

Stellar spectral classification and feature evaluation based on a random forest

Xiang-Ru Li, Yang-Tao Lin and Kai-Bin Qiu

School of Mathematical Sciences, South China Normal University, Guangzhou 510631, China; xiangru.li@gmail.com

Received 2018 December 17; accepted 2019 March 18

Abstract With the availability of multi-object spectrometers and the design and operation of some large scale sky surveys, the issue of how to deal with enormous quantities of spectral data efficiently and accurately is becoming more and more important. This work investigates the classification problem of stellar spectra under the assumption that there is no perfect absolute flux calibration, for example, when considering spectra from the Guo Shou Jing Telescope (the Large Sky Area Multi-Object Fiber Spectroscopic Telescope, LAMOST). The proposed scheme consists of the following two procedures: Firstly, a spectrum is normalized based on a 17th order polynomial fitting; secondly, a random forest (RF) is utilized to classify the stellar spectra. Experiments on four stellar spectral libraries show that the RF has good classification performance. This work also studied the spectral feature evaluation problem based on RF. The evaluation is helpful in understanding the results of the proposed stellar classification scheme and exploring its potential improvements in the future.

Key words: methods: statistical — methods: data analysis — virtual observatory tools

1 INTRODUCTION

With the development of modern telescopes, large quantities of spectra have been and are being obtained. In this massive spectrum scenario, traditional manual data processing methods and schemes with many human interventions cannot satisfy the requirements. Therefore, automatic classification is an imperative in large sky surveys and has attracted much attention (Gulati et al. 1994; von Hippel et al. 1994; Gray et al. 2009; Crowther & Walborn 2011).

Therefore, a series of schemes has been investigated for automatic classification of spectra in the last thirty years. The two most widely used schemes are template matching and artificial neural networks (ANNs). The template matching method is implemented by minimizing some metric distances or maximizing some kinds of similarity between a reference spectrum and a spectrum to be classified (Kurtz 1984; LaSala 1994; Malyuto 2002; Giridhar et al. 2006; Lee et al. 2008; Duan et al. 2009; Gray et al. 2016), for example, χ^2 minimization. The ANN method classifies a spectrum by establishing a mapping from a spectrum to its spectral type or subtype (Bailer-Jones 1997; Bailer-Jones et al. 1998; Singh et al. 1998; Weaver 2000; Bai et al. 2005; Bazarghan & Gupta 2008; Mahdi 2008; Navarro et al. 2012; Kheirdastan & Bazarghan 2016). Expert systems (Manteiga et al. 2009; Gray & Corbally 2014), support vector machines (Liu et al.

2015; Kheirdastan & Bazarghan 2016) and K-means (Qin et al. 2001; Kheirdastan & Bazarghan 2016) are also being investigated for the classification of stellar spectra.

In automatic classification of stellar spectra, a key problem is how to represent the information in a spectrum. This information representation problem is referred to as feature extraction in the machine learning community. This information representation not only affects the accuracy of a spectral classification system, its robustness to noise and calibration distortion, but also its interpretability/understandability. The interpretability means the difficulty to evaluate or identify the contribution of a specific wavelength range or spectral line in spectral classification. Good interpretability helps us understand our automatic classification scheme and its physical indications, and design an improved method by taking some physical knowledge into account.

Two typical information representation methods for a spectrum are spectral index (Malyuto et al. 1997; Lee et al. 2008; Manteiga et al. 2009; Liu et al. 2015) and principal component analysis (PCA) (Qin et al. 2001; Mahdi 2008; Kheirdastan & Bazarghan 2016). A spectral index can be an integration of spectral fluxes within a preset wavelength range, or some kind of description of a spectral line, for example, the full width at half maximum (FWHM). The greatest advantage of the spectral index methodology is its

interpretability. PCA is a kind of data compression method, and is used to obtain a compact representation by statistically minimizing the difference between some spectra and their representations. Actually, the PCA representation is a linear sum of spectral fluxes. Therefore, one typical limitation of the PCA approach is that it is difficult to evaluate/trace back the contribution of a local wavelength range of a spectrum, which is closely related to the interpretability of the computed results.

