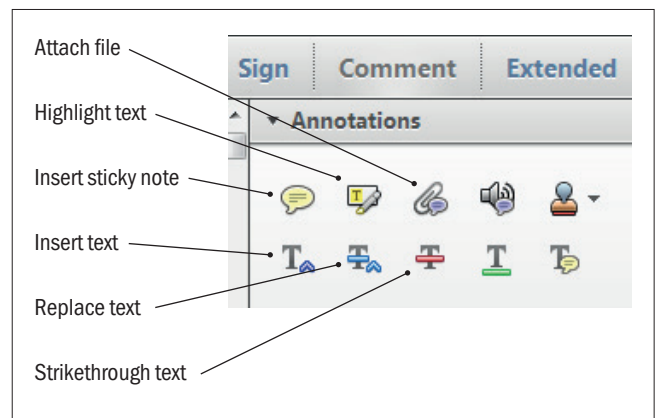


Making corrections to your proof

Please follow these instructions to mark changes or add notes to your proof. You can use Adobe Acrobat Reader (download the most recent version from <https://get.adobe.com>) or an open source PDF annotator.

The tools you need to use are contained in **Annotations** in the **Comment** toolbar. You can also right-click on the text for several options. The most useful tools have been highlighted here. If you cannot make the desired change with the tools, please insert a sticky note describing the correction.

Please ensure all changes are visible via the 'Comments List' in the annotated PDF so that your corrections are not missed.

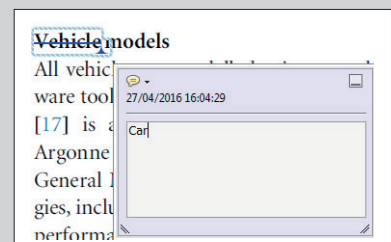


Do not attempt to directly edit the PDF file as changes will not be visible.



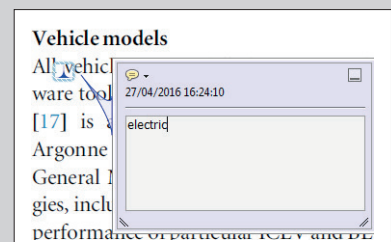
Replacing text

To replace text, highlight what you want to change then press the replace text icon, or right-click and press 'Add Note to Replace Text', then insert your text in the pop up box. Highlight the text and right click to style in bold, italic, superscript or subscript.



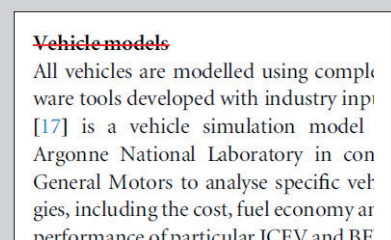
Inserting text

Place your cursor where you want to insert text, then press the insert text icon, or right-click and press 'Insert Text at Cursor', then insert your text in the pop up box. Highlight the text and right click to style in bold, italic, superscript or subscript.



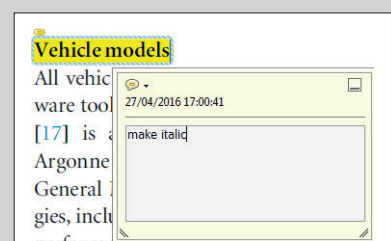
Deleting text

To delete text, highlight what you want to remove then press the strikethrough icon, or right-click and press 'Strikethrough Text'.



Highlighting text

To highlight text, with the cursor highlight the selected text then press the highlight text icon, or right-click and press 'Highlight text'. If you double click on this highlighted text you can add a comment.



QUERY FORM

JOURNAL: The Astronomical Journal

AUTHOR: Brice & Andonie

TITLE: Automated Morgan Keenan Classification of Observed Stellar Spectra Collected by the Sloan Digital Sky Survey Using a Single Classifier

ARTICLE ID: ajab40d0

Your article has been processed in line with the journal style. Your changes will be reviewed by the Production Editor, and any amendments that do not comply with journal style or grammatical correctness will not be applied and will not appear in the published article.

The layout of this article has not yet been finalized. This proof may contain columns that are not fully balanced/matched or other typographical imperfections; these issues will be resolved once the final corrections have been incorporated.

Please check that the **names of all authors as displayed in the proof are correct**, and that all **authors are linked to the correct affiliations**. Please also confirm that the correct corresponding author has been indicated. **Note that this is your last opportunity to review and amend this information before your article is published.**

If you have an Acknowledgments section, please check that the information there is complete and correct and that all relevant institutions, grant numbers, programs, and collaborators are mentioned as appropriate.

We have been provided with ORCID iDs for the authors as below. Please confirm whether the numbers are correct.
Michael J. Brice 0000-0002-1278-2404

Page 1

Q1
Please kindly define “CV.”

Page 3

Q2
Please kindly define “BOSS.”

Page 3

Q3
If possible, please kindly define “SMOTE.”

Page 6

Q4
Tables [7, 8] are not cited in order in the text. Please cite this table in sequence in the text.

Page 7

Q5
Please check the details for any journal references that do not have a link; please update with correct details and supply a CrossRef DOI or NASA ADS link if available.

Page 7

Q6

Please update if possible all references where only the arXiv preprint number is given.

Page 7

Q7

Please provide the publisher details in references [Bolón-Canedo et al. 2015, Japkowicz 2000, Oliphant 2015].

Page 7

Q8

Please provide conference details in references [Brice & Andonie 2019, Elting et al. 2008, Kohavi 1995, Zhang et al. 2008]. The format for conference citations is [Publication, Conference title, editors (Publisher City, State: Publisher Name), page number

Page 7

Q9

Please provide the initials for the author [Ailab] in reference [Powers & Ailab 2011].

