

Chocolate EDA

Manvir Kohli, Julie Song, Kelvin Wong

2022-11-18

```
## Warning: package 'kableExtra' was built under R version 4.2.2
```

Summary of the Data Set

The data set is provided by the Manhattan Chocolate Society, and was found and retrieved from the tidyTuesday data project, specifically through this link. The data set contains observations for different types of dark chocolate, including the manufacturing company, origin of the cocoa beans used to make the chocolate, the other ingredients in the chocolate, and the amount of cocoa. They have also provided a feature column that contains descriptive words relating to the characteristics of the chocolate flavor, and a final rating.

We have split the original data set into training and testing data sets using a 70%-30% split. The following data processing and EDA analysis have been performed on the training set, which contains originally contains 1771 observations with 9 features and 1 target. After data processing and analysis, we have 7 features that we will use for modelling.

Glimpsing the Data

```
## Rows: 1,771
## Columns: 10
## $ ref                <dbl> 2458, 2454, 2542, 2546, 2542, 797, 10~
## $ company_manufacturer <chr> "5150", "5150", "5150", "5150", "5150~
## $ company_location    <chr> "U.S.A.", "U.S.A.", "U.S.A.", "U.S.A.~
## $ review_date         <dbl> 2019, 2019, 2021, 2021, 2021, 2012, 2~
## $ country_of_bean_origin <chr> "Dominican Republic", "Madagascar", "~
## $ specific_bean_origin_or_bar_name <chr> "Zorzal, batch 1", "Bejofo Estate, ba~
## $ cocoa_percent       <chr> "76%", "76%", "68%", "80%", "68%", "7~
## $ ingredients         <chr> "3- B,S,C", "3- B,S,C", "3- B,S,C", "~
## $ most_memorable_characteristics <chr> "cocoa, vegetal, savory", "cocoa, bla~
## $ rating              <dbl> 3.50, 3.75, 3.00, 3.25, 3.50, 3.50, 2~
```

We have 1771 observations with 9 features and 1 target. After checking the structure and summary statistics for our data, we find the following:

- Our target variable is `rating`
- The columns `ref` and `specific_bean_origin_or_bar_name` are identifier columns and should be dropped
- The columns `company_manufacturer`, `company_location`, `country_of_bean_origin` and `ingredients` are all read as character columns but should ideally be factors (i.e. categorical columns)
- `most_memorable_characteristics` is likely a text column, containing many unique words
- `cocoa_percent` is read as a character column while it should be numeric

	Null Count
company_manufacturer	0
company_location	0
review_date	0
country_of_bean_origin	0
cocoa_percent	0
num_of_ingredients	55
ingredients	55
most_memorable_characteristics	0
rating	0

Data Processing

- We need to convert all the columns to the correct data types, but we will do this as the last step in our data processing.
- The ingredients column has two components in each cell - the number of ingredients and the actual ingredients. So we can split this column into two and have two separate features(**num_of_ingredients** and **ingredients**)
- Thus after dropping we **ref** and **specific_bean_origin_or_bar_name** we have an overall total of 8 features with 1 target.
- We also checked our columns for null values, and found that there were 55 observations with missing values in our **ingredients** and **num_of_ingredients** columns.

Exploring Categorical Columns Further: For all the factors there are many levels. We can reduce the number of levels for different factors as follows :

- For **company_location** keep only the top 10 locations and combine all other locations into “Other”
- For **country_of_bean_origin** keep only the top 25 countries and combine all other into “Other”
- For **ingredients** keep the top 5 ingredients and combine all other into “Other”
- For **company_manufacturer**, keep the top 50 manufacturers and combine all other into “Other” (50 was chosen because this categorical feature has too many levels)

Converting Data Types: Now we can convert our character columns into factors and also convert **cocoa_percent** column into a numeric column. The first few rows of our final, processed training data set is shown in Table 1.

