

doi: 10.3969/j.issn.1000-8349.0000-0000

基于样本模拟与 ResGAN 的脉冲星候选 体分类

李嘉雪¹, 李明辉¹, 卢吉光², 姜鹏², 李清云¹

(1. 贵州大学 省部共建公共大数据国家重点实验室, 贵阳 550025; 2. 中国科学院 国家天文台, 北京 100101)

摘要: 脉冲星巡天观测受限于望远镜设备灵敏度、天体引力摄动以及脉冲星本身辐射微弱等因素, 球状星团中观测到的脉冲星数量仅占其实际数量的小部分。针对已探测脉冲星数量有限的问题, 提出了一种基于 FAST 观测特征的数据增强方法, 通过模拟 FAST 系统背景噪声、射频干扰信号以及脉冲轮廓、频率-相位、时间-相位、色散量曲线等多模态信息, 生成大量新的脉冲星候选体样本, 并与部分开源数据融合, 构建了一个用于脉冲星分类的综合数据集, 为后续深度学习模型训练提供了丰富的样本支持。在分类模型搭建上, 融合半监督生成对抗网络与 ResNet-50 残差模型的优势提出了改进模型 ResGAN, 与基线模型 SGAN 相比, 改进模型在综合数据集上表现出更优的分类性能, 准确率提升了 0.9 个百分点, F1 分数提高了 1.7 个百分点。

关键词: FAST; 脉冲星候选体; 模拟; 数据集; 半监督生成对抗网络

中图分类号: P111.44 **文献标识码:** A

1 引言

脉冲星, 也是快速旋转和高度磁化的中子星, 拥有非常高而稳定的自转速度, 可归结于脉冲星的强磁场、高密度性以及脉冲星的周围相对“干净”, 受到其它天体的影响较小等原因。这种稳定性使得脉冲星在星际导航和计时方面具有重要的研究意义。脉冲星巡天任务会产生大量的脉冲星候选体, 对这些候选体进行筛选可以从大量不同类型的天文信号如快速射电暴、银河背景噪声中找到真实的脉冲星, 为后续脉冲星的物理特性分析、引力波理论验证、星际导航和计时等研究奠定基础。

早期, 主要采取人工识别和基于特征工程的传统机器学习方法, 尽管经验丰富的天文

收稿日期: 0000-00-00; 修回日期: 0000-00-00

资助项目: 贵州省科技计划项目 (黔科合基础-ZK[2023] 一般 039, 黔科合支撑 [2023] 一般 352, 黔科合平台人才-ZDSYS[2023]003)

通讯作者: 李明辉, limh@gzu.edu.cn

学家能够成功辨识出脉冲星候选体,但这种方法难以应对由 500 米球面射电望远镜 (Five-hundred-meter Aperture Spherical Radio Telescope, FAST)、平方公里阵列望远镜 (Square Kilometre Array, SKA) 等新一代高性能望远镜带来的海量数据挑战。随着深度学习技术,特别是卷积神经网络的兴起,基于自动特征学习的脉冲星分类方法逐渐成为研究的主流方向并取得了一定的进展,然而,仍有两个关键问题需要解决。

其一,以目前的天文设备观测能力,平均每个球状星团中被探测到的脉冲星数量约为 7.4 颗^[1],然而,一个星团中实际存在的脉冲星数量虽没有确切的统计数据,其实际丰度可能是已发现的几倍甚至几十倍,采取这样的数据进行训练,由于脉冲星样本稀少,限制了模型的泛化能力,使得模型在未见过的脉冲星信号上表现不佳。其二,在实际观测数据中,因存在大量的噪声和干扰等非脉冲星信号,真实脉冲星候选体样本数量和非脉冲星候选体样本数量极不均衡,正负样本比例接近 1/100000^[2],这意味着能够从几万张脉冲星候选体中找到一颗新的脉冲星都是一件幸运的事。而使用大多数标准学习算法处理不均衡数据集时,会倾向于学习多数类的特征,忽视少数类的特征表现,呈现效果易出现模型偏移的情况^[3]。

近年来,研究者们将目光聚集在如何缓解数据集不均衡的问题上。一些研究者,如刘晓飞等^[4]使用简单过采样技术对训练集中的正样本进行数据增强,并调整正负样本的比例,解决了正负样本非均衡问题;Zhang 等人^[5]结合改进的智能欠采样方法删除数据集中 86% 的负样本,有效缓和了数据集不均衡问题。Wang 等人^[6]提出了一种特别的数据增强方法,使用加权因子对多个候选体的同一张子图进行相加得到新的样本。但仅通过简单的欠采样、过采样并不能提高正样本的多样性,只能在处理已有的样本基础上,学习已有的范式,还存在过拟合的风险。针对以上问题,本文引入了一个模拟生成脉冲星候选体的创新视角,通过模拟生成未发现的但可能存在的脉冲星,可以在缓解数据集不平衡的同时极大地丰富正样本多样性。与此同时,我们在 Balakrishnan 等人^[7]的研究基础上进行改进,构建了名为 ResGAN 的半监督生成对抗网络,用于脉冲星候选体分类。生成对抗网络一般包含一个鉴别器和一个生成器,它们可以通过对抗训练不断提高各自的性能。将残差连接引入生成对抗网络,有效的缓解了卷积神经网络中的过拟合和梯度消失问题。

2 原理及方法

2.1 模拟生成脉冲星候选体

本研究通过模拟 FAST 观测数据在经过 PRESTO^[8]等软件处理后可能得到的四张特征子图,分别是频率-相位图、时间-相位图、平均脉冲轮廓、色散曲线图,得到模拟正样本。模拟过程中,充分考虑到实际情况下的 FAST 接收机带宽、系统温度等参数。从原理出发,通过随机生成不同类型的初始脉冲轮廓、随机加入射频干扰信号 (Radio Frequency Interference, RFI) 信号、进行不同程度的消色散等确保模拟生成样本的有效性和多样性。

