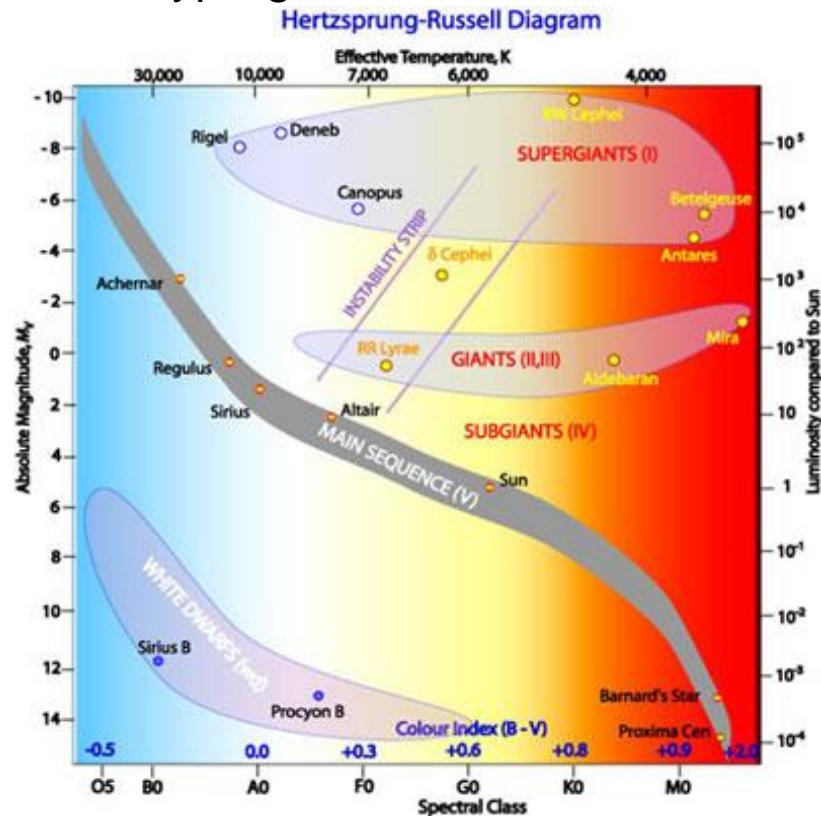# Star Classification

Crosby Pineda

# Star Cluster

Stars in the universe are grouped into various clusters, each possessing unique characteristics. These stars are categorized based on specific classifications such as Brown Dwarfs, Red Dwarfs, White Dwarfs, Main Sequence stars, Supergiant, and Hypergiants.
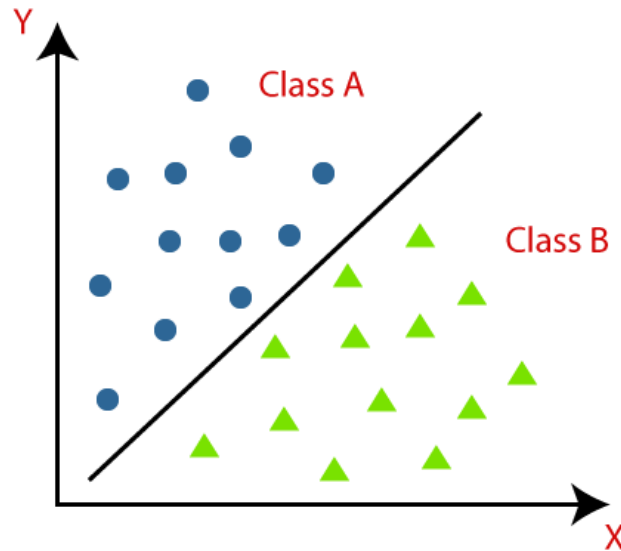


Hertzsprung-Russell Diagram

# Star Classification

It is especially because there are so many stars and clusters of them that it would be helpful to identify the type of stars, therefore the power of prediction becomes all the more needed. It is like finding a needle in a haystack of the cosmos.