Page 7

Q10

Please update reference [Yi & Pan 2010] if possible.



Automated Morgan Keenan Classification of Observed Stellar Spectra Collected by the Sloan Digital Sky Survey Using a Single Classifier*

Michael J. Brice¹ and Răzvan Andonie

Computer Science Department Central Washington University Ellensburg, WA 98926, USA
Received 2019 June 22; revised 2019 August 18; accepted 2019 August 30; published 2019 MM DD

Abstract

The classification of stellar spectra is a fundamental task in stellar astrophysics. Stellar spectra from the Sloan Digital Sky Survey are applied to standard classification **methods**, **k-nearest neighbors** and **random forest**, to automatically classify the spectra. Stellar spectra are high dimensional data and the dimensionality is reduced using astronomical knowledge because classifiers work in low dimensional space. These methods are utilized to classify the stellar spectra into a complete **Morgan Keenan** classification (spectral and luminosity) using a single classifier. The motion of stars (radial velocity) causes **machine-learning** complications through the feature matrix when classifying stellar spectra. Due to the nature of stellar classification and radial velocity, these complications cannot be corrected. However, classifiers utilizing a large set of observed stellar **spectra**, which has had **astronomical-specific** feature selection applied, performed computationally fast with extremely high accuracy.

Key words: Stellar Spectra Classification – Luminosity Classification – Morgan Keenan Classification – Automatic Classification – Redshift – Astronomical Feature Selection – K-Nearest Neighbors – Random Forest – Imbalanced Data

1. Introduction

Stellar classification is a fundamental task in stellar astrophysics. Traditionally, stellar spectra are classified by determining the wavelengths of absorption lines using wavelet transformations, statistical analysis, and using references to the Morgan Keenan (MK) **classification** scheme (Morgan et al. 1943) or they are classified by comparing the best fit of the spectra to that of templates using statistical tests (Duan et al. 2009). The traditional classification schemes require complex data transformations and analysis to identify the class of a star based on its spectrum.

The amount of astronomical data and dimensionality of said data is growing rapidly through more and more ambitious astronomical surveys. The Sloan Digital Sky Survey (SDSS) is an example of an ambitious astronomical survey with **high-quantity** and dimensional data.

Presently, SDSS is creating the most detailed **3D** maps of the **universe** ever made, with deep **multicolor** images of one-third of the sky, and spectra for more than three million astronomical objects¹ (York et al. 2000). The SDSS provides stellar spectra with observed wavelengths. The following experiments will classify stars using the SDSS data run 14 optical spectra data set.

The SDSS and other large astronomical surveys create challenging problems for a thorough and speedy analysis. As such, automated classification methods are explored. However, some classification algorithms are limited to low dimensional data, making the use of feature selection and feature extraction essential.

Radial velocity (RV) creates complications for the automated classification of stellar spectra through the **feature matrix**. The automated process for identifying RV and stellar class that the

SDSS uses is as follows (Bolton et al. 2012 and SkyServer: Redshifts, Classifications, and Velocity Dispersions²).

1. Redshift and classification templates for galaxy, quasar, and **CV** star classes are constructed by performing a rest-frame principal-component analysis (**PCA**; Shlens 2014) of training samples of a known redshift.
2. The combination of redshift and template class that yields the overall best fit (in terms of lowest reduced chi-squared) is adopted as the pipeline measurement of the redshift and classification of the spectrum.
3. The most common warning flag is set to indicate that the change in reduced chi-squared between the best and next-best redshift/classification is either less than 0.01 in an absolute sense, or less than 1% of the best model reduced chi-squared, which indicates a poorly determined redshift.

This paper proposes a novel approach to stellar classification characterized by the following:

1. **avoids** complex transformation and statistical analysis of the spectra space using machine **learning**;
2. **uses** spectra without RV **corrections**;
3. **and uses** astronomical knowledge to perform feature **selection**.

Stellar spectra are classified into a complete MK **classification** (spectral and luminosity) using a single classifier method. Astronomical knowledge is used to reduce the number of flux measurements. This results in key aspects of the spectra being preserved for classification which allows for a complete spectral and luminosity classification to be possible. However, the work conducted here deals with spectra with RV in the range of $\approx \pm 240 \text{ km s}^{-1}$.

The structure of this paper is as follows. Section 2 describes the approach to classification. Section 3 describes the

* Released on TBD.

¹ <https://www.sdss.org/>

² <https://www.sdss.org/dr12/algorithms/redshifts/>

experimental setup and the results. Section 4 provides a discussion of the results. Finally Section 5 provides the conclusions.

2. Approach to Classification

In this section, the data preprocessing, machine-learning classifiers, and feature selection are described.

2.1. Data Preprocessing

The only data preprocessing required for this approach is flux scaling using Equation (1). If a sample is known to have missing or corrupt flux measurements around the absorption lines used for feature selection, then an imputer³ method is required to fill in missing values. None of the samples in the data set used in this analysis required an imputer.

$$f_{i,\text{scaled}} = \frac{f_i - f_{\min}}{f_{\max} - f_{\min}} \quad (1)$$

where f_i is the i th flux measurement, f_{\max} and f_{\min} are the maximum and minimum flux measurements, respectively, and $f_{i,\text{scaled}}$ is the resulting scaled flux.