2.1.1 模拟背景噪声和射频干扰

脉冲星信号经过 PRESTO 消色散处理后, 无法完全消除噪声和干扰的影响。该模拟实验通过添加不同的背景噪声和随机干扰信号以更真实地反映实际的观测环境。已知的具有统计意义且影响脉冲星搜寻的噪声包括高斯噪声、泊松噪声、卡方噪声等, 而泊松噪声和卡方噪声在大数定律下近似于高斯噪声, 因此我们选择添加呈高斯分布的背景噪声, 能够充分模拟观测数据中的随机噪声特性, 并通过设置不同的脉冲信号流量以期得到不同的信噪比参数。

除去望远镜接收机本身, 脉冲星观测还受到其它多种干扰源的挑战, 如移动通信基站等地球上的人类活动设备、人造卫星等太空中的设备。特别是随着 5G 技术的普及, 电磁环境变得愈加复杂, 射频干扰信号表现出动态、多样的特征。我们模拟了两种较为广泛的干扰信号, 即频域宽带干扰和频域窄带干扰 (见图1)。

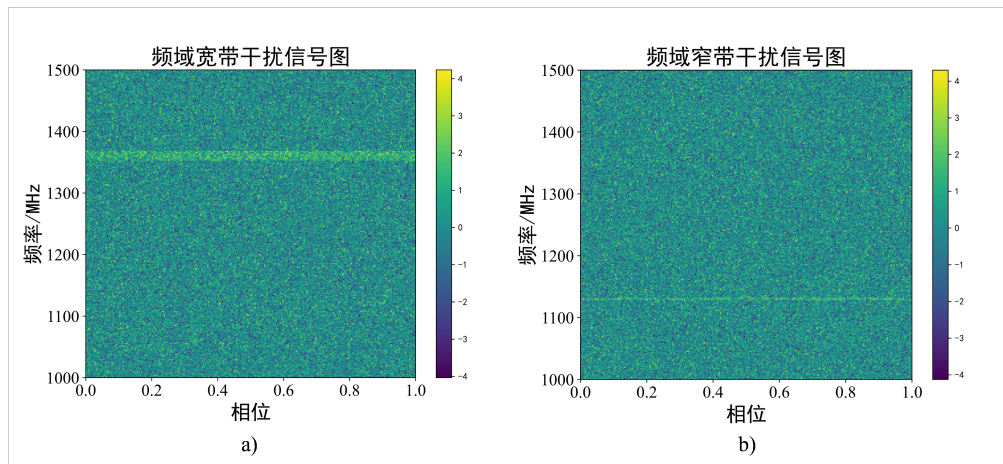


图 1 脉冲星候选体样本模拟中的频域宽带干扰和频域窄带干扰

受到干扰源的启停时间、频率漂移以及设备功率波动的影响, 干扰信号虽然集中于特定的频率范围, 但其频率位置、强度和出现时间往往具有随机性。窄带干扰影响少数频率通道, 宽带干扰影响相邻若干频率通道。为更准确的模拟 FAST 观测环境中的 RFI 的随机特性, 通过统计约 2000 个 FAST 脉冲星观测数据文件, 计算各频率通道出现 RFI 信号的概率信息, 在进行模拟实验过程中, 按照概率给相应通道增加干扰噪声, 其中, 总通道数目 4096 个, 工作频率 1.0GHz~1.5GHz, 带宽 500MHz。

2.1.2 模拟特征子图

平均脉冲轮廓的形状能够反映脉冲星发射射电波束穿过地球时的辐射强度分布^[9], 我们假设脉冲信号呈现理想的高斯分布或洛伦兹分布。其高斯分布函数如下式:

$$f(x) = A_g \cdot \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right), \quad (1)$$

其中 A_g 表示脉冲振幅, μ 表示脉冲中心位置, σ 可表示脉冲宽度。

$$f(x) = \frac{A_L}{\pi} \times \frac{\frac{\Gamma}{2}}{(x - x_0)^2 + \left(\frac{\Gamma}{2}\right)^2}, \quad (2)$$

其中 A_L 表示脉冲振幅, x_0 表示脉冲中心位置, Γ 表示脉冲宽度。首先, 随机初始化峰个数、峰相位、峰高度、峰宽度等参数, 得到初始脉冲轮廓 A 。其中, 峰个数取值范围是参考已公布的观测数据和脉冲星辐射模型的多样性所设置的, 由 Rankin 提出的得到众多研究者认可的“核 + 双锥环”模型^[10]指出: 脉冲星磁场结构复杂, 当视线从不同角度切过不同辐射区域时, 可产生不同的脉冲轮廓形态, 包括单峰、双峰、三峰和多峰, 仅切过核心时产生的单峰最常见。因此, 设定峰个数取值范围为 1 ~ 5, 不仅能涵盖常见的单峰和双峰情况, 还能捕捉到更复杂的脉冲星轮廓特征。脉冲宽度和观测频率、脉冲星类型、脉冲星辐射几何等有关, 如毫秒脉冲星的自转周期对应的辐射束扫过地球的时间较短, 脉冲绝对宽度便较窄, 脉冲星辐射束的几何结构和观测角度的变化会导致脉冲持续时间变化, 反映为脉冲宽度的变化。根据已有研究, 脉冲星的脉冲宽度通常占其周期的 3% ~ 10%, 但在某些特殊情况下, 脉冲宽度可达占 20% 或更多^[11]。因此, 我们设置脉冲宽度在周期的 2% ~ 30% 之间随机取值, 以涵盖可能存在的不同类型脉冲星的特征。脉冲高度是表征脉冲星辐射特性的重要参数, 反映脉冲星在特定观测周期内向外发射的电磁辐射强度, 我们参考了由 Wang 等^[12]公布的 FAST 数据集中绝大部分脉冲信号的流量密度, 设置初始轮廓的单脉冲流量密度取值范围为 2Jy ~ 8Jy。通常情况下, 脉冲轮廓中峰位置应该在一个相位内随机分布, 但是考虑到人工专家使用人眼筛选候选体过程中的判断依据: 正样本的主要特征在于在整个时域或者频域范围内存在一条较为明显的竖线, 因此, 为了让后期模型的训练更多关注特征图中是否存在某一条或多条竖线, 而对信号累计的位置给予较少关注度, 我们的峰位置取值范围是折叠周期的 40% ~ 60%, 多峰情况时, 只限定其中一个峰位置。

