# DATA SCIENCE FOR ECONOMISTS

ECON 220 LAB

Jafet Baca-Obando

Week 9, Handling IPUMS Data (Part 2) – 10/24/2025

# Outline

# What is IPUMS?

## Integrated Public Use Microdata Series

Operated by the University of Minnesota

## Core mission: data harmonization

IPUMS takes datasets that were originally collected with different questions, codes, and variable names and makes them consistent. Lots of recoding!

## Free access!

By providing access to detailed, anonymized individual-level data (microdata), IPUMS allows researchers to ask complex questions that can't be answered with aggregated summary tables.

How to register and extract data? Check guide on Canvas.

# Importing required libraries

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import os
path = os.getcwd()
```
Python

# Load and explore the data

```python
data = pd.read_csv('usa_00004.csv')
```
Python

| | # YEAR | # SAMPLE | # SERIAL |
|---|---|---|---|
| count | 3405809.0 | 3405809.0 | 340580 |
| mean | 2023.0 | 202301.0 | 758991.73689217 |
| std | 0.0 | 0.0 | 441473.568225181 |
| min | 2023.0 | 202301.0 | |
| 25% | 2023.0 | 202301.0 | 37238 |
| 50% | 2023.0 | 202301.0 | 75683 |
| 75% | 2023.0 | 202301.0 | 114700 |
| max | 2023.0 | 202301.0 | 151901 |

8 rows x 20 cols   10 ⌄   per page  « ‹  Page  1  of 1  › »   🔍 ▦ ▨ ⋯
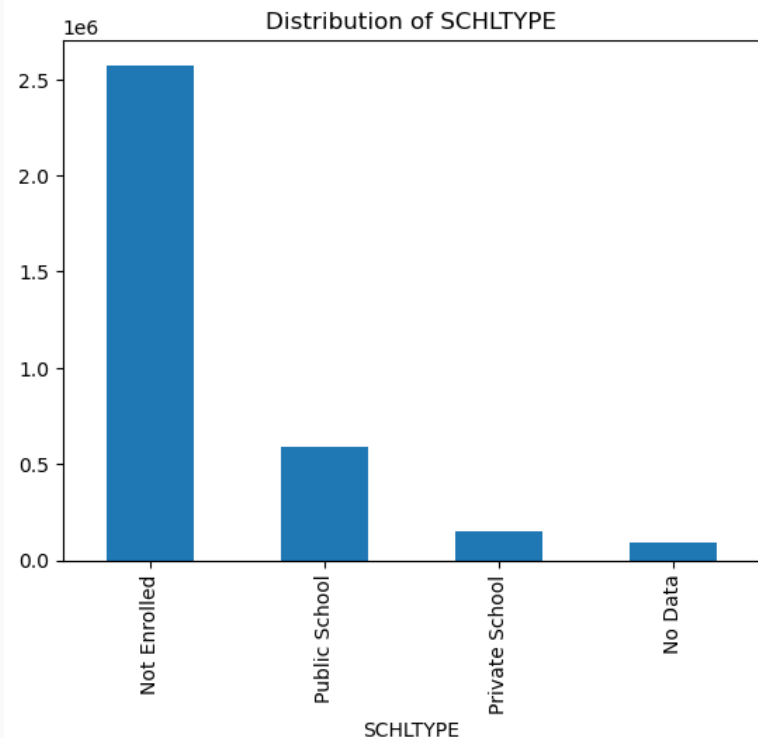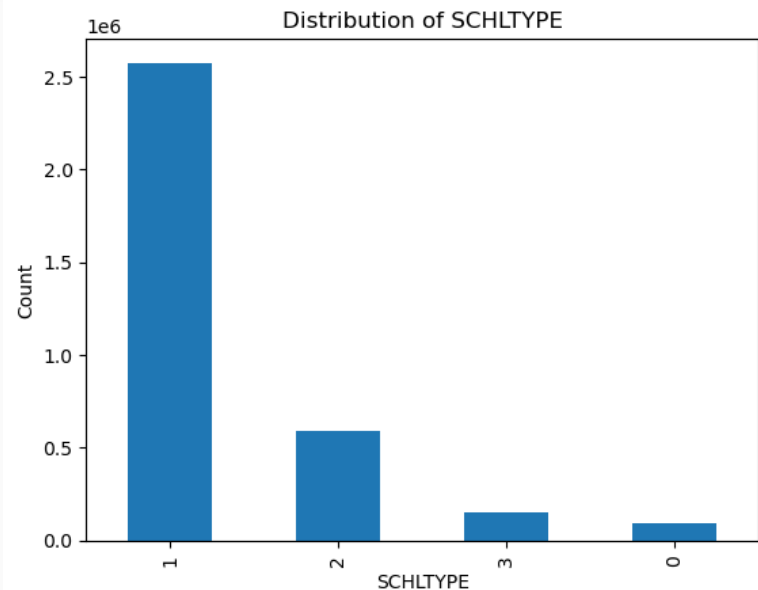
- Note the large number of observations!

# Recoding a categorical variable

- IPUMS uses numeric codes for categorical variables.
- **Example:** SCHLTYPE (School Type).
- **Need the Codebook!** Tells us: 0=N/A, 1=Not Enrolled, 2=Public, 3=Private.
