

基于深度残差网络的恒星光谱类别预测

王天翔^{1,2}, 范玉峰^{1*}, 王晓丽¹, 龙潜¹, 王传军¹

1. 中国科学院云南天文台, 云南 昆明 650011
2. 中国科学院大学, 北京 100049

摘要 近年来,随着各大光谱巡天项目的陆续实施,观测得到的天体光谱数据急剧增长。大型光谱巡天项目对光谱的自动分类和分析提出了更高的要求。本文将分类问题转化为回归问题,提出一种基于深度残差网络的光谱类别预测方法,对恒星光谱进行光谱次型预测。网络主要包括25个卷积层,1个最大池化层,1个平均池化层,全连接层以及12个残差结构。最大池化层用来筛选特征,卷积层提取特征,平均池化层用于减少模型参数,提高效率。残差结构可以防止网络退化,加深网络来提取高维抽象特征以及提高训练速度。考虑到数据有非零几率存在错误标签以及损坏数据,采用Log-Cosh作为损失函数来降低坏样本带来的负面影响。实验数据使用的是从LAMOST DR5中随机抽取的80 000条光谱,由于光谱质量等原因,每个光谱型的光谱数量不一。经过剔除坏值,流量归一化后,按7:1:2分为训练集、验证集和测试集。实验包括两个部分,第一个部分是使用数据集训练网络在光谱次型上进行类别预测,使用最大绝对误差、平均绝对误差以及标准差来比较不同形状卷积核的性能。将预测值作为横坐标,标签作为纵坐标,对测试集所有样本点使用二阶非线性拟合,得到了一条与 $y=x$ 重合的直线。证明模型可以很好的预测光谱次型。第二部分是对模型进行内部分析,使用类别激活映射的方法分别研究了模型预测A, F, G和K四种类型光谱时所关注的主要特征,赋予了模型可解释性。在文中数据集上,该方法对91.4%的光谱预测误差在0.5个光谱次型以内,预测的平均绝对误差为0.3个光谱次型。并与非参数回归、Adaboost回归树、K-Means三种方法进行同数据集比较,结果表明文中提出的方法可以很好地预测光谱次型并且速度更快,准确率更高。

关键词 恒星光谱; 光谱次型预测; 深度学习; 回归; 特征映射

中图分类号: P152 **文献标识码:** A **DOI:** 10.3964/j.issn.1000-0593(2021)05-1602-05

引言

LAMOST, 全称“大天区面积多目标光纤光谱望远镜”, 是世界上光谱获取率最高的望远镜, 可同时获得4 000个天体光谱^[1]。目前LAMOST已经发布7季数据, 在最新发布DR7中光谱数量已经高达1 448万条, 如何对海量光谱进行有效利用成为亟待解决的问题。对这些光谱进行分类是天文数据处理的重要一环。通过对恒星光谱的分类, 研究人员可以从中获取有效温度、质量和半径等物理信息, 也可以研究银河系的结构和演化过程^[2]。目前主流的恒星分类系统是MK光谱系统。每个恒星都根据其有效温度由高到低排序, 依次分为O, B, A, F, G, K和M七种光谱型, 每种光谱型又根据温度从高到低细分为0—9的次型光谱, 本文不涉及

光度型分类。

目前光谱自动分类的方法主要有三种类别, 分别是基于距离度量的方法、机器学习的方法和基于模糊逻辑知识系统的专家系统。Schierscher等^[3]将Artificial Neural Network(ANN)运用在对Sloan Digital Sky Survey(SDSS) DR7恒星光谱的分类上。Liu等^[4]对LAMOST数据使用线指数和SVM算法对恒星光谱进行MK分类。其中SVM方法对A, F和G型恒星分类效果达到90%的准确率, 对O, B, K和M型恒星只有52%的准确率。Kaushal等^[5]针对已标注数据太少, 难以训练深层神经网络分类器的问题, 提出一种半监督方法。该方法在无监督学习阶段使用自动编码器对无标签数据进行提取特征和聚类, 用有标签数据进行微调, 最后在主要光谱类别的平均准确率达到89%。在涉及光谱次型的分类模型上, Gray等^[6]提出一种专家系统, 通过直接与MK分