Table 1: Final Features and Target in the Chocolate Dataset

Company Location	Review Date	Country of Bean Origin	Amount of Cocoa (%)	Number of Ingredients	Ingredients Present	Most Memorable Characteristics	Rating (1-5)	Company Manufacturer
Other	U.S.A.	2019 Dominican Republic	0.76	3	B,S,C	cocoa, vegetal, savory	3.50	
Other	U.S.A.	2019 Madagascar	0.76	3	B,S,C	cocoa, blackberry, full body	3.75	
Other	U.S.A.	2021 Other	0.68	3	B,S,C	chewy, off, rubbery	3.00	
Other	U.S.A.	2021 Other	0.80	3	B,S,C	mildly bitter, basic cocoa, fatty	3.25	
Other	U.S.A.	2021 India	0.68	3	B,S,C	milk brownie, macadamia,chewy	3.50	
A. Morin	France	2012 Bolivia	0.70	4	B,S,C,L	vegetal, nutty	3.50	

Data Distributions

Now let us examine the distributions for each of our numerical and categorical features.

Numerical and Discrete Features The only numerical feature we have is **percent_cocoa**. The **num_of_ingredients** and **review_date** features are discrete, and our target **rating** column is also discrete, as it has values between 1 and 5 in 0.25 intervals. Figure 1 shows the distributions for these features.

