

智慧基建视角下钢材缺陷分类的机器学习算法应用研究

何海玉^{1,2}, 彭友海¹, 朱禹林^{1*}, 盧葦麟¹

¹香港珠海学院计算机系, 香港

²中铁十四局集团有限公司, 山东 济南

收稿日期: 2025年4月25日; 录用日期: 2025年5月23日; 发布日期: 2025年5月31日

摘要

本文基于美国国家标准与技术研究所和钢铁工业协会的UCI带钢缺陷公开数据库, 对智慧基建中的钢材缺陷分类问题展开研究。样本来源于北美3家钢厂2018~2021年的1941个样本, 共有27类特征。选择6种算法并采取数据预处理、特征工程以及超参数调整策略来建立高效的精准钢材缺陷智能分类方案。创新点包括: 搭建多种算法融合模型; 设计特征分类筛选与调优方案; 采用SMOTE解决不均衡样本问题; 设置完整的试验评价系统。结果表明, LightGBM和神经网络的精确率和召回率均超过96%。消融实验与参数敏感性分析证明了这些方法对于特征选取和超参数的重要性。后续的研究将扩大收集的样本数量, 并尝试结合深度学习和计算机视觉等新技术, 使模型更具有普适性、鲁棒性和更高的检测精度, 促进智慧基建更广泛的智能化发展。

关键词

智慧基建, 钢材缺陷检测, 机器学习算法, 数据预处理

Research on the Application of Machine Learning Algorithms for Steel Defect Classification from the Perspective of Smart Infrastructure

Haiyu He^{1,2}, Youhai Peng¹, Yulin Zhu^{1*}, Weilun Lu¹

¹Department of Computer Science, Hong Kong Chu Hai College, Hong Kong

²China Railway 14th Bureau Group Corporation Limited, Jinan Shandong

*通讯作者。

文章引用: 何海玉, 彭友海, 朱禹林, 盧葦麟. 智慧基建视角下钢材缺陷分类的机器学习算法应用研究[J]. 人工智能与机器人研究, 2025, 14(3): 742-753. DOI: 10.12677/airr.2025.143072

Abstract

This paper investigates the problem of steel defects classification in smart infrastructure based on the UCI Strip Steel Defects Public Database of the National Institute of Standards and Technology and the Steel Industry Association. The samples are derived from 1941 samples from 3 North American steel mills from 2018~2021, with a total of 27 types of features. Six algorithms are selected and data preprocessing, feature engineering, and hyper-parameter tuning strategies are adopted to establish an efficient and accurate intelligent classification scheme for steel defects. The innovations include: building a fusion model of multiple algorithms; designing a feature classification screening and tuning scheme; adopting SMOTE to solve the problem of unbalanced samples; and setting up a complete experimental evaluation system. The results show that the precision and recall of LightGBM and neural network are more than 96%. Ablation experiments and parameter sensitivity analysis demonstrate the importance of these methods for feature selection and hyperparameterization. Subsequent research will expand the number of collected samples and try to combine new technologies such as deep learning and computer vision to make the model more pervasive, robust, and higher detection accuracy, and promote the intelligent development of smart infrastructure more widely.

Keywords

Smart Infrastructure, Steel Defect Detection, Machine Learning Algorithm, Data Preprocessing

Copyright © 2025 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

随着智慧基建浪潮的到来，钢材作为基础设施建设中的重要材料之一，其自身质量的好坏直接影响整个工程的稳固性以及耐久性。然而在繁杂的制造过程中，钢材极易产生划伤、污渍、边缘不整齐等问题，这不仅影响了钢材本身的美观度，同时也降低了钢材的部分机械性能，存在安全隐患。现阶段，对于钢材质检的检测方式多采用最基本的人工与视觉检测方式，工作量大，效率低，并且容易出现漏检误判的情况，无法满足当下智慧基建对钢材质检的要求。近年来，伴随着人工智能的发展，机器学习的应用也为钢材缺陷检测带来了新的思路：通过不同的机器学习算法从复杂的数据中寻找规律，最终实现自动化的钢材缺陷识别和分类，减少人力，提高准确率，达到自动化的效果[1]。本课题的研究目的就是运用机器学习算法对钢铁质检行业常见的钢材缺陷类型进行智能分类。根据 UCI 机试库所提供的有关钢材缺陷数据，比较常见的机器学习算法对钢材缺陷识别的适用程度，寻找到更有效的机器学习算法来辅助完成钢材缺陷智能分类，进而推进未来智慧基建工程建设过程中的钢材智能质检工作。

2. 数据预处理与特征工程

2.1. 数据集获取与分析

本研究基于 NIST 与钢铁工业协会发布的 UCI 带钢缺陷公共数据集展开，该数据集采集自北美三家大型钢铁厂 2018~2021 年间的连铸生产线，涵盖 7 类典型表面缺陷的工业检测记录，包含 1941 个有效样

本。每个样本依据 27 种特征描述，全面涵盖钢板的几何属性、物理特性及图像相关特征，为缺陷分类提供丰富信息。研究目标是将样本精确划分到 7 类常见钢材缺陷类型中[2]。

