考虑样本类别不平衡的电网故障事件智能识别方法

卫志农1,石东明1,张明2,孙国强1,臧海祥1,沈培锋2

(1. 河海大学 能源与电气学院,江苏 南京 211100;2. 国网江苏省电力有限公司南京供电分公司,江苏 南京 210019)

摘要:电网中不同设备的故障概率存在差异,影响智能诊断技术的准确性。为解决此问题,提出了一种基于 代价敏感学习和模型自适应选择融合的电网故障事件智能识别方法。首先,利用Word2vec模型将预处理后 的电网告警信息向量化,并搭建2个双向长短期记忆网络作为基础分类器;然后,设计代价敏感损失函数,将 交叉熵损失函数与代价敏感损失函数分别应用于2个分类器中;最后,提出一种模型自适应选择融合法,融 合上述分类器,得到故障事件识别结果。实际数据测试表明,所提方法能够有效降低故障事件识别中样本类 别不平衡的影响。

0 引言

电网运行状态异常或发生故障时,监控系统将 产生大量中文文本形式的告警信息。调度人员难以 快速准确判别对应的事件类型,而基于人工智能的 故障诊断技术能通过对监控信息的推理分析实现故 障事件的自主识别^[1],有效缩短异常事件判别时间, 并提升后续事件处理效率,提高电网运行管理水平。

自然语言处理技术和机器学习的深入应用使 计算机能够学习数字化表达后的告警信息,并挖掘 海量数据中的特征,从而使电网智能告警逐渐摆脱 对人工经验的依赖[2]。同时,深度学习作为机器学 习的重要分支,通过扩展神经元层的方式构建更为 深层的神经网络,可以深入挖掘输入的电力数据中 的隐含关键特征。文献[3]构建了基于卷积神经网 络CNN(Convolution Neural Network)的电网假数据 注入攻击检测模型;文献[4]利用双向长短期记忆 网络Bi-LSTM(Bidirectional Long-Short-Term Memory network)建立了底层量测数据与电力系统暂态稳定 类别之间的非线性映射关系。上述深度学习模型具 有较好的泛化能力,但需要足量样本支撑模型训练。 电网中不同设备故障发生率存在差异,导致部分故 障样本量偏少,因此历史故障样本中存在类别不均 衡现象,不利于智能诊断系统的模型训练与参数学 习过程,影响事件识别结果。

目前,关于不平衡数据集的处理方法主要分为 数据预处理法和分类法2种。数据预处理法通过合 成或丢弃一定数量样本,降低各类别样本量的差距, 如单一的欠采样、过采样^[5-6],以及结合2种方法的混

收稿日期:2020-12-22;修回日期:2021-05-24

基金项目:国家电网公司科技项目(SGJSNJ00FCJS1800810) Project supported by the Science and Technology Project of State Grid Corporation of China(SGJSNJ00FCJS1800810)

合采样[7],该类方法改变了数据分布,一定程度上破 坏了样本特征信息。分类法能够保留样本全部初始 信息,包括代价敏感学习和集成学习。代价敏感学 习通过引入代价敏感因子,增大模型训练过程中对 少类别样本的错分代价,从而提高该类别样本的分 类可靠性。文献[8]直接将错分代价嵌入神经网络, 以降低各类别样本的平均错分代价;文献[9]提出了 一种基于代价敏感学习的决策树剪枝方法,在剪 枝阶段引入代价敏感的思想,使模型总损失值达 到最小;文献[10]通过对不同类别设置不同的代价 因子,得到总代价最小的支持向量机 SVM(Support Vector Machine)分类器,文献[11]在此基础上,将 SVM 核函数作为选取特征的标准,进一步提高了 SVM算法对不平衡数据的分类准确率。上述方法在 改善对少类别样本分类效果的同时,会影响多类别 样本的判别结果,不能有效提升模型的整体性能。 集成学习可以将多个子分类模型(下文简称子模型) 进行融合,从而得到一个整体性能较好的分类器。 Boosting、Bagging和Stacking算法^[12-13]通过不同方式 实现模型融合,但只适用于弱分类器。模型融合是 一种整合多个强分类器的集成学习方法,目前常用 的有最大值法、均值法、求和法等[14],此类方法根据 子模型计算出的各类别后验概率或结果标签,采用 特定公式进行模型融合。但这种对各类别样本分类 结果进行无差别融合的方法,原理较为简单,无法整 合子模型的优势。

针对上述方法的特点、局限性,本文以Bi-LSTM 为基础分类器,提出一种基于代价敏感学习和模型 自适应选择融合的多分类问题处理方法,在提高少 类别样本的分类精度的同时,保持对多类别样本的 准确分类。针对某市电网公司调度中心的告警信息 的测试结果表明,本文方法对于各类故障均具有良 好的判别结果,进一步验证了其在电网故障事件识 别中的优越性和可靠性。