# How?

Since this is a classification problem, different but similar to my original plan of identifying particles, I wanted to use a boosted decision tree. I also wanted to do regression, but with my data and the type of problem; it is not quite the same as a regular regression. For this reason, I decided to use the Ridge Classifier. In essence my 2 approaches are using machine learning to classify/ using the regression version of the combination of of the 2.
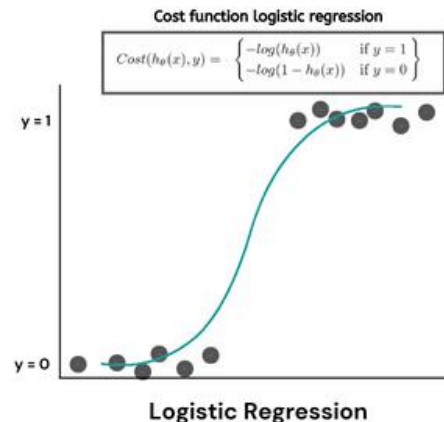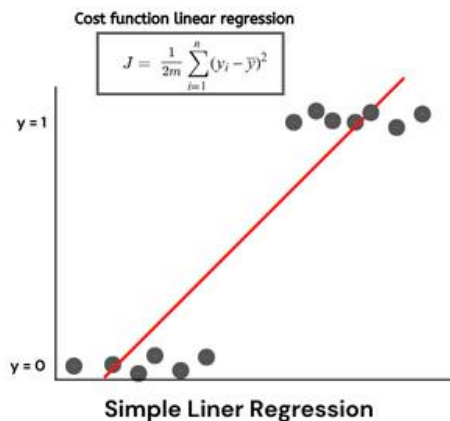
# Ridge VS BDT

Ridge Classifier is a linear model, a logistic regression model with Ridge regularization. Hence, it is closer to regression in terms of its approach. Good for linear relationships between features and classes.
BDT is good for nonlinear relationships between features and classes.
By comparing the 2 we will see which one is more accurate.

Cost function linear regression

$$J = \frac{1}{2m} \sum_{i=1}^{n} (y_i - \overline{y})^2$$

**Simple Liner Regression**

Cost function logistic regression

$$Cost(h_\theta(x), y) = \begin{cases} -log(h_\theta(x)) & \text{if } y = 1 \\ -log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

**Logistic Regression**

# XGBoost libraries

Using XGBoost I can make a Boosted Decision Tree, where many trees are trained in succession.
Using, well you guessed it, RidgeClassifier for the Ridge Classifier, to compare BDT to.

Also, an important one to make this all work is OneHotEncoder for categorical data.

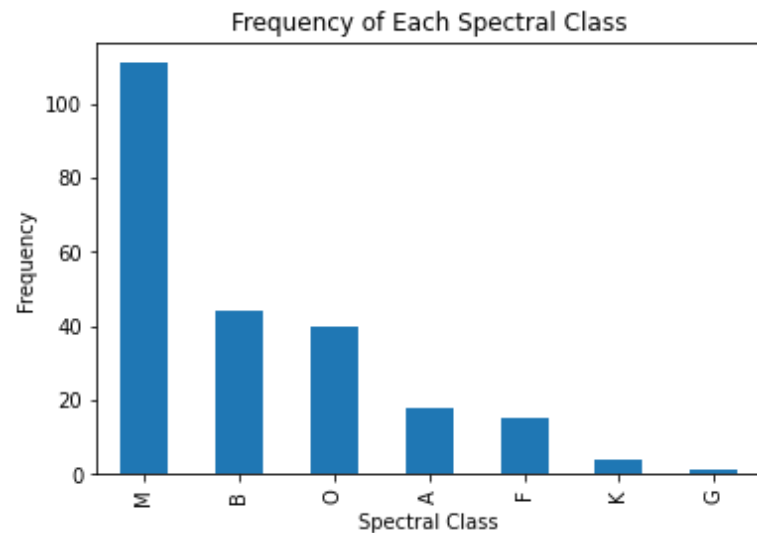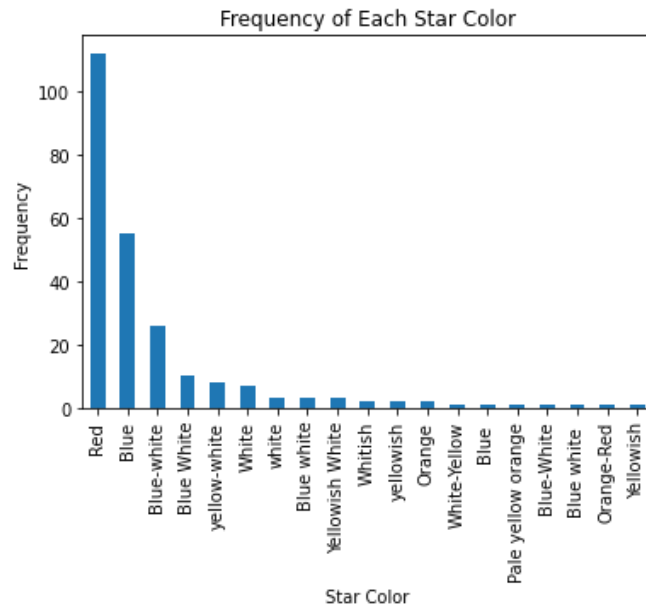| City | | Houston | Rome | Madrid | London |
|---|---|---|---|---|---|
| Houston | one-hot encoding → | 1 | 0 | 0 | 0 |
| Rome | | 0 | 1 | 0 | 0 |
| Madrid | | 0 | 0 | 1 | 0 |
| London | | 0 | 0 | 0 | 1 |

