

## Capstone 2 Milestone #1 Report

### Problem Statement

People often buy the same kind of wine over and over again. If you don't know what a wine's points value means or how origin can affect taste, you may not be as willing to try different kinds of wine. I want to build a recommendation tool that will recommend different wines that are similar to the given wine. The data will come from the following dataset:

<https://www.kaggle.com/zynicide/wine-reviews>

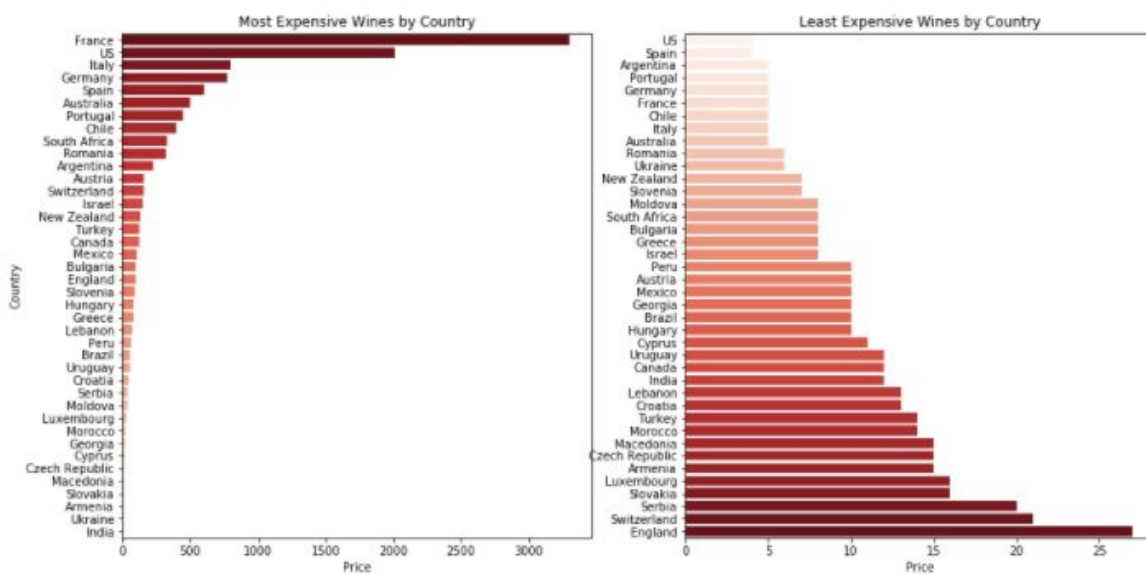
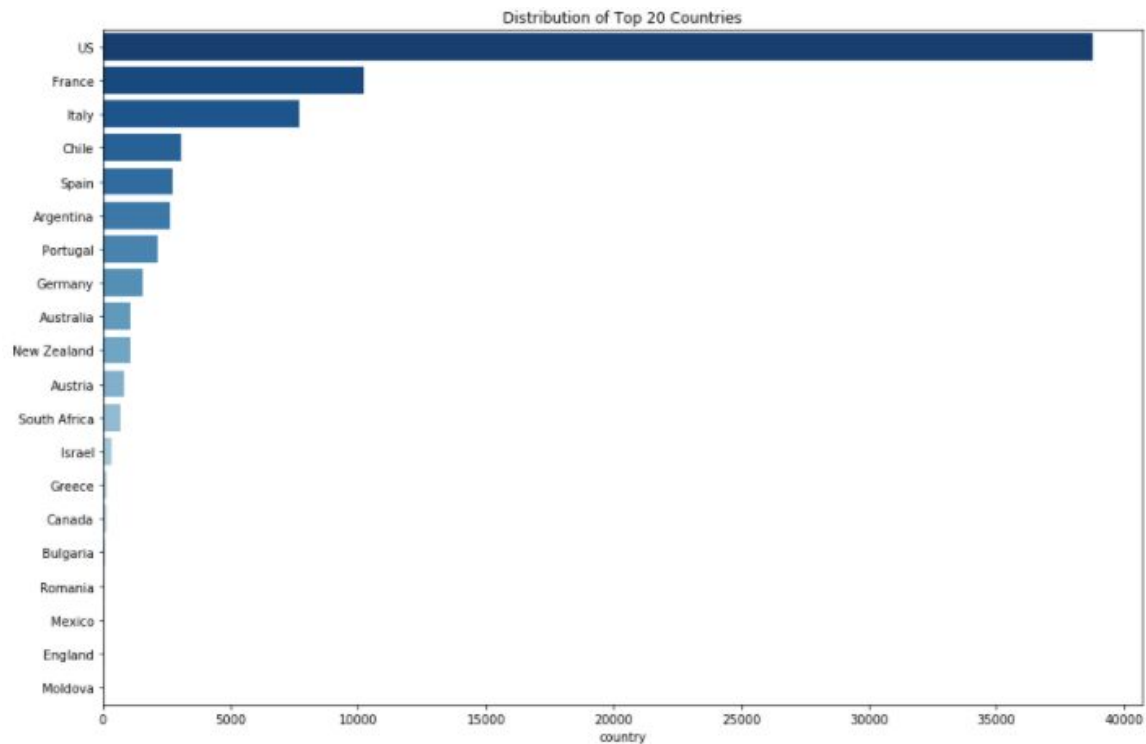
This dataset contains various attributes about a given wine, such as title, price, points, country of origin, the vineyard where the grapes were made, and a brief description of the wine from a sommelier. I will look for trends in the given data in order to provide accurate recommendations for users looking to branch out and try different types of wine.

