

ETL Report

Group 6

8/10/22

The primary focus of the data that was collected is about predicting sales information for retailers across the United States. Predicting retail sales would allow for a trickle down effect for numerous markets to more accurately predict inventory levels and plan accordingly across the timespan. We collected four separate datasets for this endeavor. The primary dataset was retrieved from the US census Bureau (Bureau, *All Sectors: Products by Industry for the U.S.: 2017* 2017) involving statistics that fall under the 44-45 retail sales code. This dataset sets the basis for how the retail industry is broken down by NAPCS and subcategorized by NAICS codes. The second dataset retrieved is an excel workbook file (Bureau, *Monthly Retail Trade Report* 2022) with numerous sheets detailing retail sales data broken down by NAICS code. With sales data spanning from the 1990's and broad range retail sales data we then chose to retrieve broad concept tables that relate to United states sales. The first data set is a macroeconomics (Dchaen, *Macroeconomics us* 2022) dataset also broken down by year and month with over 120 columns outlining macroeconomic principles across Census Bureau classifications. Our last dataset retrieved is structured similarly to the macroeconomics dataset in that it is indexed by year and month. This last data set holds consumer price index information (Avigan, *Consumer price index (CPI)* 2020). This will hereafter be referred to as CPI. With the data in place largely broken down by month we used a database composed of all the compiled data to predict future retail sales by different categorizations. Before anything could be composed all datasets needed to undergo ETL processes to be in relatable usable form.

The ET for the excel workbook files was done using python code and pandas in Visual Studio Code, the L was done using pyspark within Azure Databricks.

1. Excel workbook ETL (Bureau, *Monthly Retail Trade Report 2022*)

1. Import pandas
2. Use pandas and the read_excel method to create a dictionary of dataframes that are representative of the sheets in the excel workbook file. Ignore the top 4 rows of extraneous data.
3. Ignore the 2022 data since it is incomplete and formatted differently.
4. Drop 'Total' column since it is a calculated value
5. Drop row with adjusted [0]
6. Rename 'Unnamed: 0' and 'Unnamed: 1' as NAPCSCode and NAPCSName respectively
7. Locate all occurrences of 722 in the NAPCSCode column, write that index value to a list, select the second occurrence, and keep everything above. (not inclusive)
8. Create the adjusted sales dataframe by finding 'ADJUSTED(2)' index in the NAPCSName column, writing it to a list and selecting everything after +7 rows to remove calculated totals.
9. Drop sum rows of index 75 and 77 from the adjusted dataframe. (adjusted data frame should be complete)
10. Use the same index selector mentioned two steps previous in creating the adjusted dataframe selecting everything after -6 rows removing unwanted values.
11. Select everything after the first 7 rows in the non adjusted dataframe.
12. Drop index rows 9 and 14 to remove calculated values.(non adjusted dataframe should be complete.
13. Write the two dataframes to new respective dictionaries specifying the naming convention desired for the keys.
14. Use .melt to turn the columns, which are month and year combinations, into one column of months and years
15. Split the month and year into individual columns
16. Make sure all column data types are correct
17. Join the adjusted sales dataframe with the unadjusted sales dataframe on month and year
18. Write the superframe to csv ("Salestable_final.csv) to be used later

This data will be loaded in at the end of the document

1. Time-Series ETL

2. Using the created dictionaries of dataframes transpose each data frame with the transpose method.
3. Use the NAPCSname as the new column header
4. Drop the NAPCSname and NAPCSCode rows
5. The new index made up of the date was split by month and year. The string representation of the month was converted to a numeric integer representation and both year and month were made into new columns while keeping the index as the combination of the two.
6. These new data frames were then written to a new dictionary.
7. Combine all data frames from their respective dictionaries in a single data frame in order by date.

8. A new index was generated for each row

The other three dataset ETL processes were done using python and pyspark in the Microsoft Azure Gen10 databrick. Multiple files were manufactured to process the data sets and load the numerous information into the database. The ETL process is broken down by file in the databrick.

