# MIDTERM SKILLS EXAM
# DATA WRANGLING AND ANALYSIS

De Guzman, Jemuel Endrew C.
CPE22S3

# INTRODUCTION TO THE DATASET

# Census Income

Donated on 4/30/1996

Predict whether income exceeds $50K/yr based on census data. Also known as Adult dataset.

## Dataset Characteristics
Multivariate

## Subject Area
Social Science

## Associated Tasks
Classification

## Feature Type
Categorical, Integer

## # Instances
48842

## # Features
14

## Dataset Information

### Additional Information
Extraction was done by Barry Becker from the 1994 Census database. A set of reasonably clean records was extracted using the following conditions: ((AAGE>16) && (AGI>100) && (AFNLWGT>1)&& (HRSWK>0))
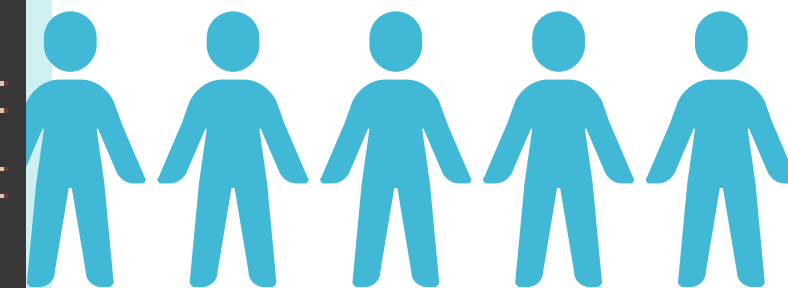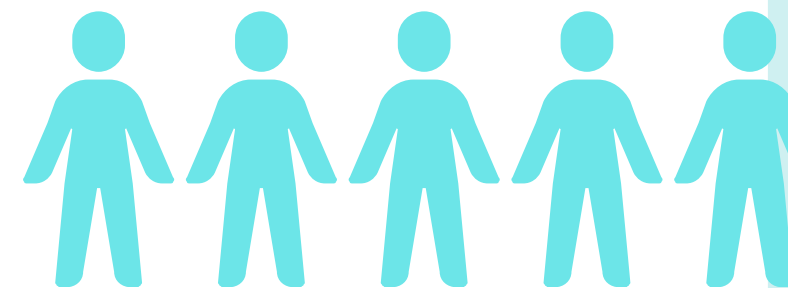
Prediction task is to determine whether a person makes over 50K a year.

SHOW LESS ^

### Has Missing Values?
Yes

```
 #   Column          Non-Null Count   Dtype
---  ------          --------------   -----
 0   age             48842 non-null   int64
 1   workclass       47879 non-null   object
 2   fnlwgt          48842 non-null   int64
 3   education       48842 non-null   object
 4   education-num   48842 non-null   int64
 5   marital-status  48842 non-null   object
 6   occupation      47876 non-null   object
 7   relationship    48842 non-null   object
 8   race            48842 non-null   object
 9   sex             48842 non-null   object
 10  capital-gain    48842 non-null   int64
 11  capital-loss    48842 non-null   int64
 12  hours-per-week  48842 non-null   int64
 13  native-country  48568 non-null   object
 14  income          48842 non-null   object
dtypes: int64(6), object(9)
```

# DATA WRANGLING

```
1 # Replacing all NaN and ? values with 'Other'
2 cols = ['workclass','occupation','native-country']
3
4 for col in cols:
5     xy[col] = xy[col].fillna('Other')
6     xy[col].replace('?','Other',inplace=True)
7
8 xy[xy.values == 'Other']
9
10 #print(xy[xy.isna().any(axis=1)]) # Shows that there are no more NaN values
11 #print('\n')
12 #print(xy.info())
```