This work studies the automatic classification of stellar spectra using a random forest (RF). An RF consists of a series of decision trees. A decision tree makes spectral classification using a small subset of fluxes on automatically selected wavelength positions. This results in the RF sometimes having the potential for superiority of interpretability. In applications, furthermore, the effects of noise and calibration imperfections can vary from spectrum to spectrum, and from wavelength to wavelength. Fortunately, the choices of effective wavelength positions are different from tree to tree in an RF. This diversity in wavelength selection makes an RF achieve good classification performance by adaptively selecting appropriate combinations of wavelength positions in the competition between the decision trees of an RF. Therefore, this work investigates the stellar spectral classification problem using the RF approach.

2 ARCHITECTURE OF THE PROPOSED SCHEME AND DATA PREPROCESSING

To reduce some negative effects arising from the incompleteness of flux calibration, we first do some preprocessing of observed spectra. Then, the stellar spectra are classified using the RF. A flowchart of the proposed scheme is presented in Figure 1.

2.1 Flux Normalization

In spectral data, the observed radiant energies from some celestial bodies with the same spectral type may vary greatly in magnitude due to detector sensitivity, brightness of the celestial body or its distance from the Earth. Some negative effects from magnitude uncertainty can be eliminated or reduced by normalizing the flux. Suppose \mathbf{x} is a spectrum, denoted as: $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, which is a vector in an n -dimensional space. The spectral flux can be normalized using the following formula (Xu et al. 2006)

$$\mathbf{y} = \frac{\mathbf{x}}{\sqrt{\sum_{i=1}^n x_i^2}}. \quad (1)$$

2.2 Continuum Normalization based on Polynomial Fitting

This paper assumes that there is no perfect absolute flux calibration. Therefore, the classification algorithms cannot be used to classify spectra firstly, but rather, the continuum normalization is. In this paper, a 17-th order polynomial fitting method is used to approximate the continuum spectrum in a stellar spectrum. Then, the continuum components computed from the spectra are removed, leaving the spectral lines. Finally, a classification algorithm, RF, is utilized for the processed spectrum.

Figures 2 and 3 show some continuum normalization results for some spectra from O, B, A, F, G, K and M spectral types. Their continua are fitted using a polynomial with order 17. The results demonstrate that the spectral line characteristics are preserved well.

3 CLASSIFYING A STELLAR SPECTRUM USING A RANDOM FOREST

An RF algorithm is an extension of the traditional decision tree. RF is implemented by combining multiple decision trees. A series of research works have verified that this combination clearly improves classification performances and increases the robustness to outliers and noise (Ho 1998; Breiman 2001). This section gives a brief introduction to the procedures involved in building an RF.

Because an RF is established by assembling a series of decision trees, the decision tree will be introduced followed by the assembly scheme.

3.1 Decision Tree

A decision tree classifier is a tree-like model. In a decision tree, there are three types of nodes: a root node, some branch nodes and some leaf nodes. If a node ‘S’ accepts signals from node ‘T’, ‘T’ is called the parent node of ‘S’ and ‘S’ is one child node of ‘T’. A root node does not have any parent node and there is a unique root node in a decision tree. Each leaf node has a parent node but does not have any child node. A branch node has one parent node and one or more child nodes. Signals can only move directly from a parent node to one of its child nodes.

A spectrum is classified by moving from root node to one leaf node. Suppose there are K leaf nodes $\{\text{leaf}_i, i = 1, \dots, K\}$ and a training set S . Based on the leaf node that a training sample can reach, the training set can be split into K subsets S_1, S_2, \dots, S_K . A subset S_k is labeled with the most frequent class in it, and is denoted by label_k . If a spectrum to be dealt with moves from root to node leaf_{k^*} , this spectrum is classified into type label_{k^*} .

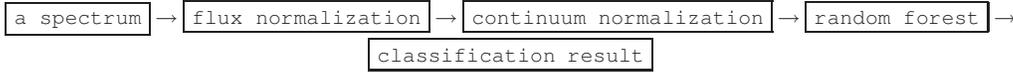


Fig. 1 A flowchart of the proposed classification scheme for stellar spectra.

To construct a decision tree, one fundamental problem is to determine which data property should be used in one parent node. For interested readers, further introduction to decision trees can be found in Quinlan (1986); Rokach & Maimon (2008).

