

# 面向不平衡数据集的矿井通风系统智能故障诊断

赵丹<sup>1,2</sup>, 沈志远<sup>1,2</sup>, 宋子豪<sup>1,2</sup>

(1. 辽宁工程技术大学 安全科学与工程学院, 辽宁 阜新 123000; 2. 辽宁工程技术大学 矿山热力灾害与防治教育部重点实验室, 辽宁 葫芦岛 125105)

**摘要:**及时准确判断故障分支的位置对保障矿井通风系统的可靠性和安全性意义重大。针对实际工况下, 矿井通风系统故障样本数据存在不平衡性, 导致传统的机器学习模型诊断能力与泛化能力差的问题, 提出了一种面向通风系统不平衡数据集的 WGAN-div-RF 故障诊断模型。以简单通风网络为例构造了不平衡比分别为 2 : 1、5 : 1、10 : 1、20 : 1 的故障数据集, 深入分析了不平衡样本集对通风系统故障诊断的影响。搭建了基于 Wasserstein 距离生成对抗网络 (WGAN-div) 对不平衡数据集进行数据增强处理, 在构建网络时创新性地加入了残差块, 提高了生成数据的质量, 实现原始样本的有效扩充。结合集成学习中的随机森林 (RF) 模型实现通风系统故障分支诊断。以东山煤矿通风系统为实验对象, 分别进行了不同数据增强模型、不同分类模型以及不同数据生成率下的故障诊断对比实验, 以多种评价指标及 t-SNE 可视化对模型有效性进行评估。结果表明: 加入残差块的 WGAN-div 模型生成数据与真实数据具有很好的相似性, 相较于 GAN 模型、WGAN 模型和 WGAN-gp 模型, WGAN-div 模型更具优越性; 应用 WGAN-div 模型进行数据增强后, 机器学习分类模型的性能提升明显; 当扩充数据集达到平衡时, 与其他集成模型及常用的矿井通风系统故障诊断 SVM 模型相比, RF 模型在  $R_e$ 、 $P_r$ 、 $G_{mean}$  和  $F_1$  指标上均占优势。

**关键词:** 矿井通风系统; 故障诊断; 不平衡数据; 生成对抗网络; 随机森林

**中图分类号:** TD72 **文献标志码:** A **文章编号:** 0253-9993(2023)11-4112-12

## Intelligent fault diagnosis of mine ventilation system for imbalanced data sets

ZHAO Dan<sup>1,2</sup>, SHEN Zhiyuan<sup>1,2</sup>, SONG Zihao<sup>1,2</sup>

(1. College of Safety Science & Engineering, Liaoning Technical University, Fuxin 123000, China; 2. Key Laboratory of Ministry of Education for Mine Thermo-motive Disaster and Prevention, Liaoning Technical University, Huludao 125105, China)

**Abstract:** It is of great significance to determine the location of fault branch timely and accurately to ensure the reliability and safety of mine ventilation system. To solve the problem that the traditional machine learning model has the poor diagnostic ability and generalization ability due to the imbalance of sample data in mine ventilation system under actual working conditions, a WGAN-div-RF fault diagnosis model is proposed. Taking a simple ventilation network as an example, the fault data sets with the imbalance ratios of 2 : 1, 5 : 1, 10 : 1, 20 : 1 are constructed, and the impact of imbalanced samples on the ventilation system fault diagnosis is analyzed indepth. The Wasserstein divergence for GANs (WGAN-div) is built, and the residual blocks are added innovatively to improve the quality of the generated data and ex-

收稿日期: 2022-12-27 修回日期: 2023-03-16 责任编辑: 王晓珍 DOI: 10.13225/j.cnki.jccs.2022.1872

基金项目: 国家自然科学基金资助项目 (52374202)

作者简介: 赵丹 (1982—), 女, 辽宁阜新人, 教授。E-mail: zhaosixiaojie\_20@163.com

通讯作者: 沈志远 (1994—), 女, 辽宁沈阳人, 博士研究生。E-mail: szy0207@foxmail.com

引用格式: 赵丹, 沈志远, 宋子豪. 面向不平衡数据集的矿井通风系统智能故障诊断[J]. 煤炭学报, 2023, 48(11): 4112-4123.

ZHAO Dan, SHEN Zhiyuan, SONG Zihao. Intelligent fault diagnosis of mine ventilation system for imbalanced data sets[J]. Journal of China Coal Society, 2023, 48(11): 4112-4123.



移动阅读

pand the original sample set. Combined with the RF model, the fault diagnosis of ventilation system is realized. Taking the ventilation system of the Dongshan Coal Mine as the experimental object, the comparative experiments are carried out respectively with different data enhancement models, different classification models, and different data generation rates. The effectiveness of the model is evaluated with various indexes and t-SNE visualization. The results show that the data generated by the WGAN-div model with residual blocks has a good similarity to the real data. Compared with GAN, WGAN, and WGAN-gp, the WGAN-div model is superior. After applying the WGAN-div model for data augmentation, the performance of the machine learning classification model is significantly improved. When the expanded data set is balanced, compared with other integrated models and the commonly used SVM model for mine ventilation system fault diagnosis, the RF model is superior in  $R_e$ ,  $P_r$ ,  $G_{mean}$  and  $F_1$  indexes.

**Key words:** mine ventilation system; fault diagnosis; imbalanced data; generate adversarial network; random forest

如何及时准确的判断故障的位置,已成为煤矿亟待解决的一个难题<sup>[1-2]</sup>。随着煤矿智能化建设的发展,应用机器学习算法实现通风系统的智能故障诊断,助力矿井通风智能化管理是研究的关键<sup>[3]</sup>。

随着大数据、工业互联网、人工智能等技术的发展,故障诊断技术在电网<sup>[4]</sup>、机械设备<sup>[5]</sup>、航空航天<sup>[6]</sup>等不同工程领域应用成熟。2018年,刘剑等<sup>[7-8]</sup>以风量作为输入特征,应用支持向量机(Support Vector Machine,SVM)算法确定了矿井通风系统故障位置及故障量,这开创了应用机器学习进行矿井通风系统故障诊断的先河,2020年应用遗传算法构建了矿井通风系统故障诊断无监督模型,无需样本参与训练,有效提升了诊断性能;HUANG等<sup>[9-11]</sup>利用卡尔曼滤波模型对矿井监测风速数据进行了预处理,并提出了基于混合编码算法的矿井通风系统无监督学习故障诊断模型,实现了故障位置和故障量的同时诊断;周启超等<sup>[12]</sup>基于改进的遗传算法对矿井通风系统故障诊断SVM模型的参数进行了优化研究,有效避免了模型易出现过拟合的问题;倪景峰等<sup>[13-14]</sup>提出了基于随机森林和决策树的通风系统故障诊断方法,并证实了随机森林模型优于决策树模型;张浪等<sup>[15]</sup>选择了SVM、神经网络和随机森林(Random Forest,RF)3种矿井通风系统故障诊断机器学习算法进行对比分析,结果表明神经网络模型具有更高的准确率;ZHAO等<sup>[16]</sup>以大明矿为研究对象,在构建的故障巷道范围库内应用改进的SVM算法进行通风系统故障诊断,缩减了故障定位的范围,提高了样本训练效率;WANG等<sup>[17]</sup>构建了基于多标签K-近邻(Multi-label K-Nearest Neighbor, ML-KNN)的机器学习模型,解决了矿井通风系统多个位置发生故障时的快速诊断问题;LIU等<sup>[18]</sup>应用4种机器学习算法:K-近邻(K-Nearest Neighbor,KNN)、多层感知机(Multilayer Perceptron,MLP)、SVM和决策树(Decision Tree,DT)对矿井通风系统故障诊断模型性能进行了充分评价,确定了KNN模型和DT模

型的优越性。虽然机器学习算法在矿井通风系统故障诊断中表现优异,但目前的矿井通风系统故障诊断模型的建立都是在数据集较为完备的前提下进行的。但是,在实际的通风系统故障情形下,完备的数据集条件是不能满足的。机器学习分类器高度依赖完备的样本集,不平衡的样本集训练出的模型通常不具有参考意义。如何在样本不平衡情况下开展故障诊断是一个严峻的挑战。机器学习领域的学者们通常从算法层面和数据层面解决不平衡数据的分类问题。文献[19]从算法层面出发构建了单分类支持向量机(One-Class SVM,OCISVM)与增量学习(Incremental Learning,IL)相结合的通风系统故障诊断模型,但是该方法依赖于特定算法,导致适用性较差。