频率-相位图是区分脉冲星候选体最为重要的特征之一, 用来展示脉冲信号在不同频率通道下的相位信息。电磁波信号在通过星际介质时, 其中的自由电子会与信号发生作用, 导致色散延时, 高频信号会比低频更先到达望远镜^[13], 观测频率与观测频率之间由色散导致的时间延迟为:

$$\Delta t \doteq 4.1488 \times 10^6 \text{ms} \times DM \times \left(\frac{1}{f_{\text{low}}^2} - \frac{1}{f_{\text{high}}^2} \right), \quad (3)$$

其中, DM (Dispersion Measure) 表示辐射传播路径上的电子柱密度, 通常以 $\text{pc} \cdot \text{cm}^{-3}$ 为单位, f_{low} 和 f_{high} 分别是最低频率和最高频率, 单位为 MHz。在模拟生成时, 取参考频率为起始频率 f_{low} , 利用上式 (3) 计算各个中心频率的通道所对应的延迟时间, 并考虑不同频率通道内, 信号的峰位置随着频率的变大朝左或者右移动。实际计算 Δt 时引入了新的变量 ΔDM , 即脉冲星真实色散与被消除的色散之差, 其表达式为:

$$\Delta DM = DM - DM_{\text{eliminated}}, \quad (4)$$

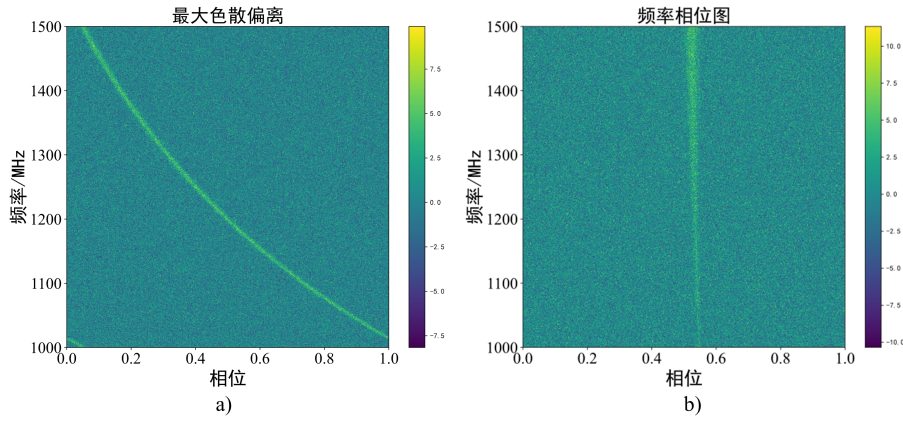
其中, $DM_{\text{eliminated}}$ 表示已被 PRESTO 消除的色散, 可理解为搜寻到的脉冲星色散。使用 ΔDM 计算色散延迟得到的频率-相位图模拟的是完成消色散处理后的情况。为验证消色散

方法的合理性, 引入了最大色散偏离 DM_{\max} , 若最高频率相对于起始频率的时间延迟 Δt 对应一整个相位延迟, 此时的 ΔDM 值即为 DM_{\max} , 图2中 a) 子图展示了未完全消除色散等于最大色散偏离时生成的频域图。另外, 我们还考虑了峰宽度和峰高度随频率变化的情况。有研究指出脉冲宽度和脉冲高度随频率变化的情况满足幂律关系^[10], 随着频率的升高, 脉冲对应的高度和宽度会随机的增大或减小, 具体关系如下:

$$W(f) \propto f^{\alpha_W}, \quad (5)$$

$$H(f) \propto f^{\alpha_H}, \quad (6)$$

其中, $W(f)$ 是观测频率 f 处的脉冲宽度, α_W 是脉冲宽度对频率的幂指数, $H(f)$ 是观测频率 f 处的脉冲高度, α_H 是脉冲高度对频率的幂指数。谱指数影响信号的能量分布和衰减特性, 我们选取的谱指数范围是 $-1 \sim 4$, 当谱指数 α_W 为 3.6、 α_H 为 2.2 时, 脉冲信号的频率-相位图如图2中 b)



注: a) 展示了未完全消除色散等于最大色散偏离时生成的频率-相位图。b) 展示了脉冲信号的宽度和强度随频率变化情况, 此时 $\alpha_W = 3.6$, $\alpha_H = 2.2$ 。

图 2 模拟不同的信号变化情况

色散量曲线反映了不同 DM 值下叠加压缩频域图得到的脉冲数据标准差与噪声标准差的比值。在模拟生成该曲线的过程中, 首先会设置一个拟合用的 DM 值, 代表当前脉冲信号的真实色散量; 然后以该值为中心, 两边各取一半 DM_{\max} 得到取值范围, 并在该范围内选取一系列点, 代表消色散操作对应的一系列 $DM_{\text{eliminated}}$; 最后, 计算不同色散量下的比值即可生成色散量曲线。脉冲信号的色散量曲线在 DM 不为零处存在峰值, 当峰值对应的 DM 值为零时, 该信号为干扰信号。

时间-相位图可反映整个观测时间内脉冲信号随时间的变化情况。在模拟脉冲星的时间-相位图时, 不仅要关注脉冲星的自转周期 P , 还需要考虑外部因素, 如天体运动和引力效应导致的加速度, 这些加速度会导致脉冲星搜寻数据中有明显观测效应的周期导数 \dot{P} 。因此,

算法引入周期导数，可有效模拟脉冲星的加速度效应。为简化模拟算法，我们忽略了一次观测中加速度存在变化的情况，若脉冲星加速度为 a ，频率导数 \dot{f} 由以下公式给出：

$$\dot{f} = \dot{f}_0 \left(1 + \frac{v_r}{c}\right)^2 + f_0 \frac{a}{c}, \quad (7)$$

其中， v_r 是脉冲星的径向速度， f_0 是脉冲星转动频率，是周期 P 的倒数， c 是光速。周期导数由以下公式给出：

$$\dot{P} = -\dot{f} \cdot P^2. \quad (8)$$

通过设置周期 P 和加速度 a 的范围，结合上式 (7)(8) 可推导出模拟实验中周期导数的随机数上限值。基于多普勒效应和加速度效应，计算每个时间步长下的信号频率和相位变化。具体相位随时间的变化公式为：

$$\varphi(t) = \varphi_0 + f \cdot t + \frac{1}{2} \dot{f} \cdot t^2, \quad (9)$$

