

# STAT 679 Final Project Report

Xiaoyang Wang

Ziang Zeng

## 1 Astronomical Challenge

Our project focuses on classifying celestial objects into stars, galaxies or quasars using their spectral characteristics. With the advancement of astronomical technology, we can obtain a large amount of data, including images and spectral information, from telescopes and large-scale photometry.

The central question of our project is: “How can we effectively use image and spectral data to accurately classify different types of stellar objects?” There are many machine learning methods and statistical models can be applied to classify the celestial objects. However, different methods and models performs differently on same data, “All models are wrong but some are useful.” Can we find a more “useful” model through combining several models together? Our solution is the voting classifier.

## 2 Data

We plan to work with the astronomy data set containing three types of data: images of the celestial objects, images of the spectrum, and the metadata of the objects. The first row in Figure 1 displays images of a galaxy, a star, and a quasar, from left to right, respectively. The second row in Figure 1 displays images of the spectrum of a galaxy, a star, and a quasar, from left to right, respectively. Table 1 provides explanations of the variables within the metadata.

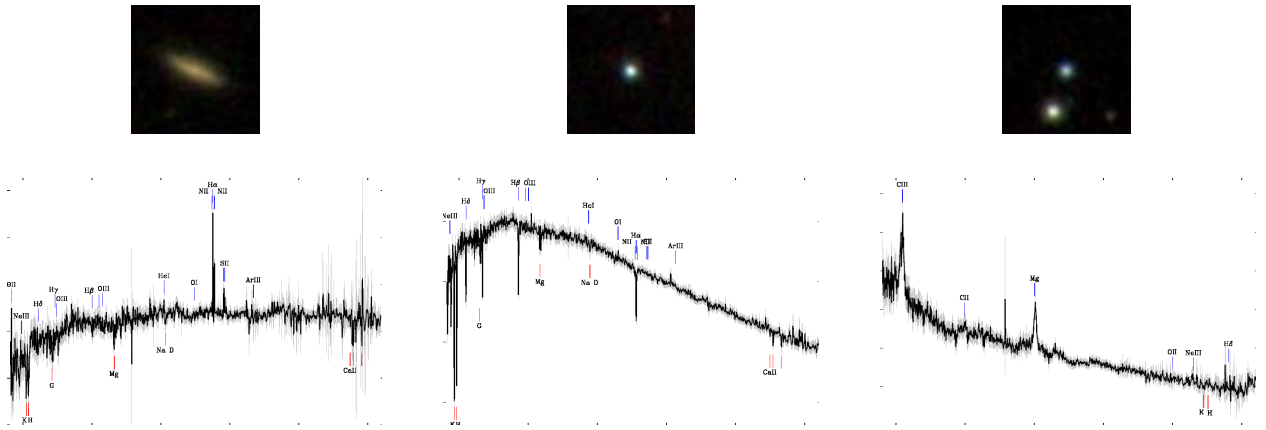


Figure 1: Images of Celestial Objects and Corresponding Spectrum Images.

They can be found in this [link](https://www.sdss.org). All the data is obtained from <https://www.sdss.org>. Moreover, the distribution of classes of the celestial objects is 33333 samples each, which is selected from the

website to make sure the dataset is balanced. Moreover, we encode the class with the following mapping: Galaxy  $\sim 0$ , Quasar  $\sim 1$ , Star  $\sim 2$ .

### 3 Exploratory Data Analysis

First, we will conduct exploratory data analysis on our dataset to better understand it and to find any possible errors. Figure 2 shows the distribution of variables that are meaningful for classification. We can see that the means of all variables across different classes are quite different. Moreover, there are no significant outliers.

