GitHub repository for files: https://github.com/knolin804/ETL_project
Project Team: Derek Hagman, Kristine Nolin, Carol Wittig
**ETL Steps:**

## EXTRACT

**FINDING DATA**

From Kaggle we downloaded the following datasets because they provided related book information and could provide data for analysis based on user ratings across three different rating groups—one a reader advisory (GoodRead) and two e-commerce sites: Amazon and Flipkart.

- "GoodReads database" (csv file) which includes over 13,000 books, title, author, isbn, avg reader rating, numbers of reviews, numbers of ratings, categories.

- Amazon () vs Flipkart (1400 entries) book prices data set(s), which contained the same set of variables: author, title, isbn, ratings count, stars (rating), and price.

## TRANSFORM

**DATA CLEANUP & ANALYSIS**

**Sources of Data:**

1) We reviewed the original csv files for their structure and columns but didn't do any do manipulation before bringing into pandas [CSV FILES IN FOLDER: amazon.csv, books.csv, flipkart.csv]

**Types of Transformation Needed:**

1) Each of us created a Jupyter notebook and imported our csv file into Panadas, creating a dataframe [JUPYTER NOTEBOOKS IN FOLDER: ETL_FINAL.ipynb (GoodReads), Final_Amazon,ipynb, and flipkart.ipynb]

2) Cleaning of the dataframes included:
   a. Removing duplicate entries (Flipkart) and also removed empty rows with missing data

   b. Formatted the text in the title column to be consistent across the three notebooks

   c. Converted rupee prices to American dollars (Flipkart and Amazon)

   d. Dropped some of the original columns that wouldn't be needed for further analysis

   e. Split the author name into first and last names. Jupyter notebooks show the details, but this was the most intensive work step, although the output doesn't reflect how long it took or how unnecessary it ended up being based on what was actually needed in the database.
      i. Additional step with dataframes needed for Flipkart because author order ws inconsistent and in Amazon, author formatting was inconsistent. GoodReads data was the cleanest, but still had challenges in that first name, middle initial, last name was not easily separated.

ii. In some instances the authors' names were split over multiple columns and these needed to be appropriately reordered. For Flipkart, Derek used the isbn value to confirm authors of the books.

iii. Saved cleaned csvs to files (amazon_df2.csv, books_cleaned.csv, and final_flipkar.csv)

## LOAD

**Final Production Database** – we selected PostGres and chose a relational database for our structure because each of our datasets were csv files in a table format with structured fields. Each contained related items in ISBN, title, and author. This type of structured and similar data, with multi-rows of similar data, and a predefined schema (column names in a table structure) lends itself to a relational database.

Setup the PostGres Database [SQL DATABASE IN FOLDER: project_queries.sql]

a) Established the PostGres database GoodReads.

FINAL TABLES USED IN PRODUCTION DATABASE [SQL TABLES IN FOLDER – project_queries.sql]

In Jupyter, created the engine to make the connection with PostGres

a) Uploaded the dataframes to create the tables in PostGres
b) In PostGres, changed the datatypes because import assigned datatypes to the columns that had to be changed before you could run queries across the tables.
   i. Tables show query for each dataset - fields included author, title, ISBN, rating.



*Figure 1: Flipkart Table*



*Figure 2: Amazon Table*



*Figure 3: GoodReads Table*

ii.    A sample query included joining tables to compare ratings across Flipkart and Amazon.

| | isbn10 character varying | flipkart_rating text | amazon_rating text |
|---|---|---|---|
| 1 | 1250127505 | 4.5 | 4.5 |
| 2 | 1250127556 | 4 | 5.0 |
| 3 | 1250127556 | 4 | 5.0 |
| 4 | 1250127556 | 4 | 5.0 |
| 5 | 1250127556 | 4 | 5.0 |
| 6 | 1260142655 | | 3.7 |
| 7 | 1405911662 | 4.4 | 3.6 |
| 8 | 1408711702 | 4.8 | 4.8 |
| 9 | 1408711702 | 4.8 | 4.8 |
| 10 | 1408711702 | 4.8 | 4.8 |
| 11 | 1408711702 | 4.8 | 4.8 |
| 12 | 1408711702 | 4.8 | 4.8 |
| 13 | 1408711702 | 4.8 | 4.8 |
| 14 | 1408711702 | 4.8 | 4.8 |
| 15 | 1408711702 | 4.8 | 4.8 |
| 16 | 1408711702 | 4.8 | 4.8 |
| 17 | 1408711702 | 4.8 | 4.8 |

Data Output   Explain   Messages   Notifications