经探索性分析发现，样本分布显著不均衡。“肮脏”类样本最多，达 680 个，占比 35%；“其他故障”516 个，占比约 26.6%；“K 型划痕”300 个，占比 15.5%；“颠簸”190 个，占比约 9.8%；“装饰缺陷”、“Z 型划痕”和“污渍”样本数量较少，分别为 80、100 和 75 个，占比 4%~5.2%。这种不均衡分布会使模型对多数类识别准确率高，但对少数类识别能力弱，影响整体分类性能。因此，后续数据预处理阶段采用 SMOTE 方法解决样本不均衡问题。

数据集整体质量较好，没有出现明显的缺失值，并且便于后期的数据处理。但在对各个特征的最大值、最小值等研究后发现，在某些特征上有些样本的值偏大或者偏小，明显高于大多数样本的情况，这种特例很有可能是特殊情况下的特殊测量值或是属于很小概率情况下产生的特例，比如有些“厚度”值偏大的几个样品会对模型对该特征的权重选取产生很大的影响。我们准备在之后的数据预处理过程中用 IQR 算法将其处理掉后再重新构建机器学习训练使用的特征组。通过对各个特征之间两两比较相关系数得出的矩阵可以看出：其中的一些特征的相关度很大，比如之前所提到过的“面积”与“长度”、“宽度”的相关度都在 0.8 左右。所以我们在之后的特征筛选过程中要将其去除多余特征。

2.2. 特征工程策略

2.2.1. 特征分类与筛选

基于物理意义和数据类型，将 27 个原始特征划分为四类：

1) 几何特征：描述钢板的物理尺寸与形状(如长度、宽度、面积)，用于捕捉与尺寸相关的缺陷(如变形、边缘不规则)。例如，当钢板出现变形时，其长度、宽度或面积等几何参数会发生改变；边缘不规则也可通过几何形状的偏离标准来识别。为检测钢板宏观尺寸和形状方面的缺陷提供关键依据，在钢板生产质量把控中，可及时发现因尺寸偏差导致的不合格产品，减少后续加工因尺寸问题造成的损失，保障产品在规格上符合生产要求。

2) 物理特征：反映材料内在属性(如厚度、反射率、张力)，关联加工过程中的力学行为与缺陷生成机制。比如，厚度不均匀可能影响钢板后续的冲压、弯折等加工性能；反射率异常可能暗示材料内部存在杂质或微观结构变化；张力的波动与缺陷生成机制相关联，可辅助判断加工过程中是否存在应力集中等问题。有助于深入理解钢材加工过程中的力学行为，从材料内在属性角度关联缺陷生成原因，为优化加工工艺、预防缺陷产生提供理论支持，提升钢材产品的质量稳定性和性能可靠性。

3) 表面特性特征：量化表面微观质量(如边缘平滑性、对称性)，直接关联表面缺陷(如划痕、凹坑)。在钢材表面质量检测中起着关键作用，直接关系到钢材表面外观质量和后续涂层、镀层等表面处理工艺的效果，确保钢材表面质量符合使用要求，提高产品的外观品质和市场竞争力。

4) 统计特征：从图像数据提取的纹理指标(如灰度共生矩阵的对比度、均匀性)，用于表征缺陷的视觉模式[3]。不同类型的缺陷在图像上会呈现出不同的纹理特征，通过分析这些统计特征，可对缺陷进行分类和识别。利用图像纹理信息实现对缺陷的视觉模式识别，为基于图像处理技术的钢材缺陷检测提供重要手段，提高缺陷检测的自动化和智能化水平，在无损检测领域具有重要应用价值。

目标变量是多分类标签，表示 7 种钢板缺陷类型分为：装饰缺陷(Pastry)、Z 型划痕(Z_Scratch)、K 型划痕(K_Scratch)、污渍(Stains)、肮脏(Dirtiness)、颠簸(Bumps)、其他故障(Other_Faults)。

对于特征筛选，采用两阶段筛选策略：

1) 相关性分析

计算特征与目标变量的 Pearson 相关系数，保留绝对值 > 0.3 的 18 个特征(如表面粗糙度、厚度等)。

相关性分析能够快速衡量特征与目标变量之间的线性关联强度。选择绝对值大于 0.3 作为筛选阈值，是因为相关系数绝对值在这个范围以上，意味着特征与目标变量有相对较强的线性关系，对目标变量的影响较大。保留这些特征有助于减少冗余信息，使后续的模型训练集中在对目标变量有显著影响的关键特征上，提升模型训练效率和性能。

2) 领域知识验证

对相关性较低但理论重要的特征(如化学成分含量, $r \approx 0.25$), 依据材料科学理论保留, 后续验证其有效性。虽然 Pearson 相关系数从线性关系角度衡量特征重要性, 但某些特征可能与目标变量存在非线性关系, 或者在材料科学领域中, 基于专业理论知识判断其对钢材缺陷形成机制、缺陷类型等有潜在影响, 即便相关性系数计算结果不高, 也不能轻易舍弃。保留这些特征并后续验证, 能避免因单纯依赖相关性分析而遗漏重要信息, 确保特征筛选过程既考虑了数据统计特性, 又结合了专业领域知识, 使特征集更加完整、合理。