3.2 Random Forest

An RF does the classification of a stellar spectrum by establishing a series of decision trees and fusing their results. Suppose S_{tr} is a training set consisting of N spectra. A novel spectral set can be generated by randomly selecting N samples from S_{tr} with replacement, and can be referred to as a bootstrap set. ‘With replacement’ means that in case of a spectrum being sampled into a bootstrap set, there is still a probability that this spectrum is sampled in the future for this bootstrap set. Therefore, some spectra in the training set may appear more than once in a bootstrap set and some other spectra probably are not present in this bootstrap set. By doing this, we can generate a number of bootstrap sets from a training set. From every bootstrap set, a decision tree is learned, and a number of decision trees is learned from these bootstrap sets, which form an RF. More about the RF can be found in Ho (1998); Breiman (2001); Hastie et al. (2008). An algorithm for building an RF is as follows:

Algorithm 3.1 RF Classifier
 Input: Training set S_{tr} , test set S_{te} , number, M , of trees
 Output: Estimated class label for every test sample
 Steps:
 1 Let $i = 1$.
 2 Generate a bootstrap set from S_{tr} and denote the set with $S_{\text{tr}}^{\text{bs}}$.
 3 Construct a decision tree from $S_{\text{tr}}^{\text{bs}}$ and denote this decision tree with tree_i .
 4 Let $i = i + 1$.
 5 Repeat steps 2, 3 and 4 M times.
 6 Estimate the class label for every test sample in S_{te} using $\{\text{tree}_i, i = 1, \dots, M\}$ and fuse the estimations from different decision trees by the corresponding majority votes as the final classification result.

4 EXPERIMENTS AND DISCUSSION

4.1 Data Sets

The proposed scheme was evaluated on two sets of stellar spectral libraries. These two spectral sets are referred to as the Jacoby, Silva and Pickles (JSP) data, from the corresponding publications listed below, and Large Sky Area Multi-Object Fiber Spectroscopic Telescope (LAMOST) data. Note, these data are described firstly in this subsection.

4.1.1 The JSP spectral set

This data set consists of 359 spectra from three representative stellar spectral libraries from Jacoby et al. (1984), Silva & Cornell (1992) and Pickles (1998). Each of the three libraries covers the spectral types from O to M.

The Jacoby spectral library has 159 spectra with a configuration of $0.14 \text{ nm pixel}^{-1}$ and a wavelength range of $351.1 - 742.8 \text{ nm}$. The Silva spectral library has 71 spectra with $0.5 \text{ nm pixel}^{-1}$ and a wavelength range of $351.0 - 893.0 \text{ nm}$. The Pickles spectral library has 129 spectra with $0.5 \text{ nm pixel}^{-1}$ and a wavelength range of $360.0 - 900.0 \text{ nm}$. In order to analyze them on the same scale, all of the spectra are resampled with a step of 0.5 nm using a linear interpolation in the wavelength range of $385.0 - 600.0 \text{ nm}$.

4.1.2 LAMOST spectral set

From Data Release 5 (DR5) published by the LAMOST project, we select 6000 stellar spectra with SNRU, SNRG, SNRR, SNRI (the signal to noise ratios of the u , g , r , i bands respectively) all higher than 20. To be consistent with the JSP spectral set, all of the LAMOST spectra are also resampled with a step of 0.5 nm using a linear interpolation on the wavelength range of $385.0 - 600.0 \text{ nm}$.

4.1.3 Spectral class representations

In the automatic classification of a stellar spectrum, a fundamental problem is how to represent the class in a computer. We represent the spectral types (O, B, A, F, G, K, M) using integers $1 \sim 7$ respectively. Ten spectral subtypes are represented by the product of the subtype number and 0.1. For example, the category of an A0 star is denoted by 3.0, and the category of an F5 star by 4.5.

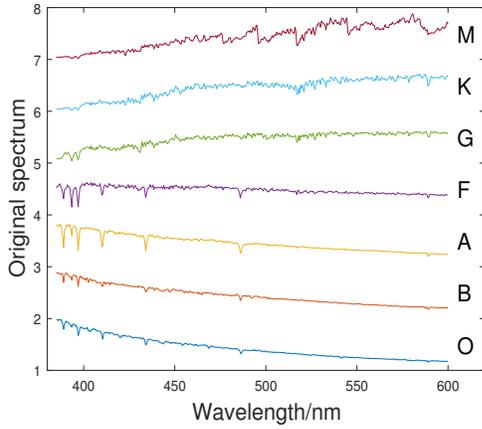


Fig. 2 Seven stellar spectra. nm: nanometer.