鉴于此,笔者从数据层面和网络体系层面开展不平衡数据集的通风系统故障诊断研究,构建了基于Wasserstein距离的生成对抗网络(Wasserstein divergence for GANs,WGAN-div),创新性地加入残差块实现原始数据增强处理,重构平衡数据集。结合集成学习中的投票机制实现通风网络分支故障诊断,确定了RF模型在通风系统故障诊断中的优越性。有效解决了实际工况下样本不平衡的故障诊断问题,为智能诊断技术真正应用到矿井提供技术支撑。

## 1 处理不平衡数据集的改进模型

### 1.1 通风系统故障数据不平衡分析

矿井通风系统实际工况下,风门、风窗等含通风构筑物的巷道,采掘工作面,主要用风巷道,通风多分支交汇点处等位置更易发生故障,产生的故障数据较多,而其他分支故障概率较低,产生的故障数据较少,各个分支产生的故障数据样本数量存在很大的差距,存在数据不平衡问题。如图1所示,不同颜色的五角星代表通风系统监测数据中的不同故障分支产生的故障样本,黄色五角星代表构筑物分支等易发生故障

巷道产生的故障样本,为多数类故障样本集合;蓝色五角星代表其他不易发生故障的分支产生的故障样本,为少数类故障样本集合。

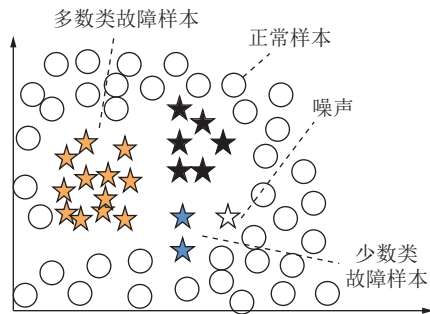


图1 数据不平衡示意

Fig.1 Schematic diagram of data imbalance

矿井通风系统故障分支不平衡数据集可以描述为

$$\begin{cases} S_{m+n} = \{X_m, Y_n\} \\ X_m = \{x_i | i = 1, 2, \dots, m\} \\ Y_n = \{y_i | i = 1, 2, \dots, n\} \end{cases} \quad (1)$$

式中,  $X_m$  为少数类故障分支数据集;  $Y_n$  为多数类故障分支数据集;  $S_{m+n}$  为通风系统故障分支不平衡数据集;  $x_i$  和  $y_i$  为各数据集中的第  $i$  个样本数据;  $m$  为少数类样本个数;  $n$  为多数类样本个数。

## 1.2 传统的 GAN 模型

生成对抗网络 (Generative Adversarial Network, GAN) 模型可以实现新样本数据的生成,从而达到调整  $X_m$  和  $Y_n$  的类间平衡度的目的。GAN 模型主要由判别器 D 和生成器 G 两部分组成,其基本结构如图 2 所示。生成器 G 将随机噪声  $z$  映射到真实样本空间生成新的数据  $\hat{x}$ ; 判别器 D 判断  $\hat{x}$  的真假即判别  $\hat{x}$  为真实数据或生成数据。2 个网络交替训练,当判别器 D 和生成器 G 达到动态平衡时,新生成的数据与真实数据具有相似特征。

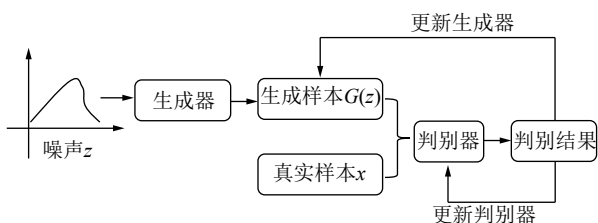


图2 GAN 模型基本结构

Fig.2 Basic structure of GAN model

GAN 模型的损失函数为

$$\min_G \max_D V(G, D) = E_{G(z) \sim P_z} (-D(G(z))) + E_{x \sim M_r} (D(x)) \quad (2)$$

其中,  $x$  为真实样本数据;  $P_z$  为随机噪声的分布;  $z$  为噪声;  $M_r$  为真实数据的分布;  $E_{G(z) \sim P_z}$  为添加噪声的期望函数;  $E_{x \sim M_r}$  为真实数据的期望函数;  $G(\cdot)$  为生成器的可微函数;  $D(\cdot)$  为判别器的可微函数。实际上,生成器 G 的损失函数相当于最小化生成数据分布和与真实数据分布之间的 JS 散度,有

$$G^* = \min_G V(G, D^*) = \min_G E_{x \sim P_G} (-D(x)) + E_{x \sim M_r} (D(x)) = \min_G 2 \text{JS}(P_r(x) || P_G(x)) - 2 \lg 2 \quad (3)$$

式中,  $P_G$  为生成数据的分布;  $G^*$ 、 $D^*$  分别为生成器损失函数和判别器损失函数的最优解;  $E_{x \sim P_G}$  为生成数据的期望函数; JS 为 JS 散度。

## 1.3 WGAN-div 模型

在 GAN 训练初期,  $P_G$  与  $M_r$  一般不会重叠,判别器 D 容易判定数据的真假,但此时,该损失函数中的 JS 散度退化为常数项  $\lg 2$ ,进而导致生成器 G 的梯度消失,无法应用梯度下降法对网络进行训练,这使得传统 GAN 模型出现训练不稳定的问题<sup>[20]</sup>。2017 年,ARJOVSKY 等<sup>[20]</sup>提出应用 Wasserstein 距离代替 JS (Jensen-Shannon) 散度以解决传统 GAN 模型梯度消失的问题,构建了基于 Wasserstein 距离的生成对抗网络 (Wasserstein GAN, WGAN) 模型。但是在 WGAN 训练过程中,通常需要保持梯度的绝对值小于某个固定值,文献[21]提出了加入惩罚因子的 GAN 模型 (Wasserstein for GANs, WGAN-gp) 模型,保证生成样本与真实样本之间满足 Lipschitz 连续,但该方案并没有理论依据。对此,文献[22]提出了不需要 Lipschitz 约束的 WGAN-div 模型,并在理论和应用上都证明了其优越性。基于前人的研究,笔者选择 WGA-div 数据增强模型,损失函数为

$$L_G = -E_{G(z) \sim P_G} [D(G(z))] \quad (4)$$

$$L_D = E_{G(z) \sim P_G} [D(G(z))] - E_{x \sim M_r} [D(x)] - k E_{\hat{x} \sim p_u} [\|\nabla D(\hat{x})\|^p] \quad (5)$$

式中,  $L_G$  为生成器损失函数;  $L_D$  为判别器损失函数;  $E_{G(z) \sim P_G}$  为生成器噪声的期望函数;  $E_{\hat{x} \sim p_u}$  为插值  $\hat{x}$  的期望函数,  $\hat{x}$  为生成样本与真实样本之间的随机插值,  $\hat{x} = \alpha x + (1 - \alpha)G(z)$ ,  $\alpha$  为系数,  $\alpha \in [0, 1]$ ;  $p_u$  为插值  $\hat{x}$  的分布;  $k$ 、 $p$  为范数的幂,根据前人研究和实验测试,设置  $k=2$ 、 $p=6$ 。

## 1.4 残差块

文献[23]针对深度神经网络训练困难问题,提出了残差学习框架,能够简化深度神经网络的训练;文献[24]应用加入残差块的生成对抗网络实现了光伏数



据的缺失值重构。鉴于此,为了防止使用深度卷积网络搭建的 WGAN-div 模型在训练过程中出现梯度消失或网络退化的问题,笔者在判别器和生成器中加入了恒等映射残差块,残差块如图 3 所示<sup>[23]</sup>。

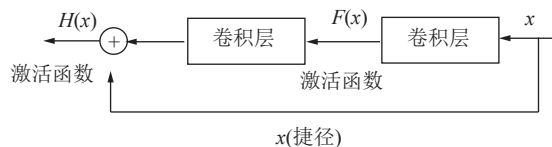


图 3 残差块示意

Fig.3 Schematic diagram of residual block

残差块以真实数据  $x$  为输入,主线径上有 2 个卷积层,其目标函数为  $H(x)$ ,定义为

$$H(x) = f(x, W) + x \quad (6)$$

其中,  $f(x, W)$  为映射函数;  $W$  为卷积层的权重。恒等映射残差块不仅可以学习  $x$  与  $H(x)$  的差别而且保证了 2 者尺寸相同。残差块的引入使得网络的训练更容易,避免了梯度消失和梯度爆炸的问题。因此,笔者采用加入了残差块的 WGAN-div 模型对通风系统监测数据不平衡样本进行数据扩充。将通风系统监测数据故障数据集中少数类样本个数由  $m$  调整到  $r = m + \sum \bar{x}_i$ , 进一步得到平衡数据集  $S' = \{X', Y_n\}$ , 其中,  $X'_i$  为平衡后的少数类样本数据集。