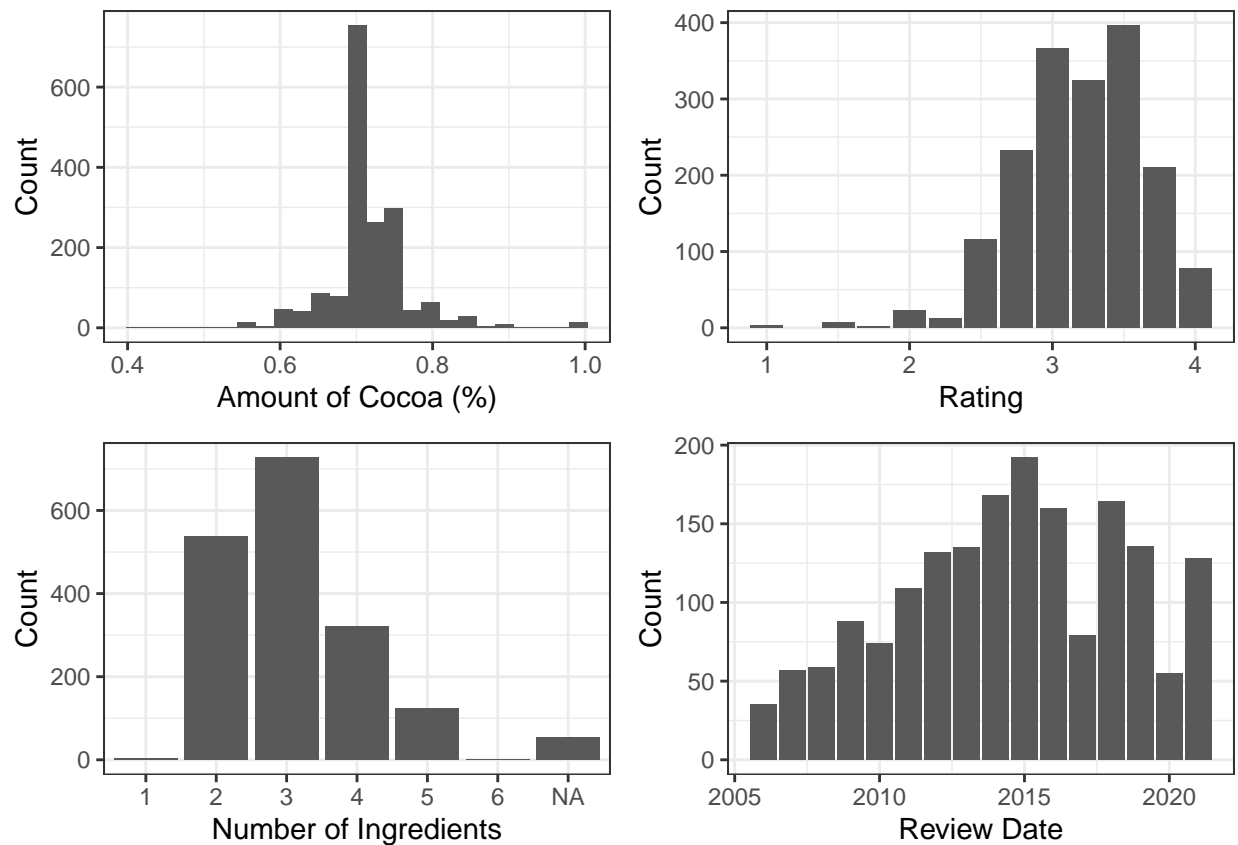


Figure 1, Distributions for numerical and discrete features in the training data set.

Categorical Features The `company_manufacturer`, `company_location`, `country_of_bean_origin`, and `ingredients` features are all categorical features. Figure 2 shows these feature distributions.

