

Homework 3 – Report

Importing Data & Data Cleaning

- After importing the necessary modules from several libraries and the data set, I renamed the columns in accordance with the codebook and replaced missing values indicated in the questionnaire with NULL.
- Then, I tried to fill the NULL values:
 - I've dropped the columns with larger than 54% of NULL values (since I want to keep "SOCIO ECONOMIC STATUS", with 53.51%, which can be an important indicator of voting behavior).
 - I've also dropped "RELIGIOUS DENOMINATION" column, since its values does not correspond to the ones shown in questionnaire.
 - For the columns with less than 10% of NULL values, I've filled them with the most common value.
 - For the other columns, I've filled them with a value randomly chosen from that column, therefore, the values with more counts will have a higher chance to be selected.

Exploratory Data Analysis

- I've divided the data set column-wise according to whether a variable is categorical, ordinal, or numerical, by checking the values they can take from the questionnaire.

Categorical Columns:

- GENDER
- MARITAL STATUS
- UNION MEMBERSHIP OF RESPONDENT
- UNION MEMBERSHIP OF OTHERS IN HOUSEHOLD
- CURRENT EMPLOYMENT STATUS
- EMPLOYMENT TYPE - PUBLIC OR PRIVATE
- INDUSTRIAL SECTOR
- LANGUAGE USUALLY SPOKEN AT HOME
- REGION OF RESIDENCE
- RACE

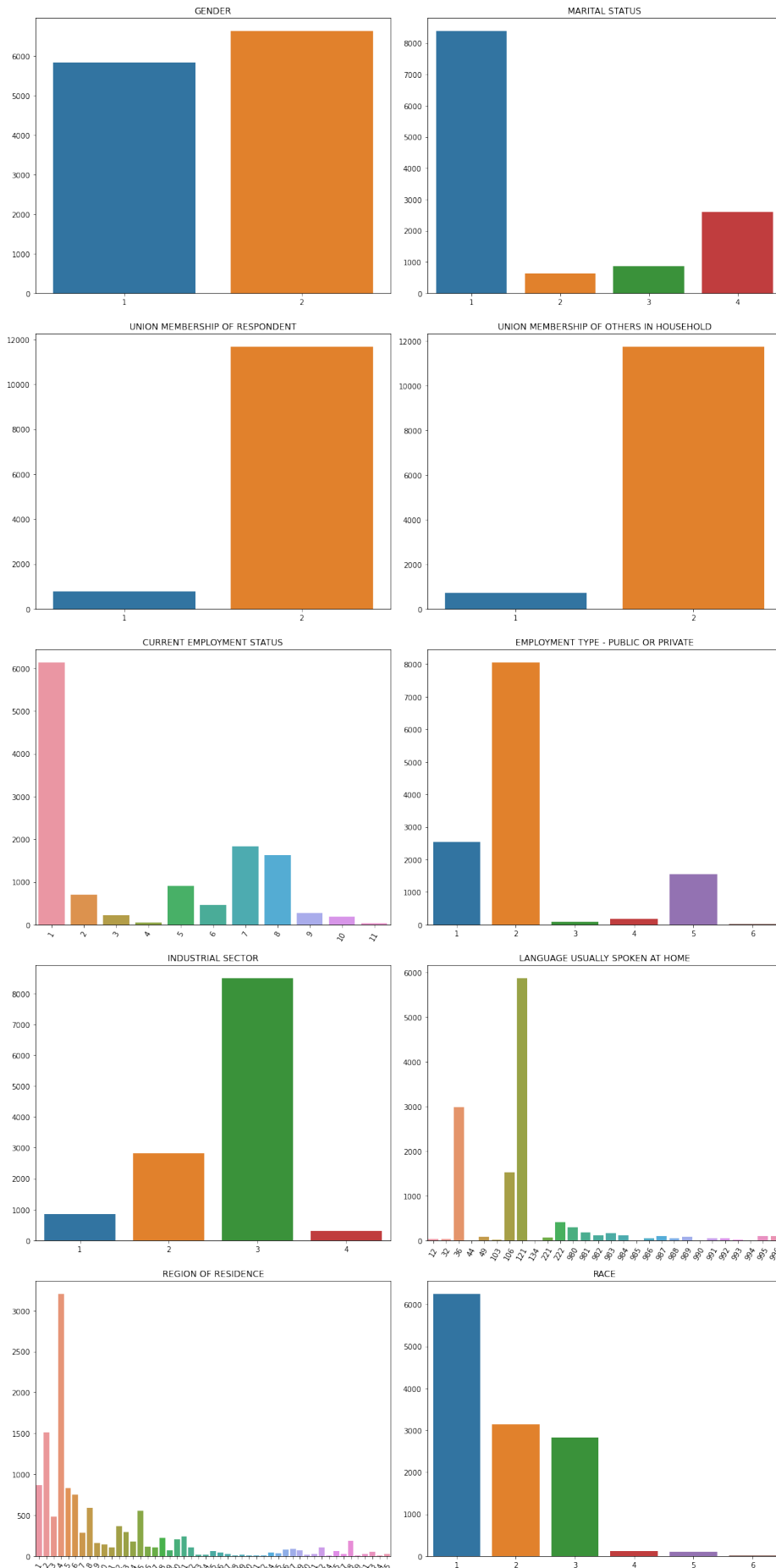
Ordinal Columns:

