

# **Which regions produces the best quality wine and what are the notes that describes them?**

## **1. Introduction**

Living in the San Francisco Bay Area, we have great proximity to the Napa Valley and Sonoma County, two of the best wine producing regions in the world. With that in mind, I got interested in which wine regions produces the best quality wines in the world.

In addition, wine is a very complex product with many different notes (flavors and aromas). For example, wines can be described as fruity, earthy, chocality, having a subtle or strong finish, or plain, among many other notes. With that in mind, I set out to see how well a computer could correctly classify a wine based on the description of the wine.

This analysis should be interested for an everyday consumer who is interested in trying and purchasing the best wines, as well as becoming more educated about what goes into each varietal.

## **2. Data Wrangling**

After reading the .csv file into python and calling .head() to get a sense for the columns and values of the dataset. My initial observation from this step include:

- Most of the columns have string objects, with the exception of the “points” and “price” columns. These have numerical values.
- The column names are well written. Everything is lowercase and columns with multiple letters are joined by an underscore, and not a space. Which is great since it’ll avoid issues in the future.
- The dataset is organized in tidy format, having one subject per column.

My next step was to look at “.info()” command to verify each columns dtypes, as well if the data has any missing values.

- Turns out the dtypes are correct. Including the “points” and “price” columns, which have int64 and float64 respectively.
- There are a lot of missing values from almost every column in this dataframe. The only columns that do not contain missing values are “description”, “points”, “title”, and “winery”.

From there, since I’m interested in comparing wine quality between countries, I tackled the missing values in country column. Given that the “winery” field has all of its values, created a

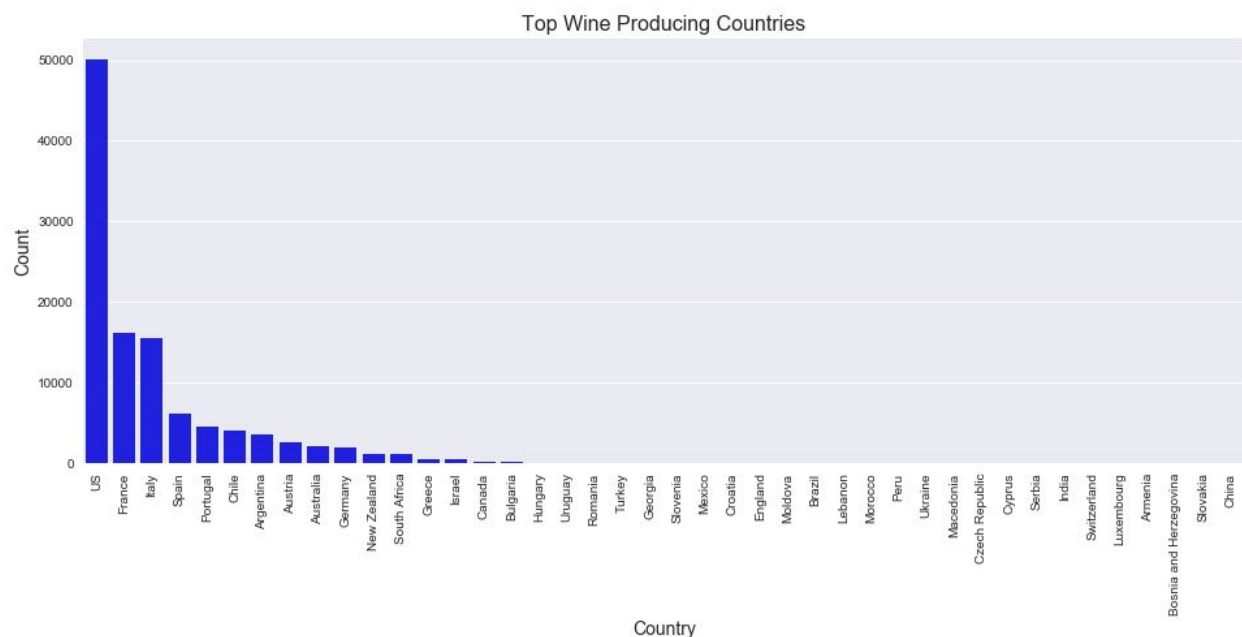
map of {winery:country} value pairs and mapped the missing country values back into the dataframe. This approach worked for 32 missing values. There are still 31 missing country values, out of over 100k. Since 31 is such a small proportion of the dataset, I dropped those rows.