1 Bi-LSTM 原理

CNN 和循环神经网络 RNN (Recurrent Neural Network)是目前应用最为成熟、广泛的2种深度学习模型。RNN考虑输入信息中的序列特征,擅长处理时序信息,Bi-LSTM 通过改进 RNN,解决了 RNN模型训练中梯度消失与梯度爆炸的问题,并结合当前输入前、后时刻的隐含信息,进一步提高了 RNN 对时序信息的挖掘能力。因此本文采用 Bi-LSTM 作为基础分类器,完成对内部具有自然时序关系的电网告警信息的处理。

Bi-LSTM的结构单元包含输入、长短期记忆网络LSTM(Long Short Term Memory network)链、输出3个部分,其中LSTM链由2个反向LSTM拼接而成,该网络结构包括输入门、遗忘门、记忆单元和输出门,具体结构见附录A图A1。

输入门对当前时刻的网络输入信息进行控制, 通过 Sigmoid 神经网络层和 tanh 层计算当前输入中 保存到记忆单元的信息,如式(1)、(2)所示。

$$\boldsymbol{i}_{t} = \boldsymbol{\sigma} \left(\boldsymbol{W}_{i} \left[\boldsymbol{h}_{t-1}, \boldsymbol{x}_{t} \right] + \boldsymbol{b}_{i} \right)$$
(1)

$$\tilde{\boldsymbol{C}}_{t} = \tanh\left(\boldsymbol{W}_{c}\left[\boldsymbol{h}_{t-1}, \boldsymbol{x}_{t}\right] + \boldsymbol{b}_{c}\right)$$
(2)

式中: i_t 、 \tilde{C}_t 分别为t时刻(当前时刻)输入门、临时记 忆单元的状态; W_i 、 W_e 分别为输入门、临时记忆单元 的权值矩阵; h_{t-1} 、 x_t 分别为t-1时刻(前一时刻)隐含 层的输入、t时刻的输入; b_i 、 b_e 分别为输入门、临时记 忆单元的偏置; $\sigma(\cdot)$ 为Sigmoid激活函数。

遗忘门保存长期重要信息,按式(3)计算*t*-1时 刻隐含层中能够保留在当前时刻记忆单元的信息。

$$f_{t} = \sigma \left(\boldsymbol{W}_{f} \left[\boldsymbol{h}_{t-1}, \boldsymbol{x}_{t} \right] + \boldsymbol{b}_{f} \right)$$
(3)

式中: f_i 为t时刻遗忘门的状态; W_f 、 b_f 分别为遗忘门的权值矩阵和偏置。

遗忘门保留序列数据的长期重要信息,输入门 临时记忆单元使得当前时刻的无用信息不进入记忆 单元,两者按式(4)共同决定记忆单元保存的信息。

$$\boldsymbol{C}_{t} = \boldsymbol{f}_{t} \odot \boldsymbol{C}_{t-1} + \boldsymbol{i}_{t} \odot \boldsymbol{\tilde{C}}_{t}$$

$$\tag{4}$$

式中:*C_t*、*C_{t-1}分别为t*时刻和*t*-1时刻记忆单元的输出值; ②表示按元素相乘。

输出门由当前时刻的输入、记忆单元和前一时 刻的隐含层确定。

$$\boldsymbol{O}_{t} = \boldsymbol{\sigma} \left(\boldsymbol{W}_{o} [\boldsymbol{h}_{t-1}, \boldsymbol{x}_{t}] + \boldsymbol{b}_{o} \right)$$
(5)

$$\boldsymbol{h}_{t} = \boldsymbol{O}_{t} \odot \tanh\left(\boldsymbol{C}_{t}\right) \tag{6}$$

式中: O_t 、 h_t 分别为t时刻输出门、LSTM的输出; W_o 、 b_o 分别为输出门的权值矩阵和偏置。

Bi-LSTM结合2个时序相反的LSTM,构成了结构单元中的LSTM链,能够同时获取当前输入前、后

时刻的特征信息,其单元结构见附录A图A2。

Bi-LSTM 输出层的输出矩阵 H_i 由正向 LSTM 的输出 \hat{h}_i 和反向 LSTM 的输出 \hat{h}_i 拼接而成。

$$\boldsymbol{H}_{t} = [\,\vec{\boldsymbol{h}}_{t}, \,\vec{\boldsymbol{h}}_{t}\,] \tag{7}$$