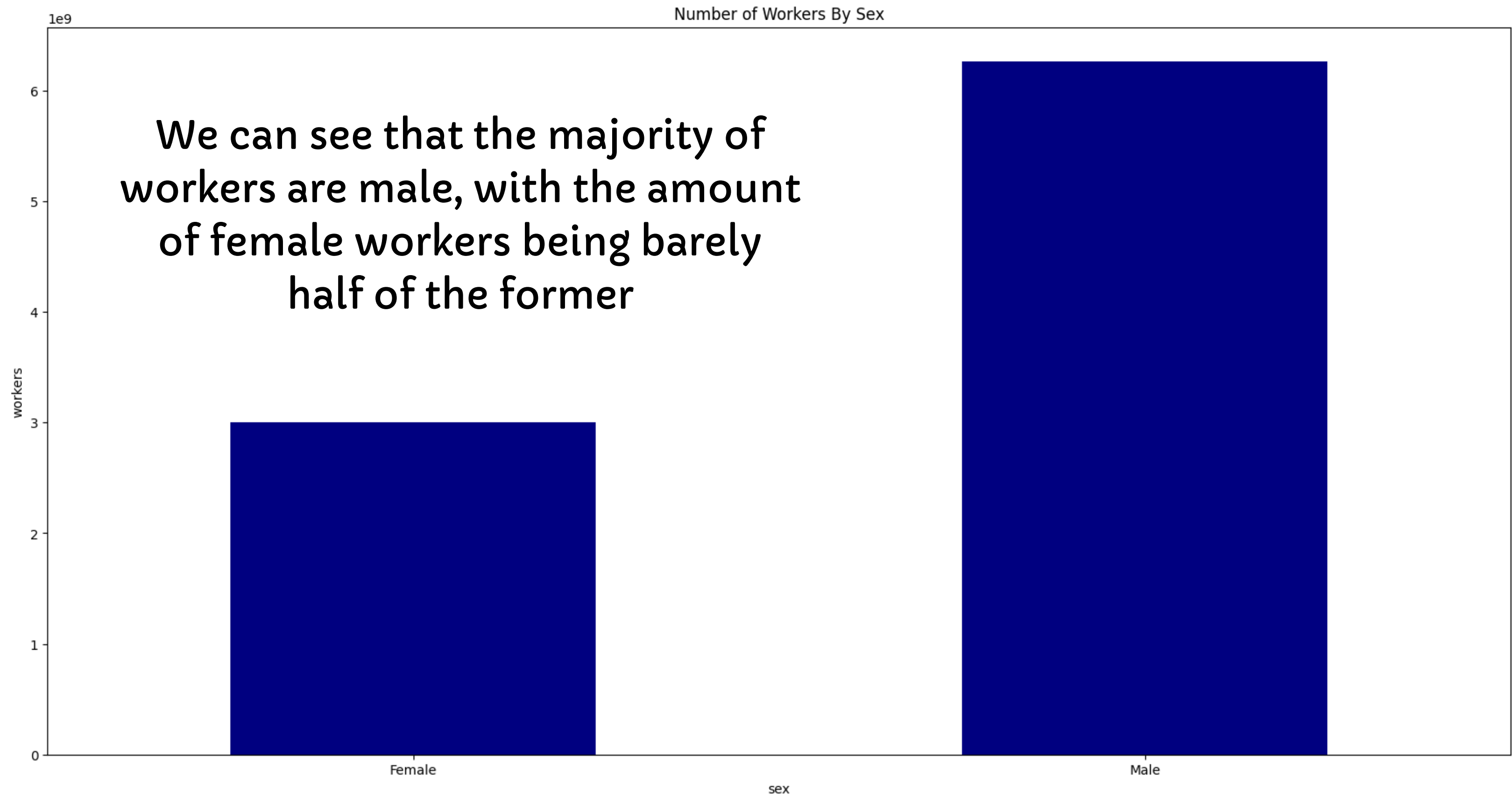| | age | workclass | fnlwgt | education | education-num | marital-status | occupation | relationship | race | sex | capital-gain | capital-loss | hours-per-week | native-country | income |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 14 | 40 | Private | 121772 | Assoc-voc | 11 | Married-civ-spouse | Craft-repair | Husband | Asian-Pac-Islander | Male | 0 | 0 | 40 | Other | >50K |
| 27 | 54 | Other | 180211 | Some-college | 10 | Married-civ-spouse | Other | Husband | Asian-Pac-Islander | Male | 0 | 0 | 60 | South | >50K |
| 27 | 54 | Other | 180211 | Some-college | 10 | Married-civ-spouse | Other | Husband | Asian-Pac-Islander | Male | 0 | 0 | 60 | South | >50K |
| 38 | 31 | Private | 84154 | Some-college | 10 | Married-civ-spouse | Sales | Husband | White | Male | 0 | 0 | 38 | Other | >50K |
| 50 | 25 | Private | 32275 | Some-college | 10 | Married-civ-spouse | Exec-managerial | Wife | Other | Female | 0 | 0 | 40 | United-States | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 48812 | 81 | Other | 26711 | Assoc-voc | 11 | Married-civ-spouse | Other | Husband | White | Male | 2936 | 0 | 20 | United-States | <=50K. |
| 48812 | 81 | Other | 26711 | Assoc-voc | 11 | Married-civ-spouse | Other | Husband | White | Male | 2936 | 0 | 20 | United-States | <=50K. |
| 48826 | 50 | Local-gov | 139347 | Masters | 14 | Married-civ-spouse | Prof-specialty | Wife | White | Female | 0 | 0 | 40 | Other | >50K. |
| 48838 | 64 | Other | 321403 | HS-grad | 9 | Widowed | Other | Other-relative | Black | Male | 0 | 0 | 40 | United-States | <=50K. |
| 48838 | 64 | Other | 321403 | HS-grad | 9 | Widowed | Other | Other-relative | Black | Male | 0 | 0 | 40 | United-States | <=50K. |

