# Extract, Transfer, Load
## *Pandas & Sqlalchemy*

## Background

We chose these data sets because the ability to join them will be straightforward given we have a primary key of zip code for the population density data table that can match the foreign key of zip code for the IRS tax return data. In addition, we contemplated what sort of analysis could be done based on the tables in this database and thought that interesting insights could be distilled relating to income, tax bracket, zip code, and the population density within a zip code.
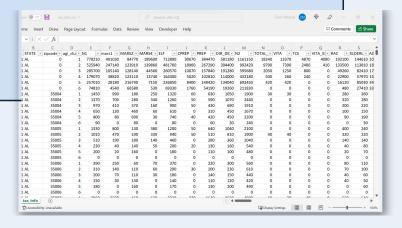
## Data Sources

**IRS Tax Returns:**
***Source***: https://www.irs.gov/e-file-providers/definition-of-adjusted-gross-income

**Population Density**
***Source:*** https://simplemaps.com/data/us-zips

## 01

## Extract

- **Imported:** Dependencies for Pandas and Sqlalchemy and CSV files into Pandas environment

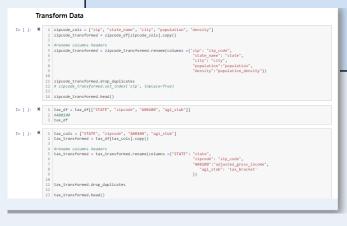- **Transformed:** read CSV files and converted them into data frame



## 02

## Transfer

**Created:** defined lists to hold values for data sources & return based on selected columns

**Identified:** found zip codes within the Tax Returns dataset not present within the Population Density dataset and add zip codes to the Population Density dataset

**Cleaned:** renamed columns and dropped duplicates and dropped index



## 03

## Load

- **Defined:** added a new database in SQL & used zip code in Population Density Primary Key linked it to zip code foreign key with IRS Tax Return

- **Created:** tables for Population Density and for IRS Tax Returns

- **Connected:** Connected SQL Database to python & insert data into Population Density and IRS Tax Returns tables using to_sql

- **Joined:** Join tables in SQL on zip codes & found Zip Codes unmapped and add them to the population density table