Course: Information Extraction, Retrieval and Integration

Unit 4: Data Integration

# Assignment Description:
# Data Integration, Bias and Fairness

Mari Carmen Suárez de Figueroa Baonza
March 2023
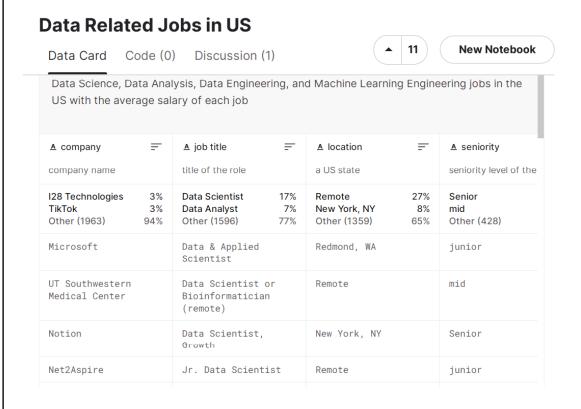Universidad Politécnica de Madrid

# General Issues

- This assignment is not the typical one.

  - We can say that it is a **testing-exploratory assignment**

- The objective of this assignment is **to identify the most important difficulties when integrating data**, taking into account **bias and fairness dimensions**.

  - Only a subsets of steps in the data integration process are going to be performed in this assignment

- **Final result** of the assignment is not an integrated dataset, but a set of lessons learned, difficulties, suggestions for improving the process, ways to proceed, among other experiences and knowledge gained.
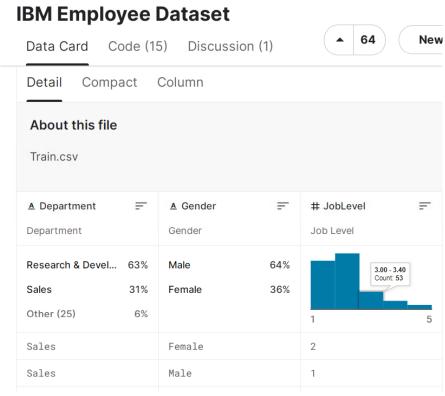
# Assignment Phases

## 1. Search for and select datasets

- **2 or more datasets** that contain overlapping data
- **Note**: Selected datasets should have some data related to attributes that can caused bias (such as gender, age, race, among others)

## 2. Identify conflicts among selected datasets

- **Data-level and schema-level** conflicts

## 3. Be aware about bias and fairness in your datasets

- Use any **tool for mitigating bias**
- Use any **tool for identifying fairness**

# 1. Search for and select datasets: Example

- **Dataset 1**: Data Related Jobs in US
  - https://www.kaggle.com/datasets/mohamedsiika/data-related-jobs-in-us
- **Dataset 2**: IBM Employee Dataset
  - https://www.kaggle.com/datasets/rohitsahoo/employee

# 1. Search for and select datasets: Dataset Repositories

- https://www.data.gov/
  - US-centric agriculture, climate, education, energy, finance, health, manufacturing data, etc.

- https://datos.gob.es/es/catalogo
  - Spanish datasets in different domains

- https://cloud.google.com/bigquery/public-data/
  - BigQuery (Google Cloud) public datasets (bikeshare, GitHub, Hacker News, Form 990 non-profits, NOAA, etc.)

- https://www.kaggle.com/datasets
  - Microsoft-owned, various (Billboard Top 100 lyrics, credit card fraud, crime in Chicago, global terrorism, world happiness, etc.)

- https://aws.amazon.com/public-datasets/
  - AWS-hosted, various (NASA, a bunch of genome stuff, Google Books n-grams, Multimedia Commons, etc.)

# 2. Identify conflicts among selected datasets: Example

- **Dataset 1**: Data Related Jobs in US
  - https://www.kaggle.com/datasets/mohamedsiika/data-related-jobs-in-us
    - job title
    - seniority
- **Dataset 2**: IBM Employee Dataset
  - https://www.kaggle.com/datasets/rohitsahoo/employee
    - Job Role
    - Job Level

# 3. Be aware about bias and fairness in your datasets: Example (Identifying Bias)

- **AI Fairness 360**
  - https://aif360.mybluemix.net/
  - https://aif360.mybluemix.net/resources#tutorials

# 3. Be aware about bias and fairness in your datasets: Example (Identifying Fairness)

- **Aequitas**

  - http://www.datasciencepublicpolicy.org/our-work/tools-guides/aequitas/

## Audit Results: Details by Fairness Measures

### Equal Parity: Failed

| What is it? | When does it matter? | Which groups failed the audit: |
|---|---|---|
| This criteria considers an attribute to have equal parity is every group is equally represented in the selected set. For example, if race (with possible values of white, black, other) has equal parity, it implies that all three races are equally represented (33% each)in the selected/intervention set. | If your desired outcome is to intervene equally on people from all races, then you care about this criteria. | **For race** (with reference group as **Caucasian**)<br>African-American with **2.55X** Disparity<br>Asian with **0.01X** Disparity<br>Other with **0.09X** Disparity<br>Native American with **0.01X** Disparity<br>Hispanic with **0.22X** Disparity<br><br>**For sex** (with reference group as **Male**)<br>Female with **0.22X** Disparity<br><br>**For age_cat** (with reference group as **25 - 45**)<br>Less than 25 with **0.52X** Disparity<br>Greater than 45 with **0.20X** Disparity |

# General Instructions

- This assignment should be performed **in groups** composed of 2/3 students.

  - Students can also decide to perform this assignment in an individual way

- **Deadline**: Wednesday 12$^{th}$ April 2023

- The **assignment delivery** should include a PDF file

  - describing the main outcomes for each assignment step: selected datasets, conflicts identified, and reports from tools for bias and fairness

  - the main decisions taken (with respect to tools for bias and fairness)

  - the lessons learned and the difficulties found

  - and any other comments and suggestions

Course: Information Extraction, Retrieval and Integration

Unit 4: Data Integration

# Assignment Description:
# Data Integration, Bias and Fairness

Mari Carmen Suárez de Figueroa Baonza
March 2023
Universidad Politécnica de Madrid