其中， φ_0 是初始相位， f 是变化着的脉冲星转动频率， \dot{f} 是频率导数。由此原理获取不同时刻脉冲信号的相位信息。然后，沿频率叠加频率-相位图生成脉冲序列，得到不同相位对应的脉冲流量，从而得到不同时刻对应的脉冲流量。最后，以整个时间段内的平均周期作为折叠参数处理整个信号序列，得到时间-相位信息，并按周期折叠得到最终积分轮廓。整个模拟算法流程见图3。

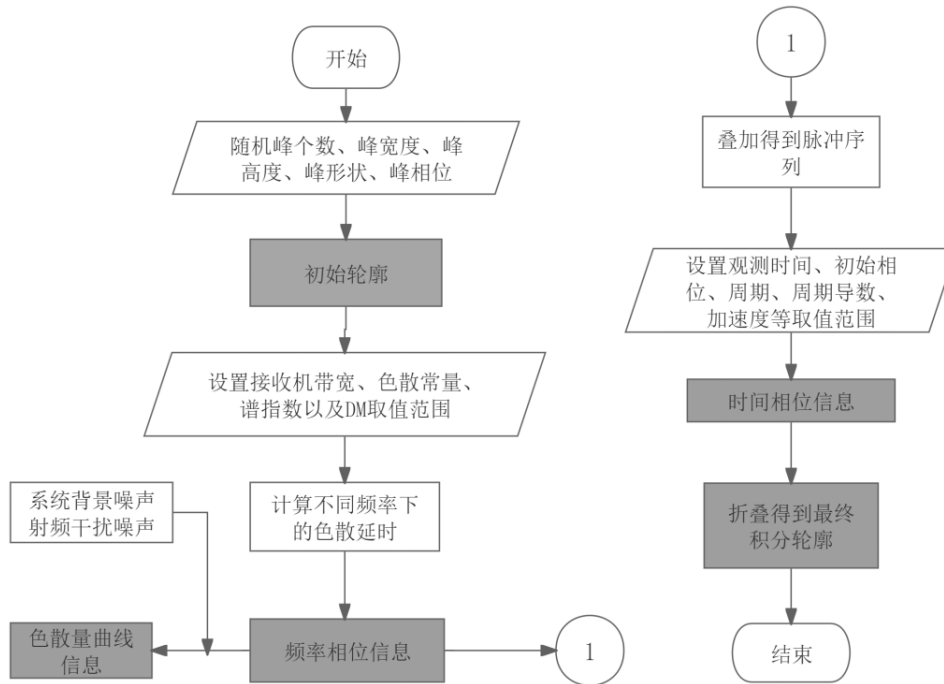
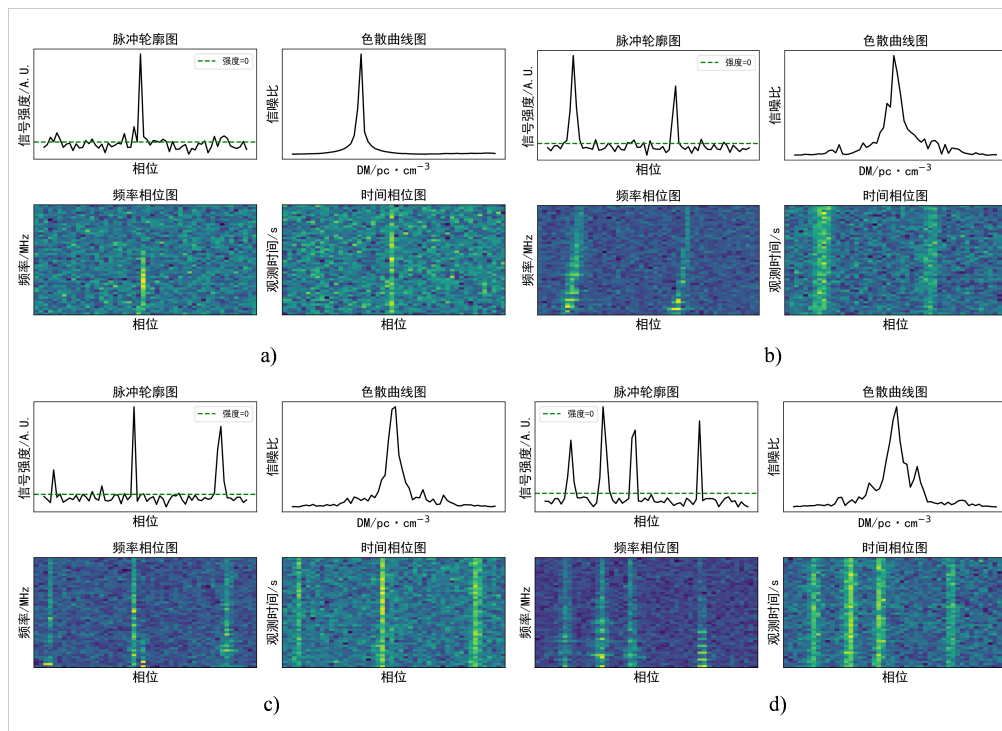