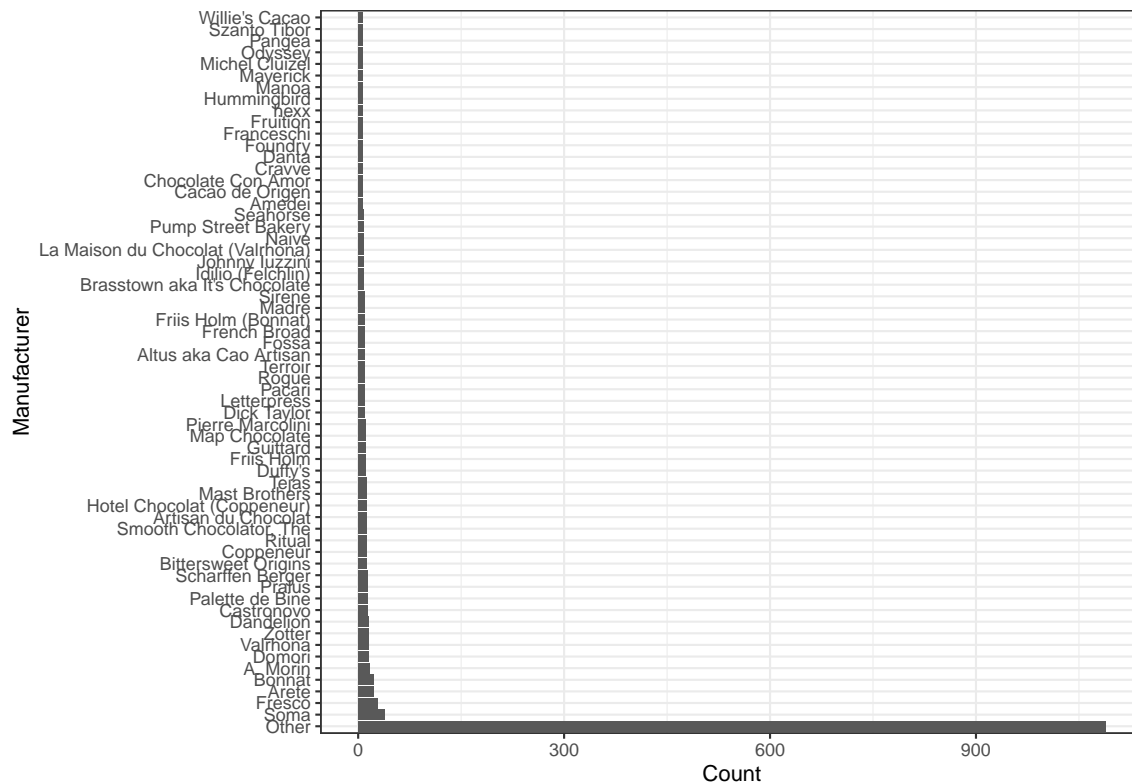
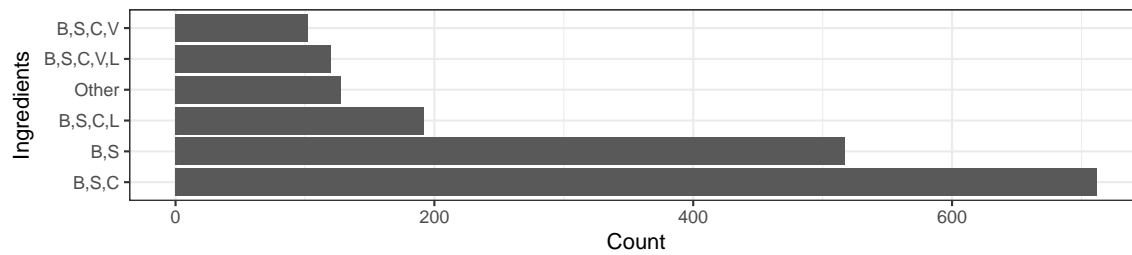
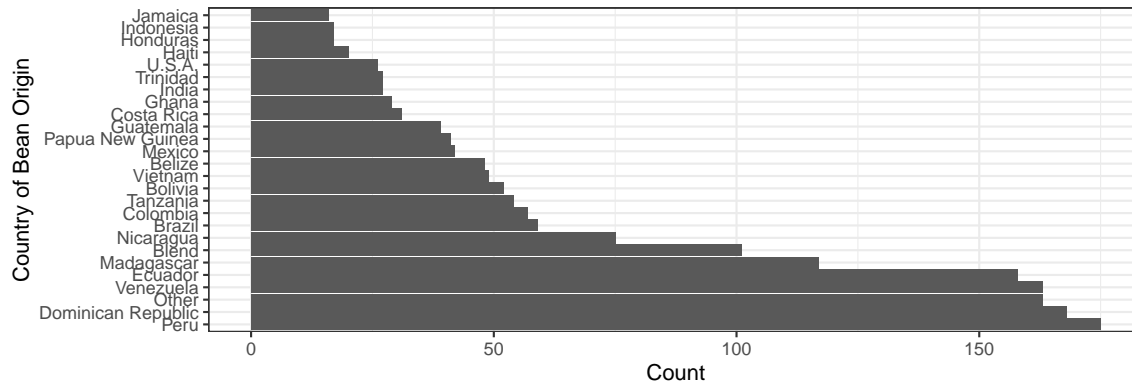
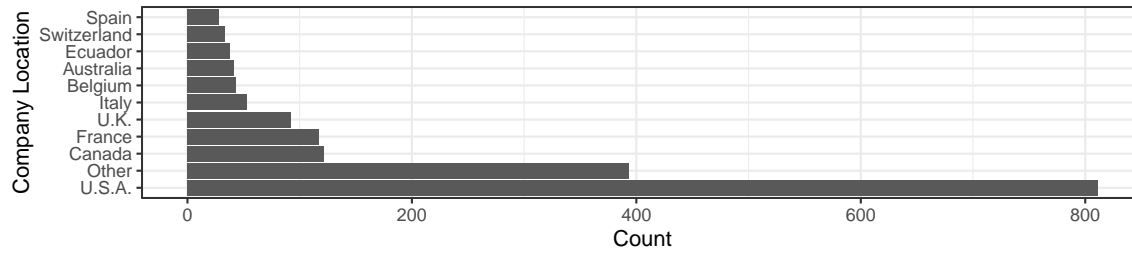


Figure 2, Distributions for categorical features in the training data set.

Based on the last plot above, it seems that there many distinct companies that manufacture chocolate in this data set, such that this feature acts more like an identifier. Therefore we can choose to drop this featurecolumn, as the values are too unique, and we would have an overwhelming Other category even if we considered the top 50 companies.

So our final dataset has the following 7 features with “rating” as our target:

Feature	Type
company_location	Factor
review_date	Numeric - Continuous
country_of_bean_origin	Factor
cocoa_percent	Numeric - Continuous
num_of_ingredients	Numeric - Discrete
ingredients	Factor
most_memorable_characteristics	Character(Text)