### Data Wrangling

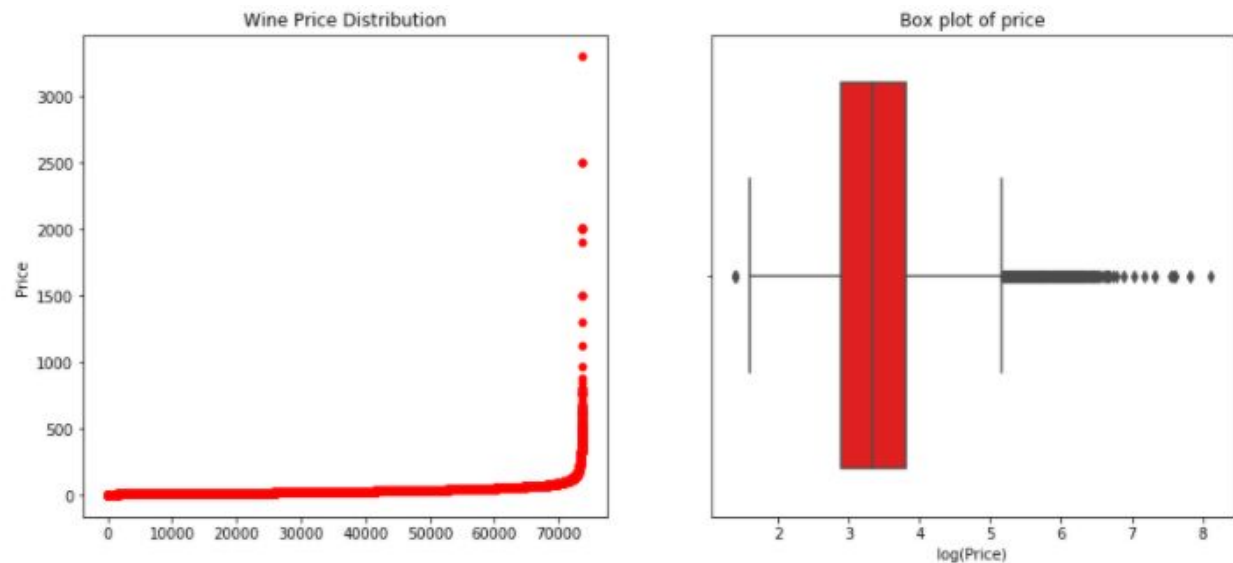
The wine data for this project is contained in two csv files, the first containing 150,000 rows and the second containing 130,000 rows. I first read in each csv file as a pandas dataframe and concatenated the two dataframes to make one large dataframe containing all of the data, which will make the analysis easier. Looking at the number of null values in each column, I removed the designation, region 1 and 2, taster name, and taster's twitter handle columns. Then, because of the duplicates as well as the null values in the title column, I decided to drop the duplicates in this column because I needed it to build the recommendation system. Other columns, such as country, province, and price, had a relatively small number of null values and contained information I was planning on using for my analysis, so I just dropped the rows containing null values in these columns rather than dropping the columns altogether. There were a large number of varieties in the dataset as well, so I decided to just focus on the top 20 for my analysis.