*H*₁经过激活函数运算后即可得到样本属于各类别的概率,默认取概率最大的类别作为计算结果。

2 基于代价敏感学习和模型自适应选择融合的电网故障识别方法

电网告警信息为中文文本形式,此类非结构化 的文本数据需要转化为结构化的数字表达,才能输 入 Bi-LSTM 模型训练学习。本文采用 Word2vec 模 型训练得到告警数据的分布式向量。Word2vec 是 一款由谷歌于 2013 年公开开源的词向量计算工 具^[15],其基本思想是通过神经网络将每个词映射成 固定维数的实数向量,所有向量构成蕴含语义信息 的词向量空间,不同词向量在该空间中的距离可以 表征词语之间的语义相似性。词向量训练完成后, 计算单条告警信息中所有词向量的平均值,得到固 定维数的故障样本句向量。

2.1 方法流程

传统 Bi-LSTM 模型更趋向于将样本判为训练集 数量多的类别,以减小损失值。本节提出一种基于 代价敏感学习和模型自适应选择融合的电网故障事 件识别方法,其能够显著降低样本类别不均衡对电 网故障事件识别结果的影响。电网故障事件识别的 流程如附录A图A3所示,具体步骤如下:

1)利用 Word2vec 模型将分词后的电网告警信 息转化为高维向量,并求均值得到告警数据句向量, 向量维度设置为300,向量化过程如图1所示;

2)构建传统深度学习模型,即采用交叉熵损失 函数的Bi-LSTM,输入故障样本进行监督训练并调 参,得到对大样本故障类别具有较好识别率的子 模型1;

3)自定义一个多分类代价敏感损失函数,代替 模型1中的交叉熵损失函数,增大模型训练过程中 对小样本的错分代价,其余过程同步骤2),得到能 够准确识别小样本故障的子模型2;

4)将每例故障样本输入子模型1、2进行判别 后,采用模型自适应选择融合方法对判别结果进行 融合,得到最终的故障事件识别结果并输出。

2.2 多分类代价敏感损失函数

传统的损失函数对所有类别的样本设置相同的 错分权重,因此少类别样本的损失易被淹没。本文 基于 Lin Tsung-yi等人提出的焦点损失函数^[16],构 建适用于多分类问题的代价敏感损失函数γ_{FL},如式 (8)所示。



图1 电网告警信息向量化过程

Fig.1 Vectorization process of power grid warning information

$$\gamma_{\rm FL} = \beta L + (1 - \beta) \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} \alpha_j \gamma_{ij} (1 - p_{ij})^2 \ln p_{ij} \qquad (8)$$

$$L = -\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} y_{ij} \ln p_{ij}$$
(9)

$$\alpha_{j} = \frac{k_{j}}{\sum_{i=1}^{m} k_{i}}, \ k_{j} = \frac{\sum_{i=1}^{N_{i}} N_{i}}{N_{j}}$$
(10)

式中:m和n分别为样本类别数和样本总数; y_{ij} 和 p_{ij} 分别为样本i属于类别j的真实概率和预测概率; $\beta \in [0,1]$,为调制因子;L为交叉熵损失函数; $\alpha_j \in [0,1]$,为权重因子,能够区分不同类别样本的错 分代价,样本量越大,该类别样本的错分代价越小, 否则错分代价越大; N_i 为属于类别j的样本的数量。

 $\gamma_{\rm FL}$ 由两部分组成,第一部分为传统交叉熵损失 函数L,第二部分为考虑类别不平衡影响的代价敏 感损失值计算。通过调制因子 β 调节两者权重, β 越 小,第二部分占比越大, $\gamma_{\rm FL}$ 对各类别样本的区分程 度越高。作为一种代价敏感损失函数, $\gamma_{\rm FL}$ 通过对各 类别样本设置不同的权重因子,提高对少类别样本 的错分代价,从而提高该类样本的分类准确性。

2.3 模型自适应选择融合方法

训练样本不平衡度较大时,γ_{EL}中少类别样本的 错分代价过大,破坏了模型对多类别样本的分类效 果。本节提出一种综合考虑召回率与准确率的模型 自适应选择融合方法,在代价敏感学习的基础上进 一步改善模型的整体分类性能。该方法首先以样本 类别为出发点,选择召回率大的子模型代表该类别 样本的分类标准,使得模型融合后能够尽可能全面 地识别出此类别样本;再结合子模型对各类别样本 的分类准确率,推理得到最终的判别结果,从而降低 模型融合后的整体误判率。该方法的流程图见附录 B图B1。以样本总数为n、样本类别为m、子模型个 数为2为例,模型融合的具体过程如下。

1)计算子模型*k*(*k*=1,2)对类别*j*(*j*=1,2,…,*m*) 样本的分类召回率*R_{ij}*,如式(11)所示。对于每个样 本类别,选择分类召回率大的子模型作为分类基准, 由此设定各类别的融合标签*σ_i*,如式(12)所示。

$$R_{kj} = \frac{\sum_{i=1}^{n} I(j = f_k(x_i), y_i = f_k(x_i))}{\sum_{i=1}^{n} I(y_i = f_k(x_i))}$$
(11)
$$\sigma_j = \begin{cases} 1 & R_{1j} \ge R_{2j} \\ 0 & R_{1j} < R_{2j} \end{cases}$$
(12)