2.2.2. 特征的类型编码与标准化

特征集确定之后, 针对不同的特征进行有针对性的编码和标准化, 使其可以适应机器学习模型所需要的输入格式。

1) 连续型特征: 长度、宽度、厚度等连续性特征用标准化方式处理, 缩放到平均值为 0、标准差为 1 的范围。标准化计算公式如下:

$$X_{std} = \frac{X - \mu}{\sigma} \quad (1)$$

其中, X_{std} 表示标准化后的特征值, X 为原始特征值, μ 和 σ 分别为该特征在训练集中的均值和标准差。这种处理方式可以消除不同特征之间由于量纲差异导致的权重偏差, 使模型在训练过程中对各个特征的权重分配更加合理, 提高模型的收敛速度和性能。像长度可能以米为单位, 厚度以毫米为单位, 不同特征的量纲(单位)不同。若不进行标准化, 在模型计算中, 量纲大的特征(比如长度数值可能较大)对结果的影响会比量纲小的特征(如厚度数值相对较小)大很多, 这就导致模型在学习时对不同特征的权重分配不合理。例如, 在计算距离等相关指标时, 量纲差异会使计算结果偏向量纲大的特征, 影响模型对真实关系的学习。

```

1 from sklearn.model_selection import train_test_split, cross_val_score
2 from sklearn.preprocessing import StandardScaler
3 from imblearn.over_sampling import SMOTE
4 import warnings
5 import pandas as pd
6 import numpy as np
7 from sklearn import metrics # 模型评价
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.ensemble import RandomForestClassifier
10 from lightgbm import LGBMClassifier
11 from sklearn.svm import SVC
12 from sklearn.decomposition import TruncatedSVD
13 from sklearn.svm import LinearSVC
14 from sklearn.decomposition import PCA
15 from sklearn.neural_network import MLPClassifier
16 from sklearn.tree import DecisionTreeClassifier
17 import matplotlib.pyplot as plt
18 import seaborn as sns
19 sns.set_theme(style="whitegrid")
20 import plotly.graph_objects as go
21 from plotly.subplots import make_subplots
22
23 warnings.filterwarnings('ignore')

```

Figure 1. Integer label

图 1. 整数标签

2) 离散型特征: 部分二值特征(如是否存在某种表面特性)直接保持 0/1 编码, 作为模型的输入。这种简单的编码方式既保留了特征的原始含义, 又方便模型进行处理。

3) 标签的类型编码: 将目标变量(即钢材缺陷类型)由文本标签转换为整数标签。我们创建了一个映射字典, 将每个类别名称(如 “Pastry”、“Z_Scratch” 等)映射为唯一的整数编码(如 0、1 等), 以便机器学习算法能够正确地识别和处理分类任务。详见图 1。

2.2.3. 特征降维与组合

为了进一步优化特征集, 提高模型的性能和训练效率, 我们采用了特征降维和组合策略[4]。

1) 特征降维: 考虑到在相关性分析中发现的一些特征之间存在较高的相关性, 我们使用主成分分析 (PCA)方法对这些特征进行降维处理。例如, 对于“长度”、“宽度”和“面积”这三个高度相关的几何特征, 我们通过 PCA 将它们转换为两个主成分, 这两个主成分能够解释原始三个特征约 95%的方差信息。这样既保留了绝大部分的信息, 又减少了特征的数量, 降低了模型的复杂度。

2) 特征组合: 在某些情况下, 单独的特征可能无法充分反映钢材缺陷的本质信息, 而多个特征的组合可能会产生更有价值的特征。例如, 我们将“厚度”和“张力”这两个物理特性特征组合成一个新特征“厚度 - 张力比”, 因为从材料力学的角度来看, 这个比值可能与钢板在加工过程中产生的内应力和变形情况有关, 进而与某些特定的缺陷类型(如“K 型划痕”)产生关联。通过这种方式, 我们构造了 3 个新的组合特征, 丰富了特征集的表达力。