收稿日期: 2020-03-15, 修订日期: 2020-07-06

基金项目: 国家自然科学基金项目(11603072, 11773074), 云南省科技厅科技入滇项目(202003AD150003)资助

作者简介: 王天翔, 1995年生, 中国科学院大学云南天文台硕士研究生 e-mail: wangtianxiang@ynao.ac.cn

* 通讯作者 e-mail: fanyf@ynao.ac.cn

$$\text{Std} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i^{\text{pred}} - \text{MAE_error})^2} \quad (6)$$

为了选择最优的卷积核形状, 本文对 4 种不同卷积核的网络在测试集上的预测结果进行对比, 结果如表 1 所示。实验表明: 网络使用 1×5 的卷积核时所得结果平均绝对误差小, 预测误差分布集中在较小值。可以取得较好的结果, 91.4% 的光谱预测误差在 0.5 个光谱次型内, 平均绝对误差降低到了 0.3 个光次谱型。

表 1 各形状卷积核实验结果
Table 1 Experimental results of convolution kernels with different shapes

卷积核大小	最大绝对误差 (光谱型)	平均绝对误差 (光谱型)	标准差
1×3	3.77	0.07	0.12
1×5	2.98	0.03	0.05
1×7	1.99	0.05	0.07
1×9	2.01	0.05	0.07

将预测值作为横坐标, 标签作为纵坐标画一个平面, 平面上的一个点代表一条光谱, 对测试集上共 16 249 个点作二阶非线性拟合, 设置置信度为 95 (如图 2 所示), 可以看出, 所得到的函数基本可以看作斜率为 1 的直线, 并且置信区间与直线基本重合, 这表示模型可以很好的预测光谱型和光谱次型。

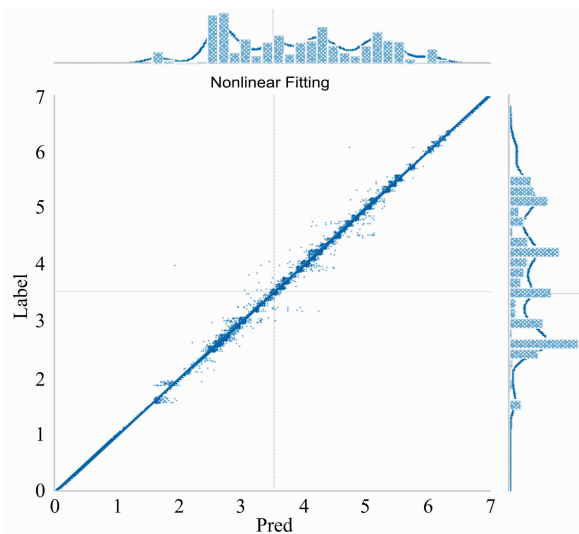


图 2 二阶非线性拟合
Fig. 2 Second-order nonlinear fitting

将文献[7-8]中使用的非参数回归、K-Means 方法, 以及 Adaboost CART 回归树算法运用在本文中的训练集和测试集上。表 2 为深度残差网络与上述三种方法的预测误差统计, 图 3 为深度残差网络与其余三种方法预测误差的分布情况。可见深度残差网络性能远优于非参数回归等方法。由于非参数回归中的核宽采用自适应方式, 取得预测样本与训练

集各个样本的最小距离, 故在大样本数据集上耗时过大。与非参数回归相比, 训练良好的深度残差网络预测速度快, 并且准确率更高, 误差更小, 更符合大数据时代光谱处理的要求。相较于 Adaboost 算法需要训练多组弱回归(分类)器, 本文的深度残差网络只需训练一个模型即可。

表 2 深度残差网络与非参数回归等方法的预测误差统计
Table 2 The statistical error of prediction by Deep residual network, Nonparametric regression, et al.

方法	最大绝对误差	平均绝对误差	标准差
深度残差网络	2.98	0.03	0.05
非参数回归	2.96	0.19	0.22
Adaboost CART Tree	3.19	0.38	0.32
K-Means	3.12	0.37	0.39

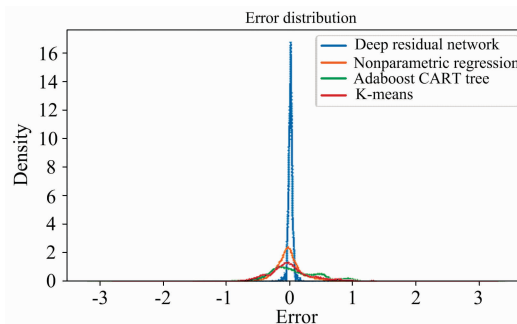


图 3 深度残差网络与非参数回归等其他方法预测误差分布情况

Fig. 3 The error distribution of prediction by Deep residual network, nonparametric regression et al.