2.2. Machine Learning and Feature Selection

The classifier methods of k -nearest neighbors (KNN) and random forest (RF) are used in this approach. KNN classifies using the k -nearest known samples. More in depth explanations of KNN can be found in Marsland (2015), Ivezic et al. (2014), and Goldberger et al. (2005). RF classifies using a forest of decision trees, where each tree votes on the classification. More in depth explanations of RF can be found in Marsland (2015), Ivezic et al. (2014), and Breiman (2001). KNN and RF were chosen because they are widely used in astronomy (Yi & Pan 2010; Ivezic et al. 2014; Bai et al. 2019), and they work well in low dimensional spaces (Breiman 2001; Goldberger et al. 2005). RF was also chosen because it demonstrated good results in earlier work found in Brice & Andonie (2019).

The difference between the work presented here and the work of other authors is feature selection. Feature selection is the act of taking a set of attributes or features and extracting or transforming the most relevant ones for classification and to reduce the number of dimensions for the input space for the classifier model (Bolón-Canedo et al. 2015). Authors Bazarghan & Gupta (2008), Almeida & Prieto (2013), and Yi & Pan (2010) do not use feature selection, rather they use the full range of wavelengths. This causes the input space for the classifier models to be very large in dimension, which makes the algorithms slow. Bazarghan & Gupta (2008) rebinned the SDSS spectra to have the same resolution as the Jacoby (Jacoby et al. 1984) spectra. One could argue that this is feature selection because they are converting and reducing the number of measurements, but a more accurate description would be that they are simply spectra fitting and not significantly reducing the number of dimensions. Almeida & Prieto (2013) use K means clustering to classify SDSS stellar spectra and do not reduce the number of dimensions. Yi & Pan (2010) utilized RF to classify stellar spectra. The authors also compared RF to neural networks (multilayer perceptron, MLP).

³ An imputer is used to fill in missing values in a feature matrix. Brice & Andonie (2019) utilize a moving average imputer for missing values in stellar spectra.

Table 1

Example of the Feature Matrix Using Two Sets of Wavelengths around Two Absorption Lines

	Wavelength: 4095.43 Å	Wavelength: 4110.55 Å	Wavelength: 4219.88 Å	Wavelength: 4235.45 Å
Spectrum 1	Flux	Flux	Flux	Flux
Spectrum 2	Flux	Flux	Flux	Flux

Note. The 34 features result in a 34 column feature matrix.

Authors Xing & Guo (2004), Bazarghan (2008), and Bailer-Jones et al. (1998) use PCA to reduce the number of flux measurements, but they maintain the shape and structure of the overall spectrum. Xing & Guo (2004) also use a wavelet transformation to reduce noisy flux measurements, but again they maintain the shape and structure of the spectrum. Others such as Zhang et al. (2008) normalized the continuum of the spectra. It is important to note that this does not reduce the dimensions of the spectra. Bai et al. (2019) uses a color space rather than spectra to classify stars. The authors use nine color bands (i.e., $g-r$, $r-i$, etc.) as their features. Elting et al. (2008) also use photometric data instead of spectra to classify stars. Schierscher & Paunzen (2011) actually reduce the spectra using a similar approach to the work in this paper by using absorption lines, but they reduce the dimensions from 2400 to 435, where this paper reduces down to 34, as explained in Section 3. However, Schierscher & Paunzen (2011) classify effective temperature ranges and not direct spectral and luminosity classes.

From an astronomical point of view, spectra contain two features: flux and wavelength. From a machine-learning point of view, spectra contain N features, where N is the number of flux measurements. In earlier work conducted in Brice & Andonie (2019), standard machine-learning feature selection methods are used which use a statistical approach to rank correlation between flux measurements and spectral classes. Then the K most correlated features (flux measurements) are taken as the input space to the classifier model. This approach does not maintain the shape or structure of the spectra. The work presented in this paper does not use machine-learning feature selection, rather astronomical knowledge of the spectra to reduce the number of dimensions in the input space.

The input space to classifier models is known as a feature matrix. Each column of the feature matrix is a unique feature/attribute/dimension of the object to be classified and the rows are the individual samples. For spectra, these features are the flux measurements, as seen in Table 1. It is important to note that the wavelength values are used as the title of the unique features, not as the features themselves.

The stellar spectra found in the SDSS data set contain on average 4617 flux measurements or in another terms the input space has 4617 dimensions. As stated above, feature selection is used to reduce these dimensions. Standard machine-learning feature selection methods used in Brice & Andonie (2019) did not work for these experiments because the spectral classes overshadow the luminosity classes. This means that statistical correlation between specific flux measurements and spectral class is stronger than the same flux measurement and luminosity class. This is apparent because the luminosity classes are based on the width of the absorption lines, which makes it difficult for individual flux measurements to be correlated to both luminosity and spectral classes. Therefore,

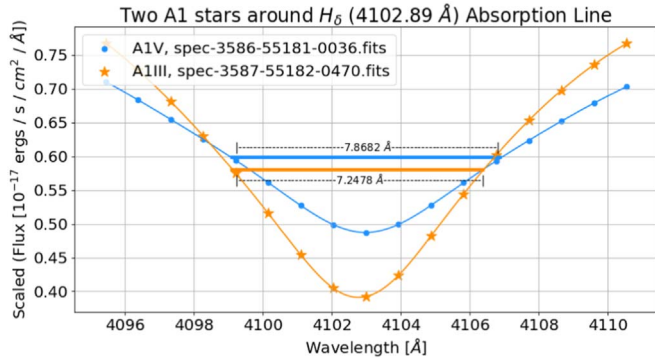


Figure 1. Example of the same spectral class with different wavelength width (FWHM) for the same absorption line for different MK classes.

flux measurements around an absorption line are used rather than the ones that are statistically correlated.