式中: $f_k(x_i)$ 为子模型k对样本 x_i 的预测标签; y_i 为样本 x_i 的真实标签; $I(\cdot)$ 为逻辑判断,括号内表达式成立时取1,否则取0。

2)对于类别 j 样本,结合 σ_j 取该类别样本分类 召回率较大的子模型,按照式(13)计算类别 j 样本 的分类准确率,将其作为准确率矩阵 Δ 的第 j 个元 素,由此得到按分类召回率大小筛选出的准确率矩 阵 Δ 如式(14)所示。

$$P_{kj} = \frac{\sum_{i=1}^{n} I(j=f_k(x_i), y_i=f_k(x_i))}{\sum_{i=1}^{n} I(j=f_k(x_i))}$$
(13)

$$\boldsymbol{\Delta} = \left[P_{\sigma_1 1} \ P_{\sigma_2 2} \ \cdots \ P_{\sigma_m m} \right] \tag{14}$$

式中: P_{k_j} 为子模型k对类别j样本的分类准确率; $P_{\sigma_j j}$ 为结合 σ_j 选取的召回率较大的子模型对类别j样本的分类准确率。

3)根据子模型分类结果,按照式(15)设置各样 本的融合标签。

$$\boldsymbol{\omega}_{ij} = \begin{cases} 1 \quad f_{\sigma_j}(\boldsymbol{x}_i) = j \\ 0 \quad f_{\sigma_j}(\boldsymbol{x}_i) \neq j \end{cases}$$
(15)

式中: ω_{ij} 为样本 x_i 对类别j的融合标签; $f_{\sigma_j}(x_i)$ 为结 合 σ_j 选取的召回率较大的子模型对样本 x_i 的分类 结果。

σ_i由式(15)计算得到,反映了能够代表类别 *j* 样本分类结果的子模型标签,若该标签对应的子模 型对样本*x*_i的分类结果与类别 *j* 一致,则将*x*_i对类别 *j*的融合标签设置为1,否则为0。在此基础上按照 式(16)计算融合后样本*x*_i属于各类别的后验概率。

$$\boldsymbol{\Pi}_{i} = [\omega_{i1}P_{\sigma_{1}1} \ \omega_{i2}P_{\sigma_{2}2} \ \cdots \ \omega_{im}P_{\sigma_{m}m}]$$
(16)
中: **\Pi**_i为由后验概率组成的矩阵,其第 *j* 列表示样

式

本x_i属于类别j的概率。Ⅱ_i中最大值对应的列索引 即模型融合的输出类别标签。Ⅱ_i=0时,取分类效 果较好的子模型的分类标签作为输出结果(默认为 子模型1)。模型融合后的输出结果表达式为:

$$f(x_i) = \begin{cases} f_1(x_i) & \boldsymbol{\Pi}_i = 0\\ j & \boldsymbol{\Pi}_i \neq 0, \ \boldsymbol{\omega}_{ij} \boldsymbol{P}_{\sigma_j j} = \max(\boldsymbol{\Pi}_i) \end{cases}$$
(17)

式中: $max(\Pi_i)$ 为 Π_i 中的最大值。

模型自适应选择融合方法依次考察子模型的召回率与准确率指标,在分析子模型分类性能的基础上进行决策,整合各子模型的分类优势,得到最终的输出结果,实现了模型的选择性融合与信息互补,同时可推广应用于2个以上子模型参与融合的场景。

3 算例分析

为验证本文方法有效性,选取某市电网调度中 心2016、2017年的历史告警信息进行算例分析。首 先根据工程需要,确定了若干种需要调控人员第一 时间重点关注的异常跳闸类事件,然后以带关键词 "分闸"的告警信息为标志,提取该信息前后一段时 间窗内的离散告警信息集合,当满足一定规则时,构 成各类标签化事件样本。从中提取9种重要故障事 件对应的样本,共得到13554例故障事件样本。从 每类故障事件样本中随机选取25例作为测试集,其 余作为训练集,并在训练过程中随机抽取训练集中 5%的样本作为验证样本,以优化模型参数。每组实 验取10次测试结果的平均值作为参考标准。故障 事件样本分布情况如表1所示。

