

**Chocolate!**

Kimberly Cable

Bellevue University

DSC 550: Data Mining

Brett Werner

August 13, 2022

## **Chocolate!**

Chocolate. Who doesn't like it? It's been around since the 19<sup>th</sup> century BCE. But what makes chocolate highly desirable? In this study, I plan to see if certain factors determine whether certain chocolate in the United States is highly delectable.

### **Why predict the best chocolate?**

Chocolatiers, bakers, and romanticists are all looking for that one thing that will set them apart from the rest. Knowing the correct combination of cocoa, butter, and/or vanilla or even a reviewer's most memorable characteristic of your chocolate could make or break your company, cookie, or even relationship.

### **Data Selection and EDA**

#### **Data**

The data I will use comes from the Manhattan Chocolate Society, Flavors of Cacao website<sup>1</sup>, and their Chocolate Bar rating table. The data includes features such as cocoa content, where it was made, bean origin, the number of ingredients, and its characteristics and rating. I will also use the USA Craft Makers table merging it with the Chocolate Bar ratings to look more closely at the US Chocolate Companies.

#### **Analysis**

In analyzing the data, I plan to see if I can predict the ratings for chocolate in the United States and see if any features play a part in good chocolate.

I will look at questions such as:

- Which states have the best chocolate?
- Which ingredients make the best chocolate?
- How much cocoa makes for good chocolate?
- What are the key characteristics of good chocolate?
- Where do the best cocoa beans originate from?

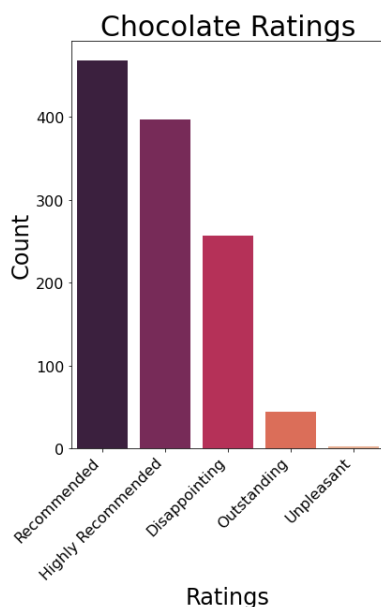
## Challenges

Some of the challenges I see are cleaning the dataset and separating out some of the columns. One major hurdle is knowing the best way to classify the ratings. Do I keep the ratings as they are? Do I group them into the 5 original rating scales or use some other grouping to best use to predict a good piece of chocolate?

## Data Cleaning

The data was scraped and moved into CSV files to begin explanatory data analysis. Some cleaning steps that needed to be taken were to standardize the spelling of some states in the Chocolate Rating table to help merge with the United State Chocolate Makers table. Once the two tables were merged, I pulled out just the United States chocolate to analyze further.

Initially looking at the data, most of the chocolate in the United States is rated “Recommended” and “Highly Recommended,” (see Figure 1). Getting the average rate per state most States are rated as “Highly Recommended” (see Figure 2). This would indicate a good chance of getting good chocolate anywhere in the United States.