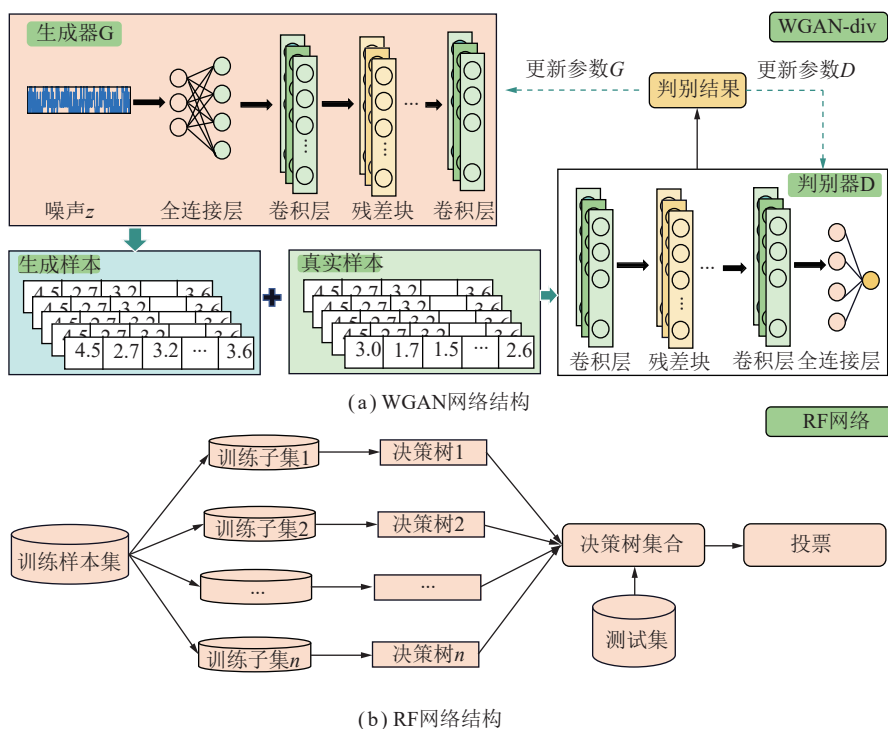


图 4 WGAN-div-RF 模型构架

Fig.4 WGAN-div-RF model architecture

(1) 由于实际工况下矿井故障样本数据获取困难,本文应用智能矿井通风仿真系统(IMVS)模拟通风系

## 2 基于 WGAN-div-RF 的通风系统故障诊断

### 2.1 RF 分类模型

随机森林作为一种典型的集成学习模型,可以处理高维数据的分类,因此笔者选择 RF 作为通风系统故障诊断多分类器。将风速数据作为 RF 分类模型的输入,将故障分支编号作为 RF 分类模型的输出。具体过程如下:对样本数据集进行 Bootstrap 采样,得到  $K_n$  个样本子集,应用子集训练出  $K_n$  个决策树,将测试数据输入  $K_n$  个决策树集合中得到  $N$  个结果,采用投票策略得到最终的分类结果为

$$F(x) = \arg \max_y \sum_{i=1}^{K_n} I(f_i(v, \theta_i) = y) \quad (7)$$

式中,  $F(x)$  为  $K_n$  个决策树投票确定的矿井通风系统故障分支;  $f_i$  为第  $i$  个决策树的分类模型;  $v$  为输入模型的特征参量,本文为风速数据;  $\theta_i$  为用于训练第  $i$  个决策树的样本子集;  $I(\cdot)$  为示性函数(分别以 1 和 0 表示集合内是否存在该数值);  $y$  为待判别的故障分支编号。

### 2.2 整体构架及流程

基于 WGAN-div-RF 的通风系统故障诊断整体构架如图 4 所示。具体流程如下:

统故障,构造通风系统故障不平衡数据集  $O$ , 将数据集划分为测试样本集  $O_{st}$  和训练样本集  $O_{in}$ 。

(2) 应用 WGAN-div 模型对不平衡的训练样本集  $O_{in}$  进行数据增强处理, 生成新的故障样本  $O_n$ , 将  $O_n$  加入到原训练样本集  $O_{in}$  中合成新的增广样本  $O_{ex}$ 。

(3) 用平衡后的增广样本集  $O_{ex}$  训练 RF 模型, 获得训练好的故障诊断模型。

(4) 将测试样本集  $O_{st}$  输入训练好的 RF 模型进行通风系统故障诊断。

### 2.3 评价指标

通风系统故障诊断多分类模型的评价通常建立在二分类混淆矩阵的基础上, 对于样本不平衡的多分类问题, 准确率指标难以实现对分类结果的准确评价, 因此, 文中增加了召回率  $R_e$ 、精确率  $P_r$ 、 $G_{mean}$ 、和  $F_1$  分数对通风故障诊断模型进行综合评价。各个指标<sup>[25]</sup>的定义如下:

$$A = \frac{\sum_{i=1}^N T_{Pi} + \sum_{i=1}^N T_{Ni}}{\sum_{i=1}^N T_{Pi} + \sum_{i=1}^N F_{Pi} + \sum_{i=1}^N T_{Ni} + \sum_{i=1}^N F_{Ni}} \quad (8)$$

$$P_r = \frac{\sum_{i=1}^N T_{Pi}}{\sum_{i=1}^N T_{Pi} + \sum_{i=1}^N F_{Pi}} \quad (9)$$

$$R_e = \frac{\sum_{i=1}^N T_{Pi}}{\sum_{i=1}^N T_{Pi} + \sum_{i=1}^N F_{Ni}} \quad (10)$$

$$F_1 = \frac{2P_r R_e}{P_r + R_e} \quad (11)$$

$$G_{mean} = \sqrt{\frac{\sum_{i=1}^N T_{Pi}}{\sum_{i=1}^N T_{Pi} + \sum_{i=1}^N F_{Ni}} \cdot \frac{\sum_{i=1}^N T_{Ni}}{\sum_{i=1}^N T_{Ni} + \sum_{i=1}^N F_{Pi}}} \quad (12)$$

式中,  $A$  为模型故障诊断准确率;  $P_r$  和  $R_e$  分别为模型的平均精确率和召回率;  $G_{mean}$  为召回率和特异度的几何平均值;  $N$  为输入模型的通风网络分支数,  $T_{Pi}$  为第  $i$  个故障分支的真正例;  $T_{Ni}$  为第  $i$  个故障分支的真负例;  $F_{Pi}$  为第  $i$  个故障分支的假正例;  $F_{Ni}$  为第  $i$  个故障分支的假负例。

### 3 不平衡数据对故障分支诊断影响实验分析

为了验证不平衡数据对通风系统故障诊断的影响, 以图 5 所示简单角联通风网络为例, 设计不同不平衡比下的故障诊断实验。该网络中分支数为 7, 节点数为 6,  $e_1$  和  $e_7$  分别为进风分支和回风分支, 调节风窗安设在  $e_4$  分支, 风机特性方程为  $1\,037.2 + 52.69q - 0.52q^2$ , 其中,  $q$  为风量。通风参数见表 1。采用智能矿井通风仿真系统 IMVS 模拟分支故障<sup>[7]</sup>(不包括源汇分支), 故障数据生成的具体方法参见文献<sup>[7]</sup>, 按照不同的不平衡比生成 4 组数据集, 构造 4 组实验方案。

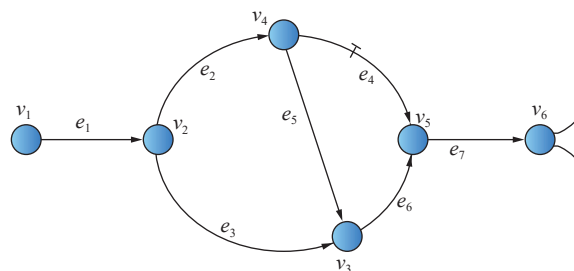


图 5 简单通风网络

Fig.5 Simple ventilation network diagram

$e_4$  分支安设了风窗, 相较于其他分支更容易发生故障, 因此通过增加  $e_4$  分支的故障次数改变不平衡比。不平衡比分别设置为 2 : 1、5 : 1、10 : 1、20 : 1,  $e_4$  分支的模拟故障次数按照不平衡比的不同分别设置为 100、250、500、1 000。为了方便比较, 实验将少数类故障样本数量设置为相同, 即除了  $e_4$  分支外,  $e_2$ 、 $e_3$ 、 $e_5$ 、 $e_6$  每个分支模拟故障 50 次, 相应的全部分支的故障样本总数分别为 300、450、700、1 200, 对应的实验

