

基于Adaboost-决策树算法的乳腺微钙化区域真假阳性检测

申楠¹, 邢素霞¹, 何湘萍², 潘子妍¹, 王瑜¹

1. 北京工商大学人工智能学院, 北京 100048; 2. 北京海淀妇幼保健院乳腺病防治中心, 北京 100080

【摘要】乳腺癌的早期症状在乳腺钼靶图像中主要表现为微钙化点,微钙化区域真假阳性检测对于乳腺癌早期筛查具有重要意义。本研究选取DDSM图像进行实验,手动截取了400个疑似钙化区域。首先提取全部区域的Haralick纹理特征和灰度游程矩阵特征建立特征集,然后使用Adaboost算法集成决策树,构建强分类器AB-DT,对400个疑似钙化区域进行分类。实验发现当集成462棵决策树时,模型分类性能最佳。最后进行10折交叉验证,AB-DT算法达到了91.75%的准确率,91.75%的敏感性,91.79%的特异性,F1指数为0.918 7。该模型在微钙化真假阳性检测上性能优越,可用于辅助乳腺微钙化点检测,具有一定的临床应用价值。

【关键词】乳腺癌;Adaboost-决策树;微钙化;Haralick纹理特征;灰度游程矩阵

【中图分类号】R318;TP301.6

【文献标志码】A

【文章编号】1005-202X(2021)08-0940-06

True- and false-positive detections of breast microcalcifications based on Adaboost-decision tree algorithm

SHEN Nan¹, XING Suxia¹, HE Xiangping², PAN Ziyang¹, WANG Yu¹

1. School of Artificial Intelligence, Beijing Technology and Business University, Beijing 100048, China; 2. Breast Disease Prevention and Control Center, Haidian Maternal and Child Health Hospital, Beijing 100080, China

Abstract: The early manifestation of breast cancer is mainly characterized by microcalcifications in mammograms. The true- and false-positive detections of microcalcifications are of great significance for the early screening of breast cancer. DDSM images were selected for the experiment, and 400 suspected calcification regions were manually intercepted. The feature set was firstly established by extracting Haralick texture features and grey-level run length matrix features of all regions; and then, Adaboost algorithm was integrated with decision tree to construct a strong classifier AB-DT for classifying 400 suspected calcification regions. It was found that the model classification performance was the best when 462 decision trees were integrated. Finally, 10-fold cross-validation was conducted, and the results revealed that the accuracy, sensitivity and specificity of AB-DT algorithm reached 91.75%, 91.75% and 91.79%, respectively, and that F1 score was 0.918 7. The proposed model has superior performance in the true- and false-positive detections of microcalcifications, and it can be used to assist the detection of breast microcalcifications, which has certain clinical application value.

Keywords: breast cancer; Adaboost-decision tree; microcalcification; Haralick texture feature; grey-level run length matrix

前言

乳腺癌作为一种具有高发病率和死亡率的恶性肿瘤,已成为威胁全球女性健康甚至生命的主要杀手。乳腺癌的早期病症在乳腺钼靶图像中主要表现为微钙化点。乳腺钼靶图像中存在较多的乳腺纤维

组织,其在图像中表现为高亮区域,而微钙化区域在图像中表现为细小的高亮区域,因此乳腺纤维组织和微钙化区域在亮度即灰度值上易混淆,增大了乳腺钼靶图像微钙化区域的检测难度。为此,很多研究在微钙化区域检测上做出了很多有意义的探索,如彭庆涛等^[1]提出基于小波分析和灰度纹理特征相结合的微钙化区域提取方法,微钙化点检出率为85%;商小宝^[2]提出一种基于旋转不变局部二值模式的早期乳腺钙化点检测方法,真阳性率为95.6%,假阳性率为5.6%;王科举^[3]利用周围区域矩阵反映射乳腺微钙化区域特征,并结合随机森林分类器,特异性达到88%,曲线下面积达到0.922 4; Karale等^[4]提

【收稿日期】2021-03-18

【基金项目】国家自然科学基金(61671028);国家重大科学研发子课题(ZLJC6 03-5-1)