*Figure 1: Chocolate Ratings for United States Chocolatiers*

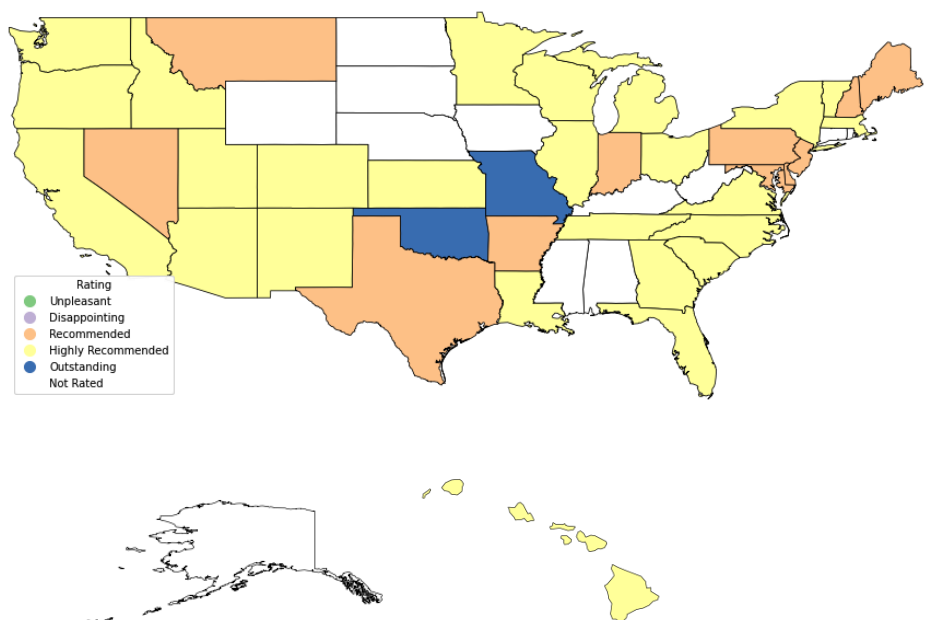


Figure 2: Average chocolate rating per State in the United States.

California has the most Chocolatiers in the United States (see Figure 3) and most of the cocoa beans originate from the Dominican Republic (see Figure 4).

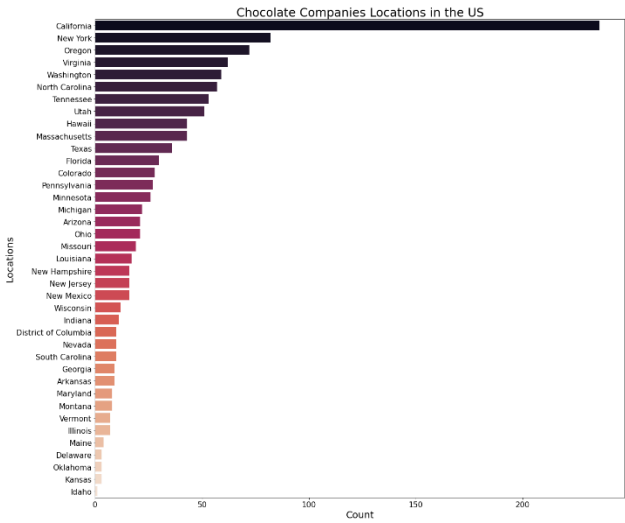


Figure 3: United States chocolate company locations.

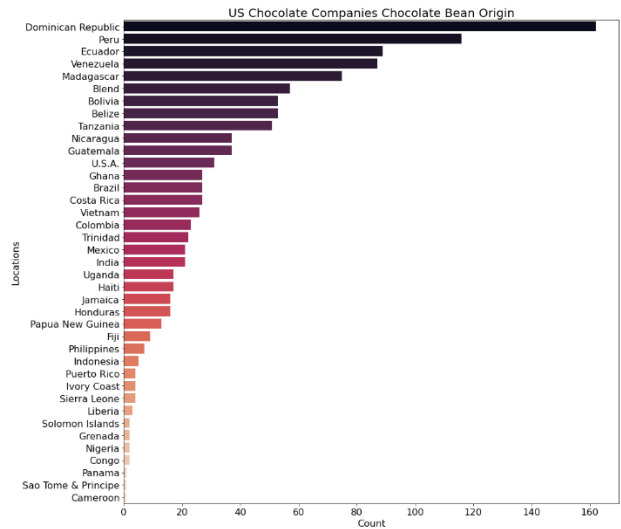


Figure 4United States cocoa bean origins.

Looking at ingredients, butter and sugar and butter, sugar and cocoa butter are the most popular ingredients (see figure 5), and cocoa percent is around 70% (see figure 6).