表 1 简单网络各分支初始参数

Table 1 Initial parameters of each branch of a simple network

分支	断面形状	断面宽度/m	断面高度/m	风量/( $m^3 \cdot s^{-1}$ )	风阻/( $N \cdot s^2 \cdot m^{-8}$ )
$e_1$	矩形	4.5	3.60	78.325	0.054
$e_2$	矩形	3.6	3.20	25.542	0.265
$e_3$	矩形	5.1	4.20	21.386	0.084
$e_4$	矩形	4.6	3.72	4.201	8.210
$e_5$	矩形	5.1	3.92	52.738	0.062
$e_6$	矩形	5.2	4.11	74.124	0.192
$e_7$	矩形	4.5	3.67	78.325	0.058

方案分别记为  $T_1$ 、 $T_2$ 、 $T_3$ 、 $T_4$ 。每一组实验均对应一个平衡数据集作为对照实验组进行对比分析。为了保证实验对比的合理性,平衡数据集的故障样本总量应与不平衡数据集保持一致即每一组实验的故障样本总数应为 300、450、700、1 200,由于平衡样本集中每一条分支的故障样本数应相同且排除源汇分支共有 5 条分支,因此,平衡数据集中 4 组实验各分支故障次数分别设置为 60、90、140、240,对应的实验方案分别记为  $D_1$ 、 $D_2$ 、 $D_3$ 、 $D_4$ 。

为严格控制相关变量,在保证故障样本量一致的同时,各个实验模型均应在最优参数下运行才具备比较意义。以最大化 F1 分数为目标进行调整,经十折交叉验证确定各实验 RF 模型最佳参数,参数定义见

表 2,参数设置见表 3。文中以风速特征作为输入,因此利用式 (13) 将通风网络解算得到的风量  $q$  转换为风速  $v$ 。

$$v = \frac{q}{l_e w_e} \quad (13)$$

其中,  $l_e$  为巷道断面高度,  $m$ ;  $w_e$  为巷道断面宽度,  $m$ 。为以  $T_1$  实验为例,其部分故障样本数据见表 4,表中  $v'_i$  为各分支风速,  $m/s$ ;  $e'_i$  为故障分支。将每一组实验数据集的 70% 划分为训练集, 30% 划分为测试集,以故障分支编号作为输出进行故障诊断实验,得到测试集的混淆矩阵如图 6 所示,横坐标表示预测故障分支编号,纵坐标表示真实故障分支编号。实验  $T_1 \sim T_4$  的综合评价指标结果如图 7 所示。

表 2 分类模型参数定义

Table 2 Definition of classification model parameters

参数	定义	参数	定义
$N'$	最大迭代次数	$\eta$	学习率
$T_d$	树的最大深度	$M_p$	节点分割所需的最小样本数
$M$	弱分类器的数量	$M_f$	叶节点所需的最小样本数
$T_f$	子样本比率	$c$	惩罚系数
$\omega$	随机采样比例	$K$	核函数

表 3 RF 模型参数

Table 3 RF model parameters

实验方案	参数设置	实验方案	参数设置
$T_1$	$M=30, T_d=25, M_p=2, M_f=1$	$D_1$	$M=10, T_d=12, M_p=1, M_f=1$
$T_2$	$M=35, T_d=30, M_p=2, M_f=1$	$D_2$	$M=20, T_d=18, M_p=1, M_f=1$
$T_3$	$M=100, T_d=50, M_p=2, M_f=1$	$D_3$	$M=50, T_d=30, M_p=1, M_f=1$
$T_4$	$M=100, T_d=50, M_p=2, M_f=1$	$D_4$	$M=50, T_d=30, M_p=1, M_f=1$

表 4  $T_1$  实验模拟故障样本集Table 4  $T_1$  simulation fault sample set

样本	$v'_1$	$v'_2$	$v'_3$	$v'_4$	$v'_5$	$v'_6$	$v'_7$	$e'_i$
1	5.326	2.182	2.765	0.234	3.614	2.418	4.581	2
2	5.235	2.112	2.452	0.211	3.871	2.563	4.113	2
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
101	5.124	1.360	3.621	2.395	5.364	4.261	3.693	4
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
200	4.590	1.329	3.262	2.151	5.238	4.517	2.581	4
201	4.020	4.670	2.152	3.240	2.69	1.266	3.657	5
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
299	4.962	2.824	2.350	1.325	3.201	2.345	3.65	6
300	4.620	3.622	2.103	2.120	2.590	2.125	2.326	6

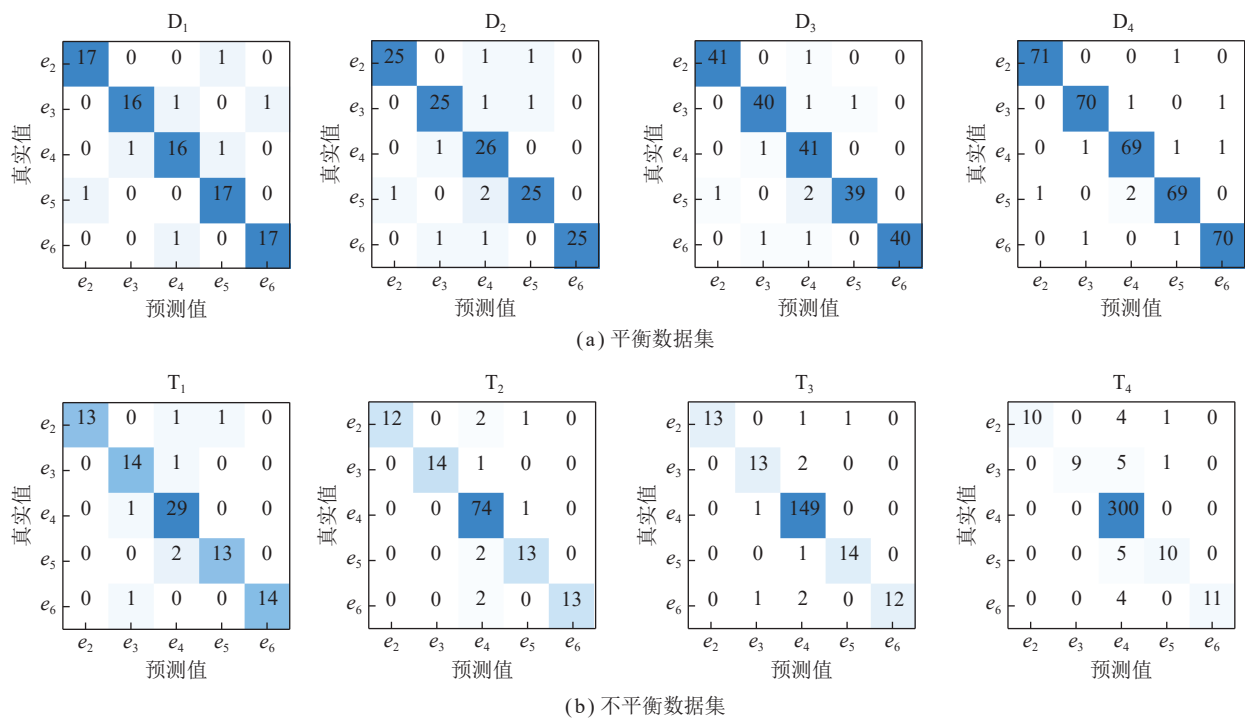


图6 简单通风网络故障诊断实验混淆矩阵

Fig.6 Confusion matrix of simple ventilation network fault diagnosis experiment

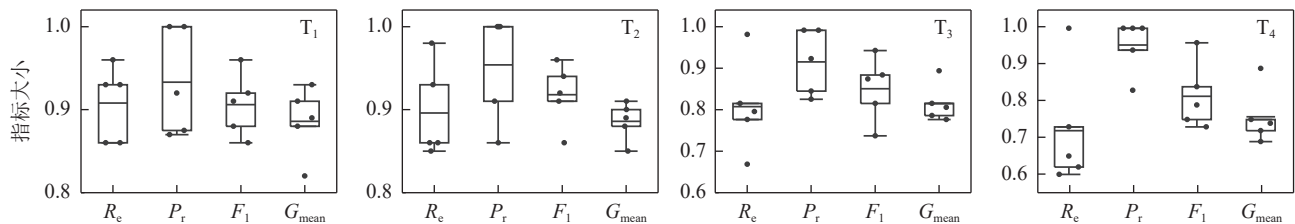


图7 简单通风网络不平衡数据集故障诊断实验评价指标

Fig.7 Experimental evaluation indexes of fault diagnosis in unbalanced data set of simple ventilation network