# Features

Deciding on the features was not too difficult. The features should just be the variables that make up different star types:
Temperature, Luminosity, Radius (size), Absolute magnitude, Star color, and Spectral Class. A total of 6 variables that I will use to try and make predictions.

# Issues

The distribution of the data is not symmetrical or even. This means there are prone to be more biases than I would like. Such as the amount of Star color descriptions or the different spectral classes.



Frequency of Each Star Color



Frequency of Each Spectral Class

# Issues

Specifically, star color descriptions seemed a bit more subjective than say the size of a star. Each one had different and varying amounts of descriptions. Not to mention it was also categorical. Such as Blue white, white, and whitish.

Blue White
White
White
White
Blue White
Yellowish White
Blue white
Yellowish White
Yellowish White
Pale yellow orange
Blue
Blue-white
Blue-white
Whitish
yellow-white
Whitish
yellow-white
yellow-white
yellow-white
yellow-white

# Data Prep

In order to handle this I had to ensure that for star color, the descriptions appeared at least twice.
Using One-Hot Encoding to convert categorical data to have unique values to be used for machine learning to handle star color and spectral class.

using aspects of linear algebra in encoding categorical features and transforming them to be put back into the feature set to be able to be used for machine learning.

# Train Test Split

To make the Train Split I made sure to use a random seed to not only replicate it, but also to randomize the distribution. While at the same time, I had to ensure that given the disruption of the data, the split had a bit of each for star color description. If I did not do this then the accuracy would be lowered if the training set did not have one of each.

# Train Test Split

I used about 80% of the data to train and 20% of the rest to test. This falls within the usual ratio. Was curious to see other ratios, and it does make both classifiers less accurate.

# BTD and RC

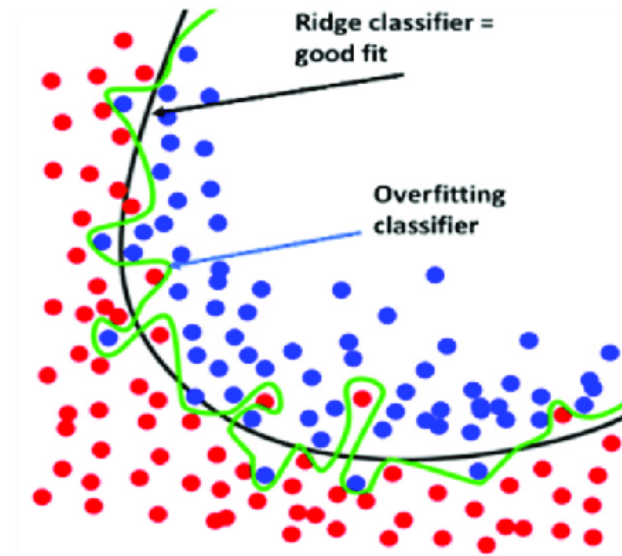After doing all of that I then used Ridgeclassifier and xgboost to do the machine leaning aspect. Most of the hard part was getting the data prep and integrating the categorical data.

# ROC Curve

As I stated in my original plan, I did want to do a ROC curve to see the accuracy. I did this for each star type to see the accuracy of each one. Which each one was 100% accurate.