2.3 模型分析

利用文献[10]中提出的类别激活映射(CAM)方法分析模型在给一条光谱预测时所关注的一些特征, 通过此分析模型可以对分类结果做出解释。将得到的 CAM 进行伪彩色变换, 拉伸, 并与光谱图像加权求和, 便可得到图 4 所示的类别特征映射图像, 其中颜色越接近红色的波段对分类越重要。实验中从 A, F, G, K 各抽取 2 条光谱画出 CAM 图像, 并在每幅图下给出了各类别的分数。对于 A 型恒星光谱, 模型关注的区域为 H 原子吸收线存在的波段, 红色精确覆盖了 Hbeta, Hgamma 和 Hdelta, 但忽视了 Halpha, 初步推断是 Halpha 较弱的原因。对于 F 型恒星光谱, 模型的关注区域为一阶 Ca 离子线存在的波段, H 原子吸收线存在的波段, 以及一阶 S 离子线存在的波段。在 F 型恒星中, 中性 H 原子谱线和一阶金属离子谱线都是比较明显的。对于 G 型恒星光谱, 模型关注区域大致在 $3\ 800 \sim 4\ 400 \text{ \AA}$ 波段, G 型星中 Ca 离子线达到了最强, 并且出现一阶 Fe 离子线与一阶 Ti 离子线, 这些谱线存在于这个波段。对于 K 型恒星光谱, 其主要以金属谱线为主, 模型主要以 Mg 线 $5\ 179 \text{ \AA}$ 附近以及 $3\ 699 \sim 4\ 390 \text{ \AA}$ 波段为判别依据。

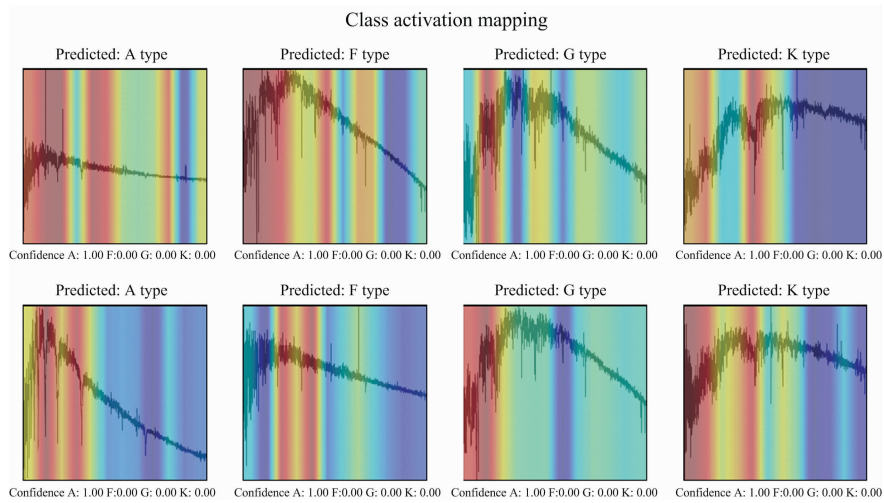


图 4 类别特征映射图 (CAM)

Fig. 4 Class activation mapping

3 结 论

光谱分类是天文数据处理的重要一环, 目前被广泛使用的模板匹配方法存在计算冗余、依赖数据质量等问题, 其他一些方法大都没有涉及光谱次型的分类。本文提出基于深度残差网络的深度学习模型来对光谱类别进行预测, 并赋予了模型可解释性。实验结果表明, 本方法在所使用的 LAMOST 数据集上可以将 91.4% 光谱预测误差保证在 0.5 个光谱次

型以内, 预测平均绝对误差为 0.3 个光谱次型。与非参数回归等方法相比有更高的准确率和预测速度。在模型分析中, 本文讨论了模型分类依据, 主要包括 Balmer 线系、金属离子谱线。对比文献[4]中线指数分类提出的, H γ , Fe 和 Mg 的组合对 O-G 分类较好, Fe, TiO₂ 和 G4300 的组合对晚期恒星分类较好, 本文 CAM 图像与文献[4]的结果基本相符, 下一步工作将通过修改模型输出维度来提高 CAM 的定位精度。