由图6(a)可知,实验 $D_1 \sim D_4$ 的平均准确率分别0.922、0.933、0.957、0.970,可以看出RF分类模型能够有效地对通风系统故障进行诊断。但值得注意的是,理想的训练样本条件是获得良好诊断结果的前提,理想的训练样本不仅意味着故障样本数据充足,还意味着故障样本数据中各个分支有着平衡的故障样本数量。然而,实际的矿井通风系统难以获得各分支故障样本均衡的数据集。由图6(b)和图7可知,实验 $T_1$ 的 $R_e$ 、 $P_r$ 、 $G_{mean}$ 和 $F_1$ 分数平均值分别为0.91、0.93、0.90、0.92;实验 $T_2$ 的 $R_e$ 、 $P_r$ 、 $G_{mean}$ 和 $F_1$ 分数平均值分别为0.89、0.95、0.88、0.91;实验 $T_3$ 的 $R_e$ 、 $P_r$ 、 $G_{mean}$ 和 $F_1$ 分数平均值分别为0.812、0.95、0.82、0.87;实验 $T_4$ 的 $R_e$ 、 $P_r$ 、 $G_{mean}$ 和 $F_1$ 分数平均值分别为0.73、0.95、0.81、0.78,可以看出随着不平衡比例的增加,除模型的精确率未发生明显变化之外,召回率、 $G_{mean}$ 和 $F_1$ 分数不断降低,由此可见不平衡数据影响了模型的整体性能,其鲁棒性降低显著,不平衡数据使得模型

出现漏判和误判的情况较多。尤其,由图7中 $T_4$ 实验可知,当不平衡比为20:1时,各故障分支中 $R_e$ 的最大值为1,最小值为0.6; $P_r$ 的最大值为1,最小值为0.83; $F_1$ 分数的最大值为0.96,最小值为0.73,各分支指标值的分布差异较大,分析认为数据不平衡易引起小析取问题,常规的机器学习分类器依据大量多数类分支( $e_4$ 分支)数据规则建立模型,而忽略了其他少样本分支的数据特点,从而导致在分类时易将其他分支故障误诊断为多数类分支( $e_4$ 分支),随着不平衡比例的增加,故障样本被误判的比例逐渐升高,这进一步说明了不平衡数据集对通风系统故障诊断模型的危害,可见研究的必要性和实用性。

## 4 生产矿井实例实验分析

### 4.1 数据准备

笔者以鸡西矿业集团东山煤矿通风系统为例进行不平衡数据故障诊断实验。实验矿井的通风方式



为对角式,该矿通风网络如图 8 所示,分支数为 96,节点数为 84,总入风量  $14\,394\text{ m}^3/\text{min}$ ,4 条进风井对应的分支编号分别为  $e_2$ 、 $e_1$ 、 $e_{23}$ 、 $e_5$ ,由南风井、西风井共同担负全矿井总回风任务,总排风量  $14\,738\text{ m}^3/\text{min}$ ,对应的分支编号分别为  $e_{54}$ 、 $e_{92}$ 。安设风门的分支编号为  $e_{47}$ 、 $e_{85}$ 、 $e_{28}$ 、 $e_{86}$ 、 $e_{48}$ 、 $e_{78}$ 、 $e_{22}$ 、 $e_7$ 、 $e_{30}$ 、 $e_{38}$ 、 $e_{29}$ 、 $e_{19}$ 、 $e_{65}$ 、 $e_{52}$ 、 $e_{84}$ 、 $e_{33}$ ;安设风窗的分支编号为  $e_{10}$ 、 $e_{83}$ 、 $e_{24}$ 、 $e_{13}$ 、 $e_{93}$ 。风机特性方程分别为:  $723.65+18.26q-0.17q^2$ 、 $614+45.2q-0.09q^2$ 。应用 IMVS 模拟分支故障(不包括源汇分支)<sup>[7]</sup>,其中风门风窗构筑物所在分

支模拟故障 200 次,其他分支模拟故障 10 次,得到 5 120 组故障样本,数据不平衡比为 20 : 1。全矿共安设了 15 台风速传感器,布设位置已在图 8 中标出(本文以矿井实际安设的传感器为基础,不考虑传感器安设数量和配置的优化问题)。将风速传感器所在分支解算得到的风量数据经式 (13) 转换为风速数据作为模型的输入,部分数据见表 5,表中  $v_i$  为各分支风速,  $\text{m/s}$ ;  $e_i$  为故障分支。将标准化处理后的故障样本数据按照 7 : 3 的比例划分为训练样本和测试样本。

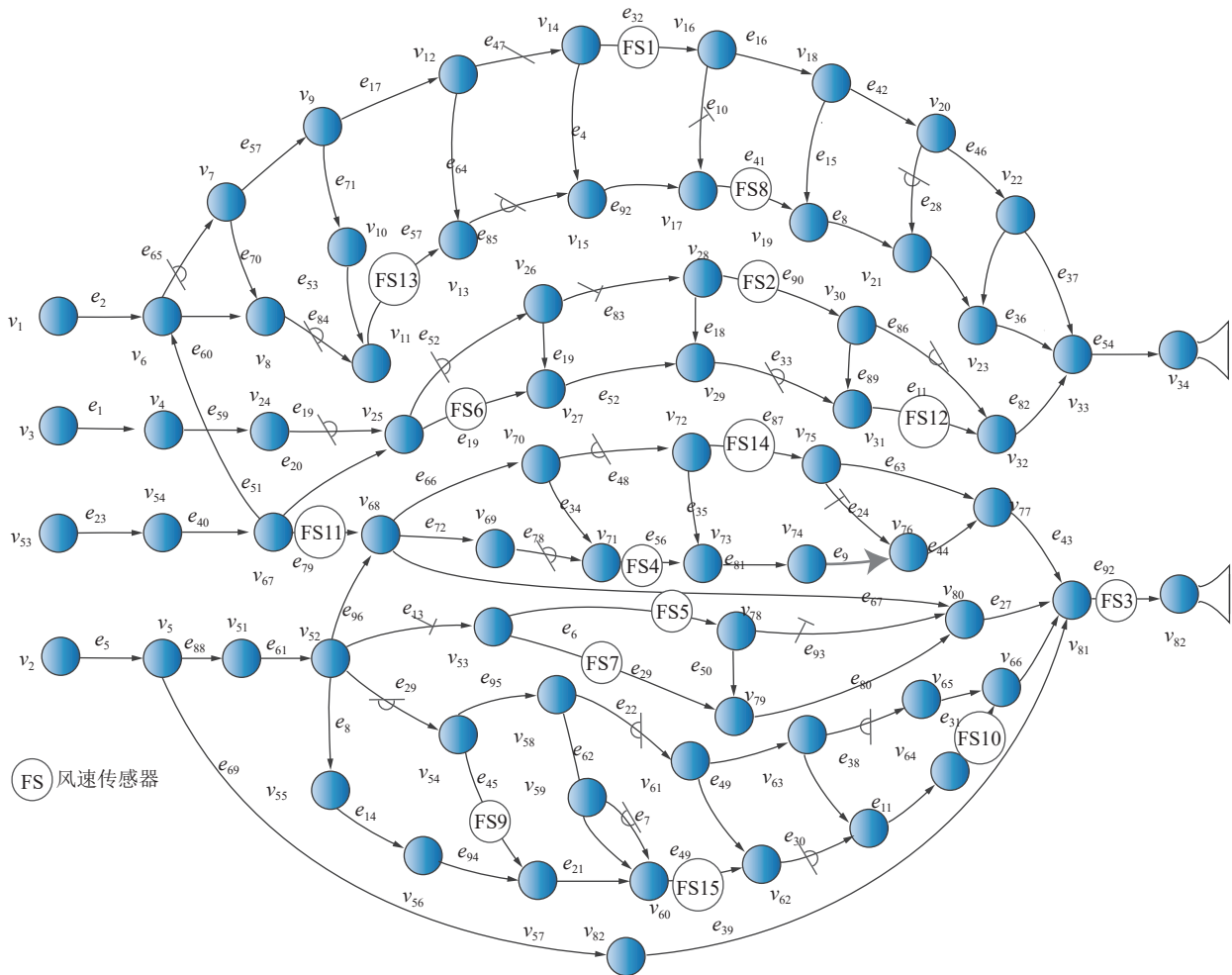


图 8 东山矿通风网络

Fig.8 Ventilation network of Dongshan coal mine

#### 4.2 WGAN-div 有效性验证

为了验证 WGAN-div 在通风系统不平衡数据处理的有效性,原始故障样本分别采用:①原始数据集  $D_{in}$ ;② GAN 模型;③ WGAN 模型;④ WGAN-gp 模型;⑤ 本文所建 WGAN-div 模型处理生成新的样本集  $O_n$ ,使得合样本集  $O_{ex}$  达到数据平衡,分类算法都选择 RF 模型。本文构建的 WGAN-div 模型生成器、判别器均包含 3 个残差块,参数设置见表 6。实验结