Since there is not one absorption line that all spectral classes share, two absorption line regions are used. The two absorption lines with rest vacuum wavelengths are H_δ (4102 Å) and Ca I (4227 Å). B and A stars have H_δ , K and M have Ca I, and F and G have both absorption lines. Spectral classes are separated because of the intensity of the flux in the regions and luminosity classes are separated because of the widths of the absorption line. Figure 1 shows how the width of the absorption line changes the flux measurements for two A-type stars. Using these two regions the feature matrix can be built. The flux measurements from the H_δ region and the flux measurements from the Ca I region are combined to create a single flux array per sample, as seen in Table 1. This feature selection must also be able to incorporate the shifting of the spectrum due to RV. Figure 2 shows that with a sufficiently sized region, RV can be incorporated.

The feature matrix can be represented as an N -dimensional hypercube. Where each feature is a dimension in this hypercube. As mentioned before, the wavelengths are used as dimension labels, where the flux measurement at that wavelength is the magnitude of the vector in that dimension. This results in each reduced spectrum being represented as an N -dimensional point. As the shape and intensity of the flux in these regions change due to spectral class, the positions of these reduced spectra change in this N -dimensional hypercube. The same happens when the width of the absorption line changes with luminosity class.

A problem arises with the feature matrix because RV causes the flux measurements to be shifted in wavelength. This is illustrated in Table 2. As mentioned earlier, each feature is unique. The problem that RV causes is that it breaks this feature uniqueness. To overcome this problem, a sufficient number of samples of the exact same class with different RV that span the range of realistic RV are required. This gives the training data set for the classifier model sufficient samples of spectral and luminosity classes at different RV. In terms of KNN, this allows for sufficient neighbors to be nearby when classifying. In terms of RF, this forces the splitting threshold of each decision tree's node to include spectra with RV.

3. Experiments

In this section, the experimental setup is described and the results are presented.

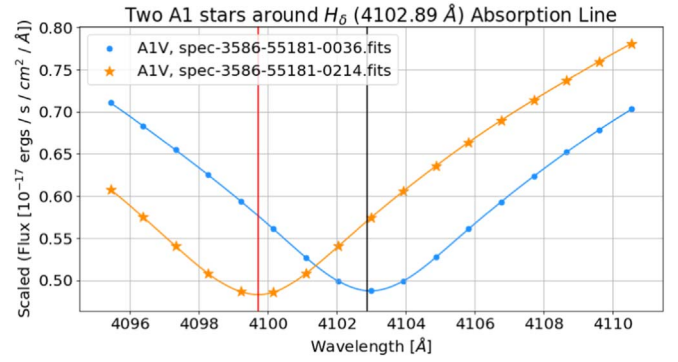


Figure 2. Example of how RV is accounted for in the feature selection window.

3.1. Experimental Setup

The data set used in these experiments comes from SDSS data run 14, which was collected using the BOSS spectrograph⁴ (Smee et al. 2013). Data run 14 contains a total of 335,844 spectra. Some of the data was rejected because it was not a spectral class of O, B, A, F, G, K, and M with a subclass of 0–9 combined with a luminosity class of I, II, III, IV, V, VII, or that the data was missing a large portion of its spectrum, similarly to the work done in Brice & Andonie (2019). The data was also preprocessed by SDSS scientists through the methods presented by Dawson et al. (2013) and Stoughton et al. (2002).

The usable data set contains 168,982 stellar spectra and 46 of the 420 class combinations. It is important to note that this is real collected data and not simulated data. The spectra are first preprocessed by scaling the flux to ensure that similar classes have similar flux measurements using Equation (1).

Then feature selection is performed using Algorithm 1, where bounds are the number of flux measurements before and after the absorption line. Bounds are set to 8 to cover the RV range of -552 to 552 km s⁻¹ and ensure that sufficient flux measurements are recorded. This results in a region of 17 flux measurements around H_δ (4102 Å) and a region of 17 flux measurements around the Ca I (4227 Å) absorption line, which is combined to form a 34-dimension feature matrix.

After the feature selection phase, the data set is split into two subsets for RV: one for RV less than 200 km s⁻¹ and one for RV greater than 200 km s⁻¹. Each subset is again divided into 10 subsets for tenfold cross validation (Kohavi 1995). One subset from each RV set is taken as testing sets respectively. The remaining subsets are combined into one training set. Due to the data set being imbalanced (Figure 3), the training set is balanced using an undersampling method (Japkowicz 2000; randomly removing samples), an oversampling method (randomly duplicating samples using SMOTE; Bowyer et al. 2011), and a hybrid method (undersampling + oversampling). Then, KNN and RF classifiers are applied to the training set and tested using the RV less than 200 km s⁻¹ test set and RV greater than 200 km s⁻¹ test set as well as a combined test set. These steps are repeated 10 times with different subsets used for testing for tenfold cross validation. The accuracy, precision, recall, and F1 score (defined below) are averaged over each cross validation. Misclassification costs for incorrectly classifying a spectrum are not explored.

⁴ Simplified BOSS spectrograph explanation: https://www.sdss.org/instruments/boss_spectrograph/.

Table 2
Example of the Feature Matrix with Observed A0 Spectra.

