

Notes On Data

Derek R Neilson

September 26, 2024

Abstract

This document contains notes on the data. The notes are intended to demonstrate how I filter and manipulate the data. The primary focus is on cleaning the data, handling missing values, and transforming data to suit the analytical objectives. These notes serve as a comprehensive guide for understanding my data processing workflow.

Chapter 1

Data Collection and Inspection

1.1 Data Collection

1.1.1 Introduction

In data analysis, the ability to effectively filter and manipulate data is crucial for extracting meaningful insights. This document outlines the methodologies and tools I employ to pre-process and analyze the dataset. The primary focus is on cleaning the data, handling missing values, and transforming data to suit the analytical objectives. These notes serve as a comprehensive guide for understanding my data processing workflow.

1.1.2 Data Collection

The data was collected from <https://fdc.nal.usda.gov/>. The dataset is 2.9GB and is labeled **Branded** and is in JSON format. I chose this dataset because it is large and one can assume that it has the most rows because it is so large.

To download the data, I used the following commands:

```
1 wget https://fdc.nal.usda.gov/fdc-datasets/FoodData_Central_branded_food_json_2024-04-18.zip
2 unzip FoodData_Central_branded_food_json_2024-04-18.zip
3 rm FoodData.Central.branded_food_json_2024-04-18.zip
```

Listing 1.1: Download, Extract, and Remove Zip File

As shown in Listing 1.1, the commands download, extract, and remove the dataset file.

It is worth noting that I am using git to track changes in the code and data. The git commands will not be shown in this document for brevity.

1.2 Data Inspection

I received the following files after extracting:

- `brandedDownload.json` - I am assuming that this is the main file
- `foundationDownload.json` - I am assuming that this is a supporting file

The first step to inspecting the data is to view it.

```
1 less brandedDownload.json # the output is too large to show here and is not useful
2 # I am going to use jq to view the data
3 jq . brandedDownload.json # this results in a segmentation fault because the file is too large
4 # I am going to use a stronger server to view the data
5 # For security reasons, the IP address and username are redacted
6 sftp -P port username@ip_address
7 put brandedDownload.json DEV/Project/Data
8 put foundationDownload.json DEV/Project/Data
9 bye
10 ssh username@ip_address -P port
```

Listing 1.2: View the Data

As shown in Listing 1.2, the file is too large to view on my local machine. I will use a stronger server to view the data. Note that there is an assumption that all commands that follow are run on the server. From here on, I will refer to the server as being the machine that I am using to view the data. To get the data on the server, I used 'sftp' to transfer the files to the server.

```
1 jq . brandedDownload.json | less # failed because the file is too large
2 jq --stream . brandedDownload.json | less # this works because it streams the data
3 jq --stream . foundationDownload.json | less
```

Listing 1.3: View the Data on the Server

After looking at the head of the data, I can see that the data is in unstructured JSON format. I will use the CSV data instead. The JSON data will not be included in this document for size reasons.

1.3 Data Collection (CSV)

```
1 rm brandedDownload.json foundationDownload.json # remove the JSON files
2 wget https://fdc.nal.usda.gov/fdc-datasets/FoodData_Central_branded_food.csv_2024-04-18.zip #
  ↳ download the CSV file
3 unzip FoodData_Central_branded_food.csv_2024-04-18.zip # unzip the file
4 mv FoodData_Central_branded_food.csv_2024-04-18/* . # move the files to the current directory
```

Listing 1.4: Download, Extract, and Remove Zip File (CSV)

