68.56 \pm 17.08$	$69.86{\pm}3.60$
Vowel	$77.17{\pm}5.13$	$76.36 \pm 17.97$	$77.37 \pm 7.70$	$38.08 \pm 19.25$	-	$45.56{\pm}5.05$	$48.28{\pm}2.84$	$48.84 \pm 3.85$
Wine	$94.12 \pm 13.16$	$95.05 \pm 12.98$	$95.06 \pm 7.70$	$95.24 \pm 13.64$	-	$94.39 \pm 13.48$	$92.16 \pm 14.45$	$95.48{\pm}29.20$
Yeast	$56.73 \pm 14.21$	$57.69 \pm 4.68$	$58.09{\pm}2.98$	$53.22 \pm 12.85$	$46.44{\pm}0.84$	$55.53 \pm 6.28$	$53.70 \pm 6.66$	$54.05 \pm 13.86$
Thyroid	$95.82 \pm 5.65$	$95.81{\pm}17.98$	$96.27{\pm}12.04$	$95.83{\pm}14.12$	$93.97{\pm}14.01$	$94.41{\pm}12.15$	$93.96{\pm}5.55$	$92.10 \pm 17.26$

从实验结果可以看出,基于 ROC 的再编码方 法在大部分数据集上能获得较好的分类效果。尽管 Re-ECOC 有时能取得更好的分类效果,但它不能实 现基分类器的选择性分类,对一些错分代价较大 的样本有一定的局限性。同时基于"一对一"编码 矩阵的 ECOC 多类分类方法在大部分数据集上的分 类性能提升优于基于随机编码的 ECOC 分类。而对 于 Glass 数据集而言,基于稀疏随机编码提高的幅 度较大,这是因为此时稀疏编码矩阵为

-1	0	0	-1	-1	-1	-1
1	0	0	0	1	0	1
1	-1	1	1	-1	-1	-1
1	1	0	0	0	-1	-1
-1	0	-1	-1	0	1	-1
-1	0	-1	-1	0	0	1

可以看出,此时的编码矩阵中"0"的数量较多,对 其采取再编码策略能够明显地提高分类性能。

为了更直观地体现再编码对码元"0"的影响, 本文引入拒分率来对编码前后的分类性能进行比较。

$$\rho = \frac{1}{l} \sum_{i=1}^{l} \frac{N_i^{\{0\}}}{N_i} \tag{8}$$

其中, N<sub>i</sub>为第 i 列编码对应的样本总数, N<sub>i</sub><sup>[0]</sup> 为其 码元"0"所应对的样本个数,即拒分样本个数, l为 分类器个数。图 7 为不同数据集在原始编码和再编 码情况下的拒分率比较。从图 7 中可以看出,再编 码之后拒分率明显下降,且基于一对一编码的拒分 率下降幅度大于随机编码。通过对"0"进行再编码 能够使三元类码字更加具体,提高了解码速度和精 度,从而提升分类效果。

## 5 结束语

在三符号 ECOC 多类分类中,码元 "0"对应 的类别样本将不参与此列基分类器的训练,因此利 用该基分类器对未参与训练的类别样本进行分类可 能会带来误差。本文提出的再编码 ECOC 方法在不 进行二次训练的前提下,通过引入"拒绝域"来实 现选择性分类,将基分类器的二值输出扩展为三值 输出,从而实现编码矩阵中的码元 "0"再编码,缩 小预测输出与类别码字的距离,提高分类正确率。 在为基分类器构造拒绝域时,一种基于 ROC 曲线 的拒绝域构造方法被应用到再编码过程中。最后实 验验证了本文方法的有效性。本文为实现多类别的 选择性分类提供了新思路。如何从降低误判风险出 发,研究有效的带拒绝域的 ECOC 多类分类方法是 下一步的工作重点。



