

Seasonal Produce

November 12th, 2020

Tigist Keneni & Fang Yu Chang

Overview

The project will aggregate data from different sources to provide users an overview of the cost of seasonal produce in the southwest USA and highest-producing countries for seasonal produce.

Sources

1. USDA Department of Agricultural <https://snaped.fns.usda.gov/seasonal-produce-guide/>
2. USDA Agricultural Marketing Service <https://www.seasonalfoodguide.org/state/california>
3. Wikipedia <https://www.wikipedia.org>

Documentation

- Datasets used and their sources.
- Types of data wrangling performed - Data cleaning, joining, filtering, and aggregating.
- The schemata used in the final production database.

Technical Report

* Extract: your original data sources and how the data was formatted (CSV, JSON, pgAdmin 4, etc).

- Our data sources were from a CSV file that was downloaded from the USDA Agricultural Marketing Service website and Wikipedia and USDA Department of Agricultural websites.
- USDA Department of Agricultural: The five produce lists were scraped from this website in order to get the fall seasonal produce lists. Reference file: "fall_produce_scraper.ipynb".
- USDA Agricultural Marketing Service: A CSV file for price in the southwest US was downloaded and cleaned and transformed using Jupyter Notebook. Reference file: "main_script.ipynb".
- After transforming, the data were imported into pgAdmin.

* Transform: required steps for cleaning and transformation of the data

- For Wikipedia data, getting the appropriate columns, renaming the columns, and adding the column in the apple_country table were performed.
- Data cleaning was done during this project by selecting specific produce that is popular in the southwest USA. More data cleaning was done when we selected data from Wikipedia.

3

- Inner data joining with the country from Wikipedia and countries were collected from USDA Marketing Services. By joining these sets of data we noticed interesting trends that helped us identify where most of the popular produce was grown at.
- We filtered data that was most important to this project and selected specific produce information to identify which countries had the most production of apples.

* Load: the final database, tables/collections, and why this was chosen.

The PgAdmin interface shows the schema.sql file with the following SQL code:

```

5 CREATE TABLE apple_country (
6   commodity TEXT,
7   country TEXT,
8   date TEXT,
9   region TEXT,
10  variety TEXT,
11  unit TEXT,
12  weighted_avg_price double precision
13 );

```

The data output table shows the following data:

index	commodity	country	date	region	variety	unit	weighted_avg_price
1	APPLE	China	11/6/2...	SOUTHW...	FUJI	3 lb bag	3.5
2	APPLE	United States	11/6/2...	SOUTHW...	FUJI	5 lb bag	3.55
3	APPLE	Poland	11/6/2...	SOUTHW...	FUJI	per pou...	1.22
4	BANANAS	India	11/6/2...	SOUTHW...	FUJI	2 lb bag	2.5
5	BANANAS	China	11/6/2...	SOUTHW...	FUJI	3 lb bag	3.99
6	BANANAS	Philippines	11/6/2...	SOUTHW...	FUJI	per pou...	1.96
7	BANANAS	Philippines	11/6/2...	SOUTHW...	GALA	3 lb bag	3.5
8	BANANAS	Philippines	11/6/2...	SOUTHW...	GALA	5 lb bag	3.55
9	BANANAS	Philippines	11/6/2...	SOUTHW...	GALA	per pou...	1.23
10	BANANAS	Philippines	11/6/2...	SOUTHW...	GALA	2 lb bag	2.62
11	BANANAS	Philippines	11/6/2...	SOUTHW...	GALA	3 lb bag	3.99
12	BANANAS	Philippines	11/6/2...	SOUTHW...	GALA	per pou...	1.52
13	BANANAS	Philippines	11/6/2...	SOUTHW...	GOLDEN D...	per pou...	0.79
14	BANANAS	Philippines	11/6/2...	SOUTHW...	GRANNY S...	3 lb bag	3.5
15	BANANAS	Philippines	11/6/2...	SOUTHW...	GRANNY S...	5 lb bag	3.55
16	BANANAS	Philippines	11/6/2...	SOUTHW...	GRANNY S...	per pou...	1.11
17	BANANAS	Philippines	11/6/2...	SOUTHW...	GRANNY S...	2 lb bag	2.62
18	BANANAS	Philippines	11/6/2...	SOUTHW...	GRANNY S...	per pou...	1.58
19	BANANAS	Philippines	11/6/2...	SOUTHW...	HONEYCRL...	per pou...	1.67
20	BANANAS	Philippines	11/6/2...	SOUTHW...	HONEYCRL...	per pou...	3.25
21	BANANAS	Philippines	11/6/2...	SOUTHW...	JONAGOLD	per pou...	1.52
22	BANANAS	Philippines	11/6/2...	SOUTHW...	RED DELIC...	per pou...	0.64
23	BANANAS	Philippines	11/6/2...	SOUTHW...	[null]	per pou...	0.53