【作者简介】申楠,在读研究生,主要从事图像处理、机器学习方面的研究,E-mail: 18210629776@163.com

【通信作者】邢素霞,博士,副教授,主要从事图像处理与嵌入式系统的研究,E-mail: xingsuxia@163.com

取乳腺微钙化区域的密度特征、形状特征、不变矩、Haralick 特征、基于方向梯度直方图的特征, 每幅图像的平均误报率为 2.59%, 并能达到 100% 的灵敏度; Suhail 等^[5]利用改进的 Fisher 线性判别方法对微钙化区域进行线性变换, 平均准确率达 96%。这些研究均为微钙化点的计算机辅助诊断的临床应用做出了贡献。

为进一步提高微钙化区域的检测准确率, 降低检测假阳性率, 本研究提出一种基于 Adaboost-决策树 (Adaboost-Decision Tree, AB-DT) 的乳腺微钙化区域真假阳性检测方法, 并采用 10 折交叉验证确定

AB-DT 模型的性能。

1 特征提取

乳腺钼靶图像为灰度图像, 钙化点的形态各异, 有管状、圆形、爆米花状、轮圈状、杆状、球状、点状和发育不良型^[6](图 1)。而乳腺组织区域中微钙化点区域相对正常区域在纹理、灰度方面差异明显^[7], 因此可提取感兴趣区域 (Region of Interest, ROI) 的纹理和灰度特征作为分类器的输入特征。本研究提取 ROI 的 Haralick 纹理特征和灰度游程矩阵特征建立特征集。

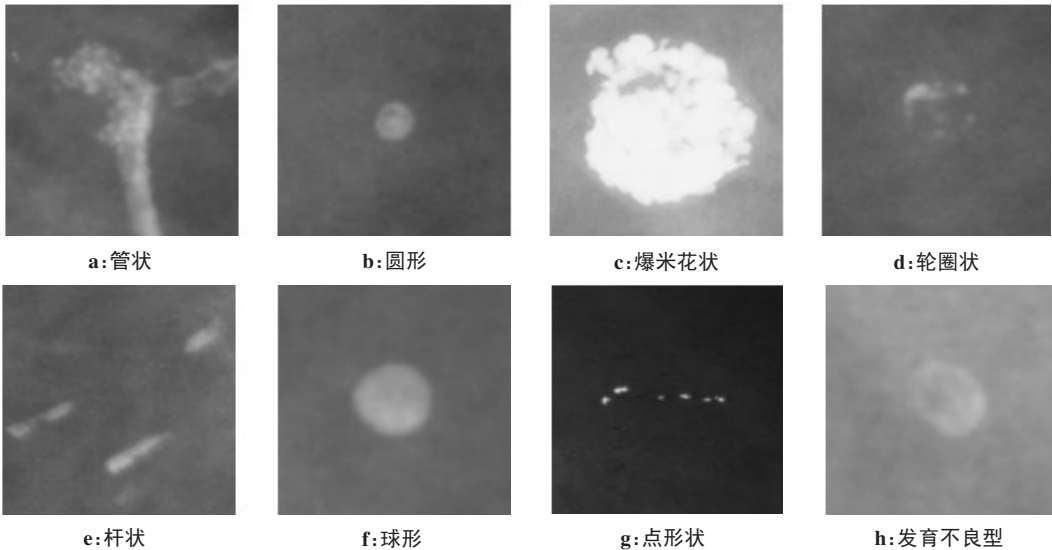


图 1 钙化点类型
Fig.1 Types of calcifications

1.1 Haralick 纹理特征

Haralick 等^[8]于 1973 年提出基于灰度共生矩阵的纹理特征统计方法。灰度共生矩阵反映了图像灰度分布关于方向、局部邻域和变化幅度的综合信息, 由于灰度共生矩阵的数据量较大, 一般不直接作为区分纹理的特征, 而是把基于它构建的一些统计量作为纹理分类特征。为简化特征提取的计算过程, 仅对 ROI 内每个像素与其 8 邻域所组成的像素对进行统计, 对应距离 $d = 1$, 且不考虑方向角度。