As shown in Listing 1.4, the commands download, extract, and remove the dataset file in CSV format. I will use the CSV data for the rest of the analysis.

```
1 import requests
2 import dotenv
3 import pandas as pd
4 import os
5 from scipy import stats
6 import numpy as np
7
8 # Load environment variables
9 dotenv.load_dotenv()
10
11 # Path to CSV
12 csv_file_path = "branded_food.csv"
13
14 # Check file existence
15 if not os.path.isfile(csv_file_path):
16     raise FileNotFoundError(f"The file {csv_file_path} does not exist.")
17
18 # Load data
19 try:
20     data = pd.read_csv(csv_file_path)
21     print("Data loaded successfully.")
22 except Exception as e:
23     print(f"An error occurred while loading the data: {e}")
24     raise
25
26 print(f"Data shape: {data.shape}")
27
28 # 1. Inspect Data Types
29 print("\n--- Data Types ---")
30 print(data.dtypes)
31
32 # Convert specific columns if necessary
33 # Example: data['price'] = pd.to_numeric(data['price'], errors='coerce')
```

```

34
35 # 2. Check for Missing Values
36 print("\n--- Missing Values ---")
37 missing_values = data.isnull().sum()
38 print(missing_values)
39 missing_percentage = (missing_values / len(data)) * 100
40 print("\n--- Percentage of Missing Values ---")
41 print(missing_percentage)
42
43 # Handle missing values (Example: Fill numerical with mean, categorical with mode)
44 numerical_cols = data.select_dtypes(include=[np.number]).columns
45 categorical_cols = data.select_dtypes(include=["object", "category"]).columns
46
47 for col in numerical_cols:
48     if data[col].isnull().sum() > 0:
49         data[col].fillna(data[col].mean(), inplace=True)
50
51 for col in categorical_cols:
52     if data[col].isnull().sum() > 0:
53         data[col].fillna(data[col].mode()[0], inplace=True)
54
55 # 3. Identify Duplicates
56 print("\n--- Duplicates ---")
57 duplicate_rows = data.duplicated().sum()
58 print(f"Number of duplicate rows: {duplicate_rows}")
59
60 if duplicate_rows > 0:
61     data = data.drop_duplicates()
62     print(f"Data shape after removing duplicates: {data.shape}")
63
64 # 4. Summary Statistics
65 print("\n--- Summary Statistics ---")
66 print(data.describe())
67
68 # 5. Detect Outliers using Z-Score
69 print("\n--- Detecting Outliers with Z-Score ---")
70 z_scores = np.abs(stats.zscore(data.select_dtypes(include=[np.number])))
71 threshold = 3
72 outliers = (z_scores > threshold).any(axis=1)
73 print(f"Number of outliers: {outliers.sum()}")
74
75 # Optionally remove outliers
76 data = data[~outliers]
77 print(f"Data shape after removing outliers: {data.shape}")
78
79 # 6. Validate Categorical Data
80 print("\n--- Categorical Data Validation ---")
81 for col in categorical_cols:
82     unique_vals = data[col].unique()
83     print(f"\nUnique values in '{col}':\n", unique_vals)
84     # Example: Standardize to lowercase
85     data[col] = data[col].str.lower()
86     # Example: Replace known inconsistencies
87     # data[col] = data[col].replace({'old_value': 'new_value'})
88
89 # 7. Ensure Data Consistency and Integrity
90 print("\n--- Data Consistency Checks ---")
91 # Example: Date consistency
92 if "manufacture_date" in data.columns and "expiry_date" in data.columns:
93     data["manufacture_date"] = pd.to_datetime(data["manufacture_date"], errors="coerce")
94     data["expiry_date"] = pd.to_datetime(data["expiry_date"], errors="coerce")
95     invalid_dates = data[data["expiry_date"] < data["manufacture_date"]]
96     print(f"Records with expiry_date before manufacture_date: {invalid_dates.shape[0]}")

```

```

97     # Remove invalid dates
98     data = data[data["expiry_date"] >= data["manufacture_date"]]
99
100 # Example: Logical numerical relationships
101 if "calories" in data.columns:
102     negative_calories = data[data["calories"] < 0].shape[0]
103     print(f"Records with negative calories: {negative_calories}")
104     # Remove negative calories
105     data = data[data["calories"] >= 0]
106
107 print("\n--- Final Data Shape ---")
108 print(data.shape)

```

Listing 1.5: Validate the Data

The code in Listing 1.5 performs the following tasks:

- Load the data from the CSV file
- Calculate the percentage of missing values per column
- Identify columns with more than 90% missing values
- Drop columns with more than 90% missing values
- Save the cleaned data to a new CSV file

The cleaned data is saved to a new CSV file called `branded_food_cleaned.csv`. The next step is to analyze the data. The reasons for dropping columns with more than 90% missing values is that they are not useful for this type of analysis and add unnecessary noise to the data.

1.4 Data Collection (API - Caloric Information)

I will use the FatSecret API to get the caloric information for each food. I will need to check the API documentation to see how to use it and what the rate limits are. I will also need to standardize the units of measurement. I will test the API with a few foods before running it on the entire dataset.