果见表 7,为 10 次运行结果的平均值 $\pm$ 标准差(最优结果加粗表示)。分析表 6 可得出:

(1) 直接采用 RF 分类模型对原始不平衡数据集进行故障分支诊断,  $A$ 、 $P$ 、 $G_{\text{mean}}$  和  $F_1$  分数都是最低。这意味着 RF 模型不能准确识别出通风系统不平衡数据集中的少数类故障样本,因此,使用原始数据集不能实现对通风系统故障分支的有效诊断。

(2) 对比原始数据集,基于 WGAN-div 数据增强



表 5 生产矿井故障样本集  
Table 5 Fault sample set in production mine

样本	$v'_1$	$v'_2$	$v'_3$	$v'_4$	$v'_5$	$v'_6$	$v'_7$	$v'_8$	$v'_9$	$v'_{10}$	$v'_{11}$	$v'_{12}$	$v'_{13}$	$v'_{14}$	$v'_{15}$	$e'_i$
1	3.6	5.6	7.6	4.2	3.9	4.1	2.6	10.4	4.7	4.2	3.2	3.8	9.6	4.2	3.2	$e_{85}$
2	3.6	5.4	7.8	4.1	3.9	4.2	2.4	10.2	4.8	4.2	3.1	3.9	9.5	4.1	3.3	$e_{85}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
201	3.5	3.6	7.9	4.1	3.8	7.1	2.6	3.6	4.7	4.3	3.2	8.0	4.3	4.2	3.1	$e_{33}$
202	3.4	3.5	7.8	4.2	3.7	7.2	2.5	3.7	4.6	4.2	3.1	7.9	4.4	4.2	3.2	$e_{33}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
5 119	6.5	5.6	7.6	4.4	3.6	4.4	2.6	0.7	4.5	4.3	3.2	3.8	0.5	4.0	3.1	$e_{92}$
5 120	6.4	5.6	7.8	4.3	3.6	4.3	2.5	0.6	4.6	4.4	3.3	3.9	0.3	4.2	3.0	$e_{92}$

表 6 WGAN-div 模型参数设置  
Table 6 WGAN-div model parameters

生成器							判别器								
名称	输入维度	激活函数	滤波器	步长	卷积核	输出维度	名称	输入维度	激活函数	滤波器	步长	卷积核	输出维度		
输入噪声	50	—	—	—	—	50	输入层	15×1	—	—	—	—	15×1		
全连接层	50	LeakyReLU	—	—	—	426	卷积层	15×1	LeakyReLU	3	1	4	15×3		
残差块1	上采样层	3×64	LeakyReLU	64	2	4	6×64	残差块1	下采样层	15×3	LeakyReLU	16	2	4	15×16
	卷积层	6×64	—	64	1	4	6×64		卷积层	15×16	—	16	1	4	15×16
残差块2	上采样层	6×64	LeakyReLU	32	2	4	12×32	残差块2	下采样层	15×16	LeakyReLU	32	2	4	12×32
	卷积层	12×32	—	32	1	4	12×32		卷积层	12×32	—	32	1	4	12×32
残差块3	上采样层	12×32	LeakyReLU	16	2	4	15×16	残差块3	下采样层	12×32	LeakyReLU	64	2	4	6×64
	卷积层	15×16	—	16	1	4	15×16		卷积层	6×64	—	64	1	4	6×64
卷积层	15×8	Tanh	1	1	4	15×1	全连接层	426	—	—	1	—	1		

表 7 不同数据增强方法的实验结果  
Table 7 Experimental results of different data enhancement methods

数据增强方法	$A$	$R_e$	$P_r$	$G_{\text{mean}}$	$F_1$
$O_{\text{in}}$	0.790±0.032	0.941±0.02	0.721±0.102	0.790±0.071	0.818±0.12
GAN	0.912±0.068	0.940±0.011	0.730±0.037	0.933±0.013	0.820±0.025
WGAN	0.920±0.021	0.943±0.017	0.724±0.02	0.934±0.024	0.813±0.01
WGAN-gp	0.931±0.018	0.936±0.017	0.845±0.054	0.947±0.02	0.892±0.065
WGAN-div	<b>0.965±0.036</b>	<b>0.962±0.104</b>	<b>0.963±0.009</b>	<b>0.961±0.041</b>	<b>0.962±0.064</b>

后,  $A$  提升了 17.5%,  $R_e$  提升了 2.1%,  $P_r$  提升了 24.2%,  $G_{\text{mean}}$  提升了 17.1%,  $F_1$  分数提升了 14.4%。由此说明, 利用 WGAN-div 模型对不平衡的故障数据进行增强, 能够有效提高原始数据的质量, 进而提高分类器的判别性能。

(3) 使用 GAN、WGAN、WGAN-gp 进行数据增强后, 虽然准确率和 G-mean 指标均增大, 分类模型对故障分支的识别能力增强, 但是在  $F_1$  分数上却没有明显的改进, 分析认为模型扩充了劣质的新故障样本, 影响了分类模型对故障分支诊断的判别。相较

于 GAN、WGAN、WGAN-gp 模型, WGAN-div 模型各项评价指标均为最高,  $A$ 、 $R_e$ 、 $P_r$ 、 $G_{\text{mean}}$  和  $F_1$  分数分别为 96.5%、96.2%、96.3%、96.1% 和 96.2%, 大幅度提高了分类模型对故障分支的识别能力, 验证了所提 WGAN-div 模型在处理不平衡数据时的优越性。

应用 t-分布随机领域嵌入 (t-Stochastic Neighbor embedding, t-SNE) 算法对 WGAN-div 模型的样本生成情况进行降维可视化分析, 图 9 展示了迭代次数  $N$  分别 0、100、200、500、800、1 000 时模型的生成样本与真实样本之间的分布情况, 图 10 展示模型损失函

数的变化情况。观察图 9、10, 随着迭代次数的增加, WAGN-div 模型的损失函数稳定收敛、逐渐平稳, 生成的新样本数据与真实数据分布逐渐交融, 生成数据与真实数据具有很好的相似性, 生成数据的质量越来越高。

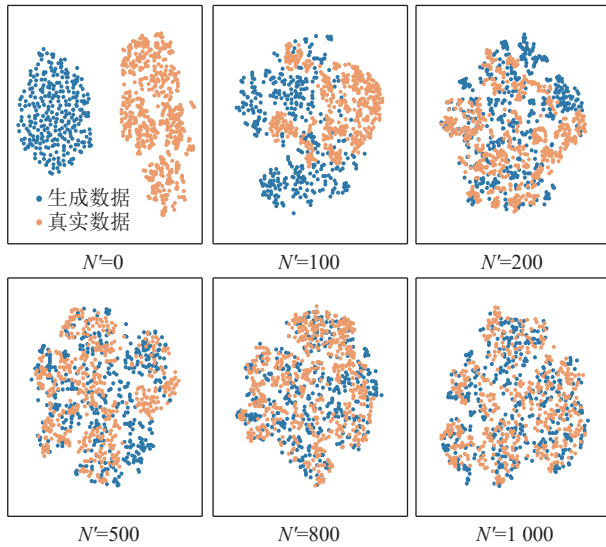


图 9 t-SNE 降维数据可视化

Fig.9 t-SNE dimension reduction data visualization

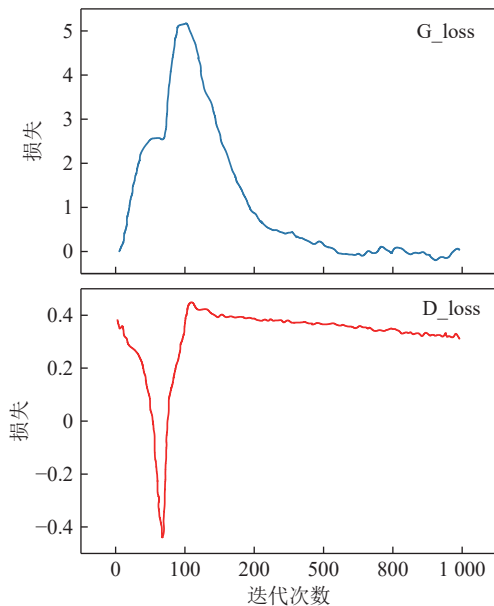


图 10 WGAN-div 损失函数

Fig.10 WGAN-div loss function

### 4.3 RF 有效性验证

为了验证 RF 模型能够更有效的对通风系统故障分支进行辨别, 原始样本经过 WGAN-div 处理后, 选用了以下经典的集成学习分类模型进行对比: 类别提升树 (CBT)、轻量梯度提升树 (LGB)、梯度提升树 (GBDT), 此外, 将文献[7]中提出的通风系统故障诊断