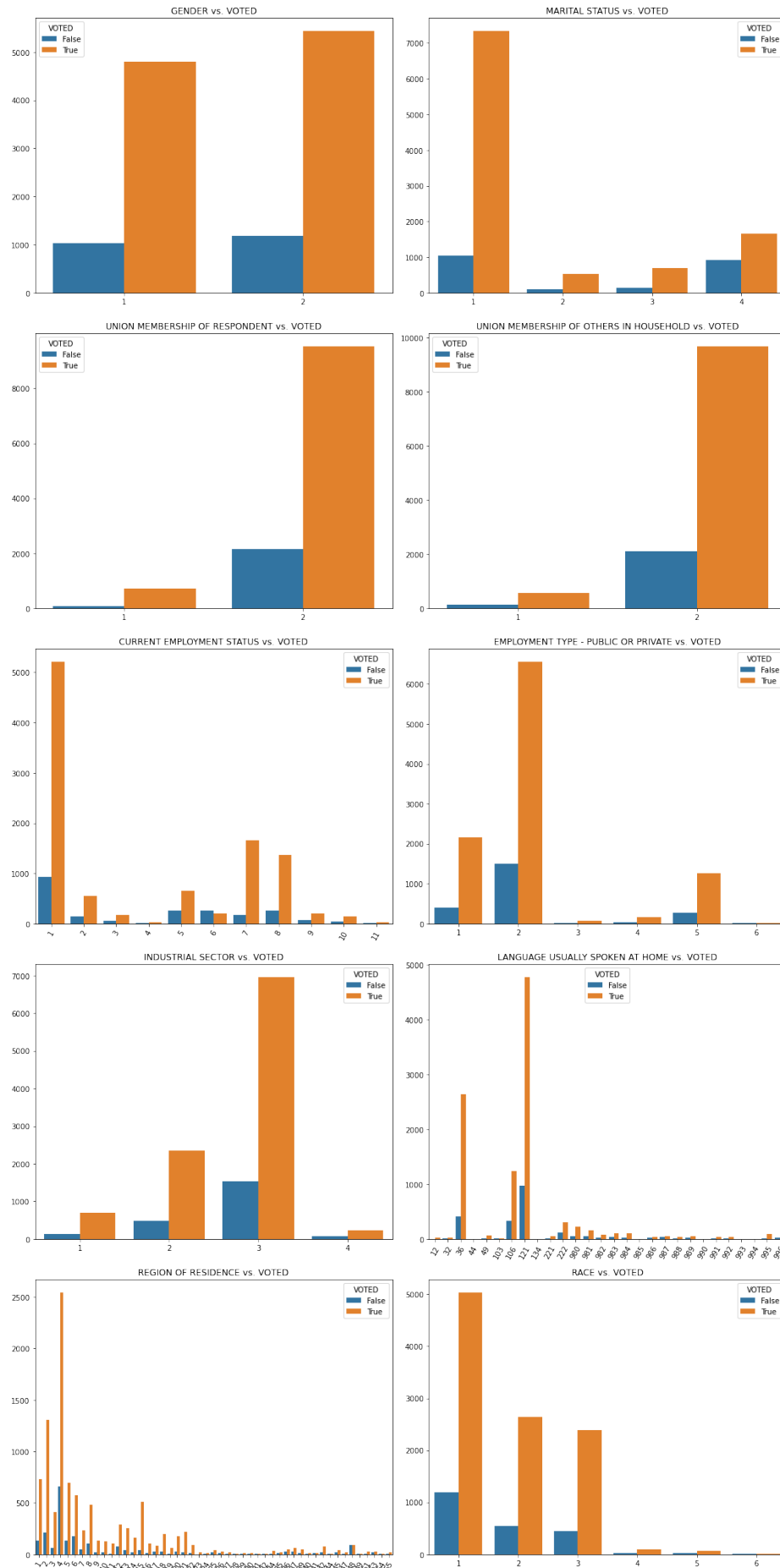
- EDUCATION
- SOCIO ECONOMIC STATUS
- HOUSEHOLD INCOME
- RELIGIOUS SERVICES ATTENDANCE
- RELIGIOSITY
- RURAL OR URBAN RESIDENCE
- NUMBER IN HOUSEHOLD IN TOTAL
- NUMBER OF CHILDREN IN HOUSEHOLD UNDER AGE 18
- NUMBER IN HOUSEHOLD UNDER AGE 6