Multi-class ROC for XGBoost Classifier

ROC curve of class 0 (area = 1.00)
ROC curve of class 1 (area = 1.00)
ROC curve of class 2 (area = 1.00)
ROC curve of class 3 (area = 1.00)
ROC curve of class 4 (area = 1.00)
ROC curve of class 5 (area = 1.00)

# Overall Accuracy

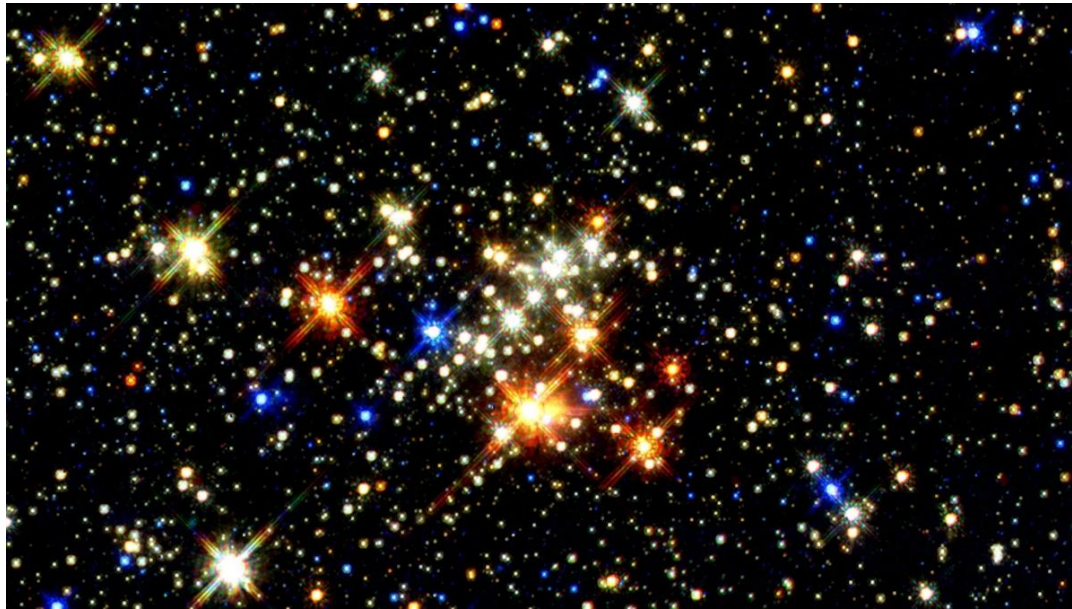The over all accuracy for my BDT was 100% as the accuracy of each class was 100%. My Ridge Classifier on the other hand was about 91%. This means in the specific circumstance that the BDT was better at predictions.

Ridge Classifier Accuracy: 0.9148936170212766

XGBoost Accuracy: 1.0

# Thoughts

While I do not think a BDT should be 100% accurate, I think it largely has to do with my sample size being rather small. It was about 200 +. Which is rather small compared to the billions of stars in the universe.

# Simplifications

My simplifications interfered with the accuracy as I eliminated star colors that do not repeat and also ensured that they were distributed between the split. This could have also been the reason for it.

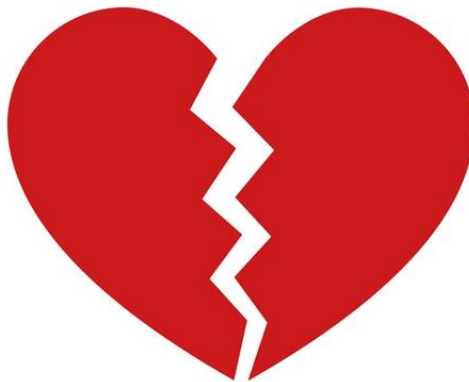Ridge Classifier

Boosted Decision Tree

# What Can we Say

We can say that perhaps from out sample size that the relationship between the features and class are not exactly linear. Physically this means that: Temperature, Luminosity, Radius (size), Absolute magnitude, Star color, and Spectral Class may not have a strictly linear relationship with the star type. After all the Ridge Classifier is effective at linear relationships, so if it was about 91% then its not quite there.



Types of Stars

Yellow Dwarf Star    Red Dwarf Star    Red Giant Star    Red Supergiant Star    Blue Giant Star    White Dwarf Star    Brown Dwarf Star

# Relationships

Rather the relationship tends to be more complex and multifaceted than just a 1 to 1. This makes sense because it required 6 variables to make a single prediction. This is why machine learning should be used for complicated relationships.

# Conclusion

This was a fun project, that made me go to areas I am not quite familiar with. We see the importance of using the appropriate type of machine learning for the problems at hand, it is not just machine learning solves all. I got my feet wet on the concepts on machine learning. There is far more on the horizon.