Next I dropped missing values from the dataframe including duplicates, missing values from the price column and variety columns. From here, the dataset is ready for exploration.

### 3. Exploratory Data Analysis (EDA)

After cleaning the dataset, I wanted to know more about the dataset. The questions and results from my findings can be seen below:

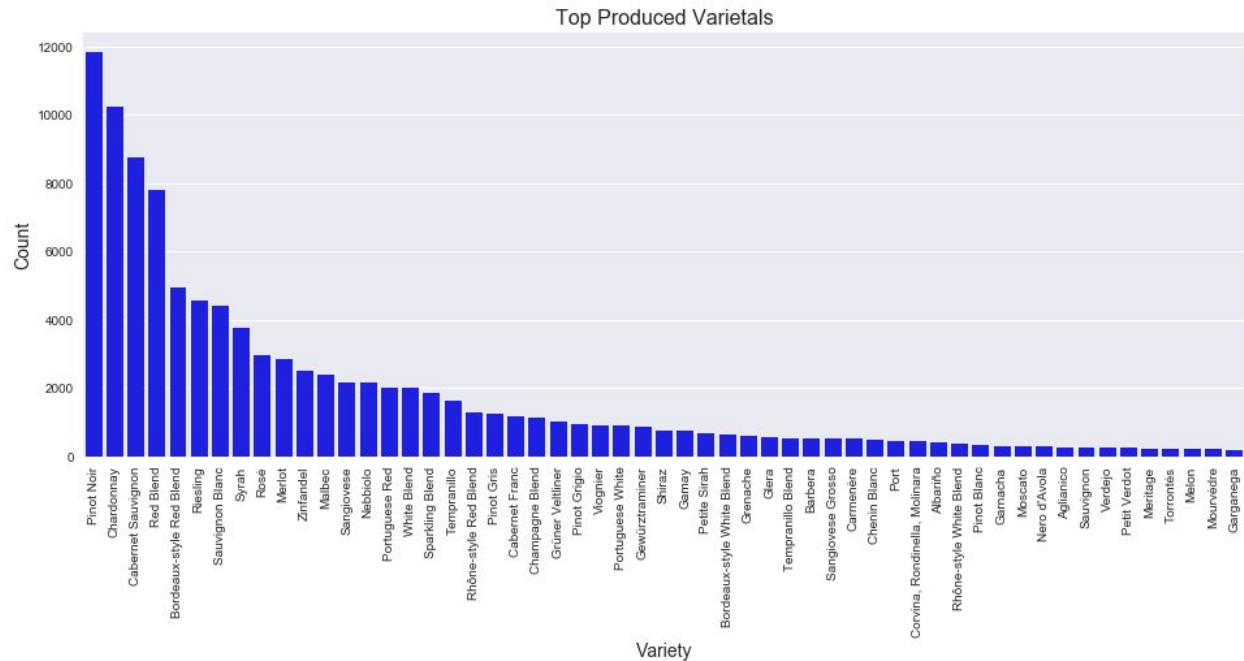
- **Which nations produce the most wine in the world?**



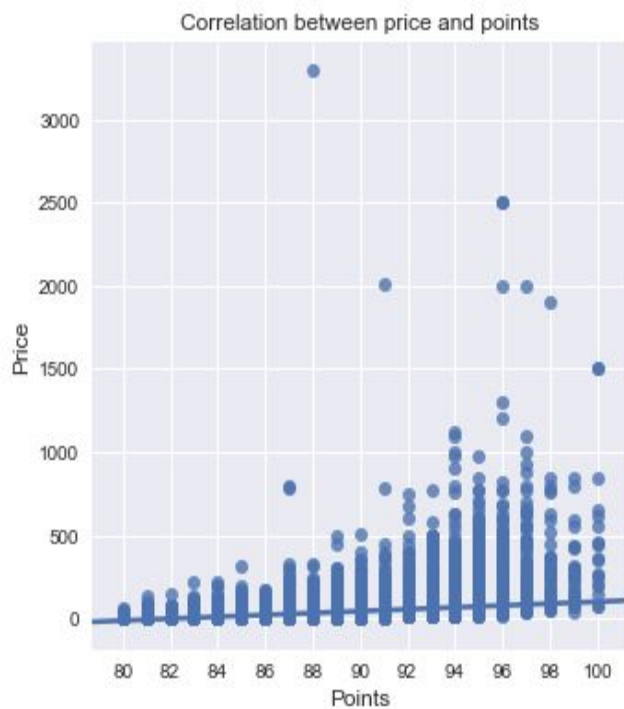
- **How many different varieties of wine are included in the dataset?**

- There are 694 different wine varieties in this dataset, including blends and champagne. In order to have a comprehensible graph, I limited the graph to varieties with at least 200 observations in the dataset.