2. Producer ETL in producer

1. Read the census dataset into a dataframe using pyspark.
2. Replace all string values in columns with Null values.
 - a. For the sales column replace D and A with Null
 - b. For the number of establishments column replace D and S with Null
3. Remove the (s) from the end of numbers in the sales column, remove the commas from the "Sales, value of shipments, or revenue of NAPCS collection code (\$1,000) (NAPCSDOL)" and "Number of establishments (ESTAB)".
 - a. This will create new columns that hold the original data but transformed. Name these "NEWSales" and "NEWEstab" respectively.
4. Drop the old columns.
 - a. "Sales, value of shipments, or revenue of NAPCS collection code (\$1,000) (NAPCSDOL)" and "Number of establishments (ESTAB)"
5. Change the Sales and Establishment columns to ints.
6. Rename the new columns to the correct names.
 - a. 'NEWSales', 'Sales, value of shipments, or revenue of NAPCS collection code (\$1,000) (NAPCSDOL)'
 - b. 'NEWEstab', 'Number of establishments (ESTAB)'

3. Census ETL (Consumer)

1. Read the census dataset into a dataframe using pyspark
2. Create census data frame excluding extraneous columns (census_df).
 - a. Chosen columns: ('2017 NAPCS collection code (NAPCS2017)', '2017 NAICS code (NAICS2017)', 'Number of establishments (ESTAB)', 'Sales, value of shipments, or revenue of NAPCS collection code (\$1,000) (NAPCSDOL)')
3. Create a specific data frame to house NAPCS name and code (napcs_df).
 - a. '2017 NAPCS collection code (NAPCS2017)', 'Meaning of NAPCS collection code (NAPCS2017_LABEL)'
 - b. Use .distinct() to prevent duplicates
4. Create a specific data frame to house NAICS name and code (naics_df).
 - a. '2017 NAICS code (NAICS2017)', 'Meaning of NAICS code (NAICS2017_LABEL)'
 - b. Use .distinct() to prevent duplicates
5. Write the dataframes to individual csvs within a blob container within Azure Datalake.

The loading is all done within a single databrick

4. Creating a Database Connection

1. To connect to the SQL database you need to define:
 - a. Database (the name of the database)
 - b. Tables (each individual table you want to read or write to needs to be defined)
 - c. User and password (these are used to give you permission to access/alter the database)
 - d. Server (where the database is located)

5. Census ETL (Final Steps)

1. Read in the csvs to dataframes
 - a. census_df, naics_id_df, napcs_id_df
2. Write the napcs and naics data to their respective tables (NAPCSTables and NAICSTable) within the SQL database
3. Read in the data from those tables within the database
 - a. This provides the primary key that was generated when each table was populated
4. Join the census_df and naics_id_df on “2017 NAICS code (NAICS2017)”
5. Drop the “2017 NAICS code (NAICS2017)” column
 - a. At this point the census data should have: 'NAICS ID', '2017 NAPCS collection code (NAPCS2017)', 'Number of establishments (ESTAB)', and 'Sales, value of shipments, or revenue of NAPCS collection code (\$1,000) (NAPCSDOL)'
6. Join that dataframe with napcs_id_df on “2017 NAPCS collection code (NAPCS2017)”
7. Drop the “2017 NAPCS collection code (NAPCS2017)” column
 - a. The final census columns should now be: 'NAICSID', 'NAPCSID', 'Number of establishments (ESTAB)', and 'Sales, value of shipments, or revenue of NAPCS collection code (\$1,000) (NAPCSDOL)'
 - b. Label this “complete_census_df”
8. Populate the census table (NAICS_NAPCS) by writing to jdbc using “complete_census_df”

6. CPI ETL

1. Load in the CPI table .csv file as a dataframe
2. Split the “DATE” column into “Year” and “Month” columns
3. Rename “USACPIALLMINMEI” column to “CPI”
4. Drop date column that was used to populate month and year column
5. Drop the leading 0 from the month column