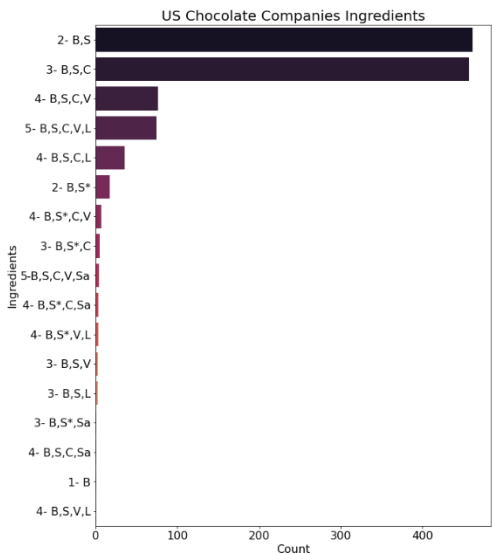


Figure 5: United States chocolate ingredient count

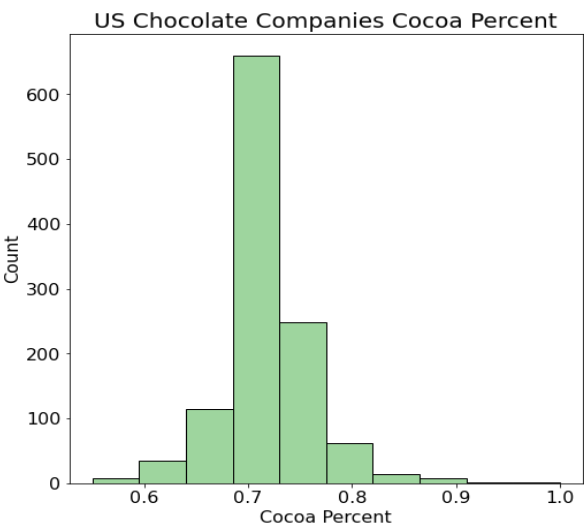


Figure 6: Cocoa percent in United State chocolate count

The top 20 most memorable characteristics in reviews used words such as “sweet”, “cocoa”, nutty”, and “roasty” (see figure 7)



not have enough samples at the lower and upper ratings to successfully predict chocolates at that level (see Table 1).

	precision	recall	f1-score	support
Above Average	0.68	0.61	0.65	83
Average	0.79	0.84	0.81	146
Below Average	0.00	0.00	0.00	1
accuracy			0.75	230
macro avg	0.49	0.48	0.49	230
weighted avg	0.75	0.75	0.75	230

Table 1: Classification Report using all features.

The precision scores indicate a good ability to correctly predict each rating for “Average” and “Above Average” ratings, but it had trouble with correctly predicting positive “Below Average” ratings as indicated by the low recall score of 0%.

Next, I looked at reducing the number of features that did not help predict the best chocolates using Principal Components Analysis. Looking at the same three models and using a grid search, I found that the Logistic Regression model predicted chocolates the best. With the reduction in features, the accuracy score was 71%. It correctly predicted a rating of “Average” 76% of the time, and “Above Average” 64%. It did not predict “Below Average” at all with a score of 0%. This indicated that we did not have enough samples at the lower rating to successfully predict chocolates at that level. (See Table 2).

	precision	recall	f1-score	support
Above Average	0.59	0.69	0.64	83
Average	0.80	0.73	0.76	146
Below Average	0.00	0.00	0.00	1
accuracy			0.71	230
macro avg	0.46	0.47	0.47	230
weighted avg	0.72	0.71	0.71	230

Table 2: Classification Report with reduced number of features.

The precision scores indicate a good ability to correctly predict each rating for “Average” and “Above Average” ratings, but it also had trouble with correctly predicting positive “Below Average” ratings as indicated by the low recall score of 0%.

### Conclusion

Rating	All Features	Using PCA
	Random Forest Classifier	Logistic Regression
	F1-Score	
Above Average	0.65	0.64
Average	0.81	0.76
Below Average	0.00	0.00
Accuracy	0.75	0.71
Weighted Average	0.75	0.71

*Table 3: Summary of Classification Report of both models*

In predicting chocolate in the United States, I believe personal taste is the true measure. Trying to predict good chocolate, and shown, can be tricky. Aside from that, looking strictly at accuracies, the Random Forest Classifier fared slightly better than using Logistic Regression, but both were above 70% (see Table 3). For me, the most important indicators were the F1 scores of each of the ratings. Using all features and the Random Forest Classifier, chocolates in the “Average” and “Above Average” ratings were correctly predicted 65 to 81% of the time. The weighted average was also 75% correctly predicting a rating. I believe using all features and the Random Forest Classifier would help initially get you good chocolate. Then, you can personally choose your tastes and preferences from the initial choices the model predicted.

The next step I would take to see if I could accurately predict good chocolate is to include all chocolates not just those made in the United States. This will give us more data to work with, especially in the bottom classification. Looking at other reviews besides the Manhattan Chocolate Society ratings will help broaden the data to different opinions.

## References

*Flavors of Cacao*. (n.d.). Retrieved from Chocolate Bar Ratings: <http://flavorsofcacao.com/>