Star Class	Wavelength: 4101.00 Å	Wavelength: 4101.94 Å	Wavelength: 4102.89 Å	Wavelength: 4103.83 Å	Wavelength: 4104.77 Å
A0	...	H _δ
A0	H _δ
A0	H _δ
A0	H _δ	...
A0	H _δ

Note. Note for the A0 stars, RV causes the H_δ absorption line to be modeled with different wavelength features.

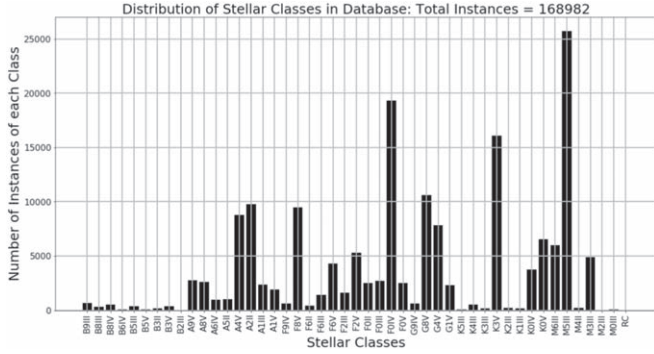


Figure 3. Distribution of classes in the data set. **RC** = remaining classes of 0 instances.

The experiments are implemented in Python, using scikit-learn (Pedregosa et al. 2011). Due to the size of the data set (16.7 GB, which is larger than the RAM used in these experiments), the Python NumPy memmap⁵ (Oliphant 2015) module was used to read very large arrays from storage rather than RAM. The experiments are performed on a personal computer with the following relevant specifications: AMD Ryzen 7 1800x 16 logical core CPU, 16 GB RAM, and 1 TB Samsung 860 EVO Solid State Drive.

KNN and RF use the scikit-learn default parameters (Pedregosa et al. 2011). Precision,⁶ recall,⁷ and F1 score⁸ are computed using functions implemented by the scikit-learn sklearn.metrics package (Pedregosa et al. 2011). Feature selection is implemented using Algorithm 1 in Python and uses the Python multiprocessing package⁹ for parallelization.

Precision is the measure of how well a predicted class compares to actual classes. (Marsland 2015). For example, in the binary class case, precision is the number of samples correctly predicted as class 1 divided by the total number of samples predicted as class 1 (Equation (2)). Recall is the measure of how well an actual class compares to predicted classes (Marsland 2015). For example, in the binary class case, recall is the number of samples correctly predicted as class 1 divided by the total number of samples that are actually class 1 (Equation (3)). F1 score is a type of harmonic mean that

combines precision and recall into a single metric (Equation (4); Marsland 2015). The closer precision, recall, and F1 score are to 1, the more accurate the model. In terms of this work, precision, recall, and F1 score are reported as an average of all the individual precision, recall, and F1 score measurements associated with each of the 46 classes respectively. Precision, recall, and F1 score are defined in the binary class case as follows:

$$\text{Precision} = \frac{\text{Total_TP}}{\text{Total_TP} + \text{Total_FP}} \quad (2)$$

$$\text{Recall} = \frac{\text{Total_TP}}{\text{Total_TP} + \text{Total_FN}} \quad (3)$$

$$F_1 = 2 \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}, \quad (4)$$

where the true positive (TP) is defined as a prediction of class 1 being correctly classified as class 1, false positive (FP) is defined as a prediction of class 1 being incorrectly classified as class 2, and false negative (FN) is defined as a sample of class 1 being incorrectly predicted as class 2. In the multiclass problem, FP is defined as a prediction of class i being incorrectly classified as class $\neq i$ and FN is defined as a sample of class i being incorrectly predicted as class $\neq i$. More information regarding precision, recall, and F1 score can be found in Powers & Ailab (2011).

4. Discussion

Tables 3 and 6 show that classification using KNN has essentially the same accuracy as RF when using hybrid and oversampling balancing. These tables demonstrate that using KNN and RF along side Algorithm 1 for feature selection are viable options for the automated classification of stellar spectra because of the high accuracy achieved. Table 3 demonstrates that using three neighbors for KNN classification performs the best for KNN. Table 6 shows that changing the number of trees used in RF does not significantly change the classification accuracy. Tables 3 and 6 show that oversampling balancing outperforms hybrid balancing.

Tables 4 and 7 show that the precision and recall is less than accuracy (Tables 3 and 6). The precision implies on average that the RF model's prediction of class i is actually 84.91% (Table 7) class i. Where recall implies on average, the RF model predicts 86.63% (Table 7) of the samples that are class i as class i. It is important to note that the test sets for all RV, $RV < 200 \text{ km s}^{-1}$, and $RV \geq 200 \text{ km s}^{-1}$ metrics are tested using the same trained model. As seen in these tables, this approach accurately classifies both $RV < 200 \text{ km s}^{-1}$ and $RV \geq 200 \text{ km s}^{-1}$. However, $RV \geq 200 \text{ km s}^{-1}$ has a higher

⁵ Memmap: <https://docs.scipy.org/doc/numpy-1.14.0/reference/generated/numpy.memmap.html>.

⁶ Precision scikit-learn: https://scikit-learn.org/stable/auto_examples/model_selection/plot_precision_recall.html.

⁷ Recall scikit-learn: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.recall_score.html#sklearn.metrics.recall_score.

⁸ F1 Score scikit-learn: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html#sklearn.metrics.f1_score.

⁹ Multiprocessing: <https://docs.python.org/3.7/library/multiprocessing.html#module-multiprocessing>.