表1 故障事件样本数量统计

1401	e i Number statist	les of fault	event sam	pies
事件类别	事件名称	训练集数量	测试集数量	总样本量
1	单相瞬时故障	4760	25	4785
2	单相永久故障	3 2 3 0	25	3 2 5 5
3	相间故障	2673	25	2698
4	电抗器 / 电容器故障	1 5 8 1	25	1606
5	所用接地变故障	346	25	371
6	主变电气量故障	317	25	342
7	主变本体重瓦斯故障	239	25	264
8	主变调压重瓦斯故障	128	25	153
9	母线故障	55	25	80
	合计	13329	225	13554

分类模型常用的评价指标有召回率、准确率、F1 值。召回率、准确率计算公式分别见式(11)、(13), 子模型*k*属于类别*j*样本的F1值的计算公式为:

$$F_{kj} = \frac{2P_{kj}R_{kj}}{P_{kj} + R_{kj}}$$
(18)

F1值是一种综合考量准确率与召回率的综合 评价指标,通常F1值越大,模型的分类性能越好。 对于多分类模型,取所有类别的F1值的期望作为该 模型的整体F1值指标。 经过测试对比,Word2vec模型和Bi-LSTM模型的参数设置情况分别见附录C表C1、C2。

3.1 基础分类器性能验证

为了验证 Bi-LSTM 在电网故障事件识别中的优 越性,设置3组对比实验,分别采用以CNN、LSTM 以及 结合 CNN 与注意力(Attention)机制的组合深度学习 模型 Attention-CNN 作为基础分类器。其中 CNN 设 置3种卷积窗口,尺寸分别为3、4、5,每种窗口的卷 积核数目为100,采用 ReLU 激活函数,其他所需参 数同附录 C表 C2;LSTM 的参数同附录 C表 C2。以 不同深度学习模型作为基础分类器,对算例进行实 验对比,得到准确率、召回率、F1值3种评价指标,结 果如图2所示。



图2 深度学习模型的评价指标对比



由图2可以看出:CNN虽然具有局部感知能力 强的特点,能够很好地处理图像信息,但在处理时序 信息时效果欠佳;Attention-CNN在CNN的基础上引 入注意力机制,能够强化局部告警信息中蕴含的关 键特征权重,以优化模型对不同的告警事件的特征 提取,但依然无法捕捉时序关联特征,导致模型总体 性能提升不大;LSTM擅长处理时序信息,电网告警 信息属于时间相关的数据,因此分类效果比CNN更 好;Bi-LSTM模型的准确率、召回率与F1值均最大, 进一步体现了Bi-LSTM基于LSTM进行的改进能够 考虑当前输入的前、后时刻的信息,优化分类效果, 作为基础分类器的性能优于其他3种对比模型。后 续实验均以Bi-LSTM模型作为基础分类器。

3.2 模型融合方法性能验证

子模型1采用交叉熵损失函数,子模型2采用由 式(8)构建的代价敏感损失函数(β=0.1)。为对比 本文的模型自适应选择融合方法(简称选择法)的实 用性,分别利用最值法、求和法对子模型进行融合。 对于每个样本,最值法取各子模型中最大后验概率 对应的类别标签作为融合结果;均值法计算所有子 模型后验概率的均值,得到融合后的后验概率,并将 最大概率对应的类别标签作为最终输出结果。子模 型与不同模型融合方法的分类召回率如表2所示, 整体评价指标对比如图3所示。

表2 子模型与融合算法的分类召回率

Table 2	Classification	recall	rate	of	submod	lels
	and fusio	n met	hode			

und fusion methods						
类别	分类召回率 / %					
	子模型1	子模型 2	最值法	求和法	选择法	
1	94.32	53.87	94.32	94.32	92.93	
2	95.93	69.21	95.93	95.93	93.68	
3	97.10	52.34	97.10	97.10	91.41	
4	99.26	93.62	99.26	99.26	99.15	
5	96.96	97.39	96.96	96.96	97.62	
6	95.79	98.87	95.79	95.79	99.17	
7	94.36	98.66	94.36	94.36	98.35	
8	92.44	99.53	92.44	92.82	99.29	
9	69.74	95.68	69.74	70.25	91.36	



图 3 子模型与模型融合方法的评价指标对比

Fig.3 Comparison of evaluation indexes among submodels and model fusion methods

对表2、图3进行分析后可得到如下结论。

1)由表2可见:由于训练样本类别的不平衡,子 模型1对样本量较大的故障事件的识别效果更好, 而对样本量小的故障事件的识别效果较差,其中对 类别9样本的分类召回率仅为69.74%;由于样本类 别不平衡度极大,子模型2中样本量大的故障事件 的权重因子很小,因此对多类别样本的召回率显著 降低,其中对类别1-3样本的分类召回率分别为 53.87%、69.21%、52.34%;而对少类别样本的召回率 明显提高,对类别9样本的分类召回率增至95.68%。