图3 模拟算法流程图

2.1.3 模拟样本的生成及其分布

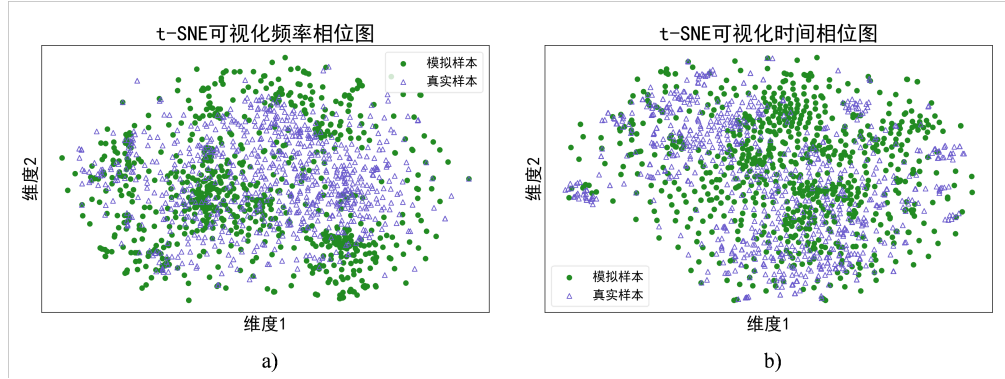
在完成脉冲星样本的模拟生成后, 我们将其与 Wang 等人^[12]公开的脉冲星候选体数据集 D_1 中的样本进行对比。数据集 D_1 包含多种类型的脉冲星信号, 在脉冲星分类任务中被广泛应用, 其中单峰脉冲星占比 83.57%, 双峰脉冲星占比 13.6%, 而在实际观测中较为稀少的多峰脉冲星仅 33 个样本, 占比 2.83%, 这限制了分类模型对稀有类型脉冲星的学习能力。另外, 该数据集中的正样本通常信噪比强度较大、消色散效果较好, 没有覆盖在实际观测中存在的低信噪比和处理过程中未完全消除色散的脉冲星候选体。而我们利用该模拟方法, 通过设置如信号强度等脉冲轮廓相关、未消除色散、谱指数等参数范围, 以及随机添加 RFI 干扰信号等处理, 生成 14853 个具有不同峰个数的脉冲星候选体, 其中, 双峰样本 2644 个, 占比提升至 17.8%, 多峰样本 1050 个, 占比提升至 7.1%。最终得到了模拟数据集 D_2 , 部分样本见图4, 相较于原始数据集 D_1 , 新数据集的四张子图不仅具有脉冲星候选体在周期性信号、频域特征等物理性质上的典型特征, 还呈现出更多样化的脉冲结构, 包括不同形状和强度的脉冲轮廓, 以及在时间和频率域上变化的细节特征。此外, 因数据是通过完全自动化的程序生成的, 数据集 D_2 在扩展性方面具有显著优势。



注: 示例中四个样本峰个数不同, 频率-相位模式均考虑设置了不同大小的谱指数, 信号强度和宽度随时间变化, b) 中频率-相位图设置了一定大小的未消除色散, 呈现较为明显的弯曲; d) 中频率-相位模式存在细微的窄带干扰。实际模拟中考虑到各种随机变量, 样本更丰富多样, 受页面限制无法逐个列出。

图 4 模拟生成的部分脉冲星候选体样本示例

为了进一步分析数据集 D_2 和 D_1 在特征空间上的分布情况，我们在两个数据集中随机抽取 786 个真实样本和 786 个模拟样本，并采用 t 分布随机邻域嵌入降维技术^[14]对抽取样本的时间-相位子图和频率-相位子图做可视化降维处理，结果见图5，模拟样本的分布较为松散，出现该现象的原因可能是模拟脉冲信号时，采用不具有周期性的高斯函数作为峰函数，若未消除色散较大带来的峰相位变化较大，时域和频域上的信号连续性会遭到破坏，需要说明的是，实际观测中也存在一些连续性遭到破坏的候选样本。总的来看，大部分模拟样本和真实样本在特征空间内可以重叠，具有很好的相似性，验证了模拟方法的可行性。



注：图中绿色的实心圆点表示数据集 D_2 中的模拟样本，紫色的空心三角表示数据集 D_1 中的真实样本，模拟样本和真实样本各有 786 个，a) 为可视化两类样本中的时间-相位图分布。b) 为可视化两类样本中的频率-相位图分布。

图 5 数据集 D_1 和 D_2 部分样本的降维可视化情况

2.2 残差半监督生成对抗网络 ResGAN

半监督生成对抗网络 (Semi-Supervised Generative Adversarial Networks, SGAN) 是基于生成对抗网络的另一种形式，与之不同的是，半监督网络的判别器不仅要区分输入数据的真假，还需要对真实数据进行分类。其中，判别器包含监督分支和非监督分支。监督模块与二分类任务有直接的关系，可以说整个半监督生成对抗模型中除了监督模块之外，其他的部分，包括非监督模块，都是服务于监督任务的分类目标。为了在脉冲星候选体分类任务上取得最佳的分类性能，我们特别选用了在图像分类任务中表现非常优秀的 ResNet-50^[15]来重新设计基线模型 (SGAN^[7]) 上监督分支的结构，得到的改进模型 (ResGAN) 网络架构见图6。

考虑到脉冲轮廓和色散量曲线是一维信号模态，我们重新设计了 ResNet 模型的一维卷积版本，专门处理一维信号数据。与传统的二维卷积不同，一维卷积更适合捕捉信号数据中的时序或空间特征。假设色散量曲线和脉冲轮廓两个模态的数据是一维向量 $\mathbf{x} \in \mathbb{R}^N$ ，其中 N 是信号的长度，卷积核 $\mathbf{w} \in \mathbb{R}^k$ 的大小为 k ，对于一维卷积操作，卷积操作可以表示为：