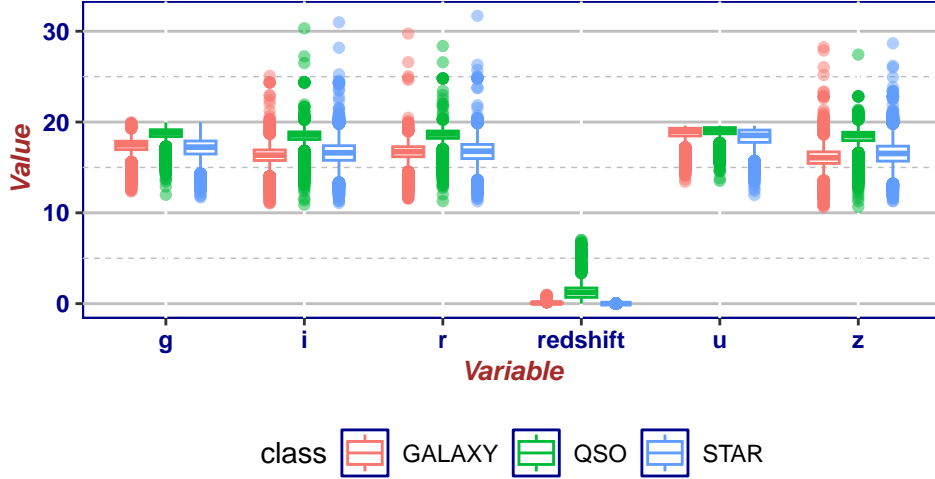


Figure 2: Boxplot

Then, Figure 3 gives the correlationship between variables in the metadata. We can see that **redshift** has strong correlationship with the class of the object. Also, the filter variables (**u**, **g**, **r**, **i**, **z**) are correlated with each other.

Next, we will check and deal with missing values in the dataset. For images of the celestial objects, there is no missing value. For images of the spectrum, we have 14115 images that are unreadable. Considering that they are hard to impute, we just ignore them and conduct the analysis on the spectrum based on the rest of the images. For the metadata, there are 1 missing value for **i** and 3 for **z**. We use regression imputation with filter variables to impute the missing values as they are correlated with each other and quite scattered.

## 4 Methods

For our three types of data, we intend to use three separate methods to build three different classification models. Then, we plan to use a voting classifier to combine three models and give our final model. In this section, we will give brief introduction to the methods that we have used.

### 4.1 kNN

The k-Nearest Neighbors (kNN) algorithm is a simple, but powerful machine learning technique that can be used for classification. At its core, kNN makes predictions about the classification of a

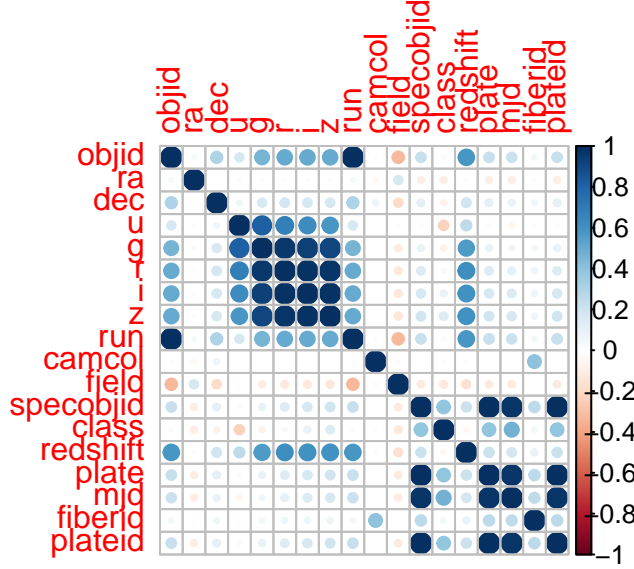


Figure 3: Correlation

data point based on the majority vote or average of its  $k$  nearest neighbors. With cross validation, we eventually selected  $k = 3$ .

## 4.2 Decision Tree

A Decision Tree is a machine learning algorithm that can be used for classification. It models decisions and their possible consequences as a tree-like structure, making it intuitive and easy to visualize. The decision-making process starts at the root node and splits the data on the feature that results in the most significant information gain (IG) or the greatest reduction in impurity (such as Gini impurity or entropy). The process continues recursively, creating decision nodes and leaf nodes. Decision nodes ask a question and branch based on the answers to those questions, leading to further splits or to leaf nodes. These nodes represent the outcome. With cross validation and consideration on the complexity of the tree, we set the maximum depth as 4, and use Gini impurity to prune the tree.