2.2.4. 单变量分析

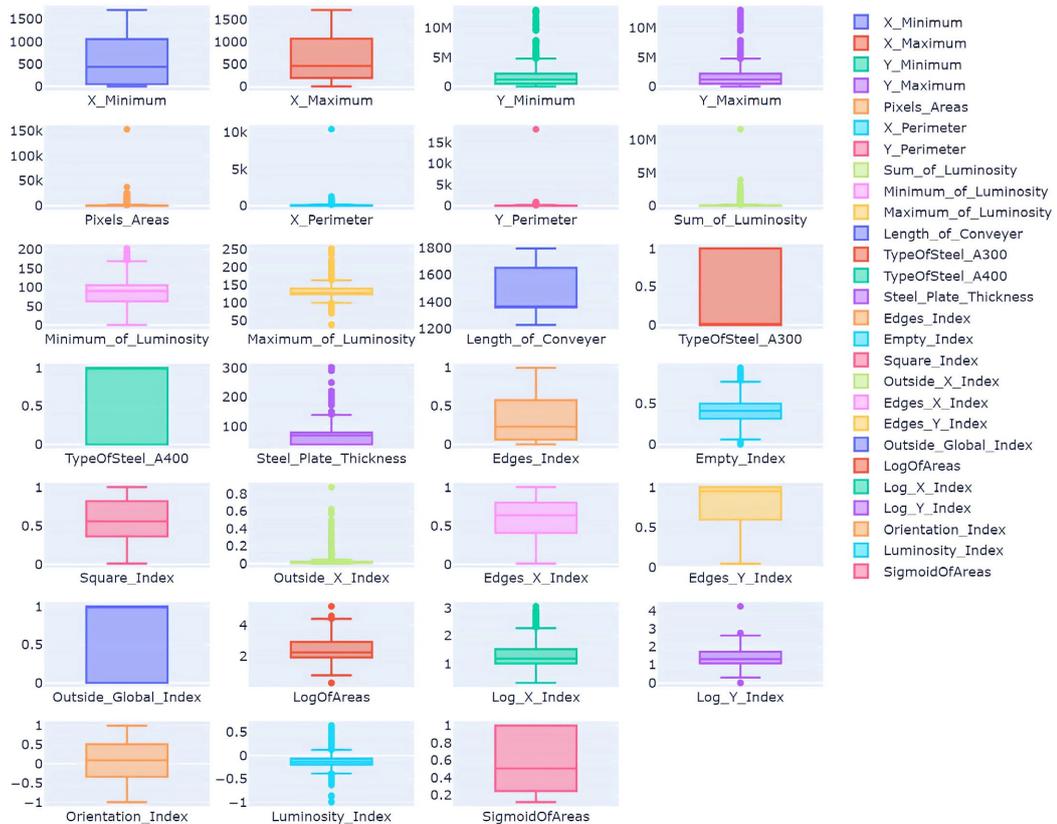


Figure 2. Box diagram
图 2. 箱型图

在数据分析过程中,我们重点关注各特征的统计特性,如均值、方差、偏度、峰度等。为直观呈现这些统计特性,选择利用箱线图来观察每个连续型特征的分布情况,尤其是 df2 数据集中“Steel_Plate_Thickness”特征的分布。从箱型图中能观察到单个特征的取值分布情况,下面绘制全部参数的取值分布箱型图,见图 2。

在异常值检测上,箱线图发挥了重要作用,它使我们发现诸如“厚度”这类特征的异常值已超出四分位数范围。针对这一情况,遂果断采用 IQR (四分位距)方法加以处理,让数据在统计分析中更具可靠性。

3. 机器学习算法的架构设计

3.1. 算法选取与原理阐述

针对智慧基建下的钢材缺陷分类任务,对算法要求较高。经过考虑各种机器学习算法的特点和应用场景后,决定采用应用逻辑回归、决策树、随机森林、SVM、LightGBM 和神经网络这六种具有代表性的算法建立钢材缺陷分类模型。通过比较各种机器学习算法的优劣,最终选择适合智慧基建下钢材缺陷分类的最优算法方案。

逻辑回归是一种广泛应用于分类任务的线性模型。尽管其名为“回归”,但通过 sigmoid 函数将线性回归结果映射到(0, 1)区间后,能够有效处理二分类问题,并可拓展至多分类场景。逻辑回归模型具有形式简单、易于理解和实现的优点,其核心在于求解以下优化问题:

$$\min_{\theta} -\frac{1}{m} \sum_{i=1}^m [y_i \ln h_{\theta}(x_i) + (1 - y_i) \ln (1 - h_{\theta}(x_i))] \quad (2)$$

其中, $h_{\theta}(x_i) = \frac{1}{1 + e^{-\theta^T x_i}}$, m 为样本数量, y_i 为样本标签。这一优化问题可通过梯度下降等优化算法进行求解。逻辑回归的优势在于其可解释性强,能够直观地体现特征对分类结果的影响,为智慧基建中钢材缺陷分类提供一种易于理解和维护的解决方案。然而,由于其本质上是线性模型,在处理复杂的非线性分类问题时可能存在局限性。

决策树模型以其直观的树形结构和强大的非线性拟合能力,在分类任务中得到了广泛应用。它通过不断地根据特征对样本进行分割,形成一系列的“如果-那么”规则。决策树的构建通常基于信息增益、信息增益率或基尼系数等指标来选择最优的分裂特征。例如,基尼系数的计算公式为:

$$G(D) = 1 - \sum_{k=1}^K \left(\frac{|C_k|}{|D|} \right)^2 \quad (3)$$

其中, D 表示数据集, C_k 表示第 k 类样本的集合, K 为类别总数。决策树的分裂过程旨在最大化分裂后的纯度,即最小化基尼系数。该算法能够处理非线性关系,且对数据的适应性强,但在处理复杂数据集时容易出现过拟合现象,导致模型的泛化能力下降。