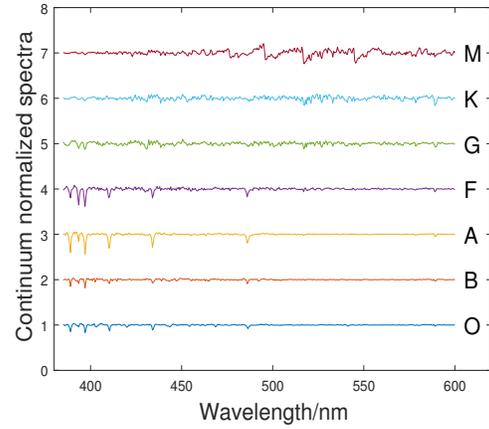


Fig. 3 The continuum normalized spectra in Fig. 2.

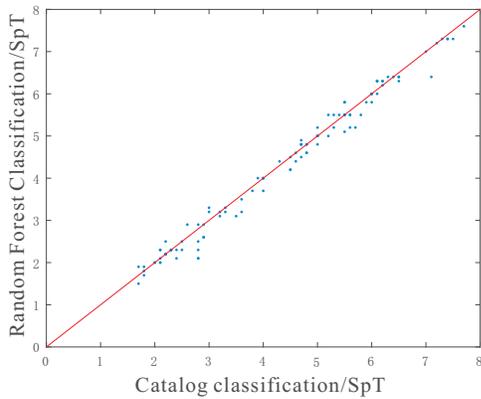


Fig. 4 Experimental results on the JSP spectral set. The vertical axis is the estimation from the proposed scheme and the horizontal axis is the reference type. SpT: spectral types.

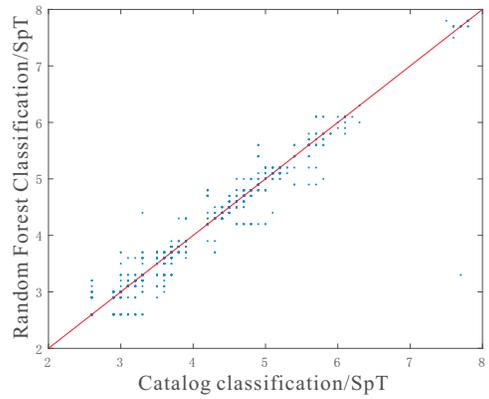


Fig. 5 Experimental results on LAMOST spectra. The vertical axis is the estimation from the proposed scheme and the horizontal axis is the reference type. SpT: spectral types.

4.2 Experimental Results

To evaluate the proposed scheme (Sects. 2 and 3), the JSP spectral set is randomly divided into two subsets, a training set and a test set. The training set consists of 70% JSP spectra and is used for learning the parameters of the RF. The other 30% of the spectra form the test set. The test set is used to evaluate the performance of the learned RF. This evaluation result may depend on the dividing of the training set and test set. To alleviate this issue and increase the objectiveness of the evaluation, we repeat the above-mentioned procedures 10 times, and take the average of ten experimental evaluations as the final result. The experimental results are presented in Figure 4. The evaluation of the LAMOST spectra set is conducted similarly and the results are displayed in Figure 5.

The results in Figure 4 show strong consistency between the estimation from the proposed scheme and the reference types. In experiments conducted on LAMOST

spectra (Fig. 5), a strong discrepancy is present in one spectrum. This discrepancy is indicated by the point in the bottom-right corner of Figure 5. This spectrum is depicted in Figure 6. Its reference type is M7 and the estimation is A3. After checking with help from Dr. Xiao Kong from the LAMOST project, we confirmed that this is a spectrum of a binary star with component types M and A.