## 4.3 Logistic Regression

Multinomial logistic regression, extends the traditional logistic regression model to handle cases where the target variable categories are more than two. Unlike binary logistic regression, which uses one binary predictor per class, multi-class logistic regression models the probabilities of the multiple classes using a softmax function, which generalizes the logistic function for multi-class problems:

$$P(Y_i = k) = \frac{e^{\beta_k \cdot X_i}}{\sum_{j=1}^3 e^{\beta_j \cdot X_i}}, i = 0, 1, 2.$$

With cross validation, we have selected the hyperparameters as follows: Regularization: L2,  $C = 1$ , where the value of  $C$  gives the strength of regularization.

## 4.4 CNN

Before we get into complex Neural Networks, we firstly try to use a simple CNN to test the performance of NN in this problem. This CNN is quite shallow (337k parameters) with two convolutional layers and a maxpooling in between, followed by three fully connected layers. The structure can be seen in Figure 4. The competitor is VGG16 which is a much deeper NN with 13 convolutional layers and 3 fully connected layers (138 million parameters).

## 4.5 Voting Classifier

We try to use two types of voting methods, soft voting and weighted hard voting. Their strategy are slightly different but the idea is similar: to combine the output of several model together to generate a more accurate one.

Suppose we have models  $\{C_1, \dots, C_n\}$ , for a given input  $x$ . Each model can have a prediction:  $y_{pred}^i|x = (y_1, y_2, y_3)$ , where one of  $y_j$  is 1, indicating the predicted class is  $j$ , while others are 0. Or a predicting probability:  $P_i(y_j|x)$  which is the predicting probability for  $x$  belong to class  $j$  for  $i_{th}$  model.

For kNN and Tree model, their prediction and predicting probability are the same which means their predicting probability is one for predicting class. For soft voting, the prediction is made by average predicting probability of all candidate models and prediction of the voting model is the class with largest predicting probability. The probabilities for voting classifier given input  $x$  is  $P(y_j|x) = \frac{1}{m} \sum_{i=1}^m P_i(y_j|x)$ . So the prediction is  $p(x) = \arg \max_{y_j} P(y_j|x)$  where  $m$  is the number of models.

For weighted hard voting, the prediction is made by sum of weighted voting for each class for all the candidate models and prediction of the voting model is the class with largest number of weighted voting. For a given input,  $C_i$  has a predict  $y^i|x : y_{k=j}^i = 1, y_{k \neq j}^i = 0$  and the prediction for weighted hard voting is  $y_{pred} = \sum_i w_i \cdot y^i|x$ . So the predict class is  $\arg \max_j y_{pred}$ .

# 5 Results

## 5.1 Metadata

For metadata, we have trained three models: kNN, Decision Tree and Logistic Regression.

For the Decision Tree model, the result is shown in Figure 5. For instance, if the scaled **redshift** is less or equal than -0.626, scaled **g** is less or equal than 1.938, then the object should be classified as quasar.

For Logistic Regression model, the result is shown in Table 2. The value of intercepts indicate the log odds of being in the respective class when all the predictor values are zero. Each coefficient for the predictors represents the change in the log odds of being in the respective class for a one-unit change in the predictor variable, holding all other predictors constant.

Figure 6 gives the confusion matrices for all three models. Combining the evaluations shown in Table 5 for those three models, we can see a success in building models through the numerical data. By using several indexes and the red shift, our simple models reach more than 96% accuracy on validation data.

## 5.2 Image of Celestial Objects

The results of this simple CNN is quite good. The cross entropy loss drop quickly after 10k iterations and get stable around 0.27. And the robustness can also be seen on its 90.2% (SGD), 91.8% (Adam) average accuracy after 10 epochs training on validation data. Which actually discourages us to apply deeper NN which is quite time and energy consuming. The training took 20 min and the cross entropy loss is under 0.2 after 3 epochs training and stable at 0.18. On the validation data set, VGG has 94.23% which is better than our simple CNN. However this high accuracy not only consume lot of time but also make it more difficult to improve the performance through the combination of models. With these facts, we plan to modify the simple CNN a little without changing its main structure, and use this simple CNN to build up a final model which is expected to have similar performance as VGG16 while consume less time.