随机森林作为一种集成学习算法,通过构建多个决策树并综合它们的预测结果来提升模型的性能。每棵决策树的训练基于对原始数据集的有放回采样(即 Bootstrap 抽样),并且在每个节点的分裂过程中,仅随机选择一部分特征进行评估。这种集成策略不仅提高了模型的准确性,还增强了其泛化能力和对抗噪声的能力。具体来说,随机森林的输出是所有决策树预测结果的多数投票或平均值。这种算法在处理高维数据和复杂分类任务时表现出色,且对数据中的异常值和缺失值具有较强的鲁棒性,但其训练过程相对耗时,且模型的可解释性相较于单棵决策树有所降低。

支持向量机(SVM)是一种基于统计学习理论的监督学习算法,其核心思想是通过寻找一个最优超平面,将不同类别的样本分开。对于线性可分问题,SVM的目标是最小化以下目标函数:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \quad (4)$$

$$\text{约束条件为: } y_i(w^T x_i + b) \geq 1, i = 1, 2, \dots, m$$

其中, w 和 b 分别是超平面的法向量和偏置项, m 为样本数量, y_i 为样本标签。对于非线性问题,SVM通过核函数(如多项式核、径向基核等)将数据映射到高维空间,使其在高维空间中线性可分。SVM的优势在于在处理数据的时候可以很好地解决维度过高的情况,并且使用了核技巧来应对复杂的非线性问题,但是它也有着计算时间过长以及解释程度弱等缺点。

LightGBM 是一种基于梯度提升决策树(GBDT)的高效算法,由微软研发,采用直方图算法和基于叶节点的分裂策略,能够显著提高模型的训练速度和效率。它是通过对原始 GBDT 在执行过程中的一些不完善的地方进行优化,用更快的直方图算法代替二元搜索找到特征的最佳值,减少了特征值的分裂点查找的时间;用更简单的叶生长方式(代价会稍微大一点)来代替老版本 GBDT 中更费时的分裂方式;解决了类不平衡的问题并且对大部分数据都适用。LightGBM 比较适用于处理大数据以及有大量变量的稀疏数据。LightGBM 有着较高的准确率和较快的模型训练速度及内存占用率,但缺点也很明显:它的参数比较繁琐需要大量的调整。

神经网络是模仿人脑的神经元结构和信息传送的方法而建立的一种网络,非常适用于拟合非线性模型。本论文所使用的 MLP 是 1 输入层 + 2 隐藏层 + 1 输出层的三层多层感知机。将输入数据经过处理之后送入输入层,中间各层使用 Relu 函数,输出端则用 Softmax 来获得多个分类的概率;神经网络通过反向传参、梯度下降来不断地迭代寻找到最优权重,使它与真实值的误差达到最小,需要依靠大量的数据样本进行训练,还需要选取合适的网络拓扑参数等等。它的优点是能够自动学习数据特征,在错综复杂的大数据集中发现内在规律,适合对非线性的数据进行回归或分类等问题,但是其缺点也十分明显:训练会受到噪声等扰动数据的影响,容易发生拟合现象。

3.2. 模型细节与优化设置

1) 针对逻辑回归模型,选择梯度下降法求解参数,并采用交叉熵损失函数作为优化目标;正则化系数 C 设置为 1.0,保证模型具有良好的稳定性和收敛速度。

2) 针对决策树模型,选取基尼系数为分裂标准,设置最大深度为 10,最小样本分裂个数为 5,防止过拟合。

3) 针对随机森林模型,建立 100 棵决策树,使用 Bagging 方式采样,每棵树分裂的特征数量设置成所有特征的平方根,减小单棵决策树的方差,提高准确度[5]。

4) 针对 SVM 算法,使用 RBF 函数解决非线性问题, C 的值设置为 1.0, γ 值设置为 0.1,对输入数据进行标准化处理。

5) LightGBM 使用叶子节点增长策略建立决策树,通过梯度提升的方法来优化目标函数,学习率为 0.1,叶子节点数为 31,使用该库内自带的处理类别不平衡问题的方式,有利于分类少数类样本。

6) 针对神经网络,创建含有输入层和两层隐藏层以及输出层组成的 MLP,其中每层隐藏层分别为 64 和 32 个神经元,并且设置每个神经元的激活函数为 ReLU 函数,最后设置一个 Softmax 函数输出层完成多分类过程,设置 Adam 更新方式改变各个权值,学习率为 0.001,设置训练批次数为 32。其他细节上做相应的优化和调整,使得算法适应了钢材缺陷分类的要求,在智慧基建中起到钢材缺陷检测作用,有