6871 rows × 15 columns

```
1 # renaming columns for clarity
2
3 xy.rename(columns={'fnlwgt':'record_count',
4                    'education-num':'education_num',
5                    'marital-status':'marital_status',
6                    'capital-gain':'capital_gain',
7                    'capital-loss':'capital_loss',
8                    'hours-per-week':'hours_per_week',
9                    'native-country':'native_country'}, inplace=True)
10
11 xy
```
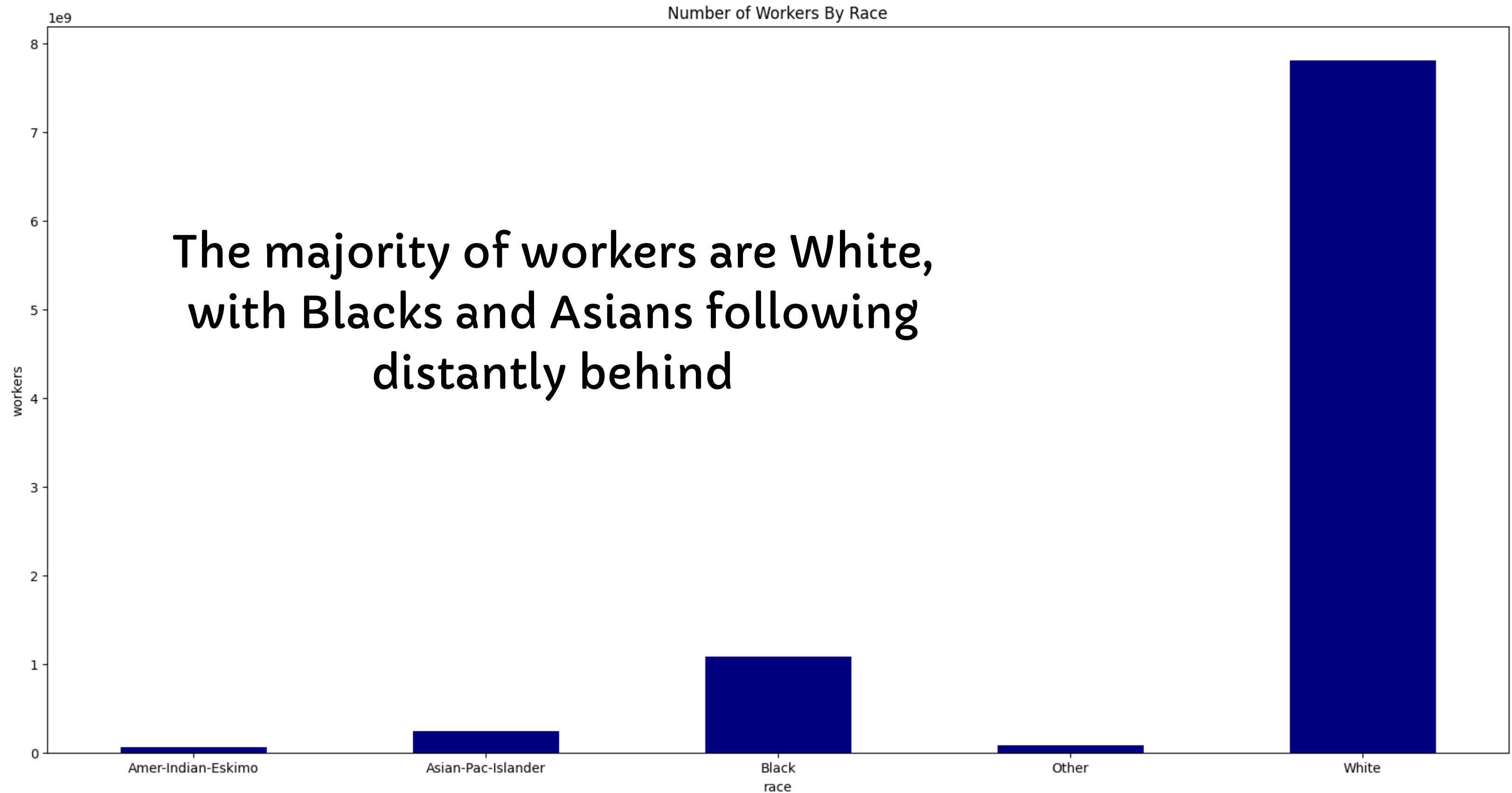
| | age | workclass | record_count | education | education_num | marital_status | occupation | relationship | race | sex | capital_gain | capital_loss | hours_per_week | native_country | income |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 39 | State-gov | 77516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family | White | Male | 2174 | 0 | 40 | United-States | <=50K |
| 1 | 50 | Self-emp-not-inc | 83311 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 13 | United-States | <=50K |
| 2 | 38 | Private | 215646 | HS-grad | 9 | Divorced | Handlers-cleaners | Not-in-family | White | Male | 0 | 0 | 40 | United-States | <=50K |
| 3 | 53 | Private | 234721 | 11th | 7 | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | 0 | 0 | 40 | United-States | <=50K |
| 4 | 28 | Private | 338409 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Wife | Black | Female | 0 | 0 | 40 | Cuba | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 48837 | 39 | Private | 215419 | Bachelors | 13 | Divorced | Prof-specialty | Not-in-family | White | Female | 0 | 0 | 36 | United-States | <=50K. |
| 48838 | 64 | Other | 321403 | HS-grad | 9 | Widowed | Other | Other-relative | Black | Male | 0 | 0 | 40 | United-States | <=50K. |
| 48839 | 38 | Private | 374983 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Husband | White | Male | 0 | 0 | 50 | United-States | <=50K. |
| 48840 | 44 | Private | 83891 | Bachelors | 13 | Divorced | Adm-clerical | Own-child | Asian-Pac-Islander | Male | 5455 | 0 | 40 | United-States | <=50K. |
| 48841 | 35 | Self-emp-inc | 182148 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 60 | United-States | >50K. |