The accuracy of CNN of images 91.73%, CNN of spectrum is 88.91%.

## 5.3 Image of Spectrum

The performance of SimpleCNN on spectrum image classification is also very good. We are able to get 0.021 cross entropy loss which is 99.14% accuracy on validation data set without missing (unreadable) images after 1467 seconds training. And if we include missing value for which the CNN randomly guess a class, the accuracy is 88.91%.

## 5.4 Voting Classifier

If we use voting strategy for only metadata model or CNN, we see a small improve of accuracy on both of them. For metadata models. the soft voting has 97.75%, weighted hard voting has 97.55% on validation data set which is better than KNN (96.80%), Tree model (97.44%) and Logistic Regression (97.00%). For CNN models, soft voting: 97.71%, weighted hard voting: 91.55% are also better than CNN celestial: 91.73% and CNN spectrum 88.91%.

The accuracy of voting classifiers are higher than any single model if we combine metadata models and CNN models together: Soft Voting:98.97%, Weighted Hard Voting: 98.63%. From the confusion matrix we can see that both of them classify star (class 2) perfectly followed by Galaxy (class 0) and Quasar (class 1).

## 6 Conclusions

The voting classifier can improve the accuracy by choosing the best predict class from all the models. And it is more robust for missing data and outliers. From Table 4 we can see that across all evaluation index, soft voting performs better than weighted hard voting. Moreover, for most of the classes, our voting model performs better than single model except for Galaxy and Quasar of image of spectrum data. CNN does a better job for these two class but much worse for Star. However, the price for better performance is more data used and more model trained. Although we use as simple model as possible for this job, the training and combining is also much more complicated for soft voting and weighted hard voting. In addition, the data we use is in a very good quality and quite sufficient for training a well performed model which is not so common in real situation. The corrupt or insufficient data may affect the accuracy of single models as well as voting classifiers.

According to these facts, our future work may focus on using less models with better voting strategies and using less data with worse data quality to test them.

## 7 References

Gao, Jialin, Jianyu Chen, Jiaqi Wei, Bin Jiang, and A-Li Luo. 2023. “Deep Multimodal Networks for m-Type Star Classification with Paired Spectrum and Photometric Image.” *Publications of the Astronomical Society of the Pacific* 135 (May): 044503. <https://doi.org/10.1088/1538-3873/acc7ca>.