SVM 模型也纳入本文的对比实验。各个模型的最优参数见表 8, 各个参数定义见表 2。

表 8 分类模型最优参数

Table 8 Optimal parameters of classification model

分类模型	参数设置
CBT	$n=300, T_d=3, \eta=0.03$
LGB	$M=100, T_d=3, \eta=0.01, T_f=1, \omega=0.8$
GBTD	$M=200, T_d=2, M_p=2, \eta=0.1, M_f=2$
SVM	$c=100, K=RBF$
RF	$M=50, T_d=12, M_p=2, M_f=1$

为了考察不同的数据生成比率下分类模型的表现是否具有明显改善, 本文将 WGAN-div 的数据生成比率分别调整为 10%、20%、50%、80%、100%。图 11 展示了基于 WGAN-div 不同数据生成比率下各分类模型的实验结果 (10 次实验的平均值), 分析如下:

(1) 从数据生成比率的角度来看, 相较于原始数据集, 数据生成比率为 10% 时, 所有模型  $R_e$ 、 $P_r$ 、 $G_{mean}$  和  $F_1$  分数平均提高了 2.7%、0.3%、1.8% 和 1.5%, 模型性能提升不明显。但是, 当数据生成比率达到 50% 时, 所有模型的  $R_e$ 、 $P_r$ 、 $G_{mean}$  和  $F_1$  分数平均提高了 19.8%、1.18%、13.4% 和 12%, 所有分类模型的性能提升明显。新数据进一步生成达到 80% 时, 模型表现的改进相对有限, 即使数据生成比率达到 100% 时, 各模型的性能达到最优, 但是相对于 50% 的生成比率模型性能提升并不优越, 因此当矿井通风系统故障分支较多时, 考虑时间成本可以将数据生成比率设置为 80%~100%。

(2) 从分类模型的角度来看, RF 模型无论是在原不平衡数据集还是增广数据集上都表现出明显的优势。当样本数据达到完全平衡时, 相较于原始数据集, RF 模型在  $R_e$ 、 $P_r$ 、 $G_{mean}$  和  $F_1$  评价指标上分别提升了 21.9%、2.7%、11.8%、11.2%。在所有的分类模型中, 传统的机器学习模型 SVM 性能要明显弱于集成学习模型, 尽管 SVM 模型在  $F_1$  指标上的表现可以通过数据增强得到显著改善, 但是其  $G_{mean}$  指标并未随着数据的平衡而明显改进, 分析认为数据增强生成的伪样本具有一定的随机性, 导致 SVM 表现不够稳定。特别地, 当扩充数据集达到平衡时, 与传统的矿井通风系统 SVM 故障诊断方法相比, RF 模型在  $R_e$ 、 $P_r$ 、 $G_{mean}$  和  $F_1$  指标上分别提高了 4.7%、2.3%、10.1%、3.5%。总的来说, 本文所提 RF 模型适用于矿井通风系统故障诊断, 当训练样本逐渐达到平衡时, RF 模型在  $A$ 、 $R_e$ 、 $P_r$ 、 $G_{mean}$  和  $F_1$  得分上的表现较其他模型更具优势。

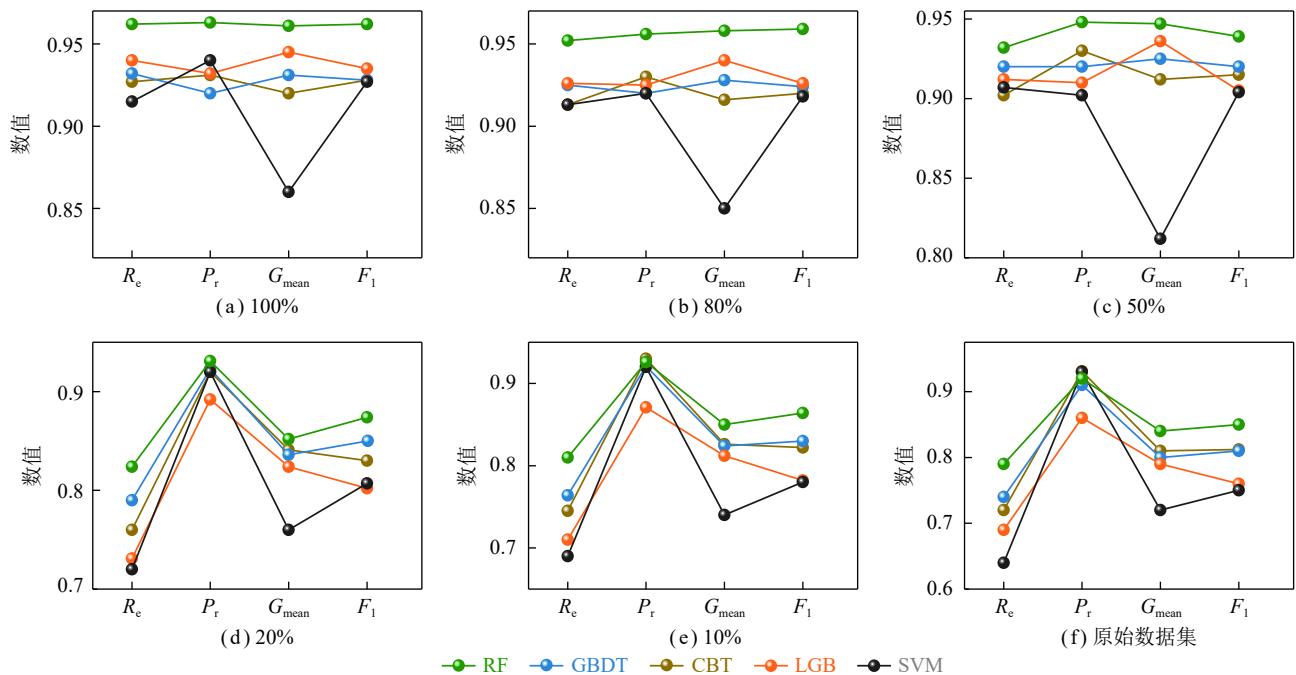


图11 不同数据生成比率下各分类模型的实验结果

Fig.11 Experimental results of different classification models at different data generation rate

课题组在东山矿进行了2次现场工业应用试验,考虑到生产安全和矿山的实际情况,通过打开关闭状态下的风门方式进行故障模拟,第1次实验在确保东山矿安全生产的前提下打开了西采区6D上左一回风巷的风门(图8中的33号分支),风门开启后采集该矿15个传感器的风速值(取风门开启后5 min内各个传感器的平均值),把15个采集到的风速作为输入值,利用预测模型对故障分支进行预测,预测结果输出为33。第2次实验打开了西采区3号上右一巷的风门(图8中的85号分支),将风门开启后采集到的15个风速传感器值输入模型,预测结果输出为85,2次试验故障分支预测结果与工业试验结果一致。

## 5 结 论

(1) 从矿井通风系统实际工况下各分支故障概率不同的角度出发,以简单的T型通风网络为例,说明了不平衡数据集对故障诊断模型的影响。建立了WGAN-div-RF故障诊断模型,有效解决了通风系统故障数据不平衡的问题,从数据层面提高了分类模型的特征提取能力,进而提高分类模型的性能。