$$\mathbf{y}_i = \sum_{j=1}^k \mathbf{x}_{i+j} \mathbf{w}_j \quad (10)$$

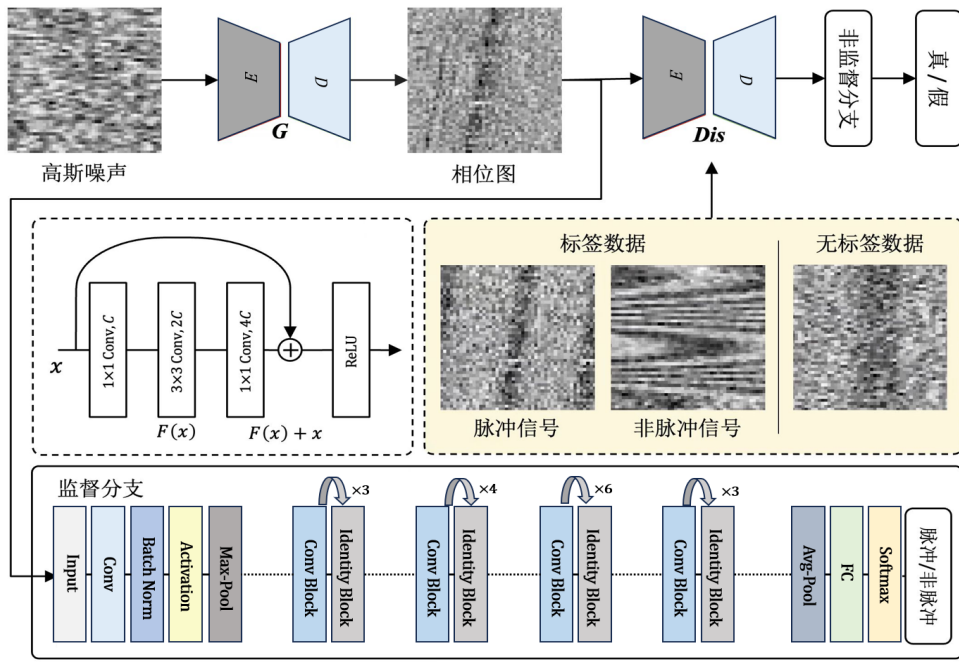


图 6 改进模型 (ResGAN) 网络架构

其中, y_i 是卷积后的输出值, x_{i+j} 是输入信号的值, w_j 是卷积核的参数。该操作通过在输入信号上沿着时间轴滑动, 提取信号的局部特征。为提取更复杂的特征, 本文将这一卷积操作扩展成多个卷积层, 并且延续残差结构形成一个更深层的架构, 使得每一层都能够学习到不同层次的特征, 多层卷积的输出可以表示为:

$$\mathbf{Y} = \text{Conv1D}(\mathbf{X}, \mathbf{W}; k, s) \quad (11)$$

其中, \mathbf{X} 是输入信号矩阵, \mathbf{W} 是卷积核, k 是卷积核的大小, s 是步长, 输出的 \mathbf{Y} 为卷积结果。在原始的 ResNet 中, 残差连接通过直接跳过一个或多个卷积层来缓解梯度消失问题, 并保持较低的计算复杂度。在一维卷积版本的 ResNet 中, 我们延续了这种经典的残差结构, 只是将每一层的卷积操作从二维变为一维。具体来说, 对于每个残差块, 残差公式输出为:

$$\mathbf{y} = \text{Conv1D}(\mathbf{x}, \mathbf{W}; k, s) + \mathbf{x} \quad (12)$$

其中, \mathbf{x} 是输入信号, \mathbf{y} 是残差块的输出。残差结构的加入可以有效避免信息丢失, 并使得信号的特征能够在多个层次中得到有效传递。

3 实验与结果

3.1 数据集和实验设置

本文所使用的综合数据集由大量脉冲星候选体构成，该数据集主要由 Wang 等人公布的 FAST 开源数据集和我们制作的私有数据集组成，私有数据集源于两个部分，一部分是利用本文数据增强算法生成的模拟候选体数据；另一部分是利用 PRESTO 处理 FAST 对 M5、M14 等星团的观测数据得到的真实候选体。

为得到合适的分类器，我们将综合数据集划分为训练集、验证集和测试集，具体分布如表 1 所示。训练集用于模型训练，验证集用于模型的超参数调优，它们的正样本全部采用模拟生成的正样本，而负样本则随机选择，样本数量和正样本保持近 1: 1 的比例；测试集用于模型的最终评估，其正样本融合了剩下的模拟样本和 FAST 开源数据集的所有正样本，负样本随机选择。将真实正样本全部放入测试集，不参与模型的训练过程，该策略旨在严格区分真正正样本和模拟正样本，从而验证模型在未见过的脉冲星候选体的识别能力。

表 1 数据集的分布情况

数据集	训练集	测试集	验证集
脉冲星	11086	3000	1930
非脉冲星	11205	7000	1865
未标记样本		10043	

本节实验在配备了高性能计算硬件的云端服务器平台上进行，其硬件环境为：15 核 Intel(R) Xeon(R) Platinum 8474C 处理器、NVIDIA GeForce RTX 4090D 显卡；软件环境为：Ubuntu20.04.4 操作系统下的 Tensorflow 2.10.0 + CUDA 11.8.0 + python3.8.17 框架。模型的训练使用了 Adam 自适应矩估计优化算法，能够自动调参数的学习率。

3.2 评价指标

在我们的数据集中，脉冲星及其谐波信号作为正样本，其它的非脉冲星作为负样本。此时脉冲星候选体分类作为二分类任务，常用的评价指标包括准确率、召回率、精准率、F1 分数，这些指标的定义如下：

$$\text{准确率} = \frac{TP + TN}{TP + TN + FP + FN} \quad (13)$$

$$\text{召回率} = \frac{TP}{TP + FN} \quad (14)$$

$$\text{精准率} = \frac{TP}{TP + FP} \quad (15)$$

$$\text{F1 分数} = 2 \times \frac{\text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}} \quad (16)$$

其中, TN(True Negative) 表示实际为负并预测为负的样本数; TP(True Positive) 表示实际为正并预测为正的样本数, 即被正确识别的脉冲星; FN(False Negative) 表示实际为正却预测为负的样本数, FP(False Positive) 表示实际为负却预测为正的样本数。准确率表示被正确分类的样本数占总样本数的比例, 100% 的准确率意味着模型正确预测了所有样本, 没有错分; 召回率表示被正确分类的正样本数占正样本总数的比例, 100% 的召回率意味着模型正确识别了所有正样本, 没有漏检; 精准率表示被正确分类的正样本数占预测为正样本数的比例, 100% 的精准率意味着模型预测的正样本全部正确, 没有误报; F1 分数是召回率和精准率的调和平均数, 100% 的 F1 分数意味着没有漏检和误报, 是一种理想状态。

3.3 实验结果对比

一个候选体样本包含四种模态的特征, 分别是色散量曲线、脉冲轮廓、频率-相位、时间-相位, 这是人类专家在对候选体进行分类时考虑的四个最重要的特征。为评估不同模态特征对分类性能的影响, 分别选择不同模态作为模型的独立输入进行训练; 为验证改进模型的有效性, 分别对基线模型 (SGAN) 和改进模型 (ResGAN) 进行训练; 为系统性评估并对比两个模型在标记数据稀缺到充足场景下的性能, 实验设置标记样本数量分别为 1%、10%、25%、100% 的对照组。上述实验结果见表 2。