Table 3
Tenfold Cross-validation Accuracy for KNN

Balance Method	RV	Accuracy (%) for K Neighbors					
		3.0	5.0	7.0	10.0	15.0	20.0
Undersampled	All	52.66	52.31	51.47	50.73	49.63	48.49
	$<200 \text{ km s}^{-1}$	53.32	52.99	52.17	51.45	50.35	49.23
	$\geq 200 \text{ km s}^{-1}$	23.97	22.94	21.18	19.38	17.98	16.67
Hybrid	All	93.02	92.63	92.31	91.85	91.26	90.73
	$<200 \text{ km s}^{-1}$	92.92	92.56	92.26	91.82	91.28	90.77
	$\geq 200 \text{ km s}^{-1}$	97.15	95.62	94.44	92.94	90.54	88.96
Oversampled	All	95.48	95.00	94.54	93.98	93.24	92.58
	$<200 \text{ km s}^{-1}$	95.46	94.99	94.56	94.02	93.32	92.69
	$\geq 200 \text{ km s}^{-1}$	96.41	95.28	93.91	92.18	89.54	87.78

Table 4
Tenfold Cross-validation Precision, Recall, and F1 Score for KNN Using Oversampling

	RV	K Neighbors					
		3.0	5.0	7.0	10.0	15.0	20.0
Precision	All	0.805493	0.795833	0.789050	0.781626	0.767124	0.756264
	$<200 \text{ km s}^{-1}$	0.803331	0.793968	0.787575	0.780606	0.766780	0.756458
	$\geq 200 \text{ km s}^{-1}$	0.926071	0.897225	0.872481	0.838610	0.793230	0.756975
Recall	All	0.835391	0.834499	0.834981	0.835475	0.833052	0.829150
	$<200 \text{ km s}^{-1}$	0.832799	0.831953	0.832601	0.833312	0.831288	0.829150
	$\geq 200 \text{ km s}^{-1}$	0.941410	0.921932	0.905449	0.879239	0.853662	0.959524
F1 Score	All	0.811617	0.805477	0.801665	0.796227	0.786297	0.776893
	$<200 \text{ km s}^{-1}$	0.809250	0.803309	0.799771	0.794651	0.785290	0.776297
	$\geq 200 \text{ km s}^{-1}$	0.989692	0.985351	0.981322	0.975739	0.967290	0.827752

Table 5
Tenfold Cross-validation Execution Times for KNN Using Oversampling for all RV

	Time in Seconds for K Neighbors					
	3.0	5.0	7.0	10.0	15.0	20.0
Feature Selection	90.89	90.89	90.89	90.89	90.89	90.89
Train	6.91	6.91	6.95	6.94	6.88	6.82
Test	3.03	3.22	3.41	3.62	3.92	4.19

accuracy, precision, and recall (Table 7) than RV $< 200 \text{ km s}^{-1}$. This could be the result of an imbalance between low RV and high RV samples and overfitting the model in the

low RV range. Future work will attempt to increase precision and recall for low RV, which includes balancing the RV range.

Tables 5 and 8 show the execution times for each experiment. These tables show that KNN performs much faster than RF. KNN has a faster train time than RF, but RF has a faster test time than KNN. For both KNN and RF, feature selection takes approximately the same amount of time, which is expected since they both use the same feature selection.

This approach takes considerably fewer steps than the one in Bolton et al. (2012) and produces excellent results. The execution times and the obtained accuracy demonstrate that, for a real application of this work, the automated classification of observed stellar spectra into a complete MK classification using a single classifier not only achieves a high accuracy but is also fast.

Table 6
Tenfold Cross-validation Accuracy for RF

Balance Method	RV	Accuracy (%) for N Trees					
		10.0	50.0	100.0	150.0	200.0	250.0
Undersampled	All	52.45	58.13	58.75	59.01	59.09	59.21
	$<200 \text{ km s}^{-1}$	53.09	58.78	59.39	59.65	59.72	59.83
	$\geq 200 \text{ km s}^{-1}$	24.62	29.79	30.85	31.48	31.77	31.91
Hybrid	All	93.83	94.44	94.48	94.48	94.50	94.49
	$<200 \text{ km s}^{-1}$	93.73	94.35	94.38	94.39	94.41	94.41
	$\geq 200 \text{ km s}^{-1}$	98.39	98.36	98.36	98.44	98.61	98.41
Oversampled	All	96.39	97.05	97.12	97.16	97.15	97.16
	$<200 \text{ km s}^{-1}$	96.34	97.01	97.09	97.13	97.12	97.12
	$\geq 200 \text{ km s}^{-1}$	98.49	98.65	98.62	98.62	98.65	98.65

Table 7
Tenfold Cross-validation Precision, Recall, and F1 Score for RF **Using** Oversampling

RV		N Trees					
		10.0	50.0	100.0	150.0	200.0	250.0
Precision	All	0.827606	0.845460	0.849466	0.849149	0.850188	0.850051
	$<200 \text{ km s}^{-1}$	0.825895	0.844077	0.848296	0.847941	0.848908	0.848765
	$\geq 200 \text{ km s}^{-1}$	0.968922	0.971134	0.967166	0.968655	0.971852	0.971852
Recall	All	0.849646	0.862420	0.866250	0.866379	0.868213	0.868580
	$<200 \text{ km s}^{-1}$	0.847837	0.861001	0.865014	0.865182	0.866917	0.867305
	$\geq 200 \text{ km s}^{-1}$	0.973081	0.973082	0.970220	0.972075	0.972433	0.972433
F1 Score	All	0.832101	0.847679	0.851134	0.851346	0.852966	0.853006
	$<200 \text{ km s}^{-1}$	0.830256	0.846204	0.849857	0.850088	0.851616	0.851663
	$\geq 200 \text{ km s}^{-1}$	0.969936	0.970727	0.967720	0.969363	0.970591	0.970591

