Classifying Celestial Bodies

Navid Mohseni, Sylvester Mensah, Ethan Corr

For our project, our research question is "From a set of scientific measurements, can we classify a celestial body as a star, galaxy, or quasar object?". We aim to use predictors like ultraviolet filtered intensity, infrared filtered intensity, and redshift of light waves (from this Kaggle <u>datatset</u> from CERN) to build 2 classification models and compare their predictions on a holdout set. The two classification models we will build are (1) a multinomial logistic classifier and (2) a multiclass support vector machine.

Methods

<u>Method 1:</u> We will use a multinomial logistic regression model to estimate the log odds of the data point belonging to each class.

$$\eta_{ik} = \log\left(rac{\pi_{ik}}{\pi_{iK}}
ight) = eta_{k0} + eta_{k1}x_1 + \dots + eta_{kp}x_p = \mathbf{x}_i'oldsymbol{eta}_k$$

<u>Method 2:</u> We will use a multi-class Support Vector Machine to create 3 decision boundaries for the 3 classes, one decision boundary for each combination of 2 classes.

$$egin{aligned} & \max & M \ _{eta ,eta , eta , eta , eta , eta , eta , eta & M \ & ext{ s.t. } & \|meta\| = 1, \ & y_i(\mathbf{x}'meta + eta _0) \geq M(1 - \epsilon_i), \ & \epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq B, \,\, i = 1, \dots, n, \end{aligned}$$

Data Exploration

The dataset contains 100k data points with 6 numerical predictors. The predictors are all characteristics of the light emitted from a celestial body. They are:

- Ultraviolet filter (in the photometric system)
- Green filter (in the photometric system)
- Red filter (in the photometric system)
- Near Infrared filter (in the photometric system)
- Infrared filter (in the photometric system)
- Redshift (based on increase in wavelength)

Sampling 1k points from each class, we performed Anderson-Darling tests on all 6 numerical predictors. For all 3 classes, none of the variables come from a normal distribution. Because of this, we avoid using Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA) due to the incompatibility of MVN distribution for each class.