48842 rows × 15 columns

# EXPLORATORY DATA ANALYSIS VISUALIZATIONS
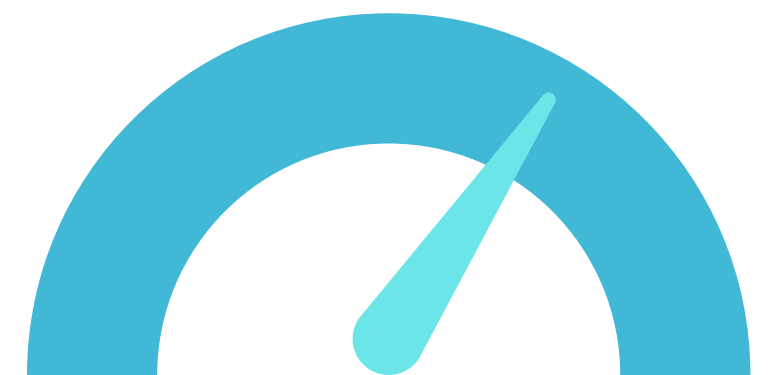
Number of Workers By Race

The majority of workers are White,
with Blacks and Asians following
distantly behind
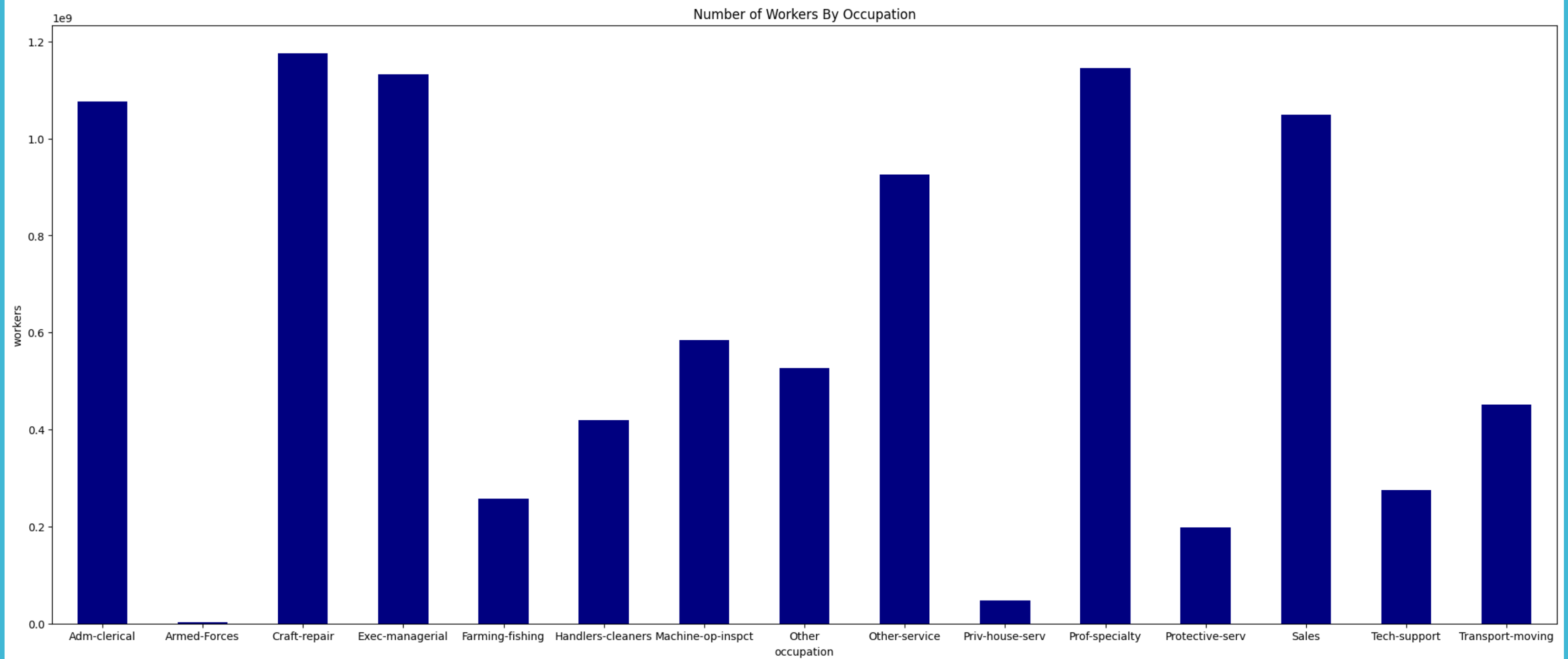
The majority of workers come from the Private sector