Table 8
Tenfold Cross-validation Execution Times for RF **Using** Oversampling for **All** RV

		Time in Seconds for N Trees					
		10.0	50.0	100.0	150.0	200.0	250.0
Feature Selection		90.89	90.89	90.89	90.89	90.89	90.89
Train		129.23	605.32	1226.75	1857.21	2448.32	2960.94
Test		0.65	2.91	6.03	9.16	11.99	14.46

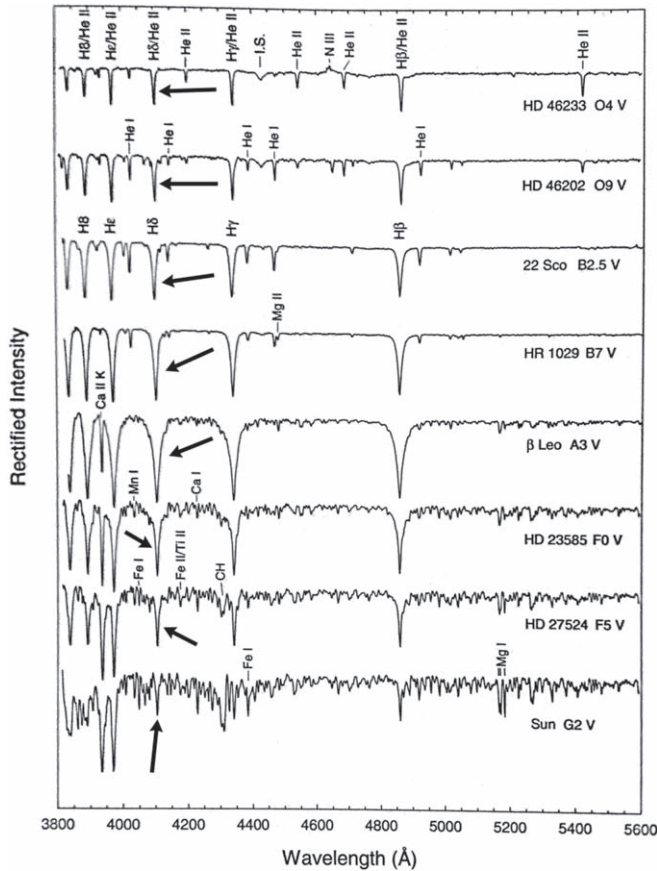


Figure 4. Sample of continuum normalized spectra from O–G type stars (Gray & Corbally 2009). The arrows point to the H_δ absorption line.

These experiments yielded an accuracy of 97.16% for RF (Table 6) using **oversampling balancing** and 150 decision trees (Table 6), which is significantly better than Xing & Guo (2004), **who report** 81.66% accuracy for just **the support vector**

machine (SVM) with no data reduction, 93.26% for **wavelet + SVM**, and 81.30% for PCA + SVM. Schierscher & Paunzen (2011) are able to achieve an 85% match in comparison with **the Sloan Extension for Galactic Understanding and Exploration (SEGUE)** Stellar Parameter Pipeline using an **artificial neural network**, but this is not comparable since they did not classify into direct spectral and luminosity classes. Schierscher & Paunzen (2011) also never addresses the fact that their data is imbalanced.

Other authors such as Bazarghan & Gupta (2008), Yi & Pan (2010), and Zhang et al. (2008) do not report metrics that are comparable to the work conducted in this paper. Bazarghan & Gupta (2008) used a **probabilistic neural network** implemented in MATLAB and they used a χ^2 value to determine classification accuracy. They make the assumption that a χ^2 value of 0.002 or lower is considered classified correctly, then they achieved a success rate of about 88% in only a few seconds. Yi & Pan (2010) find that RF performed better than the MLP with **an rms error** (RMSE) of 1.04 and **1.36**, respectively. Zhang et al. (2008) separated the classification into two classifiers. For the spectral classes, the authors used a **nonparameter** regression method. For the luminosity classes, the authors removed or normalized the continuum of the spectra and used a partial least-squared regression method. Zhang et al. (2008) used three spectra data **sources**: Silva & Cornell (1992), Pickles (1998), and Jacoby et al. (1984). They achieved a standard deviation of $\sigma = 0.7994$ for the **spectral** classes and $\sigma = 0.58159$ for the **luminosity** classes.

As described above in Section 3, these experiments deal with data collected by a real astronomical survey. As such, when an astronomical survey points their telescopes into the sky, they get the samples (classes) that they get. The experiments presented here deal with a subset of all possible class combinations. It is important to note that not all possible class combinations (O, B, A, F, G, K, and M with **subclasses** of 0–9 combined with I, II, III, IV, V, VII) are common or even found in nature. Therefore, even though this approach yielded great

Q4

results, there cannot be a claim that this approach will guarantee work for all stellar classes. There is, however, some theoretical validity to this approach.

Referencing Figure 3, **O-type** stars are the only spectral class not found in the data set. Figure 4 shows that **O-type** stars also contain the H_δ absorption line. Therefore, the missing spectral and luminosity classes are compatible with this approach because every spectral major class has at least one absorption line in the feature selection regions, which allows for the variation in width for luminosity classes.

5. Conclusion

The results shown in this paper support that accurate automatic stellar classification can be obtained using **astrophysical-specific** feature selection. Compared to previous work of other authors, there are four interesting **conclusions**. They are as follows.