利于保障工程的质量、安全及效率，下一章会继续展示其效果。

4. 实验设计与结果分析

4.1. 实验设置与流程

在实验设计方面，我们首先准备好一个 Python 运行环境，利用 scikit-learn、TensorFlow 等来搭建机器学习与深度学习的模型。在训练使用的硬件方面选择的是配备有英特尔酷睿 i7 处理器的个人电脑，内存大小为 16 G，以此来提升处理各类算法的速度和效率。而数据集的划分则按照 7:2:1 的比例进行随机划分，以便分别用来训练、验证及测试，并且确保所有的实验都有较高的统计学意义[6]。因此，我们将每一个模型都独立地训练了 5 次，在测试后求平均数作为最终模型的结果指标。然后依次按照一般机器学习项目的流程去执行：先导入并预处理好的数据集；再将该数据集传入到模型中，并初始化好相关的变量；紧接着就是模型的训练，然后在模型训练的过程中通过测试对模型参数进行调整以及超参数的调优；最后是当完成模型训练后输出模型性能结果。在模型输出的性能结果中包含了准确率、召回率、F1 值等评价矩阵。每一次循环都会根据模型输入的训练数据进行学习，在测试时用到的是训练集中的训练数据，在预测时用到的是验证集上的数据，而在得到模型预测结果的数值上使用的是测试集中的测试数据。

4.2. 结果呈现与对比

从试验结果来看，6 种算法对钢材缺陷分类的效果各有千秋。其中，LightGBM 和神经网络效果最好，精确率、召回率均超过了 96%，对于少类的识别精度较高；随机森林整体呈现一种比较平稳的状态，在大部分情况都能取得一个相对平衡的结果；逻辑回归和决策树相对来说效果差了一些，在复杂程度较强的模型之间线性关系没有前面那么好的时候没有那么好。LightGBM 在混淆矩阵中把“肮脏”、“Z 型划痕”等多数类别的分类准确度非常高，只是在这几种错误里面占了很少一部分。如把“污渍”的几个样例错分成了其他故障，这是我们后期需要继续加强对于类似问题部位的训练，提高特征的区别能力。六种机器学习方法具体性能指标数据见表 1。

Table 1. Performance metrics for learning algorithms

表 1. 学习算法的性能指标

算法	精确率	召回率	F1 值	训练时间(s)
逻辑回归	0.89	0.87	0.88	2.3
决策树	0.85	0.83	0.84	1.8
随机森林	0.93	0.91	0.92	15.6
SVM	0.90	0.88	0.89	32.4
LightGBM	0.97	0.96	0.96	5.7
神经网络	0.96	0.96	0.96	28.3

在钢材缺陷分类研究中，对比 6 种机器学习算法的性能指标(精确率、召回率、F1 值、训练时间)，可以发现 LightGBM 和神经网络在精确率和召回率上表现突出，均超 96%，但也存在参数调整繁琐、易过拟合等问题。以下从多个方面对这些模型的优势与局限性展开深入探讨。

1) 逻辑回归

优势：形式简单、可解释性强，能直观体现特征对分类结果的影响，便于理解和维护，计算复杂度

低，训练速度快，在处理线性可分问题时效率高。

局限性：本质是线性模型，难以处理复杂非线性分类问题，在钢材缺陷数据特征复杂、非线性关系多的情况下，分类性能受限，对复杂缺陷类型识别能力弱。

2) 决策树

优势：树形结构直观，构建基于信息增益等指标，能处理非线性关系，对数据适应性强，可自动处理特征间的交互作用，无需对数据进行复杂预处理。

局限性：处理复杂数据集易过拟合，导致泛化能力下降，对噪声数据敏感，可能因噪声产生错误分裂，影响模型准确性，且不适合处理高维数据，计算量会随维度增加而剧增。

3) 随机森林

优势：通过集成多个决策树，综合预测提升准确性，泛化能力和抗噪声能力强，对高维数据和复杂分类任务表现出色，对异常值和缺失值有较强的鲁棒性。

局限性：训练过程相对耗时，构建多棵决策树需大量计算资源和时间，模型可解释性较单棵决策树降低，难以直观理解每个特征对分类结果的具体影响。

4) 支持向量机(SVM)

优势：能有效解决维度过高问题，通过核函数将数据映射到高维空间处理非线性问题，在小样本、非线性及高维模式识别中表现良好，理论完善，有坚实数学基础。

局限性：计算时间长，尤其是处理大规模数据集时，核函数参数选择和调优困难，不同核函数及参数设置对模型性能影响大，解释程度弱，难以直观理解分类决策过程。

5) LightGBM

优势：基于梯度提升决策树优化，采用直方图算法和叶节点分裂策略，训练速度快、效率高，内存占用率低，适合处理大数据和大量变量的稀疏数据，准确率高，对类别不平衡问题有较好处理能力。

局限性：参数繁琐，需大量调整才能达到最佳性能，调参过程复杂且耗时，对使用者专业知识和经验要求高，模型可解释性相对较弱，难以清晰解释决策过程。

6) 神经网络(MLP)

优势：能自动学习数据特征，在复杂大数据集中发现内在规律，适合处理非线性数据的回归和分类问题，具有强大的非线性拟合能力，理论上可逼近任何连续函数。

局限性：训练受噪声等扰动数据影响大，易过拟合，需大量数据样本训练，数据不足时性能不佳，网络拓扑参数选取困难，不同参数设置对模型性能影响大，训练时间长，计算资源消耗大。