## 8 Appendix

### 8.1 Figures

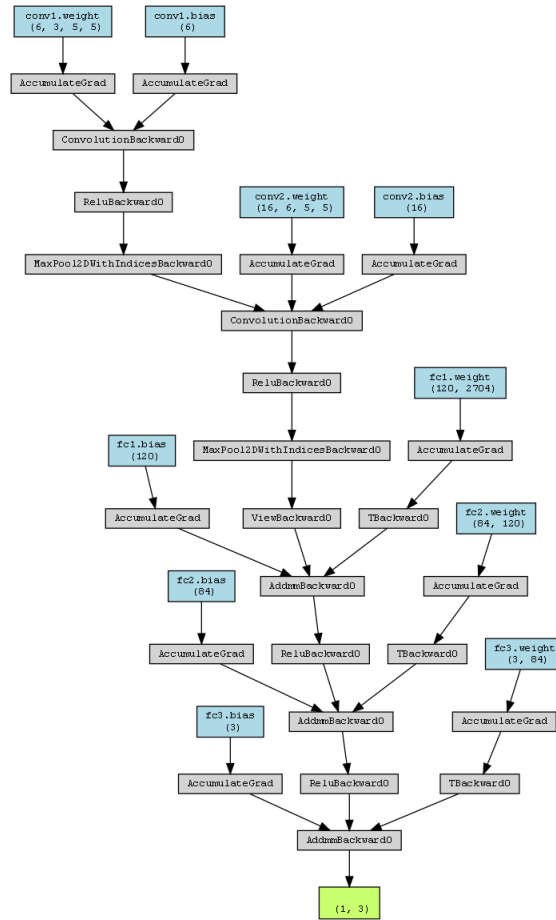


Figure 4: Structure of CNN

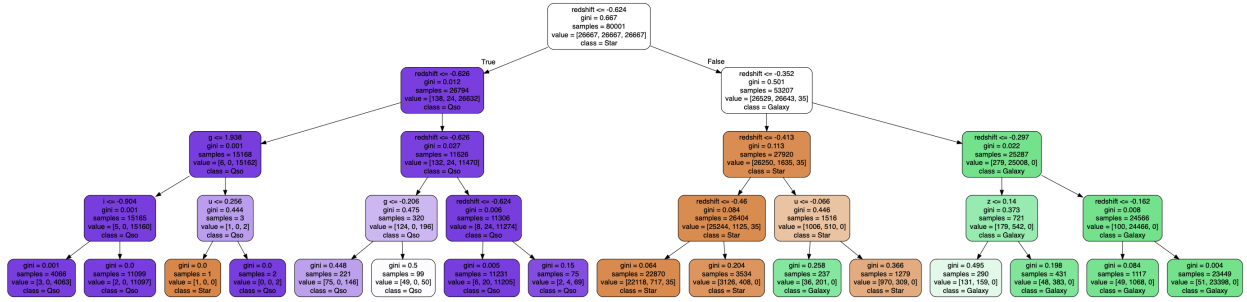


Figure 5: Result of Decision Tree

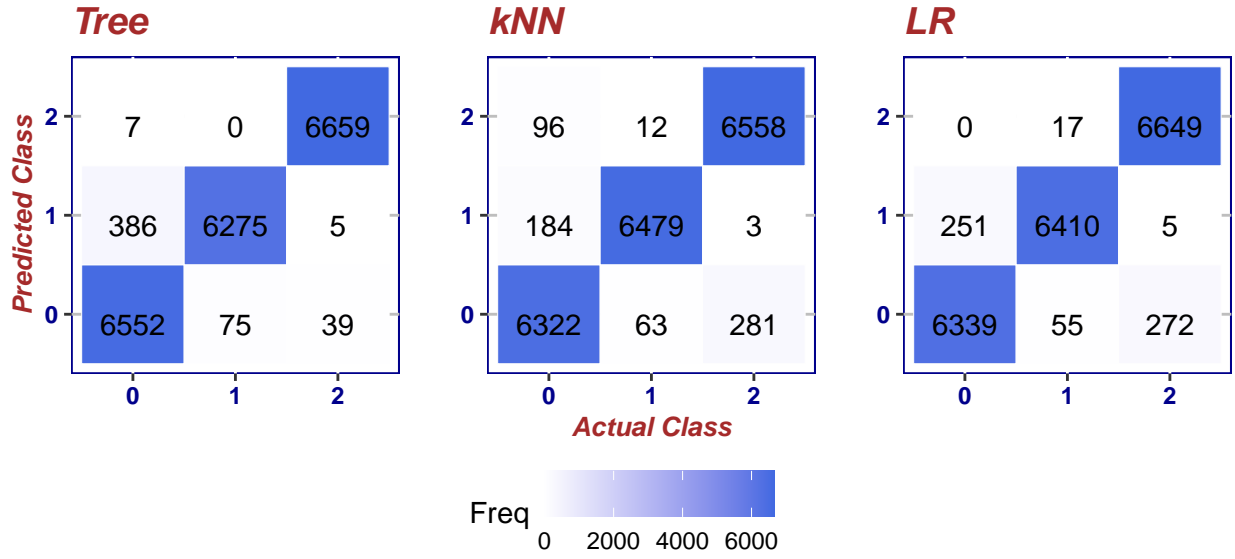


Figure 6: Confusion Matrices of Metadata Models

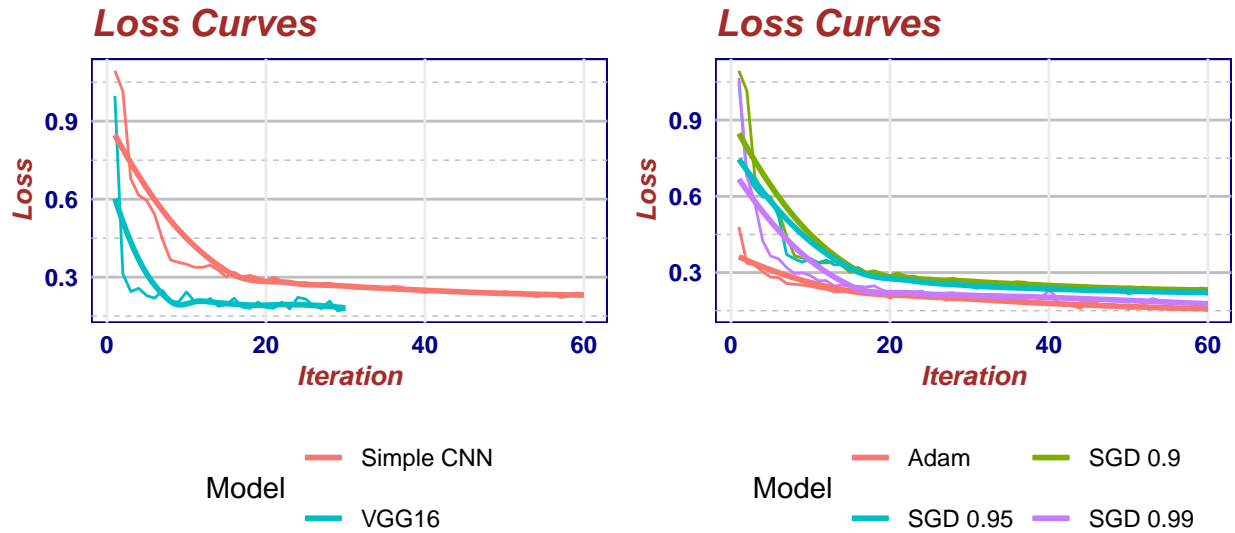


Figure 7: Loss Curve Plot for Different CNN



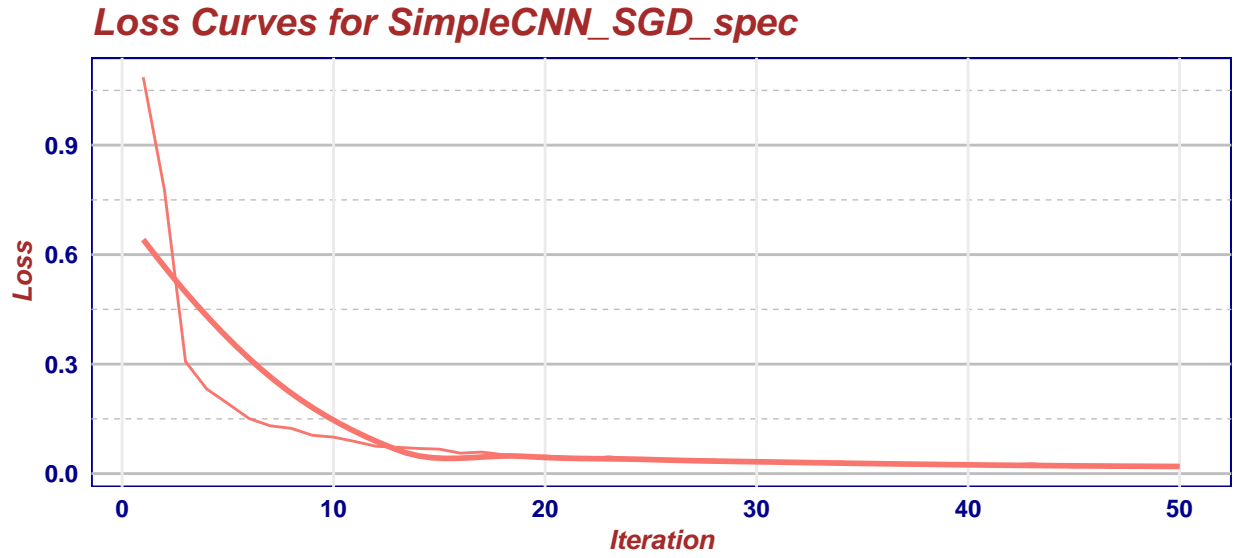


Figure 8: Loss Curve

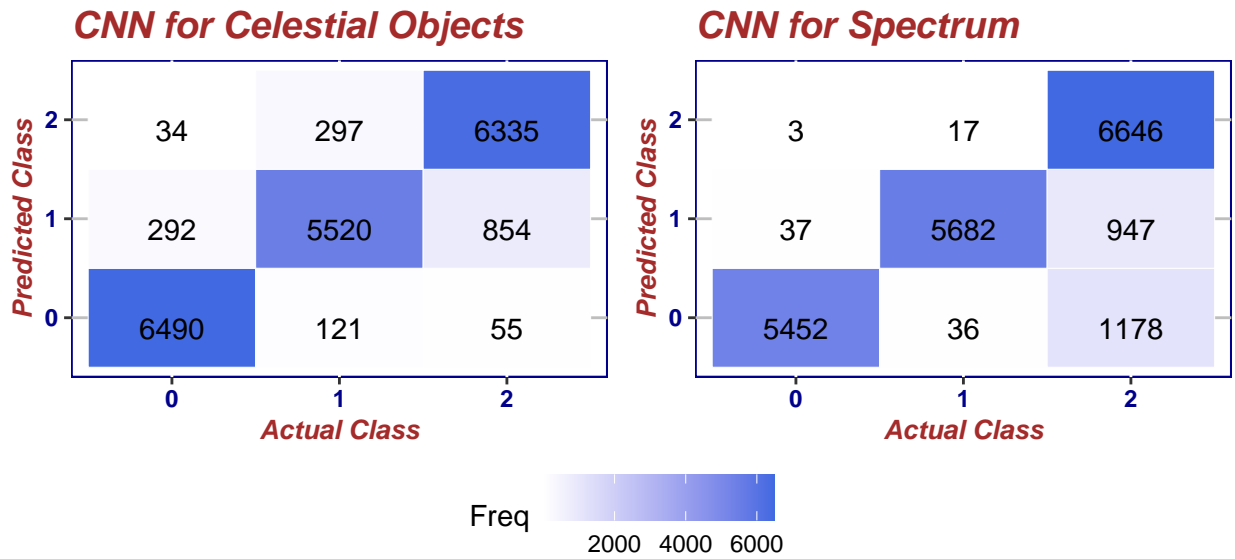


Figure 9: Confusion Matrices of CNN Models

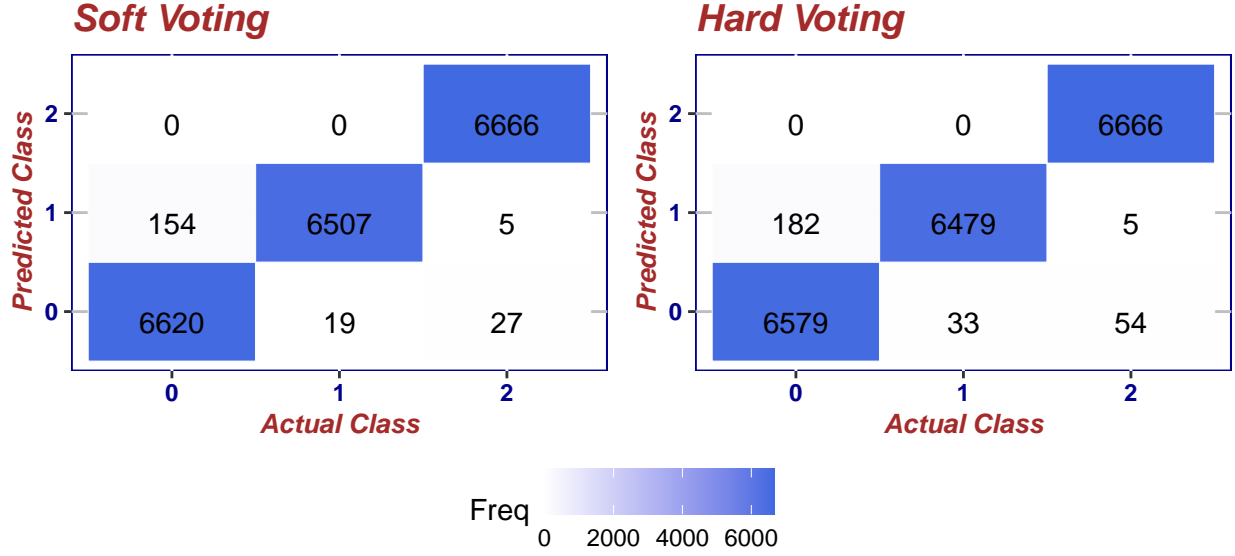


Figure 10: Confusion Matrices of Voting Classifier