Number of Workers By Educational Attainment Level

The majority of workers are high school graduates, while most have either graduated, dropped out, or never reached college.
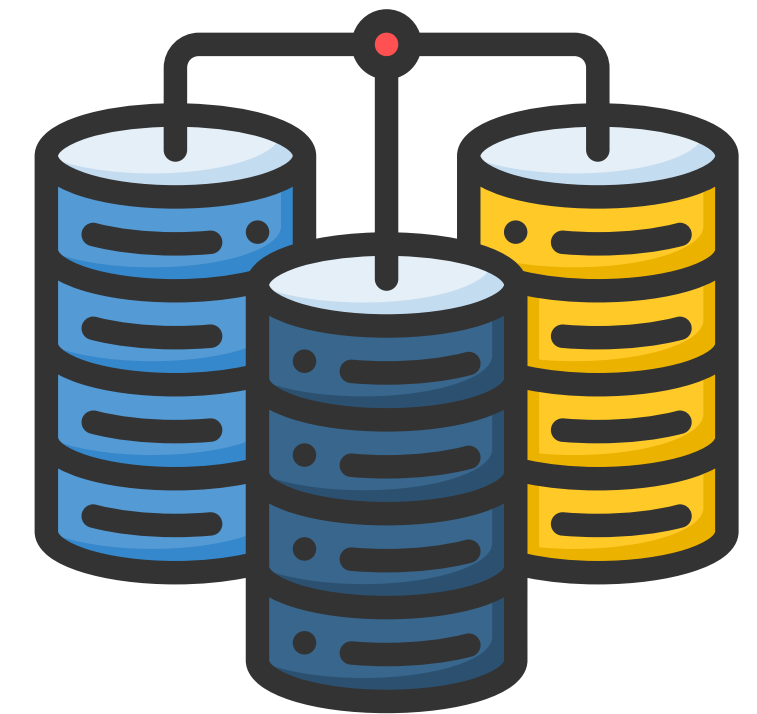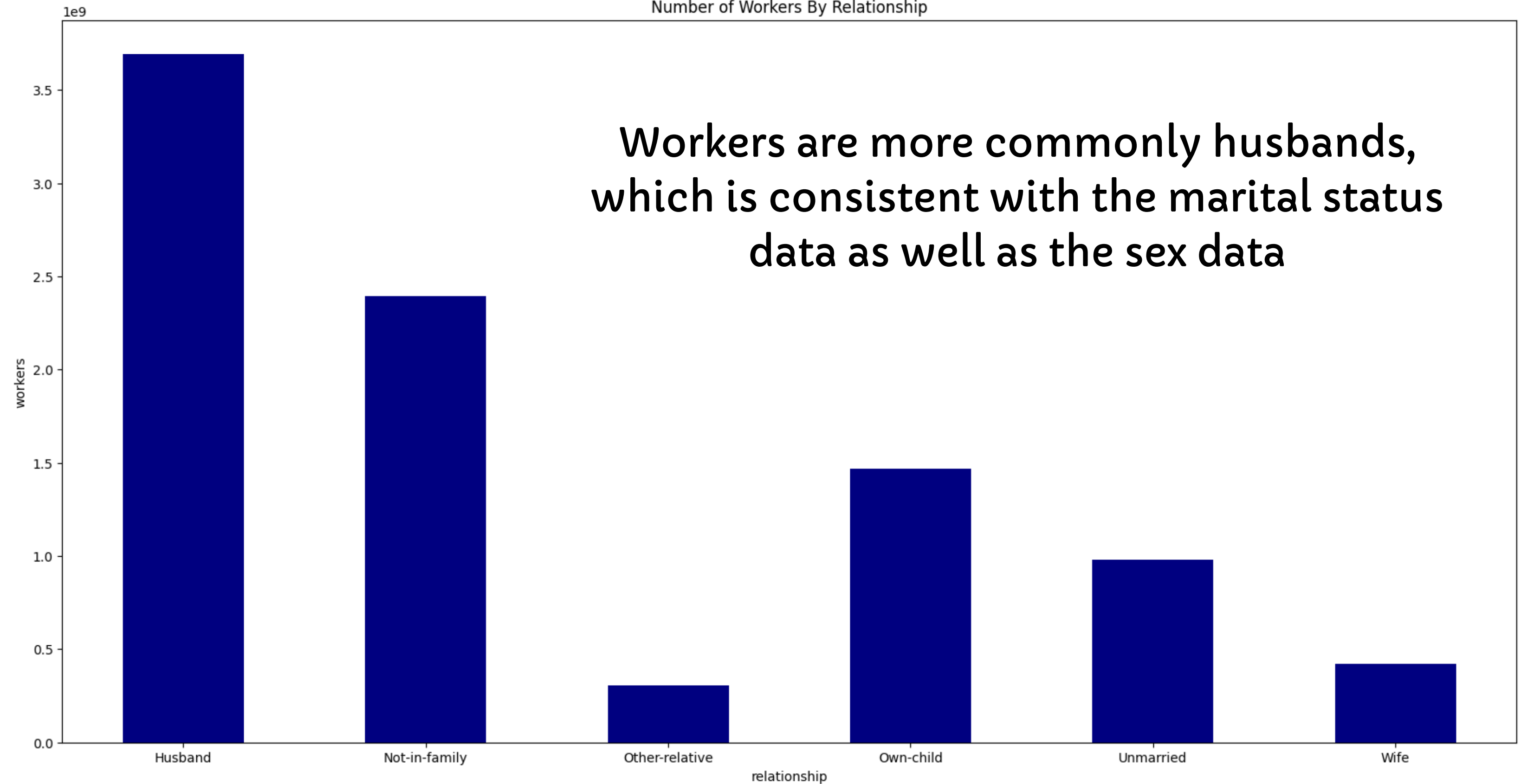
Number of Workers By Marital Status