4.3. 参数敏感性分析

为了评估模型对关键参数的敏感性，分别对随机森林、SVM 和 LightGBM 进行了调参实验。

对于随机森林，树的数量($n_estimators$)是一个关键参数。实验发现，树数从 10 增加到 100 时，模型的精确率和召回率稳步提升，超过 100 后并没有明显增加，并且训练时间花费也会更多，因此可设置为 100。

在对 SVM 进行参数敏感性分析时，重点考察了惩罚系数 C 和核函数参数 γ 。较大的 C 值虽然提高了模型对训练数据的拟合程度，但也容易导致过拟合。通过实验，将 C 设置为 1.0、 γ 设置为 0.1 时，模型在准确率和泛化能力之间取得了最佳平衡。

LightGBM 对于学习率来说，它控制的是每一颗 XGBOOST 树的权重变化幅度；数值越小说明最后收敛就会越精准，但是也会更慢。叶节点数目也是同样的道理，设置学习率为 0.1，叶节点数目为 31 较好。

三种算法的关键参数敏感性分析结果见表 2。

Table 2. Results of sensitivity analysis of key parameters
表 2. 关键参数敏感性分析结果

算法	参数	取值范围	最优值	对应精确率	对应召回率
随机森林	n_estimators	10~200	100	0.93	0.91
SVM	C	0.1~10	1.0	0.90	0.88
SVM	gamma	0.01~1	0.1	0.90	0.88
LightGBM	学习率	0.01~0.3	0.1	0.97	0.96
LightGBM	叶节点数	10~50	31	0.97	0.96

4.4. 消融实验

为验证不同特征对模型性能的影响,设计了消融实验,分别移除几何特征、表面特性特征和统计特征。

实验结果显示,移除几何特征后(长度、宽度、面积),模型性能平均下降约 5%,说明几何特征对分类有一定贡献,尤其是在区分与形状相关的缺陷类型时较为重要;移除表面特性特征(边缘平滑性),性能下降显著,平均达 12%,表明表面特性如边缘平滑性、对称性等是模型识别钢材缺陷的关键依据,对分类任务影响最大;移除统计特征(灰度对比度),性能下降约 8%,反映出统计特征在细分缺陷类别、提升模型精准度方面有一定价值。

从实验数据可以看出,表面特性在钢材缺陷分类中起核心作用,几何特征和统计特征也分别从不同角度为模型提供重要信息。这一结论凸显了特征工程中筛选和保留关键特征的重要性,后续实际应用中应重点确保表面特性数据的质量和完整性,同时综合考虑几何与统计特征,以维持模型对各类钢材缺陷的高效识别能力。

5. 智慧基建的应用前景与价值

5.1. 实际应用场景展望

智慧基建领域中,钢材缺陷分类技术的智能化应用前景广阔。在建桥、建筑工地施工等场景里,钢材是关键材料,其质量直接关乎工程的成败。对于建筑工地而言,现场布置钢材检测仪器,让工作人员手持设备就能检测钢材缺陷,一旦发现缺陷可立刻更换新材料,既能保障工期,又能预防结构安全隐患。桥梁结构健康监测方面,桥梁作为交通建设的基石,其钢材长期处于露天环境,易受腐蚀和疲劳作用影响。利用机器学习算法,配合高清摄像头与传感器实时收集钢材表面图像并检测分类,能及时发现损伤并修复,节省维修费用、延长使用寿命、确保安全运营。工业园区厂房装配时,钢材质量影响着厂房结构安全与使用寿命。智能检测系统可确保钢材符合质量标准,避免缺陷带来的安全隐患。机器学习算法赋能的钢结构装配式生产智能检测系统,与生产流水线结合,能快速检测、分拣钢材,反馈生产问题,提升生产过程的智能化、自动化水平。在钢材仓储与物流管理环节,智能检测系统可自动化检测、分类钢材,确保入库钢材质量,避免缺陷钢材流通引发风险,提高仓储物流效率,降低管理成本。

5.2. 经济与社会效益分析

从经济效益方面看,过去人工检测费时费力还容易出错,现在智能检测设备能快速批量检测钢材,大大节省了时间和人力。比如在中型建筑项目里,原本需要几名钢结构专业检测员花几天才能完成的钢

材缺陷检测工作，用上智能设备后，几个小时就能搞定。而且精准的检测能让企业在钢材刚出厂时就发现问题，及时处理，降低因不合格产品带来的损耗和废品率，提高材料利用率，减少成本。这对大规模基建行业来说是个大优势，能帮企业省下不少开支。