利用灰度共生矩阵提取 19 个特征参数, 分别是^[9]: 能量、熵、对比度、相异性、相关性、自相关系数、突出聚类、阴暗聚类、差熵、差方差、同质性、相关信息度 1、相关信息度 2、逆差距、最大概率、和平均、和熵、和方差以及平方和, 定义如表 1 所示, 表 2 为所用到的变量定义。

1.2 灰度游程矩阵

灰度游程矩阵统计图像中某个方向上同灰度值 g 连续的像素长度 l 的出现频率, 是一个二维统计矩

表 1 Haralick 纹理特征的定义
Tab.1 Definition of Haralick texture features

特征参数	特征定义
短游程优势	$\frac{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} \left[p(i,j \theta) / j^2 \right]}{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} p(i,j \theta)}$
长游程优势	$\frac{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} j^2 p(i,j \theta)}{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} p(i,j \theta)}$
长游程不均匀性	$\frac{\sum_{j=1}^{N_g} \left[\sum_{i=1}^{N_g} p(i,j \theta) \right]^2}{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} p(i,j \theta)}$
游程百分比	$\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} \left[\frac{p(i,j \theta)}{N_p} \right]$
灰度不均匀性	$\frac{\sum_{i=1}^{N_g} \left[\sum_{j=1}^{N_g} p(i,j \theta) \right]^2}{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} p(i,j \theta)}$
低灰度级游程优势	$\frac{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} \left[p(i,j \theta) / i^2 \right]}{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} p(i,j \theta)}$
高灰度级游程优势	$\frac{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} i^2 p(i,j \theta)}{\sum_{i=1}^{N_g} \sum_{j=1}^{N_g} p(i,j \theta)}$

表2 Haralick 纹理特征的变量描述
Tab.2 Variable description of Haralick texture features

特征参数	特征定义	特征参数	特征定义
能量	$\sum_{i=1}^N \sum_{j=1}^N [p(i,j)]^2$	同质性	$\sum_{i=1}^N \sum_{j=1}^N \frac{p(i,j)}{1+(i-j)^2}$
熵	$-\sum_{i=1}^N \sum_{j=1}^N p(i,j) \log p(i,j)$	相关信息度1	$\frac{HXY - HXY1}{\max(HX, HY)}$
对比度	$\sum_{i=1}^N \sum_{j=1}^N (i-j)^2 p(i,j)$	相关信息度2	$\sqrt{1 - e^{-2(HXY2 - HXY)}}$
相异性	$\sum_{i=1}^N \sum_{j=1}^N i-j \cdot p(i,j)$	逆差距	$\sum_{i=1}^N \sum_{j=1}^N \frac{p(i,j)}{1+ i-j }, i \neq j$
相关性	$\sum_{i=1}^N \sum_{j=1}^N \left(\frac{i - \mu_x}{\sigma_x} \right) \left(\frac{j - \mu_y}{\sigma_y} \right) p(i,j)$	最大概率	$\max_{ij} p(i,j)$
自相关系数	$\sum_{i=1}^N \sum_{j=1}^N (i \cdot j) p(i,j)$	和平均	$\sum_{k=2}^{2N} k p_{x+y}(k)$
突出聚类	$\sum_{i=1}^N \sum_{j=1}^N (i+j-2\mu)^3 p(i,j)$	和熵	$-\sum_{k=2}^{2N} p_{x+y}(k) \log [p_{x+y}(k)]$
阴暗聚类	$\sum_{i=1}^N \sum_{j=1}^N (i+j-2\mu)^4 p(i,j)$	和方差	$\sum_{k=2}^{2N} (k - \mu_{x+y})^2 p_{x+y}(k)$
差熵	$-\sum_{k=0}^{N-1} p_{x-y}(k) \cdot \log [p_{x-y}(k)]$	平方和	$\sum_{i=1}^N \sum_{j=1}^N (i-\mu)^2 p(i,j)$
差方差	$\sum_{k=0}^{N-1} (k - \mu_{x-y})^2 p_{x-y}(k)$		