Most workers are either married or have never been married, while some are from dysfunctional families

Number of Workers By Occupation

It can be said from the data that the distribution between white-collar and blue-collar workers is fairly even. While it does seem like the white-collar jobs are of higher count based on the graph, the same can be said for the blue-collar jobs, but they are more distributed into more columns while the latter are condensed into lesser columns.

Number of Workers By Native Country

The majority of workers are from the United States, which is consistent with the race data where only a small amount of workers are from a different race not native to the US

Number of Workers by Age

The small number of workers from the minority age group as well as from the 40s and above age group shows how middle-aged workers are more favored over the previous 2 age groups.
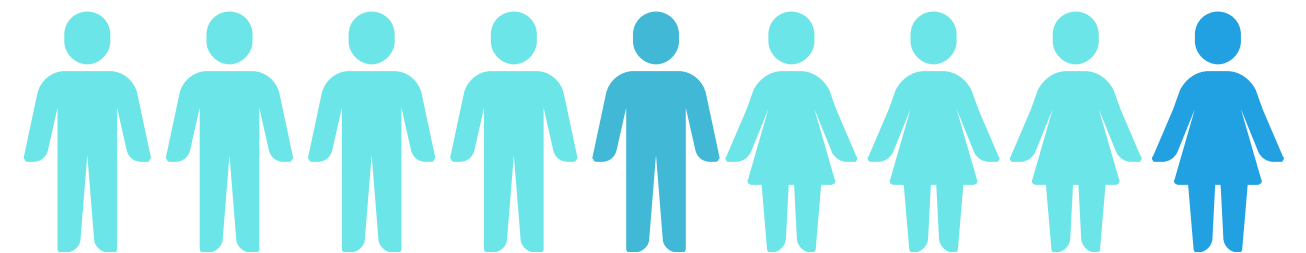
Work Hours per Week by Age

However, when the work hours are compared, only the minority have lower work hours. It can be inferred that these workers are probably interns or part-time workers, explaining the shorter work hours
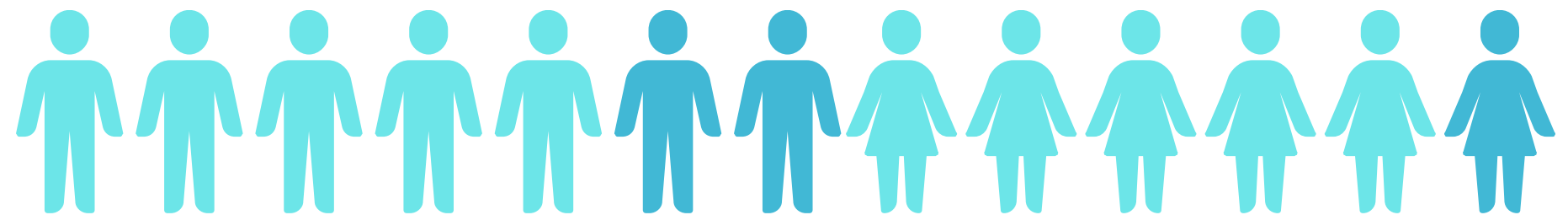
CORRELATION ANALYSIS