1. A high level of accuracy can be obtained by considering only flux measurements at wavelengths near the H_δ and Ca I absorption lines.
2. A complete MK classification can be identified using a single classifier with a high level of accuracy for stars with small RV (the range of $\approx \pm 240 \text{ km s}^{-1}$). However, for stars with large RV (outside the range of $\approx \pm 240 \text{ km s}^{-1}$), increasing the value of bounds in Algorithm 1 should compensate for large RV, but at the cost of an increase to the computational time because of the increase in dimensions.
3. Correcting for RV is not necessary because of a sufficient distribution of samples with different RV.
4. Aside from flux scaling, any additional spectrum **preprocessing** after recombining and **rebinning** the spectra as presented by Dawson et al. (2013) and Stoughton et al. (2002) is unnecessary for SDSS stellar spectra.

Therefore, this new approach for the automatic classification of stellar spectra is feasible, useful, and accurate. Future work will be conducted with the concepts of this approach to automatically classify stars and large redshift objects such as galaxies and quasars. Future work will also include building the training set with an equal distribution of RV and experiment using SMOTE or similar algorithms to build a training set using a small number of real samples to simulate how this approach could work for a new spectroscopic survey.

ORCID iDs

Michael J. Brice  <https://orcid.org/0000-0002-1278-2404>

References

- Almeida, J. S., & Prieto, C. A. 2013, **ApJ**, **763**, 50
- Bai, Y., Liu, J., Wang, S., & Yang, F. 2019, **AJ**, **157**, 9
- Bailer-Jones, C. A. L., Irwin, M., & von Hippel, T. 1998, **MNRAS**, **298**, 361
- Bazarghan, M. 2008, arXiv:0804.2742
- Bazarghan, M., & Gupta, R. 2008, **Ap&SS**, **315**, 201
- Bolón-Canedo, V., Sánchez-Marroño, N., & Alonso-Betanzos, A. 2015, Feature Selection for High-Dimensional Data (Springer International Publishing)
- Bolton, A. S., et al. 2012, **AJ**, **144**, 1
- Bowyer, K. W., Chawla, N. V., Hall, L. O., & Kegelmeyer, W. P. 2011, CoRR, abs/1106.1813, arXiv:1106.1813
- Breiman, L. 2001, **Machine Learning**, **45**, 5
- Brice, M., & Andonie, R. 2019, in Proc. Int. Joint Conf. Neural Networks (IEEE)
- Dawson, K. S., et al. 2013, **AJ**, **145**, 1
- Duan, F.-Q., Liu, R., Guo, P., Zhou, M.-Q., & Wu, F.-C. 2009, **RAA**, **9**, 341
- Elting, C., Bailer-Jones, C. A. L., & Smith, K. W. 2008, in AIP Conf. Proc. **1082**, 9
- Goldberger, J., Hinton, G. E., Roweis, S. T., & Salakhutdinov, R. R. 2005, in Advances in Neural Information Processing Systems 17, ed. L. K. Saul, Y. Weiss, & L. Bottou (Cambridge, MA: MIT Press), 513
- Gray, R. O., & Corbally, J. C. 2009, Stellar Spectral Classification (Princeton, NJ: Princeton Univ. Press)
- Ivezic, Z., Connolly, A., VanderPlas, J., & Gray, A. 2014, Statistics, Data Mining, and Machine Learning in Astronomy: A Practice Python Guide for the Analysis of Survey Data (Princeton, NJ: Princeton Univ. Press)
- Jacoby, G. H., Hunter, D. A., & Christian, C. A. 1984, **ApJS**, **56**, 257
- Japkowicz, N. 2000, Papers from the AAAI Workshop **Tech. Rep.** WS-00-05 (AAAI Press), 10
- Kohavi, R. 1995, in Proc. 14th Int. Joint Conf. Artificial Intelligence, Vol. 2, IJCAI'95 (San Francisco, CA, USA: Morgan Kaufmann Publishers Inc), 1137
- Marsland, S. 2015, Machine Learning: An Algorithmic Perspective (2nd ed.; Boca Raton, FL: CRC Press)
- Morgan, W. W., Keenan, P. C., & Kellman, E. 1943, An Atlas of Stellar Spectra, with an Outline of Spectral Classification (Chicago, IL: Univ. Chicago Press)
- Oliphant, T. E. 2015, A Guide to NumPy (Continuum Press)
- Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, Journal of Machine Learning Research, **12**, 2825
- Pickles, A. J. 1998, **PASP**, **110**, 863
- Powers, D., & Ailab 2011, J. Mach. Learn. Technol, **2**, 2229
- Schierscher, F., & Paunzen, E. 2011, **AN**, **332**, 597
- Shlens, J. 2014, arXiv:1404.1100
- Silva, D. R., & Cornell, M. E. 1992, **ApJS**, **81**, 865
- Smee, S. A., et al. 2013, **AJ**, **146**, 1
- Stoughton, C., et al. 2002, **AJ**, **123**, 485
- Xing, F., & Guo, P. 2004, in Advances in Neural Networks), ed. F.-L. Yin, J. Wang, & C. Guo (Berlin: Springer), 616
- Yi, Z., & Pan, J. 2010, IEEE, 3129
- York, D. G., et al. 2000, **AJ**, **120**, 1579
- Zhang, J. N., Luo, A. L., & Tu, L. P. 2008, in Inter. Congress on Image and Signal Processing, 249, CISP 2008

Q5

Q6

Q7

Q8

Q9

Q10