

AAS-PROVIDED PDF • OPEN ACCESS

## The Application of Machine Learning to Quasar and Seyfert Classification

To cite this article: Vivek Abraham *et al* 2024 *Res. Notes AAS* **8** 46

Manuscript version: AAS-Provided PDF

This AAS-Provided PDF is © 2024 The Author(s). Published by the American Astronomical Society.



Original content from this work may be used under the terms of the Creative Commons Attribution 4.0 licence. Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

Everyone is permitted to use all or part of the original content in this article, provided that they adhere to all the terms of the licence  
<https://creativecommons.org/licenses/by/4.0>

Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions may be required.

View the [article online](#) for updates and enhancements.

1  
2  
3 DRAFT VERSION FEBRUARY 2, 2024  
Typeset using L<sup>A</sup>T<sub>E</sub>X default style in AASTeX631

## The Application of Machine Learning to Quasar and Seyfert Classification

4 VIVEK ABRAHAM,<sup>1</sup> JOEL DEVILLE,<sup>1</sup> AND GARV KINARIWALA<sup>1</sup>

5 <sup>1</sup>*University of Texas at Austin*

### 6 ABSTRACT

7 Machine learning can be utilized to classify spectra flagged as Active Galactic Nuclei (AGN) belonging  
8 to Seyferts or Quasars, expediting data collection and aiding in analyzing the AGN types. While  
9 many properties of Seyferts and Quasars can be used as feature points in training a machine learning  
10 model, one relatively available property with high information density is the spectra of the AGN types.  
11 This paper aims to describe the training and results of a K-Nearest Neighbors (KNN) and a Dense  
12 Neural Network (DNN) machine learning model built to classify AGNs as Seyfert type 1s, Seyfert type  
13 2s, or Quasars.

### 14 1. INTRODUCTION

15 Since their discovery, two types of Active Galactic Nuclei (AGN), Seyferts and Quasars, have gained significant  
16 traction in the scientific community. Seyferts appear as very luminous celestial objects, with their total radiation  
17 output rivaling the amount emitted by the entirety of the constituent stars in their host galaxies. Seyfert objects  
18 can be further classified into two major classes based on the relative line widths of their emission spectra: Seyfert 1s  
19 and Seyfert 2s. While both Seyfert types have broad line widths relative to non-AGN spectra, Seyfert 1s have very  
20 broad permitted lines on the order of  $10^4$  kilometers per second, which Seyfert 2s lack. Like Seyferts, Quasars emit  
21 tremendous amounts of radiation; however, while Seyferts emit radiation comparable in intensity to their host galaxy,  
22 quasars outshine them immensely as they put off more light at further distances generally. Both Seyferts and Quasars  
23 have been attributed to be hosts of supermassive black holes in the center of galaxies, making them hot topics of  
24 scientific research. Understanding the processes behind Seyfert and Quasars will help expand our understanding of  
25 galaxy morphology, star formation, and the composition of the early universe. It is useful to consider other studies,  
26 which used similar machine learning methods, such as Ma et al. (2019) and Cavaudi et al. (2013). By looking at  
27 these studies and our own, we believe machine learning is crucial to accelerate the classification of AGNs based on  
28 the spectral values measured. We seek to fit a model that is able to accurately classify Seyferts and Quasars across  
29 varying redshifts, which proves difficult normally. This model will focus on doing spectral classifications of Quasars  
30 versus Seyferts.

### 31 2. DATA COLLECTION

32 Our study primarily utilized data from the Sloan Digital Sky Survey (SDSS) Data Release 18 (Kollmeier et al. 2019)  
33 and the Veron Catalog of Quasars & AGN (VeronCAT), 13th edition (Véron-Cetty & Véron 2010). We specifically  
34 retrieved data using SQL queries over data from SDSS-V and the Science Archive Server for FITS files. We were able  
35 to interface with the SDSS data through the Catalog Archive Server (CAS). CAS provided a vast array of objects  
36 with detailed spectra, notably from its SpecObj subset which includes about 5.1 million objects. However, as SDSS's  
37 classification of AGNs is limited, we incorporated the VeronCAT Catalog for more nuanced AGN classifications. We  
38 utilized TopCAT to cross-match a batch SQL query from SDSS with a TAP query from VeronCAT using a maximum

akka.viv@gmail.com

joeldeville114@gmail.com

garvkinariwala@gmail.com

39 error margin of 8 arc-seconds. After cross-matching, we had 61,270 data points, out of which approximately 35,000  
 40 spectra FITS files were available from the SDSS Science Archive Server.

41 We then processed this data with a Python program; This program extracted spectral information and downloaded  
 42 spectra files. The Python program created a CSV file from this data collected, which we later used with Pandas to  
 43 serve as the input to our model. The final analysis, involving trend identification and visualization, was performed  
 44 using Python libraries such as Pandas and Matplotlib. This approach enabled us to efficiently categorize and analyze  
 45 the data, particularly focusing on AGN sub-classifications and their properties.