# Sex vs Income

When comparing the income ranges for each gender, we can see that the majority of males earn less than 50K, as well as for females but at a lesser volume. This is consistent with the sex data shown previously
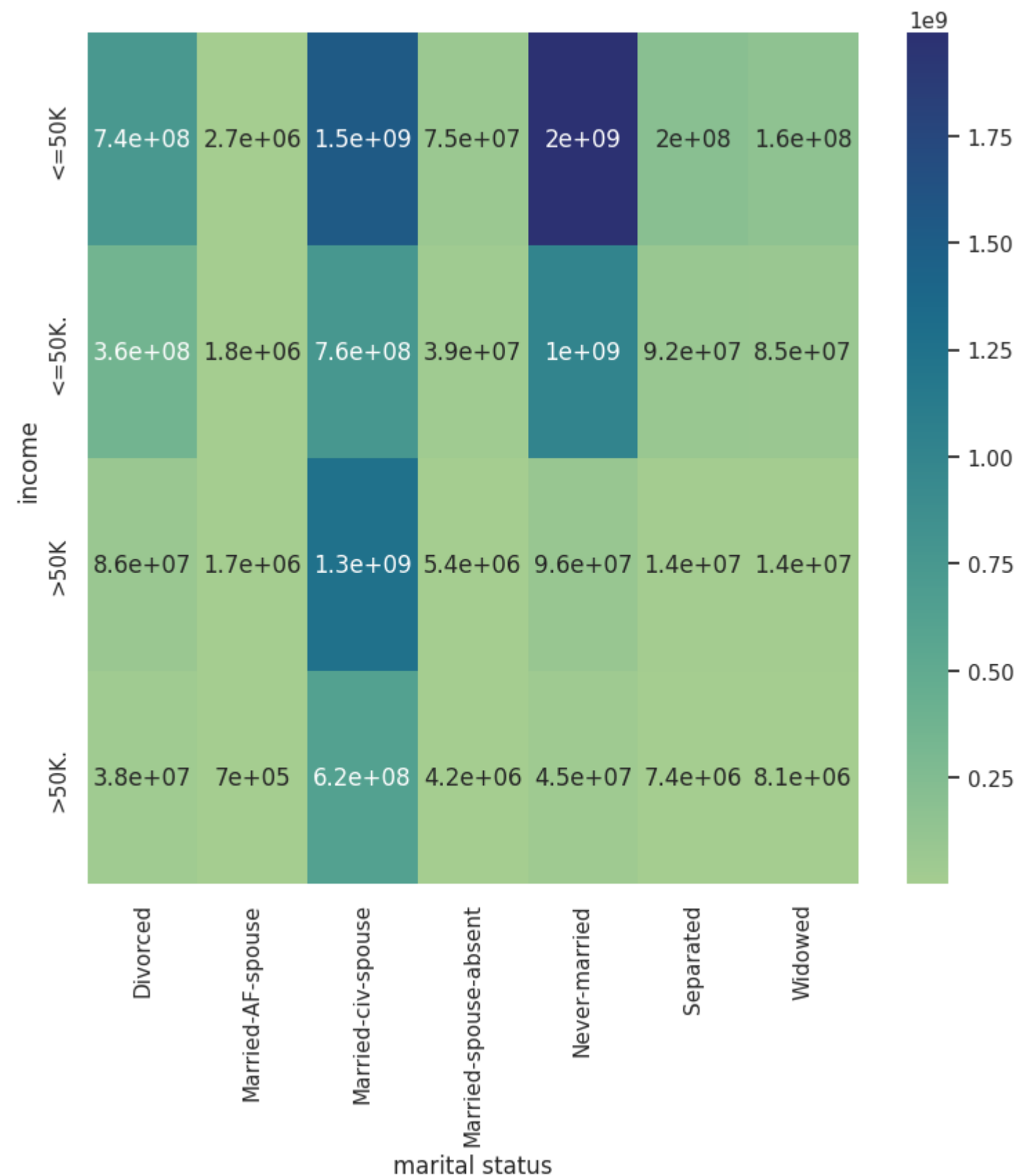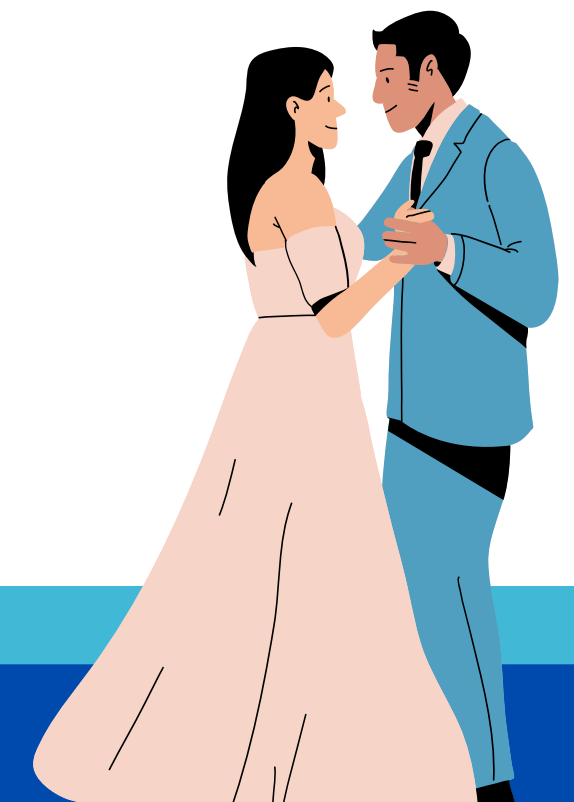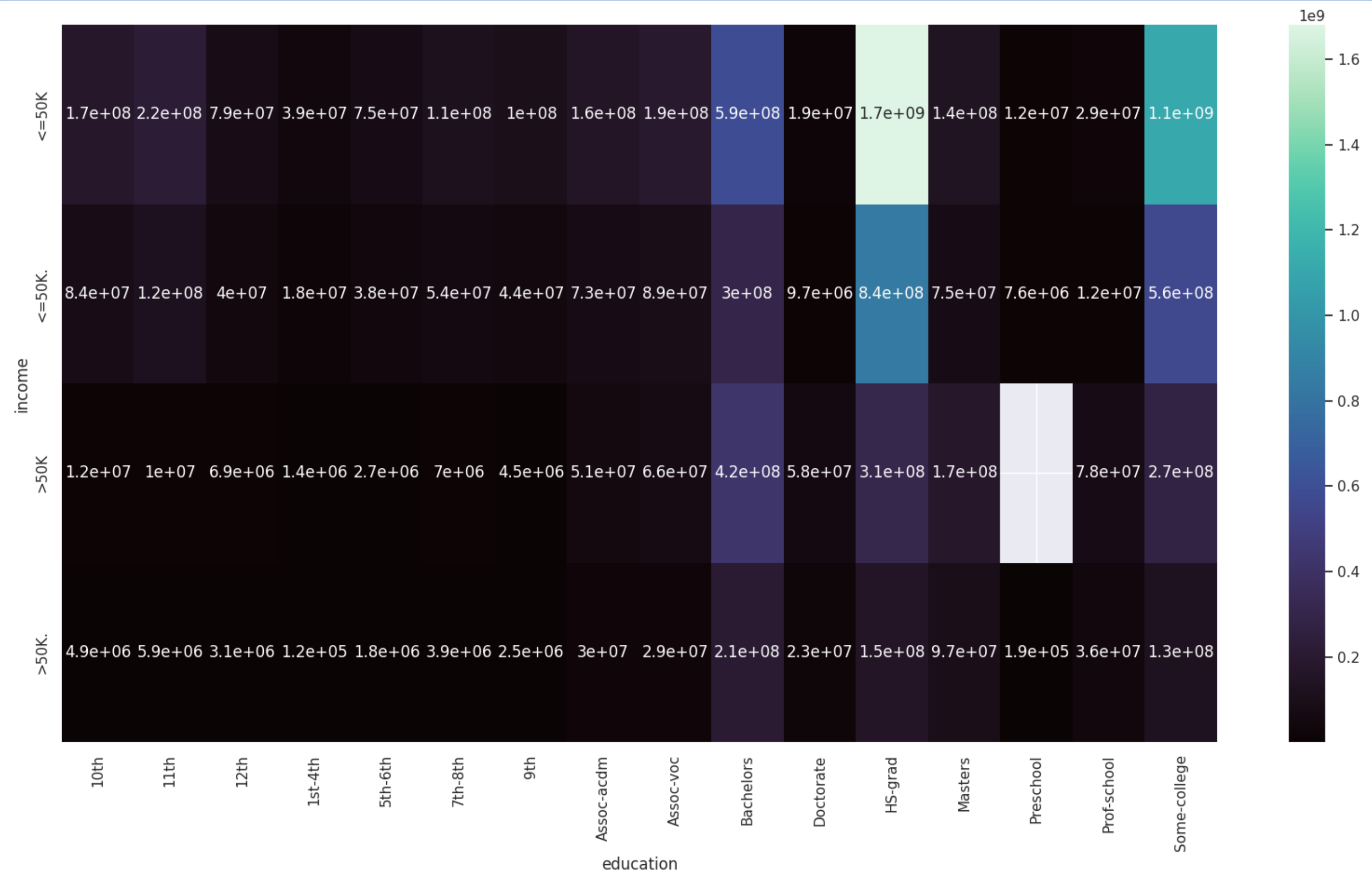
# Marital Status vs Income