7. Macro ETL

1. Load in the CPI table .csv file as a dataframe
 - a. Only select these columns: 'observation_date', 'RPI', 'RETAIL', 'USTRADE', 'USWTRADE'
2. Rename column “observation_date” to “Date”
3. Split the “DATE” column into “Year” and “Month” columns
 - a. Ignore the day information from the “DATE” column, it is unnecessary
4. Drop the “DATE” column

8. Merge MACRO and CPI

1. Merge the Macro dataframe and the CPI dataframe on “Month” and “Year” (merged_and_ordered_df)

2. Change the type of “Month” and “Year” to int
3. Order the dataframe in ascending order by both “Year” and “Month”

9. DateTable

1. Create a dataframe (time_df) that only contains the “Year” and “Month” from the previous dataframe
2. Populate the DateTable by writing to jdbc using time_df
3. Read in the DateTable and label the data frame (dateid_df)

10. MacroTable

1. Join the “merged_and_ordered_df” dataframe with “dateid_df” on “Year” and “Month”
2. Use .na.drop(subset = ‘DateID’) to drop all rows with null DateIDs
3. Drop the “Year” and “Month” columns
4. Drop all columns with null values for everything but “DateID”
 - a. Call this dataframe “macro_table_df”
5. Populate the MacroTable by writing to jdbc using “macro_table_df”

11. Excel workbook ETL (Final Steps)

1. Load in the “Salestable_final.csv” into a dataframe called “sales_df
2. Create a dataframe titled append_NAICS_df that takes the NAICScode and NAICSname from “sales_df”
 - a. Use .distinct() to ensure there are no duplicates
 - b. Rename the columns '2017 NAICS code (NAICS2017)' and 'Meaning of NAICS code (NAICS2017_LABEL)' respectively
 - c. Append the NAICSTable by writing to jdbc and using “append_NAICS_df”
3. Drop the ‘_c0’ column from “sales_df”, it is unnecessary
4. Inner join sales_df with dateid_df on “Year and “Month”
 - a. Label this “cleaned_sales_df”
 - b. This provides the DateID
 - c. Drop “Year” and “Month”
5. Read in the updated NAICSTable as a data frame
6. Create a new data frame that joins cleaned_sales_df with the NAICSTable
 - a. This provides the NAICSID
7. Drop all NAICS related columns from cleaned_sales_df except for the NAICSID
8. Replace all string values within the “Adjusted Sales” and “Unadjusted Sales” columns with null values
9. Populate the SalesTable by writing to jdbc

The steps provided took an ample amount of time across the entirety of our team. By dividing the datasets into subtasks and splitting the work it was able to be completed in a timely manner. Through these processes we were able to produce a working database that is normalized with no repeated data and easily traversed structure.

Sources:

Bureau, U. S. C. (2017). *All Sectors: Products by Industry for the U.S.: 2017*. Explore census data. Retrieved August 3, 2022, from <https://data.census.gov/cedsci/table?q=ECNNAPCSPRD2017.EC1700NAPCSPRDIND&n=N0600.44&tid=ECNNAPCSPRD2017.EC1700NAPCSPRDIND&hidePreview=true>

Avigan, A. (2020, February 16). *Consumer price index (CPI)*. Kaggle. Retrieved August 4, 2022, from <https://www.kaggle.com/datasets/aavigan/consumer-price-index-usa-all-items>

Dchaen. (2022, May 11). *Macroeconomics us*. Kaggle. Retrieved August 4, 2022, from https://www.kaggle.com/datasets/denychaen/us-macro?select=US_MACRO110522.csv

Bureau, U. S. C. (2022, July 15). *Monthly Retail Trade Report*. United States Census Bureau. Retrieved August 10, 2022, from <https://www.census.gov/retail/index.html>