表 2 不同模态和不同训练样本比例下 SGAN 和 ResGAN 方法对比

模态	样本百分比	SGAN			ResGAN		
		精准率 /(%)	召回率 /(%)	F1 分数 /(%)	精准率 /(%)	召回率 /(%)	F1 分数 /(%)
色散量曲线	1%	87.569	89.948	88.743	86.451	90.38	88.372
	10%	89.235	94.428	91.758	89.356	93.58	91.419
	25%	91.716	93.149	92.427	95.784	91.652	93.672
	100%	96.565	96.332	96.448	96.158	96.819	96.487
脉冲轮廓	1%	58.056	69.571	63.294	58.765	68.409	63.221
	10%	73.561	65.464	69.276	68.476	68.901	68.688
	25%	73.215	69.383	71.248	72.472	70.68	71.565
	100%	75.838	81.345	78.495	77.203	81.675	79.376
频率-相位	1%	64.975	86.283	74.128	64.829	88.68	74.902
	10%	88.832	90.029	89.426	89.348	91.085	90.208
	25%	89.252	94.533	91.817	90.032	93.904	91.927
	100%	92.537	98.691	95.515	93.687	97.691	95.647
时间-相位	1%	57.642	83.353	68.153	57.581	85.755	68.899
	10%	68.008	90.054	77.494	69.707	89.034	78.194
	25%	74.046	93.241	82.542	82.058	92.414	86.929
	100%	91.672	96.899	94.213	95.559	95.844	95.701

注: 为了便于比较, 表格中表现更好的数据项已加粗处理。

结合表 2 中 F1 分数可知, 针对色散量曲线和脉冲轮廓这两种一维信号模态, 在训练样本较少的情况下, 我们提出的 ResGAN 模型对其分类效果没有提升, 而随着样本数量的增加, ResGAN 逐渐展示出了一定的优势, 例如, 在脉冲轮廓模态的 100% 样本下, ResGAN 的 F1 值为 0.794, 相较于基线模型的 0.785 大约提升了 0.9 个百分点。在频率-相位和时间-相位模态这两种二维图像模态中, ResGAN 模型则展现了明显的优势, 尤其在训练样本量充足的时候。以时间-相位模态为例, 当训练样本占比 1% 时, 改进模型的 F1 值在基线模型基础上提升了 0.74 个百分点, 这表明在极少训练样本的情况下, ResGAN 已经能够从图像模态中挖掘出关键特征进行分类。总体而言, ResGAN 模型除了在样本量占比 10% 及以下时对一些一维信号模态的分类略弱于基线模型, 在其它样本比例下对候选体四种模态的分类表现几乎都更好。这一趋势也进一步验证了对 ResGAN 架构的设计假设: 即其强大的图像处理能力使其在图像信号处理任务中表现更好, 而在一维信号任务中仍有优化空间。

与此同时, 从实验结果可看出, 若分别使用候选体的四种模态特征对脉冲星进行分类, 分类的效果有所不同。其中, 色散量曲线是表现最好的特征, 几乎所有比例条件下的分类情况都比其它三种模态好, 这可能是因为大量正样本在 DM 值不为 0 处有明显的峰值, 而绝大部分负样本的色散量曲线表现出随机波动, 没有显著的局部最大值或者最小值, 数据点波动性小。另外, 脉冲星的脉冲轮廓因峰的个数、形态、宽度、高度等不同表现为形态各异, 缺乏统一的特征, 又受到仪器噪声、射频干扰等各种因素影响, 难以靠单一模态对候选体进行分类, 故分类效果是四个模态中最差的。频率-相位和时间-相位两个模态表现效果较好。在单一模态下, 脉冲星信号和非脉冲星信号的特征可能重叠, 导致分类效果不佳。例如, 某些射频干扰在时间-相位模态上可能与脉冲星信号非常相似, 而某些微弱的脉冲星信号, 或者脉冲频率与 RFI 信号类似的脉冲星信号, 它们的时间-相位图或频率-相位图与非脉冲星信号特别相似。综上, 我们认为综合考虑四个模态特征对脉冲星候选体进行分类是必要的, 表 3 是在样本比例为 100% 的时候, 两种组合模型在综合数据集上的分类效果对比。

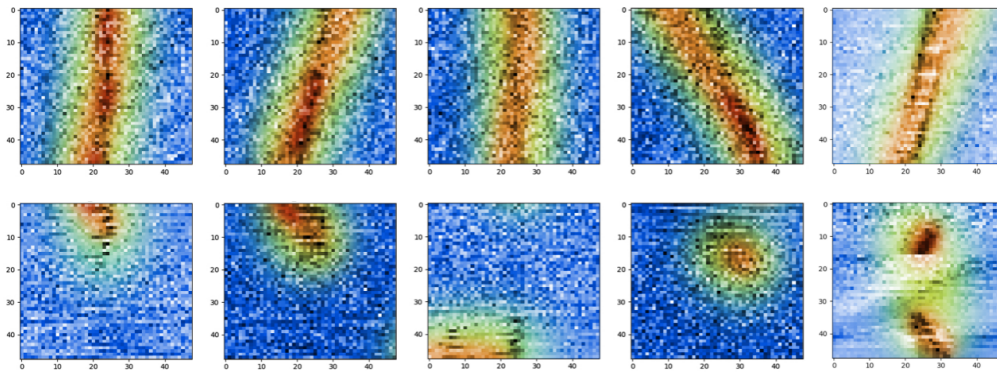
表 3 SGAN 和 ResGAN 在综合数据集上的表现

模型	准确率	精准率	召回率	F1 分数
SGAN	0.962	0.961	0.965	0.963
ResGAN	0.971	0.978	0.982	0.980

从表 3 中的结果可以看出, 组合四个模态特征的分类模型较单个模态而言在综合数据集上的表现有了明显的提升。与此同时, 通过与基线模型相比, 可以看出新模型在各个指标上的表现都更好, 展现了其在脉冲星分类任务上的优越性能。具体而言, ResGAN 在各项指标上均有明显提升, 其准确率达到了 0.971, 较基线模型提高了 0.009; F1 分数作为精准率和召回率之间的调和平均值, 达到了 0.98。