阵^[10]。设一个 $M \times N$ 大小图像的统计矩阵为 $p(i,j|\theta)$,即在 θ 方向构建的游程矩阵在 (i,j) 坐标下的位置, N_g 为最大像素值, N_r 为不同的像素沿方向 θ 的行走距离^[11]。基于灰度游程矩阵可以计算 7 个纹理特征参数^[12]:短游程优势、长游程优势、长游程不均匀性、游程百分比、灰度不均匀性、低灰度级游程优势和高灰度级游程优势,详见表 3。

2 AB-DT 算法

2.1 Adaboost 算法

Schapire 等^[13]对 Boosting 算法进行改进得到 Adaboost 算法。Adaboost 算法运用迭代的思想,在使用样本训练集的过程中,挑选其中的关键分类特征,增加前一轮被错误分类的样本的权重,减小被正确分类的样本的权重,重复多次,逐步训练各弱分类器,并采用加权多数表决的方法调整各弱分类器的权重,最终筛选出权重系数最小的弱分类器构造成一个强分类器。

Adaboost 算法具有很强的适应性和灵活性。弱分类器可以与大多数的分类器兼容,如决策树、支持向量机、朴素贝叶斯以及 K 最近邻等算法,可根据实

际应用组合分类器,以获得最佳的分类识别效果^[14]。

2.2 决策树

决策树是一种重要的机器学习和数据挖掘算法^[15]。由于决策树算法易于理解和实现,对噪声数据具有良好的鲁棒性,与此同时具有很好的预测性能,因此被广泛用于各种实际领域。决策树是一种分类树结构的预测模型,描述了对象属性与对象值之间映射属性关系。学习时,根据最小化损失函数的原则,利用训练集数据建立决策树模型;预测时,利用决策树模型对新的数据进行分类。其中,单层决策树(也称决策树桩)是一种较为典型的简单决策树,只基于单个特征来做决策,仅有一次分裂的过程^[16],因此处理数据非常迅速、简单易行、实时性好,十分适合作为弱分类器。

3 实验过程及结果

3.1 数据来源

DDSM (Digital Database for Screening Mammography)是由南佛罗里达大学于 1999 年提供数字化乳腺图像库,该数据库是一个高分辨率的乳腺钼靶图像标准数据库,含 2 620 个病例^[17]。每个病

表3 基于灰度游程矩阵的特征定义
Tab.3 Feature definition based on grey-level run length matrix

符号	定义
$x(ij)$	非标准化灰度共生矩阵中的元素 ij
N	表示图像中的离散强度级数
$p(ij)$	$\frac{x(ij)}{\sum_{i=1}^N \sum_{j=1}^N x(ij)}$, 是任意距离和方向的灰度共生矩阵
$p_x(i)$	$\sum_{j=1}^N p(ij)$, 表示边缘行概率
$p_y(i)$	$\sum_{i=1}^N p(ij)$, 表示边缘列概率
μ_x	表示边缘行概率 p_x 的平均值
μ_y	表示边缘列概率 p_y 的平均值
σ_x	表示边缘行概率 p_x 的标准差
σ_y	表示边缘列概率 p_y 的标准差
H	$-\sum_{i=1}^N \sum_{j=1}^N p(ij) \cdot \log p(ij)$, 表示灰度共生矩阵 $p(ij)$ 的熵
HX	$-\sum_{i=1}^N p_x(i) \cdot \log p_x(i)$, 表示边缘行概率 p_x 的熵
HY	$-\sum_{i=1}^N p_y(i) \cdot \log p_y(i)$, 表示边缘列概率 p_y 的熵
$p_{x+y}(k)$	$\sum_{i=1}^N \sum_{j=1}^N p(ij), i+j=k, k=2,3,\cdots,2N$
$p_{x-y}(k)$	$\sum_{i=1}^N \sum_{j=1}^N p(ij), i-j =k, k=0,1,\cdots,N-1$