- **Recoding**: replace numbers with meaningful labels.
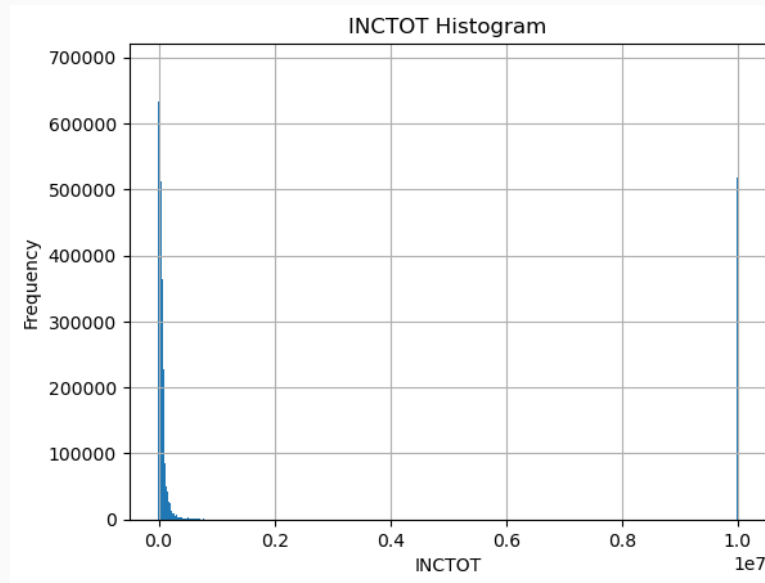
# Recoding a categorical variable

# Outlier detection and treatment

- **Outliers:** Extreme values that might distort analysis.
- Idea:
  - Use a histogram (.hist()) to see the distribution.
  - Observe asymmetry and potential extreme high/low values (including negative).
- Treatment:
  - Interquartile range (IQR) =Q3 (75th percentile) - Q1 (25th percentile)
  - **Define Bounds:** Lower = Q1 - 1.5*IQR, Upper = Q3 + 1.5*IQR.
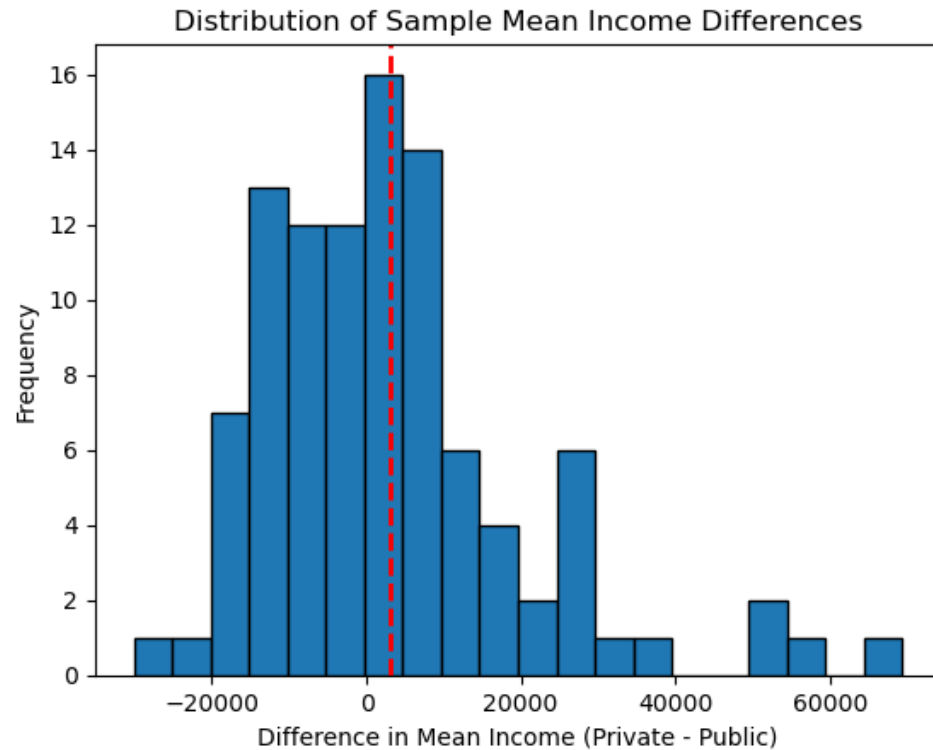  - **Filter:** Keep data only within these bounds.

# Outlier detection and treatment

- IQR might not catch all problematic values (e.g., 0 or negative income).
- **Domain knowledge**: Is zero/negative income plausible/useful for this analysis?
- **Filter**: Remove observations with INCTOT <= 0.

# Outlier detection and treatment

# Random sample



Distribution of Sample Mean Income Differences

# Recap

- **Check the Codebook**: It's essential for understanding IPUMS numeric codes.

- **Recode Variables**: Use .replace() to make your data readable (e.g., 2 → "Public School").

- **Find Outliers**: Use .hist() to visualize data and spot extreme values.

- **Sample for Speed**: Use .sample() for quick analysis on large datasets, but be aware of sampling variability.

# To-do list

- **Complete Data Exercise 6**

  - Upload Jupyter notebook (.ipynb file) and HTML file on **October 26**

- **Complete Data Exercise 7**

  - Upload Jupyter notebook (.ipynb file) and HTML file on **November 2**