## Exploratory Analysis

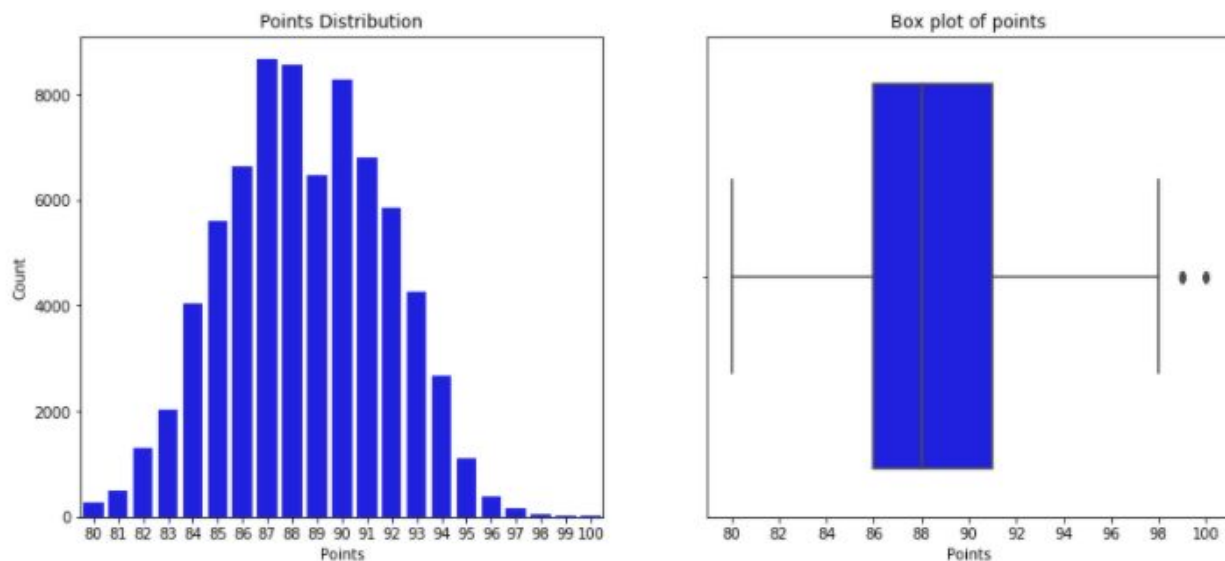
Out of the 50 countries represented in the dataset, the majority of the wine is produced by the US, followed by France and Italy. There is also a similar group of countries atop the list of most expensive wine sold and the least expensive wine sold, leading me to believe there is a great range of wines sold in countries like the US, France, Portugal, Argentina, and Italy. I noticed there was one country named “US-France”, which is likely a wine that is produced in the US and sold in France or vice-versa.