## 8.2 Tables

Table 1: Metadata of the celestial objects

Variables	Explanations
objid	Object Identifier
ra	Right Ascension angle (at J2000 epoch)
dec	Declination angle (at J2000 epoch)
u	Ultraviolet filter
g	Green filter
r	Red filter
i	Near Infrared filter
z	Infrared filter
run	Run Number
rerun	Rerun Number
camcol	Camera column
field	Field number
specobjid	Unique ID used for optical spectroscopic objects
class	Object class
redshift	Redshift value based on the increase in wavelength
plate	Plate
mjd	Modified Julian Date
fiberid	fiber ID
plateid	Plate ID

Table 2: Coefficients of Logistic Regression

	Intercept	u	g	r	i	z	redshift
Galaxy	15.10	1.11	-1.70	-0.15	0.61	-0.02	23.36
Qso	16.81	-2.88	5.21	0.80	-1.22	-2.14	32.51
Star	-31.91	1.77	-3.51	-0.64	0.61	2.16	-55.86

Table 3: CNN Model Comparison

Name	Accuracy	Training.Time
SimpleCNN	91.68%	635 s
VGG16	94.11%	1281 s
Res18	94.89%	1324 s

Table 4: Simple CNN Tuning Comparison

Optimizer	Accuray	Training.Time
SGD_mom0.9	91.68%	458 s
SGD_mom0.95	92.84%	463 s
SGD_mom0.99	93.78%	462 s
Adam	93.91%	463 s

Table 5: Evaluation of Models

Data		M	M	M	IC	IS	M+IC+IS	M+IC+IS
Model		kNN	DT	LR	CNN	CNN	SVC	HVC
Accuracy		0.968	0.9744	0.97	0.9173	0.8891	0.9897	0.9863
Precision	Galaxy	0.9576	0.9434	0.9619	0.9522	0.9927	0.9773	0.9731
	Qso	0.9886	0.9882	0.9889	0.9296	0.9908	0.9971	0.9949
	Star	0.9585	0.9934	0.96	0.8745	0.7577	0.9952	0.9912
Recall	Galaxy	0.9484	0.9829	0.9509	0.9736	0.8179	0.9931	0.9869
	Qso	0.9719	0.9413	0.9616	0.8281	0.8524	0.9761	0.9719
	Star	0.9838	0.9989	0.9974	0.9503	0.997	1	1
F1	Galaxy	0.953	0.9628	0.9564	0.9628	0.8969	0.9851	0.98
	Qso	0.9802	0.9642	0.9751	0.8759	0.9164	0.9865	0.9833
	Star	0.971	0.9962	0.9784	0.9109	0.861	0.9976	0.9956

*Note:*

M: Metadata. IC: Image of Celestial Objects. IS: Image of Spectrum.