To quantitatively evaluate the performance of the proposed scheme, we use four measures — mean of the squared difference (MSD), mean of the absolute difference (MAD), mean of difference (MD) and accuracy of spectral type (AST). Suppose $S = \{(\mathbf{x}^i, y_i), i = 1, \dots, m\}$ is a set of spectra with their associated type labels; \hat{y}_i is the estimation of y_i , where m is an integer representing the number of spectra in S . On S , MSD, MAD and MD are defined as follows:

$$\text{MSD}(S) = \sqrt{\frac{\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2}{m}}, \quad (2)$$

Table 1 Quantitative Performance Evaluation

Spectral set	AST	MSD	MAD	MD
JPS spectra	0.9537	0.2151	0.1481	0.0574
LAMOST spectra	0.9377	0.1954	0.0787	0.0067

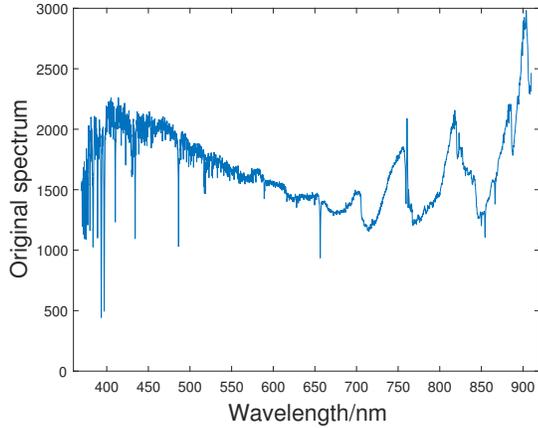


Fig. 6 The spectrum with the most notable classification inconsistency in Fig. 5.

$$\text{MAD}(S) = \frac{\sum_{i=1}^m |y^{(i)} - \hat{y}^{(i)}|}{m}, \quad (3)$$

$$\text{MD}(S) = \frac{\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})}{m}. \quad (4)$$

Suppose there are n spectra in S whose estimated spectral types are consistent with their reference value. The corresponding AST is defined as

$$\text{AST}(S) = n/m. \quad (5)$$

Some quantitative evaluation results are presented in Table 1.

4.3 Effects of Spectrum Preprocessing

In the proposed scheme (Fig. 1), two essential procedures are flux normalization and continuum normalization. There are more or fewer deviations and distortions in observed spectra. Therefore, these two preprocessing procedures evidently improve the spectral classification performance on both JPS data and LAMOST data (Table 2). In particular, the JPS spectra are observed with multiple telescopes and calibrated using multiple pipelines, and a greater variety of calibration deviation and distortions exists in them. Therefore, much larger performance improvement is observed on JPS spectra than on LAMOST spectra (Table 2).

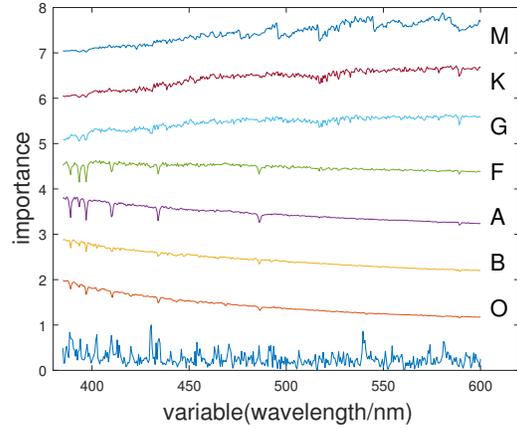


Fig. 7 The importance scores of stellar spectrum features on the JPS spectral data. The above seven curves are computed based on some spectra with spectral types O~M respectively; the bottom curve indicates the effectiveness of the spectral fluxes.

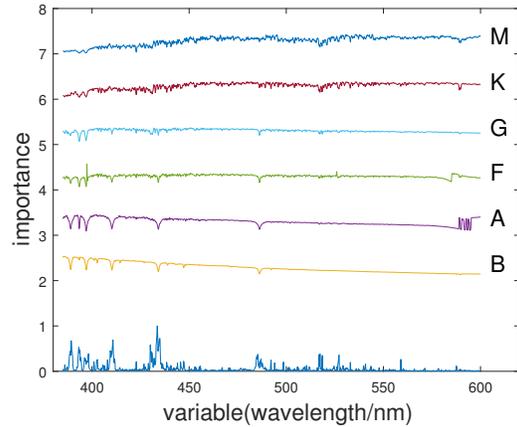


Fig. 8 The importance scores of stellar spectrum features on the LAMOST data. The above six curves are computed from spectra with spectral types B~M respectively; the bottom curve indicates the effectiveness of the spectral fluxes.