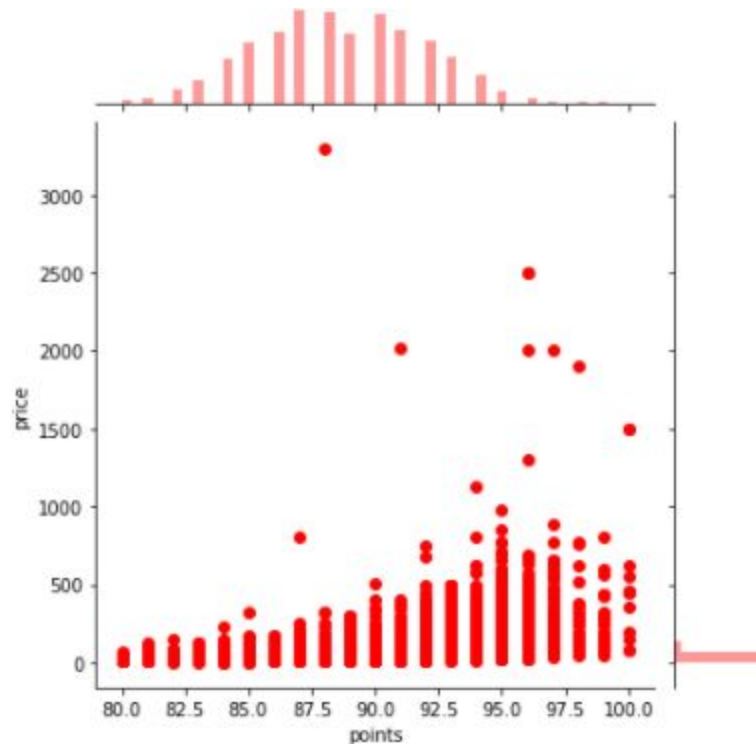
Looking at the distribution of wine prices, there appears to be some upper outliers, so I wrote a simple function to output the number of outliers in a given column. In the price column, there are no lower outliers and 751 higher outliers out of 73,691 total observations. This only represents 1.0296% of the data. Also, none of these outliers appear to be from erroneous data, so I decided to leave them in. I also decided to use a logarithmic scale because the range in prices was high.



The points distribution appears to be normally distributed with a mean of about 88. There are no lower outliers and only 35 upper outliers, representing 0.0475% of total observations. It appears that any wine given a rating of 98 or higher was deemed an outlier in this dataset. For the same reason as the price outliers, I will not exclude them from the analysis.

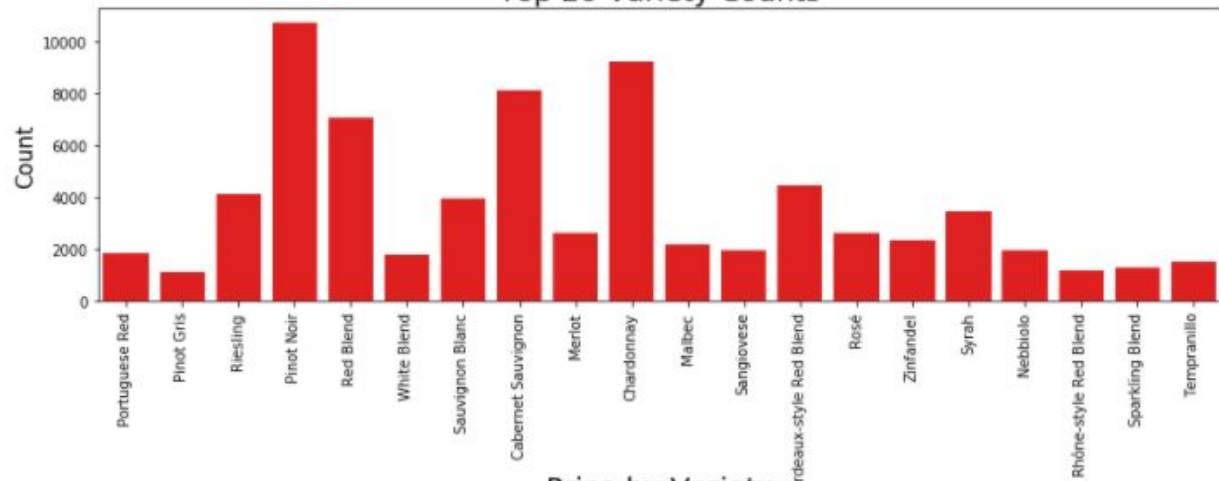


I also wanted to see if there was a correlation between price and points, so I made a joint plot with the points data on the x-axis and the price data on the y-axis. Points and price are generally positively correlated, though it appears there are some relatively affordable wines that are highly rated. Most interestingly, the highest priced wine in the dataset at \$3,300 was only given a rating of 88, which is around the mean of the points data. It lends credence to the possibility that not all expensive wine is worth the price.