例包含右侧头尾位、左侧头尾位、右侧侧斜位、左侧侧斜位这4个视图的图像(图2)。每个视图的标注文件包含医师手动标注的病灶区域、病灶类型、病灶等级等相关信息。

将来自 DDSM 的图像转换格式后,从中手动截取 400 幅 128×128 像素的 ROI,其中 200 幅为医生标注的钙化区域图像,另外 200 幅为疑似钙化区域图像。如图 3 所示,图 3a 为乳腺钼靶的原始图像,图 3b 为用 OVERLAY 文件生成的金标准图像,图 3c 和图 3d 分别为从乳腺钼靶图像中截取的疑似钙化区域和微钙化区域图像。

3.2 AB-DT 参数的确定

乳腺微钙化区域检测方法流程如图 4 所示。提取图像的 Haralick 纹理特征及灰度游程矩阵特征,26 维特征组成特征集来量化乳腺的微钙化区域。决策树作为弱分类器,通过 Adaboost 集成算法生成强分类器 AB-DT。

在 AB-DT 算法中,决策树的数量 k 影响着分类器的性能,当 k 值较小时,AB-DT 算法的分类误差较大,算法分类性能差;而 AB-DT 算法的复杂性与 k 值成正比, k 值越大,算法的复杂度越高,运行时间也越长。在 400 幅 ROI 图像中,随机选取 80% 的 ROI 图像作为训练集,剩余的 20% 作为测试集。通过 MATLAB 中的集成学习工具箱,搭建 AB-DT 模型。

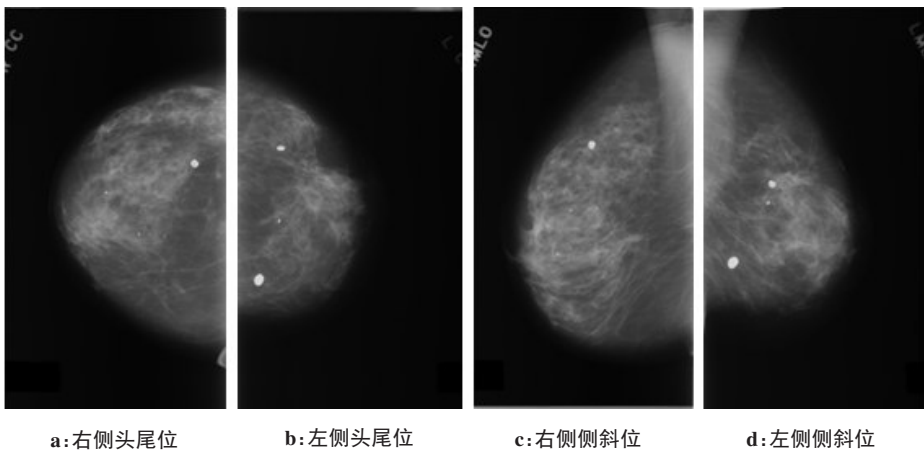


图2 DDSM 图像
Fig.2 Images from digital database for screening mammography

设置决策树的数量为 500,绘制决策树的数量与分类错误率的关系曲线,如图 5 所示。横坐标表示决策树的数量,纵坐标表示分类错误率。
由图 5 可知,当 $k \geq 484$ 时,分类错误率趋于稳定,图像存在多个分类错误率最低点,综合分类准确率和运算复杂度,决策树的数量 k 应设置为 462(即第一个分类错误率最低点的横坐标)。

3.3 评价标准

为验证和量化此分类算法的效果,采用准确率(Accuracy)、敏感性(Sensitivity)、特异性(Specificity)、阳性预测值(Positive Predictive Value, PPV)、阴性预测值(Negative Predictive Value, NPV)及 F1 分数(F1-score)对训练的模型进行评价,其定义如下所示。

