

# ETL Project

## Team Members:

Alexandra Taft  
Anthony English  
Arlette Varela  
Nathan Bolt

---

## Extract:Transform:Load - Automating Jobs Away

August 2020

## OVERVIEW/OUTLINE

The datasets we will be using for our project combines automation probability data with a breakdown of the number of jobs by salary in each occupation by state (within the US). For the purposes of our project we will filtering down to the detailed level of class titles and assigning them an occupational group and then region based on US census categorizations. We will be loading into PgAdmin as a SQL database.

---

## DATASET(S):

Kaggle Datasets: Occupation, Salary, and the Likelihood of Automation  
[https://www.kaggle.com/andrewmvd/occupation-salary-and-likelihood-of-automation?select=automation\\_data\\_by\\_state.csv](https://www.kaggle.com/andrewmvd/occupation-salary-and-likelihood-of-automation?select=automation_data_by_state.csv)

[https://www.kaggle.com/andrewmvd/occupation-salary-and-likelihood-of-automation?select=automation\\_data\\_by\\_state.csv](https://www.kaggle.com/andrewmvd/occupation-salary-and-likelihood-of-automation?select=automation_data_by_state.csv)



Occupation, Salary and Likelihood of Automation • kaggle.com

---

## TASK BREAKDOWN:

- ✓ Set up Github Repository @Arlette V
  - ✓ Each member create a branch
- ✓ Search for data @doc @
- ✓ Set up notebooks (Jupyter, VSC, etc.) @doc
- ✓ Create ERD
  - ✓ Export ERD code into PgAdmin
- ✓ Bring CSV files into Jupyter
  - ✓ Organize and clean dataframes, this will include:
    - ✓ automation
    - ✓ occupation\_broad
    - ✓ occupation\_detail
    - ✓ occupation\_major
    - ✓ occupation\_minor
- ✓ Set up PgAdmin
  - ✓ Create tables/schemas
  - ✓ Make sure tables are good in Jupyter