The PgAdmin interface shows the schema.sql file with the following SQL code:

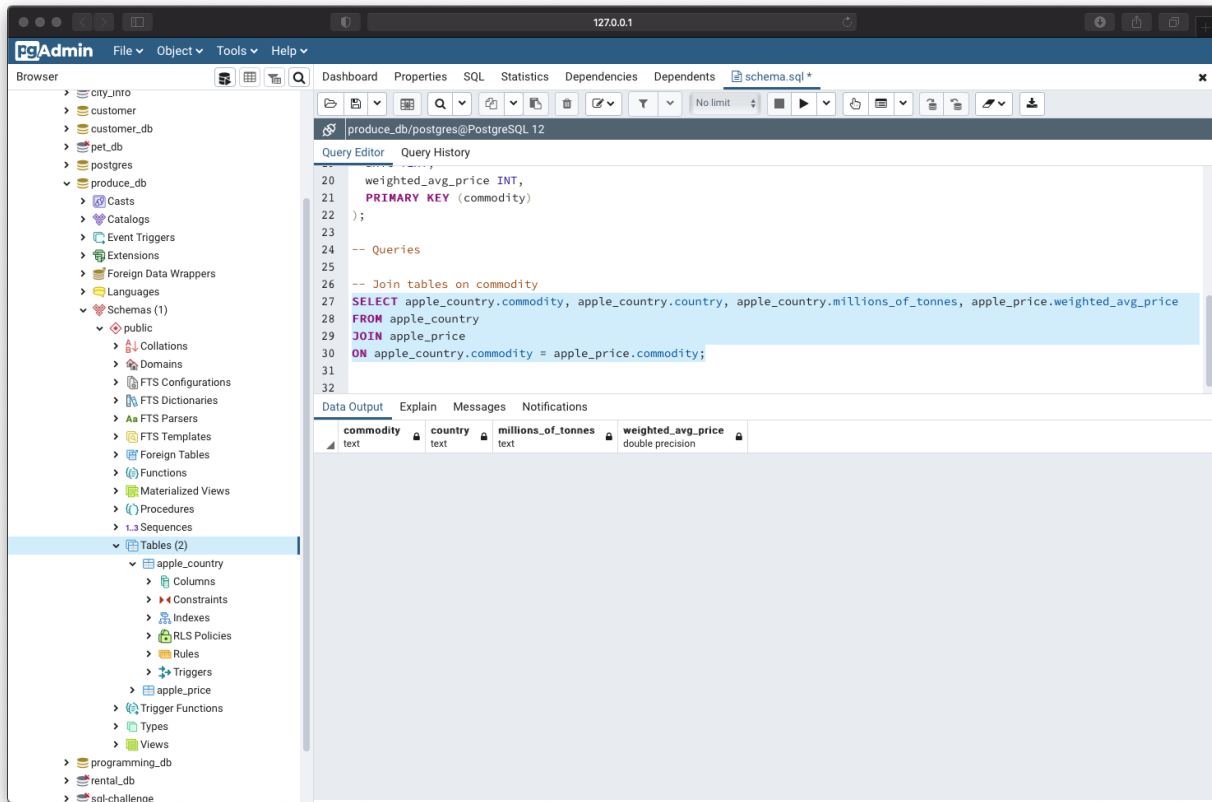
```

1 -- Create tables
2
3 SELECT * FROM apple_country
4
5 CREATE TABLE apple_price (
6   commodity TEXT,
7   country TEXT,
8   millions_of_tonnes INT,
9   PRIMARY KEY (commodity)
10 );
11
12 SELECT * FROM apple_price
13
14 CREATE TABLE apple_price (

```

The data output table shows the following data:

index	commodity	country	millions_of_tonnes
1	APPLE	China	39.2
2	APPLE	United States	4.7
3	APPLE	Poland	4.0
4	BANANAS	India	30.5
5	BANANAS	China	11.2
6	BANANAS	Philippines	6.1



Summary

- These final schemas identify relevant information that gave us insight into the question we started with, how much do popular seasonal fruits cost?

- From this data analysis, we were able to discover some of the most popular fruits and decided to focus on aggregating data from different sources. After looking at the data, we had 5 fruits to work with. We then decided to aggregate relevant data for one specific fruit which was apples.
- Some of our takeaways from working on this project was to ask a specific question and to keep iterating on the project scope as we discovered more information from the data that was available.