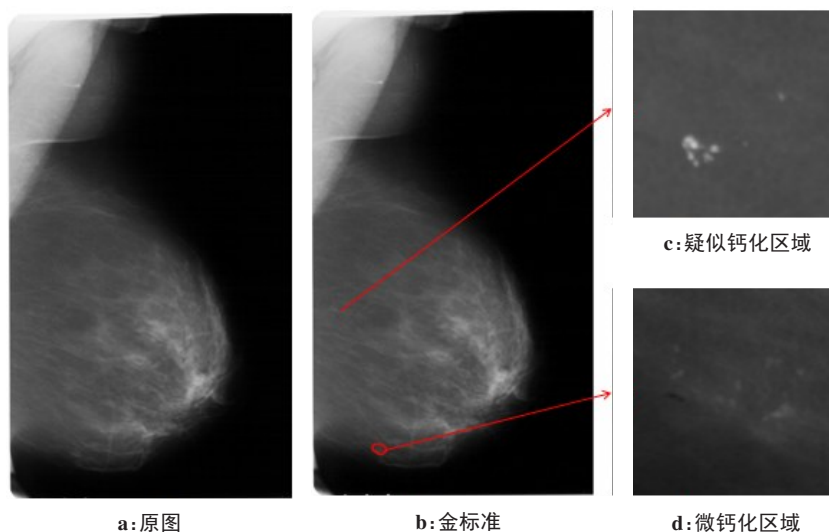


图3 乳腺钼靶图像、对应的金标准图像及其所含的疑似钙化区域和微钙化区域放大图
Fig.3 Mammography, the corresponding gold standard image and the enlarged images of suspected calcification region and microcalcification region

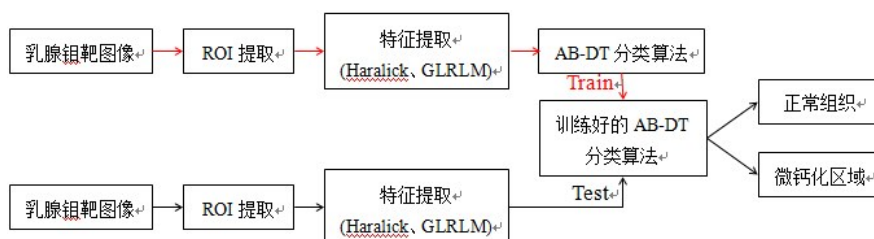


图4 乳腺微钙化区域检测方法流程图
Fig.4 Flowchart of breast microcalcifications detection

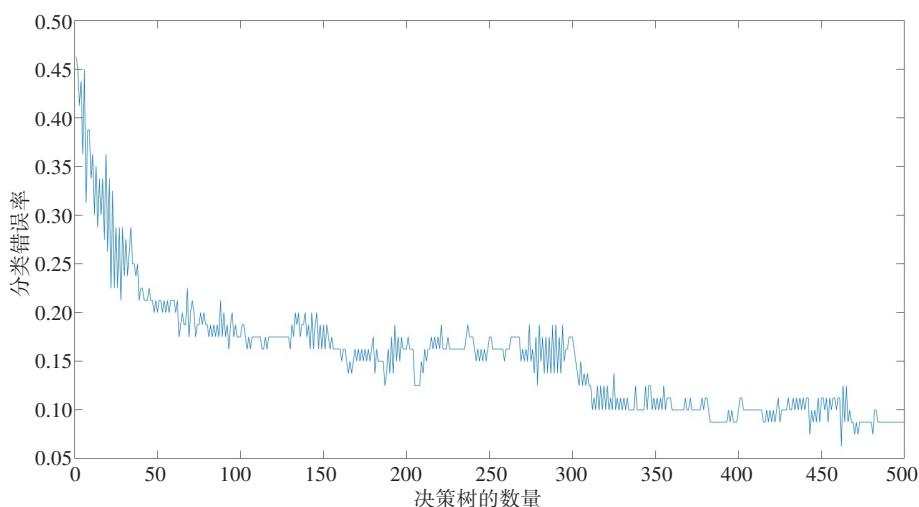


图5 决策树数量与分类错误率的关系曲线
Fig.5 Relationship between the number of decision trees and classification error rate

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \times 100\% \quad (1)$$

$$\text{Sensitivity} = \text{Recall} = \frac{TP}{TP + FN} \times 100\% \quad (2)$$

$$\text{Specificity} = \frac{TN}{FP + TN} \times 100\% \quad (3)$$

$$\text{PPV} = \text{Precision} = \frac{TP}{TP + FP} \times 100\% \quad (4)$$

$$\text{NPV} = \frac{TN}{TN + FN} \times 100\% \quad (5)$$

$$\text{F1 - score} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}} \quad (6)$$