Table 2 Effects of Spectrum Preprocessing

IE	FN	CN	AST	MSD	MAD	MD
(a) On JPS spectra						
1	no	no	0.4630	1.9975	1.2120	1.0954
2	yes	no	0.8333	0.2287	0.1620	0.0250
3	no	yes	0.4907	2.0000	1.2231	1.1472
4	yes	yes	0.9537	0.2151	0.1481	0.0574
(b) On LAMOST spectra						
5	no	no	0.8391	0.3025	0.1642	0.0175
6	yes	no	0.9289	0.2242	0.0723	0.0067
7	no	yes	0.9091	0.2152	0.0999	0.0043
8	yes	yes	0.9377	0.1954	0.0787	0.0067

IE: index of an experiment; FN: flux normalization; CN: continuum normalization.

4.4 Comparisons with Related Works in the Literature

Zhang et al. (2009) studied the classification of stellar spectra using a non-parametric regression method with a continuum spectrum normalization on the JSP data, and achieved an accuracy of MSD=0.3226 and MAD=0.2554. Kheirdastan & Bazarghan (2016) obtained three accuracies of MSD=1.39, 1.53 and 1.64 using an ANN, as well as the SVM and K-means methods combined with PCA on some spectra from SEGUE-2 (Yanny et al. 2009) and SEGUE-1 from the Sloan Digital Sky Survey (SDSS) III. On the JSP spectra, Liu et al. (2017) achieved an accuracy of MSD=0.2214 and MAD=0.1632 using the non-parametric regression method. The experimental results in Table 2 show that the RF approach enables good performance on both the JSP spectra and LAMOST spectra.

5 SPECTRAL FEATURE EVALUATIONS

The evaluation of variable effectiveness is a fundamental procedure to understand the potential physical interpretations and study more effective schemes. The RF algorithm estimates the importance of a variable by looking at how much prediction error increases in the case of one variable being permuted with all others left unchanged. Conventional calculation methods of variable importance measure (VIM) in an RF are divided into two types: One is based on the Gini index and the other is the Out-of-Bag (OOB) data error rate. The score statistics for the variable X_j are denoted by $VIM_j^{(Gini)}$ and $VIM_j^{(OOB)}$ respectively. Interested readers are referred to Breiman (2002) for their definitions. In the existing literature on RF, the $VIM_j^{(OOB)}$ score statistic is more extensive than the $VIM_j^{(Gini)}$ score statistic. Therefore, this article ranks the importance of variables based on the $VIM_j^{(OOB)}$ score statistic.

The evaluation results are presented in Figures 7 and 8. In Figure 7, the eight curves from the bottom to top are the importance scores of spectral features at every wavelength computed from the JSP spectral data with spectral types from O to M respectively. The seven curves in Figure 8 from the bottom to top are the importance scores of the spectral features at every wavelength computed from the LAMOST spectral data with spectral types from B to M respectively, showing the relationship between the important spectral variables and spectral lines of each type. The results from these figures signify that for the spectral data from different systems, the important features selected by the RF are approximately similar, which indicates that the RF selects the important spectral lines from each type of spectrum as the basis for classification. The evaluation is helpful in understanding the results of the proposed stellar

classification scheme and exploring its potential improvements in the future.

6 CONCLUSIONS

Although there is a series of research papers in the literature on the automatic classification of stellar spectra, the performance of automatic classification is still being improved, especially for some spectra without flux calibration or with only relative flux calibration, for example, spectra from LAMOST.

This work proposed a stellar spectral classification scheme based on an RF, and experimental results demonstrate its superiority on real spectral data. The characteristics of this work are a comprehensive investigation of the effects from flux normalization and continuum normalization. This work also studied the evaluation of spectral features. This evaluation is helpful in understanding potential physical interpretations and designing more effective schemes.

Acknowledgements This work is supported by the National Natural Science Foundation of China (Grant Nos: 61273248 and 61075033), the Natural Science Foundation of Guangdong Province (2014A030313425 and S2011010003348), China Scholarship Council (201706755006) and the Joint Research Fund in Astronomy (U1531242) under cooperative agreement between the National Natural Science Foundation of China and Chinese Academy of Sciences.