(2) 故障诊断实验以及t-SNE可视化结果表明,加入残差块的WGAN-div模型能够生成高质量的新数据实现对样本集的扩充,WGAN-div模型的 $A$ 、 $R_e$ 、 $P_r$ 、 $G_{mean}$ 和 $F_1$ 分数分别为96.5%、96.2%、96.3%、96.1%和96.2%,相较于其他数据增强模型在处理不平衡数据时更具优越性。

(3) 针对通风系统故障诊断高维多分类问题,结合集成学习中的投票机制对通风网络分支进行分类,所得结果要优于传统的SVM模型,其中RF模型在不同数据生成比率上各评价指标得分较其他集成模型更具优势。

## 参考文献(References):

- [1] 刘剑. 矿井智能通风关键科学问题综述[J]. 煤矿安全, 2020, 51(10): 108-111, 117.  
LIU Jian. Overview on key scientific and technical issues of mine intelligent ventilation[J]. Safety in Coal Mines, 2020, 51(10): 108-111, 117.
- [2] 范京道, 李川, 闫振国. 融合5G技术生态的智能煤矿总体架构及核心场景[J]. 煤炭学报, 2020, 45(6): 1949-1958.  
FAN Jingdao, LI Chuan, YAN Zhenguo. Overall architecture and core scenario of a smart coal mine incorporating 5G technology ecology[J]. Journal of China Coal Society, 2020, 45(6): 1949-1958.
- [3] JIA J, JIA P, LI Z. Theoretical study on stability of mine ventilation network based on sensitivity analysis[J]. Energy Science & Engineering, 2020, 8(8): 2823-2830.
- [4] ELSISI M, TRAN M, MAHMOUD K, et al. Effective IoT-based deep learning platform for online fault diagnosis of power transformers against cyberattacks and data uncertainties[J]. Measurement, 2022, 190: 110686.
- [5] 孟宗, 关阳, 潘作舟, 等. 基于二次数据增强和深度卷积的滚动轴承故障诊断研究[J]. 机械工程学报, 2021, 57(23): 106-115.  
MENG Zong, GUAN Yang, PAN Zuozhou, et al. Fault diagnosis of rolling bearing based on secondary data enhancement and deep convolutional network[J]. Journal of Mechanical Engineering, 2021, 57(23): 106-115.

- [6] 戴金玲, 许爱强, 申江江, 等. 基于 OCKELM 与增量学习的在线故障检测方法[J]. 航空学报, 2022, 43(3): 378–389.  
DAI Jinling, XU Aiqiang, SHEN Jiangjiang, et al. Online fault detection method based on kernel incremental learning and OCKELM[J]. Acta Aeronautica et Astronautica Sinica, 2022, 43(3): 378–389.
- [7] 刘剑, 郭欣, 邓立军, 等. 基于风量特征的矿井通风系统阻变型单故障源诊断[J]. 煤炭学报, 2018, 43(1): 143–149.  
LIU Jian, GUO Xin, DENG Lijun, et al. Resistance variant single fault source diagnosis of mine ventilation system based on air volume characteristic[J]. Journal of China Coal Society, 2018, 43(1): 143–149.
- [8] 刘剑, 尹昌胜, 黄德, 等. 矿井通风阻变型故障复合特征无监督机器学习模型[J]. 煤炭学报, 2020, 45(9): 3157–3165.  
LIU Jian, YIN Changsheng, HUANG De, et al. Unsupervised machine learning model for resistant variant fault diagnosis of mine ventilation system with composite features[J]. Journal of China Coal Society, 2020, 45(9): 3157–3165.
- [9] HUANG D, LIU J, DENG L, et al. An adaptive kalman filter for online monitoring of mine wind speed[J]. Archives of Mining Sciences, 2019, 64(4): 813–827.
- [10] HUANG D, LIU L, DENG L. A hybrid-encoding adaptive evolutionary strategy algorithm for windage alteration fault diagnosis[J]. Process Safety and Environmental Protection, 2020, 136: 242–252.
- [11] 黄德, 刘剑, 刘永, 等. 矿井通风阻变故障观测特征组合选择试验研究[J]. 煤炭学报, 2021, 46(12): 3922–3933.  
HUANG De, LIU Jian, LIU Yong, et al. Experimental research on combination selection of observation feature of resistance variation fault in mine ventilation[J]. Journal of China Coal Society, 2021, 46(12): 3922–3933.
- [12] 周启超, 刘剑, 刘丽. 基于 SVM 的通风系统故障诊断惩罚系数与核函数系数优化研究[J]. 中国安全生产科学技术, 2019, 15(4): 45–51.  
ZHOU Qichao, LIU Jian, LIU Li. Research on fault diagnosis penalty coefficient and kernel function coefficient optimization of ventilation system based on SVM[J]. Journal of Safety Science and Technology, 2019, 15(4): 45–51.
- [13] 倪景峰, 李振, 乐晓瑞, 等. 基于随机森林的阻变型通风网络故障诊断方法[J]. 中国安全生产科学技术, 2022, 18(4): 34–39.  
NI Jingfeng, LI Zhen, LE Xiaorui, et al. Resistance variant fault diagnosis method of ventilation network based on random forest[J]. Journal of Safety Science and Technology, 2022, 18(4): 34–39.
- [14] 倪景峰, 乐晓瑞, 常立峰, 等. 基于决策树的矿井通风阻变型故障诊断及传感器优化布置[J]. 中国安全生产科学技术, 2021, 17(2): 34–39.  
NI Jingfeng, LE Xiaorui, CHANG Lifeng, et al. Resistance variant fault diagnosis and optimized layout of sensors for mine ventilation based on decision tree[J]. Journal of Safety Science and Technology, 2021, 17(2): 34–39.
- [15] 张浪, 张迎辉, 张逸斌, 等. 基于机器学习的通风网络故障诊断方法研究[J]. 工矿自动化, 2022, 48(3): 91–98.  
ZHANG Lang, ZHANG Yinghui, ZHANG Yibin, et al. Research on fault diagnosis method of ventilation network based on machine learning[J]. Journal of Mine Automation, 2022, 48(3): 91–98.
- [16] ZHAO D, SHEN Z. Study on roadway fault diagnosis of the mine ventilation system based on improved SVM[J]. Mining, Metallurgy & Exploration, 2022, 39(3): 983–992.
- [17] WANG D, LIU J, DENG L, et al. Intelligent diagnosis of resistance variant multiple fault locations of mine ventilation system based on ML-KNN[J]. PloS One, 2022, 17(9): e0275437.
- [18] LIU L, LIU J, ZHOU Q, et al. Machine learning algorithm selection for windage alteration fault diagnosis of mine ventilation system[J]. Advanced Engineering Informatics, 2022, 53: 101666.
- [19] 赵丹, 沈志远, 刘晓青. 基于 OCISVM 的矿井通风系统在线故障诊断[J]. 中国安全科学学报, 2022, 32(10): 76–82.  
ZHAO Dan, SHEN Zhiyuan, LIU Xiaoqing. Online fault diagnosis of mine ventilation system based on OCISVM[J]. China Safety Science Journal, 2022, 32(10): 76–82.
- [20] ARJOVSKY M, CHINTALA S, BOTTOU L. Wasserstein GAN adversarial networks[C]//Proceedings of the 34th International Conference on Machine Learning (ICML). Australia: PMLR 70, 2017: 214–223.
- [21] GULRAJANI I, AHMED F, ARJOVSKY M, et al. Improved training of wasserstein GANs[C]//Proceedings of the 31st International Conference on Neural Information Processing Systems. United States: Curran Associates Inc, 2017: 5767–5777.
- [22] WU J, HUANG Z, THOMA J, et al. Wasserstein divergence for GANs[C]//Proceedings of the 15th European Conference on Computer Vision. Germany: Springer, 2018: 653–668.
- [23] HE K, ZHANG X, REN S, et al. Deep residual learning for image recognition[C]//Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition. United States: IEEE, 2016: 770–778.
- [24] 殷豪, 丁伟锋, 陈顺, 等. 基于生成对抗网络和纵横交叉粒子群算法的光伏数据缺失重构方法[J]. 电网技术, 2022, 46(4): 1372–1381.  
YIN Hao, DING Weifeng, CHEN Shun, et al. Reconstruction method for missing data in photovoltaic based on generative adversarial network and crisscross particle swarm optimization algorithm[J]. Power System Technology, 2022, 46(4): 1372–1381.
- [25] MARINA S, GUY L. A systematic analysis of performance measures for classification tasks[J]. Information Processing & Management, 2009, 45(4): 427–437.