设阳性代表微钙化区域,阴性表示非钙化区域。

式中,TP(True Positive, 真阳性)表示测试集中被正确分类的微钙化区域样本个数;FP(False Positive, 假阳性)表示测试集中被错误分类的微钙化区域样本个数;TN(True Negative, 真阴性)表示测试集中被正确分类的正常组织样本个数;FN(False Negative, 假阴性)表示测试集中被错误预测的正常组织样本个数。

3.4 实验结果及分析

采用10折交叉验证来验证所使用的AB-DT模型的性能,将样本集中所有的样本数据随机分成10组,选择其中9组数据作为训练样本,训练出分类模型,最后1组样本数据作为测试数据,验证训练的模型的准确率。

将决策树的数量设置为462,测试结果显示模型分类准确率为91.75%,敏感性为91.75%,特异性为91.79%,阳性预测值为92.35%,阴性预测值为91.56%,F1分数为0.918 7。实验结果表明此模型具有较强的学习能力和泛化能力,并且具有较高的预测精度。

为进一步验证AB-DT算法的鲁棒性与有效性,将本文算法与其他文献中所提方法进行性能比较。Cai等^[18]提出基于卷积神经网络的乳腺钙化区域分类算法,该算法在中山大学附属肿瘤医院和南方医科大学附属南海医院的数据库上取得88.59%的准确率,89.32%的精准率,86.89%的特异性和88.43%的敏感性;王科举^[3]提出基于周围区域矩阵的微钙化区域检测算法,在MIAS数据集上得到88.84%的准确率,90.00%的敏感性和88.80%的特异性;蔡雅丽等^[19]利用局部二值模式和灰度共生矩阵特征进行乳腺钙化检测,支持向量机、随机森林和Adaboost算法均可较好地区分正常样本和钙化样本,检测的准确率分别为90.0%、81.5%、87.5%,敏感性分别为83.60%、78.40%、79.30%,特异性分别为91.10%、86.00%、87.00%;Chakravarthy等^[20]提出萤火虫算法进行微钙化区域的检测,在MIAS数据库中的实验准确率为86.75%,敏感性为90.08%,特异性为83.42%,阳性预测值为86.78%,阴性预测值为90.21%,F1分数为0.874 2。比较每种算法的准确率、敏感性和特异性,本文提出的AB-DT算法在微钙化区域检测上具有更好的准确率、敏感性、特异性。

4 结 语

本研究针对微钙化点检查精度不高的问题,提出一种乳腺钼靶图像微钙化真假阳性检测的方法。首先提取ROI的Haralick纹理特征和灰度游程矩阵特征,然后结合Adaboost算法和决策树算法,构建强分类器AB-DT对区域进行分类,将微钙化区域和正常组织分离开,并通过10折交叉验证,验证了该分类方法的有效性,分类正确率高达91.75%。本研究提出的方法在辅助乳腺微钙化点检测中具有一定的临床应用价值。