The majority of workers who are either never married or have married a civilian earn more or less than 50K. It is more distributed for married civilians in comparison to the ones who have never married.
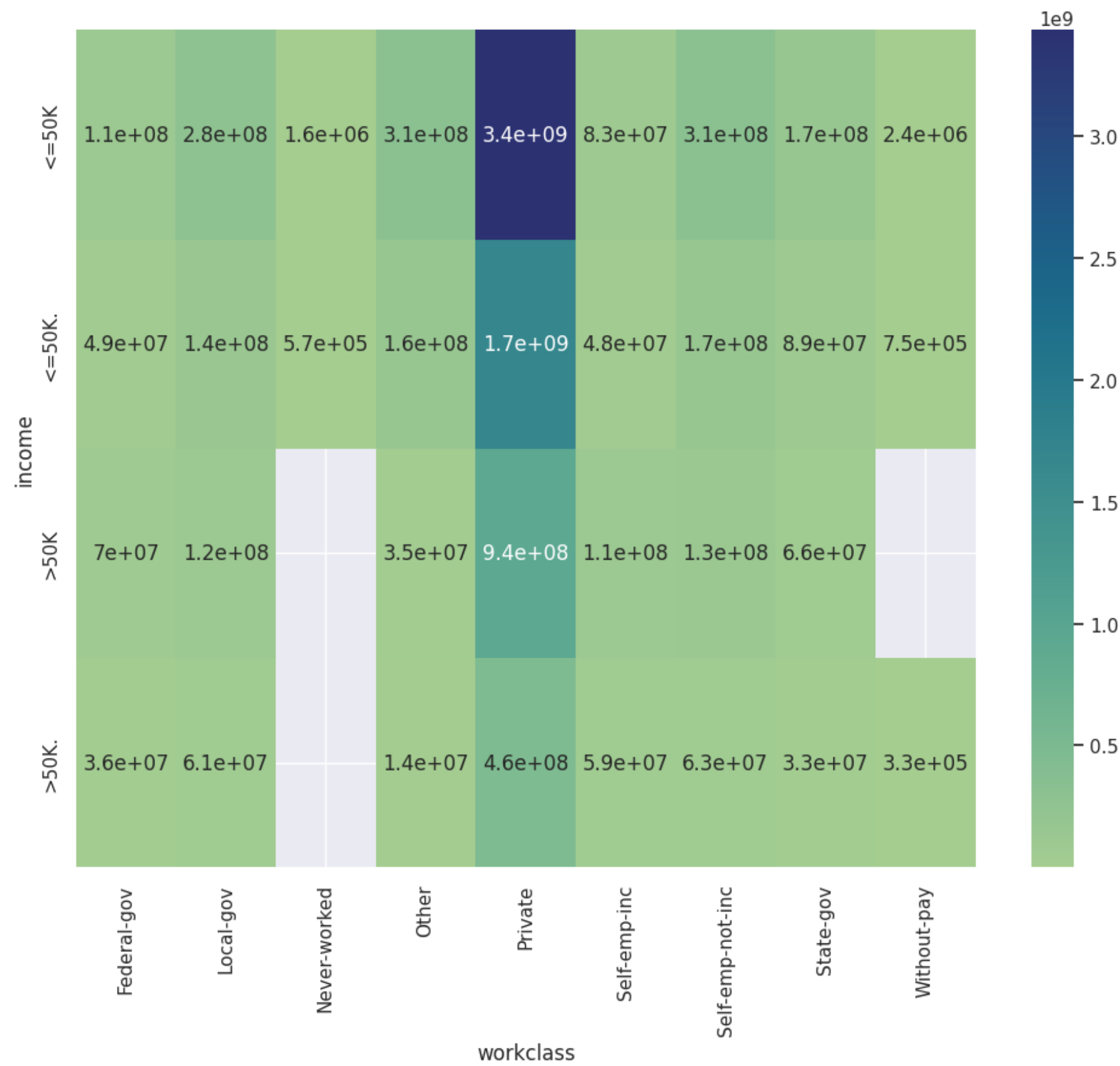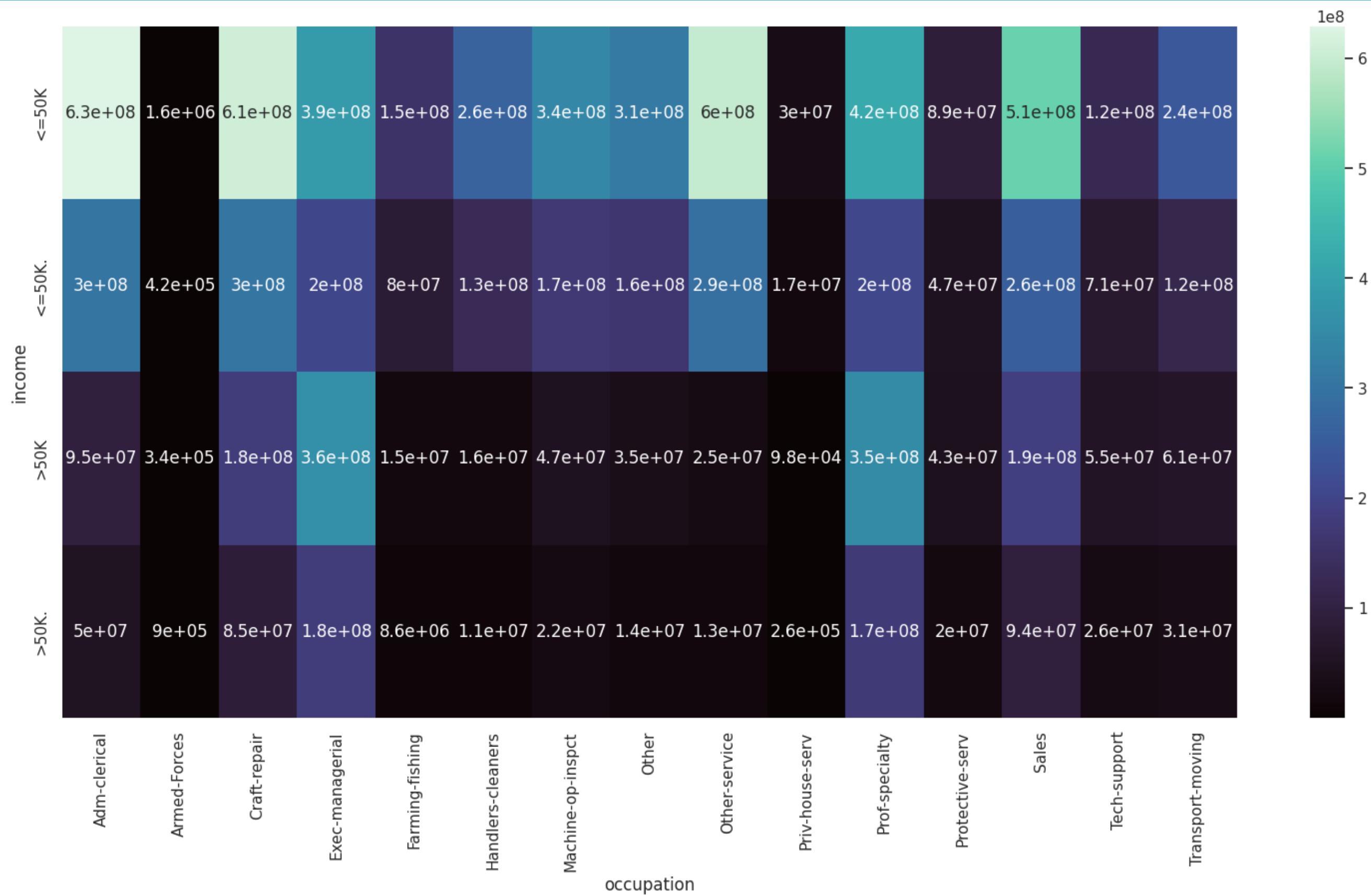
# Educational Attainment vs Income

It can be said that a majority of those who are earning less than 50K are those who have dropped out of college or are high-school graduates. This is also consistent with the education data, which shows that these 2 are the majority of workers in the US.

Workclass vs Income

# Conclusion

From the dataset, I was able to clean up some missing values as well as rearrange the dataset to make it easier to analyze. As I was analyzing the dataset, I was able to determine that the record_count column (previously the fnlwgt column) was going to be the most important data in the dataset. This is because it is the most consistent quantity present in the dataset. Most columns in the dataset, as well as the income column, are mostly qualitative, which doesn't give much space for any computations and accurate analysis. With the record_count column, I was able to visualize and compare the amount of people represented by the records per specific demographic, as well as correlate the income category in specific columns through the use of the record_count column. I think that being able to analyze and identify how to work with a dataset is key to being able to interpret it. While I was only able to visualize multiple relationships in the dataset, I think that I have done enough of it that some form of interpretation and conclusion can be drawn from some of them combined.

# THANK YOU FOR LISTENING!