2)结合表2和图3可以看出:最值法、求和法单 纯从子模型预测的后验概率出发,不能对子模型的 性能进行分析,因此无法有效结合各子模型学习到 的信息,导致整体分类结果无明显改善;模型自适应 选择融合方法,综合考虑了子模型的召回率与准确 率指标,对于每个样本均能够灵活地选择子模型 的预测结果,从而保留子模型的优势性能,实现信 息互补,在保证多类别样本的分类效果的同时,有效 增强了对少类别样本的识别能力,准确率、召回率、 F1值相比子模型均有进一步的提升,分别达到了 95.97%、95.78%、95.74%。

3.3 整体性能验证

使用基于 Python 的 imblearn 工具包设置 4 组实 验,对比分析本文方法在整体性能上的优越性与 可靠性。在进行模型训练前,4 组实验分别采用 少数类别样本合成技术 SMOTE(Synthetic Minority

Oversampling TEchnique)^[17]、Borderline-SMOTE方法 (kind='borderline-1')^[18]、SMOTE 与编辑最近邻混合 采样方法(SMOTE-ENN)^[19]和 SMOTE-Tomek^[20]混合 采样方法按默认参数处理训练样本,依次记为方法 1—4。4种对比方法与本文算法的分类召回率见表 3,整体评价指标对比见图4。

	表3	对比方法和本文方法的分类召回率	×
--	----	-----------------	---

 Table 3
 Classification recall rate of comparison methods and proposed method

		-	-		
类别		分	类召回率/	/ %	
	方法1	方法2	方法3	方法4	本文方法
1	87.22	89.27	83.55	88.52	92.93
2	93.35	92.67	84.36	93.25	93.68
3	94.73	93.96	89.55	94.96	91.41
4	98.27	98.88	98.27	98.42	99.15
5	97.22	97.33	98.23	98.34	97.62
6	96.45	94.36	94.51	95.63	99.17
7	83.67	86.45	91.73	91.52	98.35
8	91.42	87.06	91.24	89.03	99.29
9	91.43	83.85	94.36	93.83	91.36



图 4 对比方法和本文方法的评价指标对比 Fig.4 Comparison of evaluation indexes among comparison methods and proposed method

综合表3和图4可以看出:

1)与其他考虑样本类别不平衡的对比方法相比,本文方法的3种评价指标均为最大,达到了95%以上,对各类故障事件的分类召回率也均在90%以上;

2)对于样本类别不平衡度较大的数据集,过采 样算法易合成噪点数据,破坏样本分布信息;混合采 样算法中欠采样的引入会丢失部分样本特征,破坏 模型对多数类样本的识别效果;

3)本文方法不改变样本初始分布,保留全部特征信息,在提高少数类样本的分类召回率的同时,有效维持了多数类样本的分类召回率并提高了其分类 准确率,因此整体故障识别效果得到了显著提高。

4 工程实际应用

以2018年8月17日"温比亚"台风过境当天所 截取的某信息密集时段内监控信息作为对象,验证 本文方法的实际应用效果。

当天13:27-13:31时段共产生了4146条告警

信息,系统从告警信息中提取出7项事故跳闸事件, 并通过本文方法在0.5 s内得到故障事件识别结果, 包括线路单相瞬时故障、单相永久故障、相间故障以 及一项历史样本极少的母线故障实例,经过验证,识 别结果均正确,其中母线故障事件识别结果如表4 所示。虽然在线应用样本量少,但是本文方法表现 出较高的识别准确率,并正确识别出一项发生概率 极低的母线故障事件,具有良好的工程应用价值。

表4 母线故障实例识别结果

Table 4 Recognition result of instance of bus fault

关键告警信息	事件识别结果
××市××变110 kV 事故总动作 ××市××变全站事故总动作 ××市××变110 kV.2号主变中压侧702分闸 ××市××变702开关控制回路断线动作 ××市××变220 kV第一套母线保护装置异常动作 ××市××变220 kV第二套母线保护装置异常动作 ××市××变702开关事故总动作 ××市××变752开关控制回路断线动作 ××市××变110 kV 母差出口跳副母动作 ××市××变110 kV 事故总动作	母线故障

5 结论

本文针对电网故障事件中的样本类别不平衡现 象,提出一种基于代价敏感学习和模型自适应选择 融合的多分类问题处理方法,实现了电网告警事件 的智能识别。基于对某市电网公司调度中心告警历 史信息的实验测试,所得结论如下:

1)通过本文构建的多分类代价敏感损失函数, 在损失函数中引入代价敏感因子,增大了少数类电 网故障事件的错分代价,优化模型对该类样本的特 征学习能力,从而改善模型对少数类电网故障事件 的识别性能;