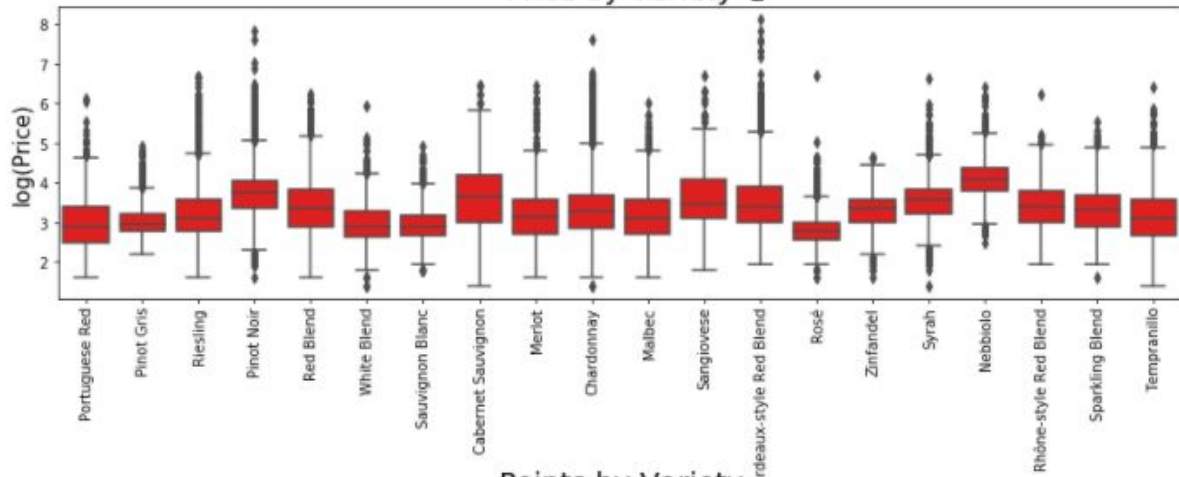


Without doing any machine learning, I thought the easiest way to make wine recommendations in practice was by recommending the same or similar variety, so I wanted to see a breakdown of the top 20 varieties in the dataset, along with boxplots of the price and points grouped by variety. The top 5 varieties in the dataset are pinot noir, chardonnay, cabernet sauvignon, red blends, and Bordeaux-style red blends. Once again, I used a logarithmic scale for the prices in each variety since there the range in prices is high, even when grouped by variety. The 5 most expensive varieties on average are nebbiolo, cabernet sauvignon, Bordeaux-style red blends, pinot noir, and sangiovese. The 5 highest-rated varieties on average are nebbiolo, riesling, pinot noir, syrah, and Rhone-style red blends.