【参考文献】

- [1] 彭庆涛,吴水才,高宏建,等. 基于小波分析和灰度纹理特征的乳腺X线图像微钙化点区域的提取[J]. 北京生物医学工程, 2015, 34(5): 462-467.
PENG Q T, WU S C, GAO H J, et al. Extraction of the regions of microcalcification in breast X-ray image based on Wavelet analysis and image texture feature[J]. Beijing Biomedical Engineering, 2015, 34(5): 462-467.
- [2] 商小宝. 医学图像降噪处理及计算机辅助诊断[D]. 杭州: 浙江大学, 2018.
SHANG X B. Medical image noise reduction processing and computer-aided diagnosis[D]. Hangzhou: Zhejiang University, 2018.
- [3] 王科举. 基于周围区域矩阵反映射特征的乳腺钙化点区域检测[D]. 兰州: 兰州大学, 2017.
WANG K J. Calcification regions extraction based on inverse mapping features of surrounding region matrix [D]. Lanzhou: Lanzhou University, 2017.
- [4] KARALE V A, EBENEZER J P, CHAKRABORTY J, et al. A screening CAD tool for the detection of microcalcification clusters in mammograms[J]. J Digit Imaging, 2019, 32(5): 728-745.
- [5] SUHAIL Z, DENTON E R, ZWIGGELAAR R. Classification of microcalcification in mammograms using scalable linear Fisher discriminant analysis[J]. Med Biol Eng Comput, 2018, 56(8): 1475-1485.
- [6] 郭亚南. 乳腺钼靶X线病灶检测研究[D]. 兰州: 兰州大学, 2019.
GUO Y N. Research on lesion detection for mammograms [D]. Lanzhou: Lanzhou University, 2019.
- [7] 顾广娟. 乳腺X线影像微钙化簇检测技术研究[D]. 哈尔滨: 哈尔滨理工大学, 2009.
GU G J. Research of microcalcification cluster detection technology in mammograms [D]. Harbin: Harbin University of Science and Technology, 2009.
- [8] HARALICK R M, SHANMUGAM K, DINSTEN I. Textural features for image classification[J]. Studies in Media and Communication, 1973, SMC-3(6): 610-621.
- [9] 贾立男. 数字病理图像特征提取算法的研究与实现[D]. 保定: 河北大学, 2020.
JIA L N. Research and implementation of digital pathological image feature extraction algorithm[D]. Baoding: Hebei University, 2020.
- [10] DASARATHY B V, HOLDER E B. Image characterizations based on joint gray level-run length distributions[J]. Pattern Recogn Lett, 1991, 12(8): 497-502.
- [11] 宋炎. 基于深度学习与影像组学的乳腺钙化诊断方法[D]. 广州: 华南理工大学, 2018.
SONG Y. Breast calcification classification by deep-learning and radiomics descriptors [D]. Guangzhou: South China University of Technology, 2018.
- [12] 李曦. 基于图像处理的HIFU治疗中生物组织变性识别方法研究[D]. 长沙: 湖南师范大学, 2020.
LI X. Study on the recognition of denatured biological tissue during HIFU therapy based on image processing [D]. Changsha: Hunan Normal University, 2020.
- [13] SCHAPIRE R E. A brief introduction to boosting [C]//Sixteenth International Joint Conference on Artificial Intelligence. Morgan Kaufmann Publishers Inc, 1999: 1401-1406.
- [14] 余勇. 基于多智能体样本交换的分散式集成学习方法研究[D]. 重庆: 西南大学, 2020.
YU Y. Research on decentralized ensemble learning based on sample exchange among multiple agents [D]. Chongqing: Southwest University, 2020.
- [15] PANCHALKAR A R, DOYE D D. A novel approach to build accurate and diverse decision tree forest[J]. Evolutionary Intell, 2021(1): 1-15.
- [16] 刘万安. 基于深度学习算法的高原地区云雪分类[D]. 南京: 南京信息工程大学, 2019.
LIU W A. Cloud and snow classification in plateau area based on deep learning algorithms [D]. Nanjing: Nanjing University of Information Science and Technology, 2019.
- [17] HEATH M, BOWYER K, KOPANS D, et al. Current status of the digital database for screening mammography [J]. Digit Mammography, 1998: 457-460.
- [18] CAI H, HUANG Q, RONG W, et al. Breast microcalcification diagnosis using deep convolutional neural network from digital mammograms[J]. Comput Math Meth Med, 2019, 2019: 1-10.
- [19] 蔡雅丽, 蔡盛, 施敏敏, 等. 计算机辅助系统在乳腺钙化性病变X线摄影诊断中的应用[J]. 中国医学影像学杂志, 2019, 27(12): 910-913.
CAI Y L, CAI S, SHI M M, et al. Application of computer aided system in X-ray diagnosis of breast calcification lesions [J]. Chinese Journal of Medical Imaging, 2019, 27(12): 910-913.
- [20] CHAKRAVARTHY S R, RAJAGURU H. Detection and classification of microcalcification from digital mammograms with firefly algorithm, extreme learning machine and non-linear regression models: a comparison[J]. Int J Imaging Sys Technol, 2019, 30(2): 126-146.

(编辑:谭斯允)