2)综合考虑召回率与准确率的模型自适应选择 融合方法,对2个具有不同性能特点的模型进行融 合,结合子模型的优势,实现了模型的信息集成与优 势互补,在保留对多数类故障事件识别能力的基础 上,提高了少数类故障事件的识别率,得到整体效果 更好的电网故障识别模型。

后续可考虑将规则推理方法与深度学习进行深度结合,提高电网中人工智能模块的可靠性,同时进一步扩展可识别事件的类型。

附录见本刊网络版(http://www.epae.cn)。

参考文献:

孙国强,沈培锋,赵扬,等.融合知识库和深度学习的电网监控告警事件智能识别[J].电力自动化设备,2020,40(4):40-47.
 SUN Guoqiang, SHEN Peifeng, ZHAO Yang, et al. Intelligent recognition of power grid monitoring alarm event combining knowledge base and deep learning[J]. Electric Power Automation Equipment,2020,40(4):40-47.

- [2] 汪崔洋,江全元,唐雅洁,等. 基于告警信号文本挖掘的电力调度故障诊断[J]. 电力自动化设备,2019,39(4):126-132.
 WANG Cuiyang, JIANG Quanyuan, TANG Yajie, et al. Fault diagnosis of power dispatching based on alarm signal text mining[J]. Electric Power Automation Equipment,2019,39(4): 126-132.
- [3] 李元诚,曾婧. 基于改进卷积神经网络的电网假数据注入攻击 检测方法[J]. 电力系统自动化,2019,43(20):97-104.
 LI Yuancheng,ZENG Jing. Detection method of false data injection attack on power grid based on improved convolutional neural network[J]. Automation of Electric Power Systems, 2019,43(20):97-104.
- [4] 孙黎霞,白景涛,周照宇,等. 基于双向长短期记忆网络的电力 系统暂态稳定评估[J]. 电力系统自动化,2020,44(13):64-72.
 SUN Lixia, BAI Jingtao, ZHOU Zhaoyu, et al. Transient stability assessment of power system based on bi-directional long-shortterm memory network[J]. Automation of Electric Power Systems, 2020,44(13):64-72.
- [5] 谢桦,陈俊星,赵宇明,等.基于SMOTE和决策树算法的电力 变压器状态评估知识获取方法[J].电力自动化设备,2020,40 (2):137-142.

XIE Hua, CHEN Junxing, ZHAO Yuming, et al. Knowledge acquisition method of power transformer condition assessment based on SMOTE and decision tree algorithm[J]. Electric Power Automation Equipment, 2020, 40(2):137-142.

[6] 李雅欣,侯慧娟,胥明凯,等.基于策略梯度和生成式对抗网络的变压器油色谱案例扩充方法[J].电力自动化设备,2020,40 (12):211-218.

LI Yaxin, HOU Huijuan, XU Mingkai, et al. Oil chromatogram case generation method of transformer based on policy gradient and generative adversarial networks [J]. Electric Power Automation Equipment, 2020, 40(12):211-218.

[7] 谈林涛,李军良,任昺,等. 基于RB-XGBoost算法的智能电网 调度控制系统健康度评价模型[J]. 电力自动化设备,2020,40 (2):189-198.

TAN Lintao, LI Junliang, REN Bing, et al. Health evaluation model of smart grid dispatch and control system based on RB-XGBoost algorithm[J]. Electric Power Automation Equipment, 2020, 40(2):189-198.

- [8] ZHANG Zhongliang, LUO Xinggang, GARCÍA S, et al. Costsensitive back-propagation neural networks with binarization techniques in addressing multi-class problems and non-competent classifiers[J]. Applied Soft Computing, 2017, 56:357-367.
- [9] BRADFORD J P, KUNZ C, KOHAVI R, et al. Pruning decision trees with misclassification costs [C]//Proceedings of the 10th European Conference on Machine Learning. [S.I.]:Springer-Verlag, 1998:131-136.
- [10] 黄海松,魏建安,康佩栋.基于不平衡数据样本特性的新型过 采样SVM分类算法[J].控制与决策,2018,33(9):1549-1558.
 HUANG Haisong,WEI Jian'an,KANG Peidong. New over-sampling SVM classification algorithm based on unbalanced data sample characteristics[J]. Control and Decision, 2018, 33(9): 1549-1558.
- [11] DHAR S, CHERKASSKY V. Development and evaluation of cost-sensitive universum-SVM[J]. IEEE Transactions on Cybernetics, 2015, 45(4): 806-818.
- [12] 刘波,秦川,鞠平,等. 基于XGBoost与Stacking模型融合的短期母线负荷预测[J]. 电力自动化设备,2020,40(3):147-153.
 LIU Bo,QIN Chuan,JU Ping, et al. Short-term bus load fore-casting based on XGBoost and Stacking model fusion[J]. Electric Power Automation Equipment,2020,40(3):147-153.
- [13] 潘国兵,龚明波,贺民,等.基于Stacking模型融合的专变用户 电费回收风险识别方法[J].电力自动化设备,2021,41(1):

152-160.