## 参考文献

- DIETTERICH T G and BAKIRI G. Solving multi-class learning problems via error-correcting output codes[J]. Journal of Artificial Intelligence Research, 1995, 34(2): 263–286. doi: 10.1613/jair.105.
- [2] PHYO K S, JIAN G W, and EAM K T. Facial age range estimation with extreme learning machines[J]. *Neurocomputing*, 2015, 149A: 364–372. doi: 10.1016/ j.neucom.2014.03.074.
- [3] ELIF D Ü. ECG beats classification using multiclass support vector machines with error correcting output codes[J]. *Digital Signal Processing*, 2007, 45(17): 675–684. doi: 10.1016/j.dsp. 2006.11.009.
- [4] SERGIO E, DAVID M, ELOI P, et al. Online error correcting output codes[J]. Pattern Recognition Letters, 2011, 32(3): 458-467. doi: 10.1016/j.patrec.2010.11.005.
- [5] ERIN L A, ROBERT E S, YORAM S, et al. Reducing multiclass to binary: a unifying approach for margin classifiers[J]. Journal of Machine Learning Research, 2000, 39(1): 113–141. doi: 10.1162/15324430152733133.
- [6] SERGIO E, ORIOL P, and PETIA R. Separability of ternary error-correcting output codes[J]. *Pattern Recognition Letters*, 2009, 30(5): 285–297. doi: 10.1016/j.patrec.2008.10.002.
- [7] SERGIO E, DAVID M J T, ORIOL P, et al. Subclass problem-dependent design for error-correcting output codes
  [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2008, 30(6): 1041–1054. doi: 10.1109/TPAMI. 2008.38.
- [8] 周进登,王晓丹,周红健.基于混淆矩阵的自适应纠错输出编码多类分类方法[J].系统工程与电子技术,2012,34(7):220-226.doi:10.3969/j.issn.1001-506X.2012.07.38.
  ZHOU Jindeng, WANG Xiaodan, and ZHOU Hongjian. Multiclass classification of adaptive error-correcting output codes based on confusion matrix[J]. Systems Engineering and Electronics, 2012, 34(7): 220-226.doi: 10.3969/j.issn.1001-506X.2012.07.38.
- [9] WANG Y, CHEN S C, and XUE H. Can under-exploited structure of original-classes help ECOC-based multi-class classification?[J]. *Neurocomputing*, 2012, 89: 158–167. doi: 10.1016/j.neucom.2012.02.035.
- [10] SERGIO E, ORIOL P, and PETIA R. Re-coding ECOCs without Re-training[J]. Pattern Recognition Letters, 2010, 31(7): 555–562. doi: 10.1016/j.patrec.2009.12.002.
- [11] MIGUEL A B, SERGIO E, XAVIER B, et al. On the design of an ECOC-compliant genetic algorithm[J]. Pattern

Recognition, 2014, 47(2): 865–884. doi: 10.1016/j.patcog. 2013.06.019.

- [12] FRANCESCO C, ORIOL P, and PETIA R. ECOC-DRF: discriminative random fields based on error correcting output codes[J]. *Pattern Recognition*, 2014, 47(6): 2193–2204. doi: 10.1016/j.patcog.2013.12.007.
- [13] MIKEL G, ALBERTO F, EDURNE B, et al. DRCW-OVO: distance-based relative competence weighting combination for one-vs-one strategy in multi-class problems[J]. Pattern Recognition, 2015, 48(1): 28–52. doi: 10.1016/j.patcog. 2014.07.023.
- [14] LEI L, WANG X D, LUO X, et al. Hierarchical error-correcting output codes based on SVDD[J]. Pattern Analysis and Applications, 19(1): 163–171. doi: 10.1007/ s10044-015-0455-5.
- [15] TADEUSZ P. On the use of ROC analysis for the optimization of abstaining classifiers[J]. Machine Learning, 2007, 68(2): 137–169. doi: 10.1007/s10994-007-5013-y.
- [16] ZHOU J D, and WANG X D. Research on the unbiased

probability estimation of error-correcting output coding[J]. Pattern Recognition, 2011, 44(7): 1552–1565. doi: 10.1016/ j.patcog.2010.12.020.

- [17] ZHOU J D, YUN Y, ZHANG J M, et al. Constructing ECOC based on confusion matrix for multiclass learning problems[J]. Science China Information Sciences, 2016, 59(1): 1–14. doi: 10.10071/s11432-015-5321-y.
- [18] 邹洪侠,秦峰. 二类分类器的 ROC 曲线生成算法[J]. 计算机 技术与发展, 2009, 19(6): 109-112.
  ZOU Hongxia and QIN Feng, Algorithm for generating ROC curve of two-classifier[J]. Computer Technology and Development, 2009, 19(6): 109-112.
- 雷 蕾: 女,1988年生,博士生,研究方向为机器学习、多类分类.
- 王晓丹: 女,1966年生,教授,博士生导师,研究方向为机器学 习、计算机视觉处理、信息融合.
- 罗 玺: 男, 1988年生, 硕士, 讲师, 研究方向为智能信息处理.