从社会效益方面看，精准检测钢材缺陷对保障工程质量、确保安全很重要。建筑、桥梁这些工程的质量和钢材息息相关。要是钢材有严重问题，很可能会引发安全事故。所以用智能化手段检测钢材，提前发现问题并解决，能避免后续施工中出现事故。拿桥梁建设来说，用智能化检测系统可以发现桥体钢板表面的微小裂痕，防止裂痕扩大引发大桥垮塌，保障人们出行安全。这样做还能延长基础设施的使用寿命，减少后期维修和更换的次数，降低社会成本，对社会发展很有帮助。

5.3. 未来研究方向

未来研究将致力于提升钢材缺陷分类技术的广度与深度，全力增强模型泛化能力。一方面，积极拓展数据收集范围，突破地域与材料种类的限制。当前数据虽源自北美三家钢厂的 1941 个样本，涵盖七类典型表面缺陷，但不足以勾勒全球钢材缺陷全貌。不同地区因生产工艺、原料质量差异，钢材缺陷特征千差万别。纳入欧洲、亚洲等地钢厂数据，以及合金钢、碳素钢等不同钢种和罕见缺陷数据，将极大丰富模型的学习素材，使其更能适应复杂多变的实际情况，全方位提升识别能力。

在深度拓展上，深度学习与计算机视觉技术将成为研究重点。深度学习在图像识别领域已取得瞩目成就，计算机视觉技术更是能直接处理钢材表面图像。其中，卷积神经网络(CNN)凭借多层卷积和池化操作，可自动提取钢材表面缺陷的纹理、形状特征，摆脱传统机器学习算法手工提取特征的束缚，检测精度大幅跃升。同时，图像增强、去噪等预处理技术为模型输送更优质图像数据，进一步夯实精准检测的根基。不仅如此，循环神经网络(RNN)及其变体(LSTM, GRU)的独特优势也将被挖掘，借助其对时间序列数据的强大处理能力，分析钢材生产环节连续监测数据，有望实现缺陷的超前预测与早期发现，为智慧基建筑牢防线。

且看物联网技术，正徐徐拉开实时监测系统的建设帷幕。在钢材的生产、仓储、使用全流程中部署传感器与摄像头，表面图像、温度、应力等关键数据即刻被捕获，并借物联网之翼疾速飞抵云端服务器。在那里，机器学习与深度学习模型实时开启分析处理，一旦发现钢材缺陷，警报瞬时拉响，关键信息火速送达。建筑工人手持设备接收检测反馈，得以迅速更换问题钢材；桥梁健康监测系统则能抢占先机，察觉损伤并及时修复，为桥梁安全运营保驾护航。而实时监测系统与企业生产管理系统的深度融合，将彻底革新生产管理模式，智能化生产管理不再是遥不可及的梦想，生产效率与质量控制水平将实现质的飞跃，智慧基建的智能化宏图也将徐徐铺展至全新高度。

6. 结论

在本研究中，我们针对智慧基建领域中的钢材缺陷分类问题，系统探索了机器学习算法的应用，取得了一系列创新成果。研究的创新点主要体现在：1) 构建了多算法融合的钢材缺陷分类模型，综合逻辑回归、决策树等算法优势，提供多样化解决方案；2) 提出特征分类筛选与优化方法，通过相关性分析等技术挖掘关键特征，提升分类性能；3) 将 SMOTE 方法与机器学习算法创新结合，有效解决样本不均衡问题，增强模型对少数类样本的识别能力；4) 设计全面实验评估体系，通过参数敏感性分析和消融实验，为模型优化提供关键指导。

尽管取得了积极成果，但研究仍存在不足。未来我们将拓展数据收集范围，涵盖更广泛的钢材类型和缺陷形态，提升模型的普适性和鲁棒性。同时，探索结合深度学习和计算机视觉技术，实现更高精度的缺陷检测与分类，并结合物联网技术构建实时监测系统，推动智慧基建向更高智能化的水平发展。

参考文献

- [1] 廖晓群, 李丹, 徐清钊. 基于深度学习的钢板表面缺陷检测研究综述[J/OL]. 计算机测量与控制: 1-14. <http://kns.cnki.net/kcms/detail/11.4762.tp.20250116.1232.002.html>, 2025-04-21.
- [2] 马磊, 李晔, 王宇翔. YOLOv8-FD:YOLOv8 改进的钢板表面缺陷检测方法[J]. 计算机工程与应用, 2024, 60(24): 211-221.
- [3] 郑贵君, 邹伯昌, 马瑞. 基于多模态与自适应特征融合的钢材表面缺陷检测[J]. 物联网技术, 2025, 15(7): 20-26+31.
- [4] 王瑞仪, 徐沛航, 任杰. 基于深度学习的钢材表面缺陷检测算法研究[J]. 电脑编程技巧与维护, 2025(2): 129-131+144.
- [5] Saberironaghi, A., Ren, J. and El-Gindy, M. (2023) Defect Detection Methods for Industrial Products Using Deep Learning Techniques: A Review. *Algorithms*, **16**, Article 95. <https://doi.org/10.3390/a16020095>
- [6] 李键, 李华, 胡翔坤, 等. 基于深度学习的表面缺陷检测技术研究进展[J]. 计算机集成制造系统, 2024, 30(3): 774-790.