References

- [1] Luo A-li, Zhao Yongheng, Zhao G, et al. RAA (Research in Astronomy and Astrophysics), 2015, 15: 1095.
- [2] Yi Zhenping, Pan Jingchang. Image and Signal Processing (CISP). 3rd International Congress, 2010.
- [3] Schierscher F, Paurzen E. Astronomische Nachrichten, 2011, 332(6): 597.
- [4] Liu Chao, Cui Wenyuan, Zhang Bo, et al. Research in Astronomy and Astrophysics, 2015, 15: 1137.
- [5] Kaushal Sharma, Ajit Kembhavi, Aniruddha Kembhavi, et al. Monthly Notices of the Royal Astronomical Society, 2020, 491: 2280.
- [6] Gray R O, Corbally C J. The Astronomical Journal, 2014, 147(4): 80.
- [7] LIU Rong, QIAO Xue-jun, ZHANG Jian-nan, et al(刘 蓉, 乔学军, 张健楠, 等). Spectroscopy and Spectral Analysis(光谱学与光谱分析), 2017, 37(5): 1553.
- [8] Kheirdastan S, Bazarghan M. Astrophysics and Space Science, 2016, 361(9): 304.
- [9] He K, Zhang X, Ren S, et al. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016. 770.
- [10] Zhou Bolei, Aditya Khosla, Agata Lapedriza, et al. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016. 2921.

Prediction of Stellar Spectrum Categories Based on Deep Residual Network

WANG Tian-xiang^{1, 2}, FAN Yu-feng^{1*}, WANG Xiao-li¹, LONG Qian¹, WANG Chuan-jun¹

1. Yunnan Observatories, Chinese Academy of Sciences, Kunming 650011, China

2. University of Chinese Academy of Sciences, Beijing 100049, China

Abstract In recent years, the spectral data of celestial bodies observed have achieved a dramatic increase thanks to the successful implementation of various projects of spectral sky survey. Therefore, higher requirements for the automatic classification and analysis of spectrum are proposed for large-scale projects of spectral sky survey. The classification problem is transformed into a regression one in this paper, and a method of spectral category regression based on the residual depth network is put forward to conduct a prediction of MK spectral subtype on stellar spectrum. The network is mainly composed of 25 convolution layers, 1 maximum pooling layer, 1 average pooling layer, full connection layer and 12 residual structures. The maximum pooling layer is used to filter features, the convolution layer to extract features, and the average pooling layer to reduce parameters and improve efficiency. The residual structure can prevent the degradation of the network, extract high-dimensional abstract features by deepening the network and improve training speed. Considering the non-zero probability of data with false labels and corrupted data, Log-Cosh is adopted as a loss function in this paper to reduce the negative impact of bad samples. 80 000 spectra that are randomly selected from LAMOST DR5 are used as the experimental data. The spectra are divided into the training set, verification set and test set according to the proportion of 7 : 1 : 2 after eliminating the bad value and normalizing the flow. The experiment includes two parts. In the first part, the network is adopted to carry out a prediction on the spectral subtype, and the maximum absolute error, the average absolute error and the standard deviation are used to compare the performance of convolution kernels with different shapes. The predicted value is taken as the abscissa and the label as the ordinate, and the second-order nonlinear fitting is used for all sample points in the test set, a straight line that is coincident with $y=x$ is obtained, proving that the model can predict the spectral subtype well. The second part is concerning the internal analysis of the model. The main characteristics of the model in predicting four types of spectra, A, F, G, K, are mainly explored with the method of category activation mapping, thus endowing the model with interpretability. In the text data set, 91.4% of the spectral prediction errors of this method are within 0.5 spectral subtypes, and the average absolute error of the prediction is 0.3 spectral subtypes. It is shown that the method proposed in this paper can better predict spectral subtypes with faster speed and higher accuracy according to the comparison of the same data set with nonparametric regression, Adaboost regression tree and K-means.

Keywords Stellar spectrum; MK classification; Deep learning; Regression; Feature mapping

(Received Mar. 15, 2020; accepted Jul. 6, 2020)

* Corresponding author