### 46 3. MODEL PARAMETERS

47 For the dense neural network, we used Tensorflow as the machine learning library. The training parameters were  
 48 approximately 3900 measured flux values for different wavelengths, extracted from the FITS file, as well as the estimated  
 49 redshift for the object. The training target was the category of the object (Quasar, Seyfert 1, or Seyfert 2). We used  
 50 one-hot encoding for our three classification classes. This converted our categorical data of classifications into numerical  
 51 values to be used as input and output for the dense neural network. For the hidden layers, we decided to use the  
 52 ReLu activation function, which was found to be the best general activation function for dense neural networks in  
 53 most cases (Bai 2022). The softmax activation function was utilized for the output layer for an easy conversion into  
 54 labels. We decided to test two different hidden layer patterns: a 256-32 layer pattern and a 1024-256-64-16 layer  
 55 pattern. For both patterns, we trained for 200 maximum epochs with early stopping enabled with a patience of 40.  
 56 Finally, we tested our model both with L2 regularization implemented and without regularization, creating a total of  
 57 four model architectures for the neural network; To reduce variation in our results, we tested each architecture three  
 58 times with different random states and averaged the results together. For the KNN model, we tested four different  
 59 nearest-neighbor sizes: 5 neighbors, 10 neighbors, 20 neighbors, and 50 neighbors. For each size, we tested a uniform  
 60 and an inverse-to-distance weight function for prediction. Similar to the neural network test, we ran each subtest three  
 61 times with different random states. To calculate accuracies, we used the Scikit-learn accuracy score metric.

KNN Results				
	5 Neighbors	10 Neighbors	20 Neighbors	50 Neighbors
Distance	91.39%	91.76%	91.41%	89.88%
Uniform	91.76%	91.40%	91.00%	90.12%
Neural Network Results				
	Small		Large	
No Regularization	93.16%		93.98%	
Regularization	92.18%		92.36%	
Class Accuracy Breakdown				
	Quasars	Seyfert 1s	Seyfert 2s	
10 Neighbor Distance KNN	98.89%	56.85%	79.06%	
Non-regularized Large DNN	98.73%	71.97%	76.78%	

62 **Table 1.** Comparison of Model Accuracies

### 63 4. RESULTS

64 The majority of the differences in predictions were from the KNN mispredicting true Seyfert 1s as Quasars, which  
 65 the neural network was better at predicting. Another larger source of differences was with one model predicting an  
 66 AGN as a Seyfert 1 and the other model predicting an AGN as a Seyfert 2. Both models generally agreed on Quasar  
 67 predictions; moreover, an AGN classified as a Quasar by one model was seldom classified as a Seyfert 2 by the other  
 model.

68 The kNN model tended to perform better with correctly predicting Seyfert 1s with lower redshifts relative to  
 69 predicting Seyfert 1s with higher redshifts. This is likely due to Quasar objects tending to have a higher redshift

70 value than Seyferts, allowing the model to use the redshift values to distinguish Quasars from Seyferts; moreover, high  
71 redshifts compress the spectra of objects, further aiding with differentiation. We hypothesize that the neural network  
72 performed better with predicting even high-redshift Seyfert 1s due to the neural network suffering less from the very  
73 high dimensionality of our training set (University 2018). The KNN model performed significantly worse than our  
74 dense neural network when dealing with Seyfert 1s. However, both the KNN model and the neural network performed  
75 worse on intermediate Seyfert subclassifications than their respective Quasar classifications; Subclassifications which  
76 gave our models issues were Seyfert 1.2s (which we grouped as a Seyfert 1), Seyfert 1.8s (which we grouped as a Seyfert  
77 2), and especially Seyfert 1.5s (which we grouped as a Seyfert 1). Attempting to group Seyfert 1.5s as Seyfert 2s instead  
78 did not seem to improve the accuracy of our models, with the large, non-regularized neural network performing at  
79 93.68% accuracy. A further breakdown on the best-performing models and information about other models can be  
80 found in Table 1.

81 We hypothesize that the greater degree of mispredictions for these Seyfert subclasses is due to the spectra of  
82 intermediate Seyferts being more ambiguous, containing features that are not distinctly similar to those of either  
83 true Seyfert 1s or Seyfert 2s. Our results corroborate with the unification scheme for Active Galactic Nuclei, which  
84 suggests that all AGN are formed through the same cosmic phenomenon (Spinoglio & Fernandez-Ontiveros 2019). The  
85 unification scheme theory implies that AGN features occupy a continuous range instead of distinct types, which means  
86 that some nuclei can't be easily classified by a machine learning model. It is worthwhile to note this theory is still a  
87 topic under debate, and the current inability to accurately classify intermediate classes of objects does lend support  
88 to the notion of AGN properties being on a spectrum.

89 We want to thank Dr. Mitra, Dr. Gebhardt, our peer mentors, and all members of the Geometry of Space stream in  
90 the Freshman Research Initiative at the University of Texas at Austin for all the help they gave along the way.

## REFERENCES

91 Bai, Y. 2022, RELU-Function and Derived Function  
92 Review, ,  
93 Cavuoti, S., Brescia, M., D'Abrusco, R., Longo, G., &  
94 Paolillo, M. 2013, Photometric Classification of Emission  
95 Line Galaxies with Machine-Learning Methods, ,  
96 Kollmeier, J., Anderson, S. F., Blanc, G. A., et al. 2019,  
97 SDSS-V Pioneering Panoptic Spectroscopy, ,  
98 Ma, Z., Xu, H., Zhu, J., et al. 2019, A Machine Learning  
99 Based Morphological Classification of 14,245 Radio  
100 AGNs Selected from the Best-Heckman Sample, ,  
101 Spinoglio, L., & Fernandez-Ontiveros, J. A. 2019, AGN  
102 Types and Unification Model, ,  
103 University, C. 2018, Lecture 2: k-nearest neighbors / Curse  
104 of Dimensionality, ,  
105 Véron-Cetty, M. P., & Véron, P. 2010, A catalogue of  
106 quasars and active nuclei: 13th edition, , ,  
107 doi:10.1051/0004-6361/201014188