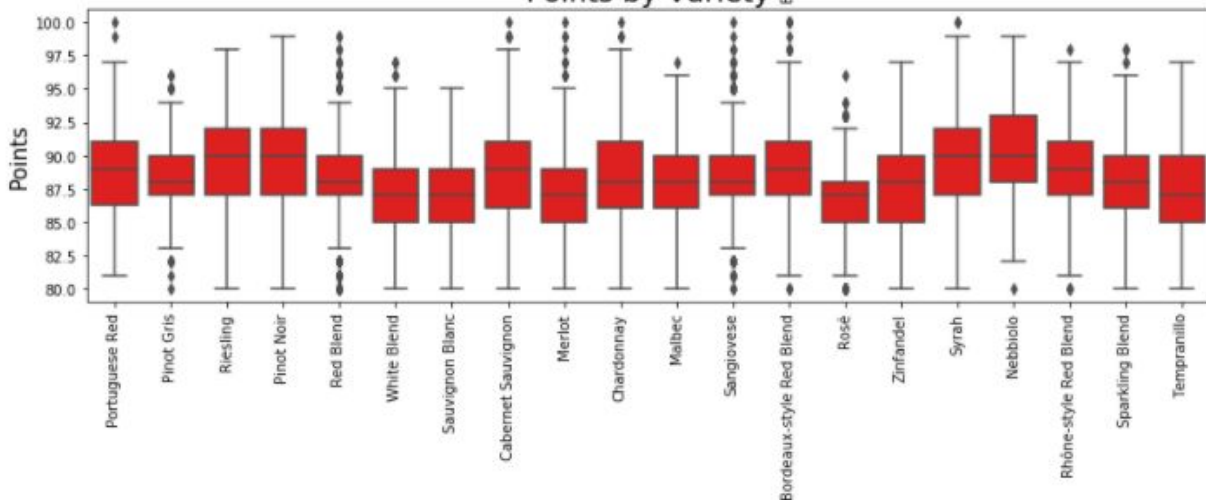
### Top 20 Variety Counts



### Price by Variety

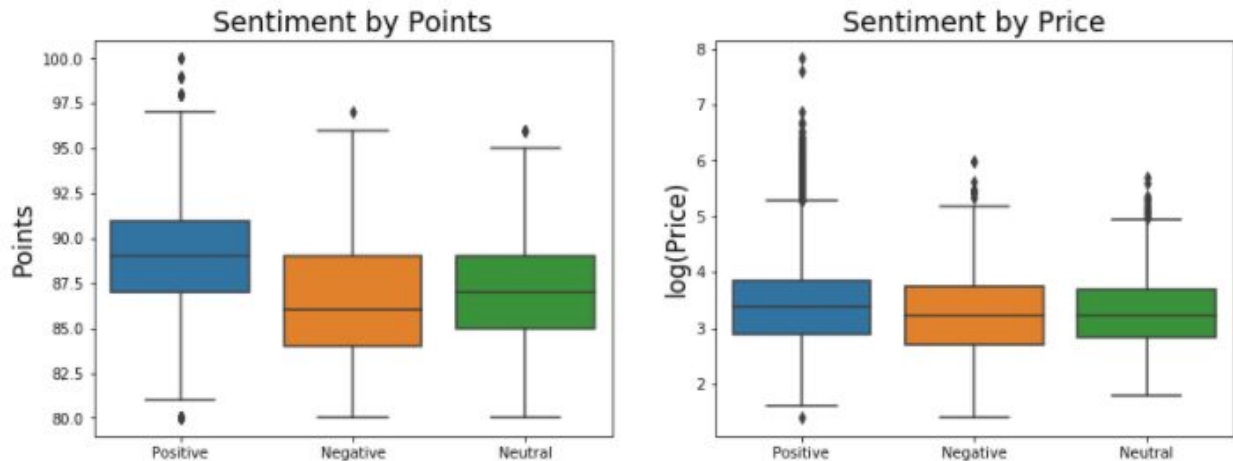


### Points by Variety









There appears to be a more positive sentiment among the higher rated wines and a more negative sentiment among lower rated wines, though there is less variation on sentiment when grouping by price.

## Inferential Statistics

In order to test for statistical independence between the price and points columns, I first wrote functions to assign categories to each wine based on its points rating, as well as categories for price range. This allowed me to summarize the data in a contingency table shown below.

| price_range | 1-30  | 100+ | 31-60 | 61-100 |
|-------------|-------|------|-------|--------|
| quality     |       |      |       |        |
| bad         | 12037 | 26   | 1505  | 184    |
| good        | 4019  | 1721 | 10356 | 4584   |
| great       | 7     | 324  | 100   | 219    |
| ok          | 25327 | 385  | 10941 | 1956   |

Then, I performed a chi-squared test with the `scipy.stats` package. The test statistic was very large at 25,185.857, which was much larger than the critical value of 16.919. The p-value was also less than the chosen alpha, so I rejected the null hypothesis that these two variables are independent.