References

- Bai, L., Guo, P., & Hu, Z.-Y. 2005, *ChJAA* (Chin. J. Astron. Astrophys.), 5, 203
- Bailer-Jones, C. A. L. 1997, *PASP*, 109, 932
- Bailer-Jones, C. A. L., Irwin, M., & von Hippel, T. 1998, *MNRAS*, 298, 361
- Bazarghan, M., & Gupta, R. 2008, *Ap&SS*, 315, 201
- Breiman, L. 2001, *Machine Learning*, 45, 5
- Breiman, L. 2002, *Manual on Setting Up, Using, and Understanding Random Forests v3.1*, Statistics Department University of California Berkeley, CA, USA
- Crowther, P. A., & Walborn, N. R. 2011, *MNRAS*, 416, 1311
- Duan, F.-Q., Liu, R., Guo, P., Zhou, M.-Q., & Wu, F.-C. 2009, *RAA* (Research in Astronomy and Astrophysics), 9, 341
- Giridhar, S., Muneer, S., & Goswami, A. 2006, *Mem. Soc. Astron. Italiana*, 77, 1130
- Gray, R. O., Corbally, C. J., & Burgasser, A. J. 2009, *Stellar Spectral Classification* (Princeton, NJ: Princeton Univ. Press)
- Gray, R. O., & Corbally, C. J. 2014, *AJ*, 147, 80
- Gray, R. O., Corbally, C. J., De Cat, P., et al. 2016, *AJ*, 151, 13
- Gulati, R. K., Gupta, R., Gothoskar, P., & Khobragade, S. 1994, *ApJ*, 426, 340

- Hastie, T., Tibshirani, R., Friedman, J. 2008, *The Elements of Statistical Learning* (2nd ed., Springer)
- Ho, T.K. 1998, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20, 832
- Jacoby, G. H., Hunter, D. A., & Christian, C. A. 1984, *ApJS*, 56, 257
- Kheirdastan, S., & Bazarghan, M. 2016, *Ap&SS*, 361, 304
- Kurtz, M. J. 1984, in *The MK Process and Stellar Classification*, ed. R. F. Garrison, 136
- LaSala, J. 1994, in *Astronomical Society of the Pacific Conference Series*, 60, *The MK Process at 50 Years: A Powerful Tool for Astrophysical Insight*, ed. C. J. Corbally, R. O. Gray, & R. F. Garrison, 312
- Lee, Y. S., Beers, T. C., Sivarani, T., et al. 2008, *AJ*, 136, 2022
- Liu, C., Cui, W.-Y., Zhang, B., et al. 2015, *RAA (Research in Astronomy and Astrophysics)*, 15, 1137
- Liu, R., Qiao, X. J., Zhang, J. N., & Duan, F. Q. 2017, *Spectroscopy and Spectral Analysis*, 37, 1555
- Mahdi, B. 2008, *Bulletin of the Astronomical Society of India*, 36, 1
- Malyuto, V., Oestreicher, M. O., & Schmidt-Kaler, T. 1997, *MNRAS*, 286, 500
- Malyuto, V. 2002, *New Astron.*, 7, 461
- Manteiga, M., Carricajo, I., Rodríguez, A., Dafonte, C., & Arcay, B. 2009, *AJ*, 137, 3245
- Navarro, S. G., Corradi, R. L. M., & Mampaso, A. 2012, *A&A*, 538, A76
- Pickles, A. J. 1998, *PASP*, 110, 863
- Qin, D.M., Hu, Z.Y., & Zhao, Y.H. 2001, In *Object Detection, Classification, and Tracking Technologies*, 4554, 268, *International Society for Optics and Photonics*
- Quinlan, J. R. 1986, *Machine Learning*, 1, 81
- Rokach, L., & Maimon, O. 2008, *Data Mining with Decision Trees: Theory and Applications* (World Scientific Pub Co Inc.)
- Silva, D. R., & Cornell, M. E. 1992, *ApJS*, 81, 865
- Singh, H. P., Gulati, R. K., & Gupta, R. 1998, *MNRAS*, 295, 312
- von Hippel, T., Storrie-Lombardi, L. J., Storrie-Lombardi, M. C., & Irwin, M. J. 1994, *MNRAS*, 269, 97
- Weaver, W. B. 2000, *ApJ*, 541, 298
- Xu, X., Wu, F.-C., & Hu, Z.-Y. 2006, *Spectroscopy & Spectral Analysis*, 26, 182
- Yanny, B., Rockosi, C., Newberg, H. J., et al. 2009, *AJ*, 137, 4377
- Zhang, J. N., Zhao, Y. H., & Liu, R. 2009 *Spectroscopy and Spectral Analysis*, 29, 3424.