PAN Guobing, GONG Mingbo, HE Min, et al. Identification method of electricity charge recovery risk of specialized transformer user based on Stacking model fusion [J]. Electric Power Automation Equipment, 2021, 41(1):152-160.

- [14] 韩笑,王新迎,韩帅,等.基于不均衡数据集成学习的大型电力 变压器状态评价方法[J].电网技术,2021,45(1):107-114.
 HAN Xiao,WANG Xinying,HAN Shuai, et al. Ensemble learning method for large-scale power transformer status evaluation based on imbalanced data[J]. Power System Technology, 2021,45(1):107-114.
- [15] MIKOLOV T,SUTSKEVER I,CHEN K,et al. Distributed representations of words and phrases and their compositionality[J]. Advances in Neural Information Processing Systems, 2013, 26 (5):3111-3119.
- [16] LIN T Y, GOYAL P, GIRSHICK R, et al. Focal loss for dense object detection [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2020, 42(2):318-327.
- [17] CHAWLA N V, BOWYER K W, HALL L O, et al. SMOTE: synthetic minority over-sampling technique[J]. Journal of Artificial Intelligence Research, 2002, 16:321-357.
- [18] HAN Hui,WANG Wenyuan,MAO Binghuan. Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning

[C]//International Conference on Intelligent Computing. Hefei, China:Springer-Verlag, 2005:878-887.

- [19] LU Tao, HUANG Youpeng, ZHAO Wen, et al. The metering automation system based intrusion detection using random forest classifier with SMOTE+ENN[C]//2019 IEEE 7th International Conference on Computer Science and Network Technology(ICCSNT). Dalian, China; IEEE, 2019; 370-374.
- [20] ZENG M, ZOU B J, WEI F R, et al. Effective prediction of three common diseases by combining SMOTE with Tomek links technique for imbalanced medical data[C]//2016 IEEE International Conference of Online Analysis and Computing Science(ICOACS). Chongqing, China: IEEE, 2016: 225-228.

作者简介:



卫志农(1962—),男,江苏江阴人,教 授,博士,主要研究方向为电力系统运行分 析与控制、输配电系统自动化等(E-mail: wzn ni@263.net):

石东明(1998—),男,安徽安庆人,硕 士研究生,主要研究方向为电力系统数据挖 掘(**E-mail**:1540805438@qq.com)。

卫志农

(编辑 任思思)

Intelligent identification method of power grid fault events considering sample classification imbalance

WEI Zhinong¹, SHI Dongming¹, ZHANG Ming², SUN Guoqiang¹, ZANG Haixiang¹, SHEN Peifeng² (1. College of Energy and Electrical Engineering, Hohai University, Nanjing 211100, China;

2. Nanjing Power Supply Company of State Grid Jiangsu Electric Power Co., Ltd., Nanjing 210019, China)

Abstract: In order to solve the problem that the difference of the failure probability of different equipments in the power grid affects the accuracy of fault intelligent diagnosis technology, an intelligent identification method of power grid fault events based on cost-sensitive learning and model adaptive selection fusion is proposed. Firstly, the Word2vec model is used to vectorize the pre-processed power grid alarm information, and two bidirectional long-short-term memory networks are established as basic classification models. Then, the cost-sensitive loss function is designed. The cross-entropy loss function and cost-sensitive loss function are respectively applied to the two classification models. Finally, a model adaptive selection fusion method is proposed to fuse the above classification models, so as to obtain the identification results of fault events. Actual data test shows that the proposed method can effectively reduce the impact of sample classification imbalance in the fault event identification.

Key words: identification of power grid fault events; deep learning; classification imbalance; cost-sensitive learning; model fusion

附录 A



Fig.A3 Flowchart of power grid fault events identification

附录 B



附录 C

Table C1 Parameters of Word2vec model						
模型参数	参数值	参数含义				
训练算法	0		CBOW 算法			
窗口大小	5	中心词与预测词在一条信息中的最大距离				
最小词频	5	词》	词频少于设置次数的词语会被丢弃			
训练加速策略	1	hierarchical softmax				
词向量维度	300	每个词的向量维度				
表 C2 Bi-LSTM 模型参数						
Table C2 Parameters of Bi-LSTM model						
参数名称	参	爹数值	参数名称	参数值		
输入词向量维	度	300	全连接层输出维度	9		
隐含层数量		1	学习率	0.001		
隐含层节点数 6		64	优化函数	adam		
全连接层丢弃率 0		0.5	训练批次	6		
全连接层激活函数 softmax 单次训练样本数量 12			128			

表 C1 Word2vec 模型参数