3.4 定性分析

对监督分支的卷积层及其对应的特征图进行了可视化 (见图7), 并深入分析了 ResGAN 模型的监督分支在推理脉冲星信号的 ROI (Region of Interest) 区域的表现。以频率-相位模



注: 图中第一排是频率-相位模态, 第二排是时间-相位模态, 坐标值对应原始图像的像素位置。

图 7 Grad-CAM 方法下的 ROI 区域效果

态和时间-相位模态为例, 深度网络能够有效地捕捉脉冲信号的内在规律。具体来说, 脉冲信号的集中区域往往位于图像的中央, 且形态特征呈现为竖形的条带状结构, 表明网络识别到了脉冲信号的特征模式。在不同的模态图中, 网络则表现出对多个 ROI 区域的关注, 尤其是在一些特定的时间或相位点, 这些区域的信号特征更加突出。通过这些细致的特征分析, 可以推测, ResGAN 在训练过程中成功地学习到了分类推理的关键信号特征, 进一步验证了其在复杂模式识别任务中的有效性和鲁棒性。

4 总结与展望

在利用深度学习对脉冲星候选体进行分类识别的过程中, 为满足提高脉冲星候选体样本多样性、缓解数据集正负样本不均衡的需求, 我们提出了一种模拟生成脉冲星候选体的新方法, 通过模拟系统背景噪声、RFI 干扰信号, 设置不同的信号类型、脉冲宽度、脉冲高度、脉冲周期、周期导数等物理参数, 生成了大量可能存在的脉冲星样本并在一定程度上缓解了类别不均衡问题。考虑到当前搜寻脉冲星使用的首要筛选标准还是“觉得它可能是一个脉冲星”, 那么通过人为模拟生成样本来进行筛选的模式, 相当于把专家认为是脉冲星的样本做为正样本, 这样的正样本中高信噪比的部分可以认为至少已经覆盖了目前已发现的脉冲星, 使用 T-SNE 降维技术对模拟生成样本和真实样本的分布情况可视化结果也能辅助说明模拟方法的有效性, 这为脉冲星候选体识别的数据增强步骤提供了新的思路和方法。除此之外, 改进模型 ResGAN 在不同样本比例下对不同模态的特征子图分类以及组合四个模态的分类结果显示, 在少量样本情况下, 新模型对一维信号模态的识别较基线模型没有明显提升, 但在样本充足的情况下, 新模型对四种模态的识别均优于基线模型。

下一步研究拟通过参考不同类型脉冲星的物理性质, 设置更具体的参数进行模拟实验, 得到更具有代表性的脉冲星候选体, 期望利用我们模拟产生的脉冲星样本, 训练并应用于更

多的分类模型。另外,在本实验的基础上继续优化,尝试设置不同数量的未标记样本,探究未标记样本数据量对模型性能的影响。

参考文献:

- [1] Yin D J, Zhang L Y, Li B D, et al. RAA, 2023, 23(5): 055012.
- [2] Morello V, Barr E D, Bailes M, et al. MNRAS, 2014, 443(2): 1651-1662.
- [3] Wang Y, Zhang Z, Liu Y, et al. 2023 International Conference on Power, Electrical Engineering, Electronics and Control (PEEEEC). IEEE, 2023: 38-40.
- [4] 刘晓飞, 劳保强, 安涛, 等. 天文学报, 2021, 62(2): 20.
- [5] Zhang S C, Kong X C, Zhou Y Y, et al. RAA, 2021, 21(10): 257.
- [6] Wang Y C, Li M T, Pan Z C, et al. RAA, 2019, 19(9): 133.
- [7] Balakrishnan V, Champion D, Barr E, et al. MNRAS, 2021, 505(1): 1180-1194.
- [8] Ransom S. Astrophysics Source Code Library. <https://ui.adsabs.harvard.edu/abs/2011ascl.soft07017R/>.
- [9] 卢吉光. 现代物理知识, 2022, 34(3): 4-12.
- [10] Rankin J M. ApJ, Part 1, 1983, 274: 333-368.
- [11] Slee O B, Bobra A D, Alurkar S K. Australian Journal of Physics, 1987, 40(4): 557-586.
- [12] Wang H F, Zhu W W, Guo P, et al. Science China Physics, Mechanics & Astronomy, 2019, 62: 1-10.
- [13] Condon J J, Ransom S M. Essential radio astronomy. Princeton University Press, 2016.
- [14] Van der Maaten L, Hinton G. Journal of Machine Learning Research, 2008, 9(11).
- [15] He K, Zhang X, Ren S, et al. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016: 770-778.

Pulsar Candidate Classification Based on Sample Simulation and ResGAN

LI Jia-xue¹, LI Ming-hui¹, LU Ji-guang², JIANG Peng², LI Qing-yun¹

(1. State Key Laboratory of Public Big Data, Guizhou University, Guiyang 550025, China; 2. National Astronomical Observatories, Chinese Academy of Sciences, Beijing 100101, China)

Abstract: In the process of pulsar survey observations, considering the constraints of telescopic sensitivity, gravitational perturbations from celestial bodies, and the intrinsically faint radiation emitted by pulsars, the number of detected pulsars in globular clusters represents only a fraction of their actual population. To address the limited number of detected pulsars, this paper proposes a data augmentation method based on the observational characteristics of the Five-hundred-meter Aperture Spherical Radio Telescope (FAST). This method simulates the background noise and radio frequency interference of the FAST system, along with multimodal information including pulse profiles, frequency-phase, time-phase, and dispersion

measure curve, to generate a substantial number of new pulsar candidate samples. By integrating these samples with same openly available data, we have constructed a comprehensive dataset for pulsar classification, which offers improved sample support for the training of subsequent deep learning models. For the classification model architecture, we developed ResGAN, an innovative framework that combines the strengths of semi-supervised generative adversarial networks with the ResNet-50. In comparison to the baseline model (SGAN), ResGAN exhibits superior classification performance on the comprehensive dataset, achieving a 0.9 percentage point improvement in accuracy and a 1.7 percentage point enhancement in the F1 score.

Key words: FAST; pulsar candidates; simulation; dataset; semi-supervised generative adversarial network