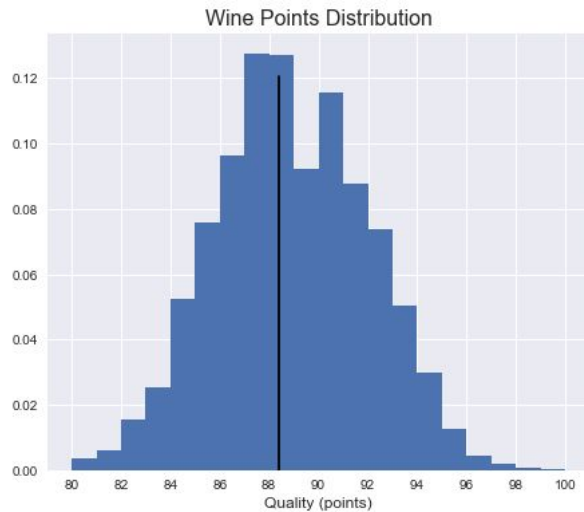
- **What are the top varieties?**



- What is the relationship between wine price and points (quality)?



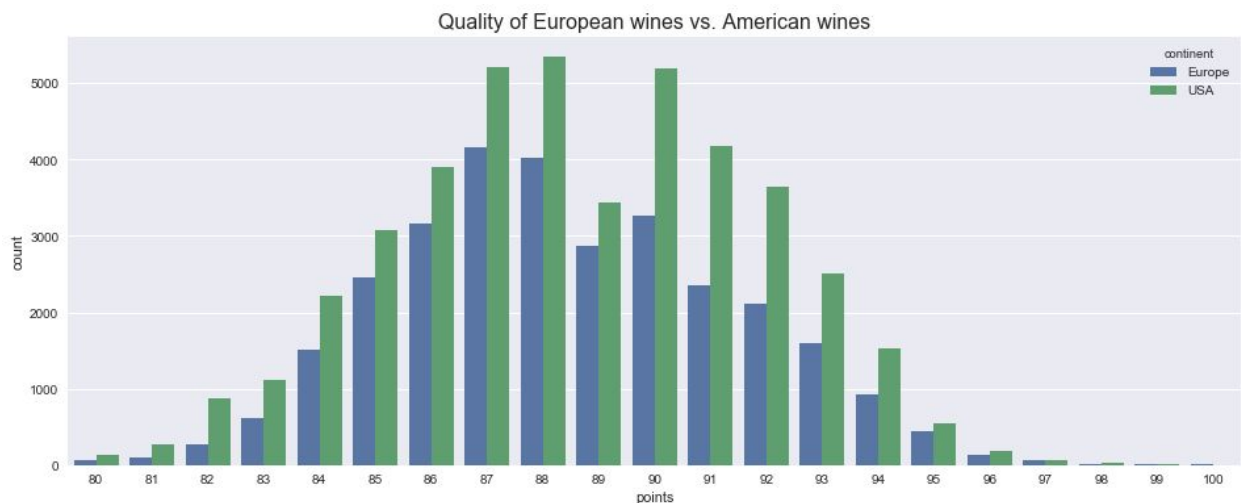
- What type of distribution do points and price have?
  - The black vertical bars in the graph below are the averages.



## 4. Statistics

Given all the findings above, I set out to see if there's a significant difference in wine quality (points) between the United States and France, Italy, Spain, and Portugal, the next four wine producing countries. From here on out our refer to those four European nations simply as Europe.

I put together an initial graph to visualize the data:



You can see that in absolute terms, the United States produces a greater quality of wine when compared to Europe. This, however, is likely because the United States has more wine compared when compared to Europe.

In order to avoid this issue, I compared means with the null hypothesis that there is no difference between the quality of wines from the United States versus the wines from Europe.

As can be seen from the calculations and low P-value shown below, we can reject the null hypothesis and confirm that on average, wines from the United States have greater quality than the wines from Europe.

## **5. Machine Learning**

To be included....