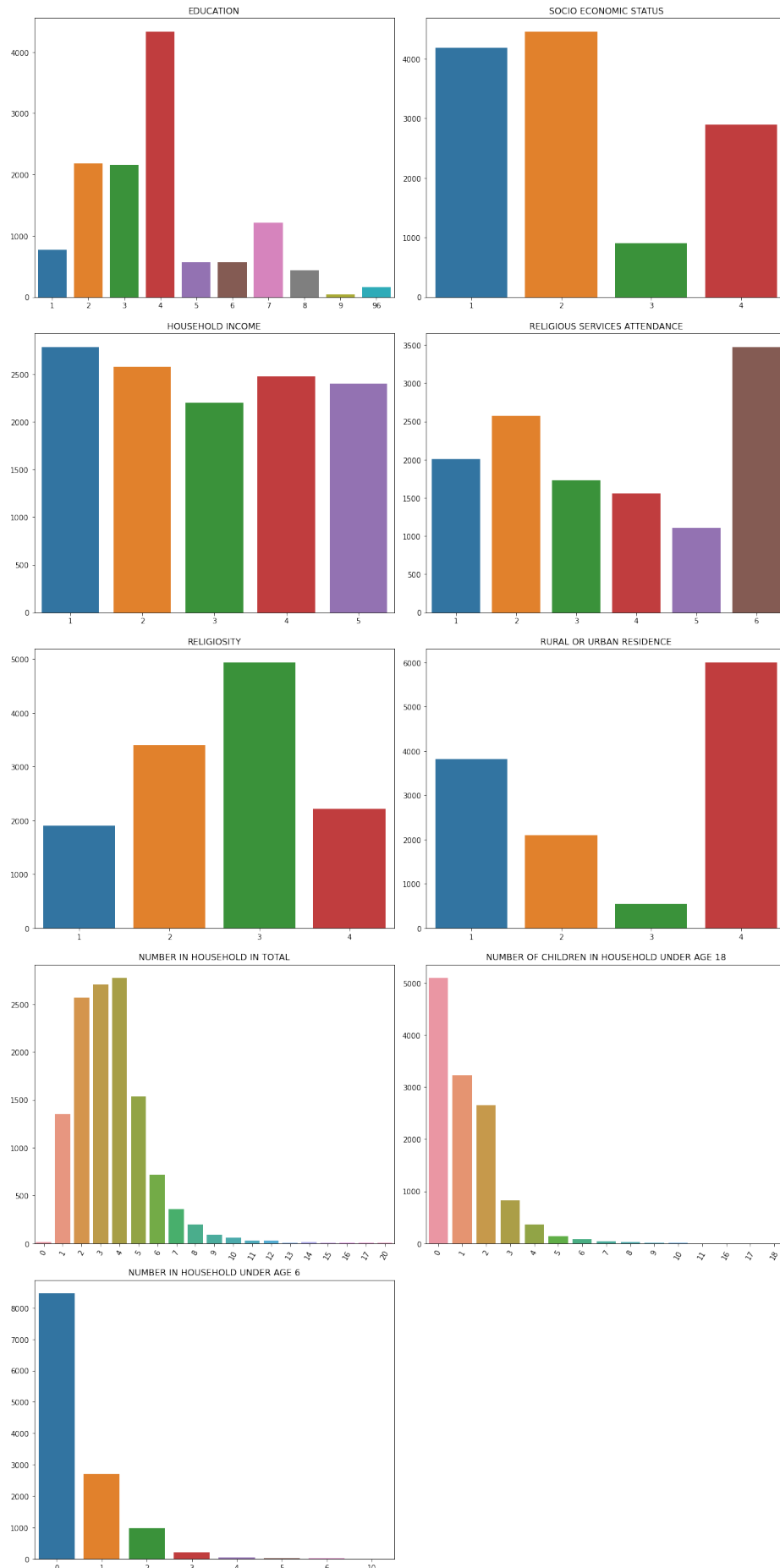
Numerical Columns:

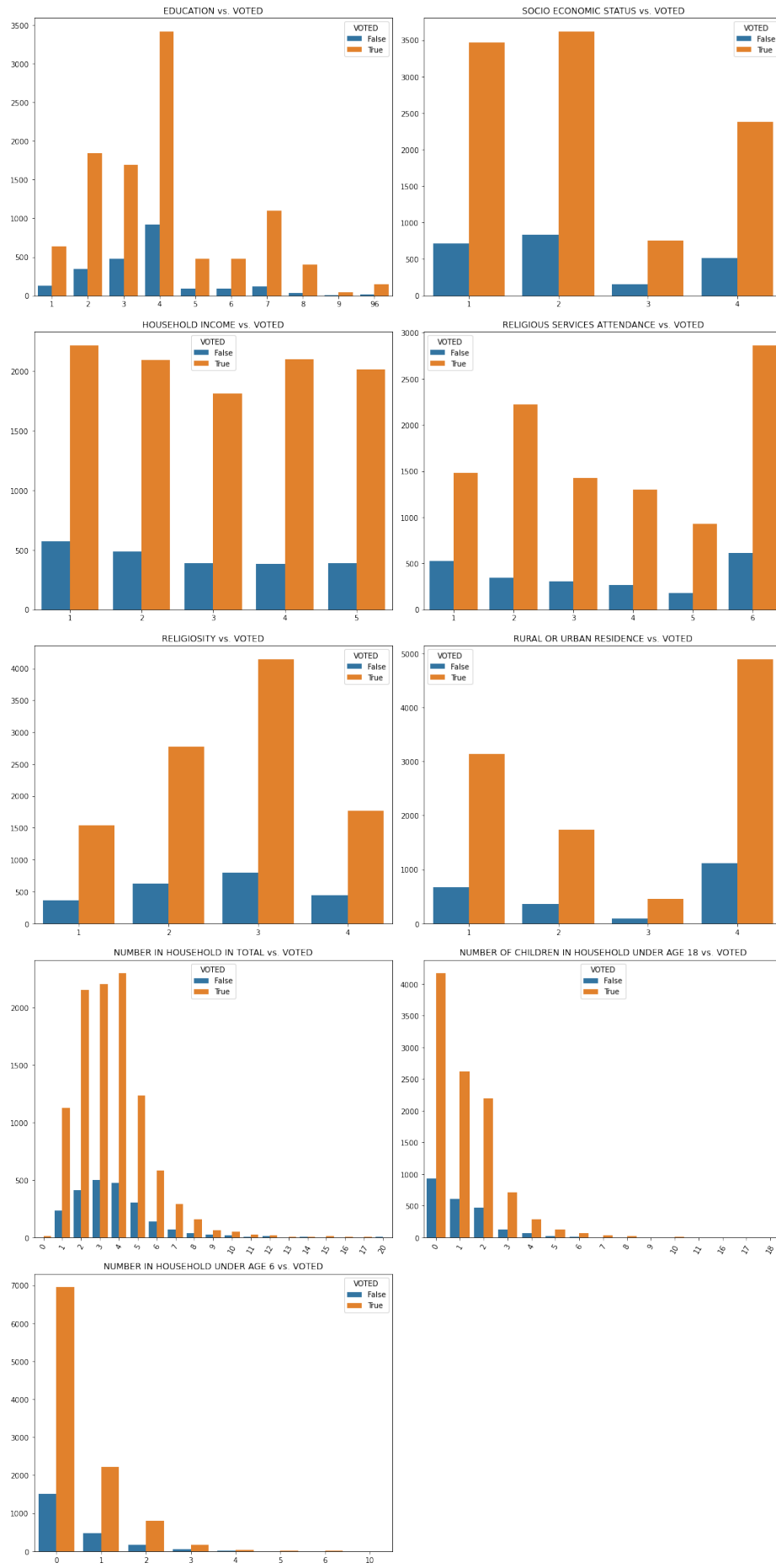
- AGE

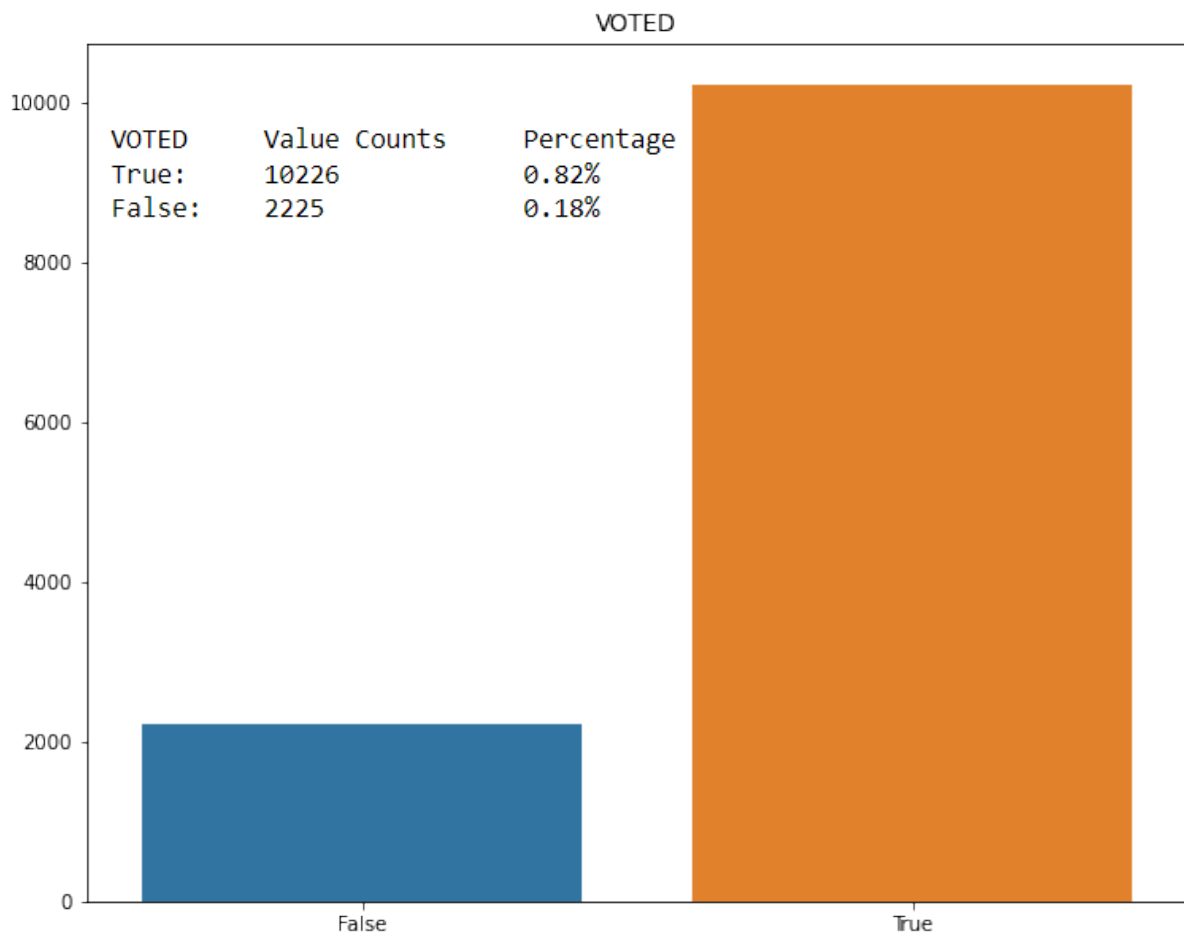
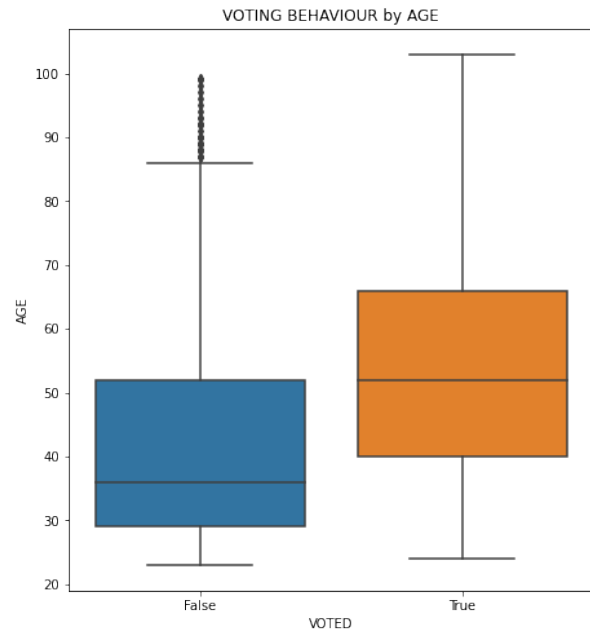
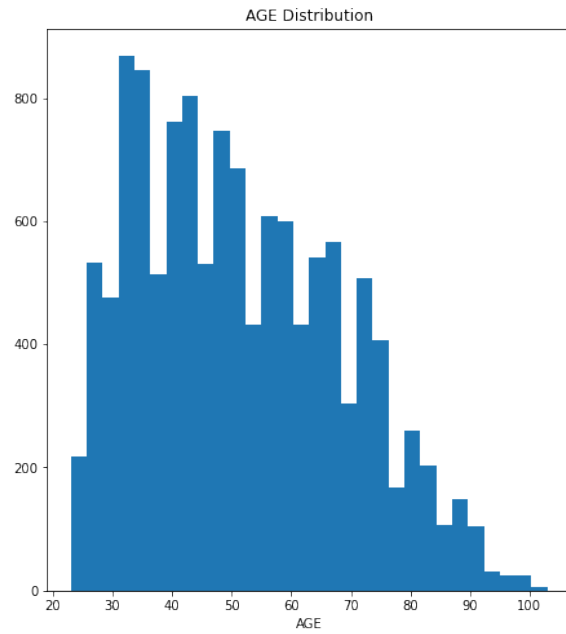
- Then, I've visualized each variable, by looking at the counts in categorical and ordinal variables, and distribution in the numerical variable. The plots in the following pages demonstrate these counts and distributions.
- I've also looked at the counts of the target variable, "VOTED", which indicated that the target variable is indeed imbalanced. Therefore, I've decided to look not only at the accuracy metric but also precision, recall, F1 score and ROC AUC metrics.
- Lastly, I have created a heatmap of the correlation between each of the variables. "AGE" variables seems to be the most correlated (0.26) with the target variable "VOTED", while "MARITAL STATUS" is the second mostly correlated with a negative correlation score (-0.23).

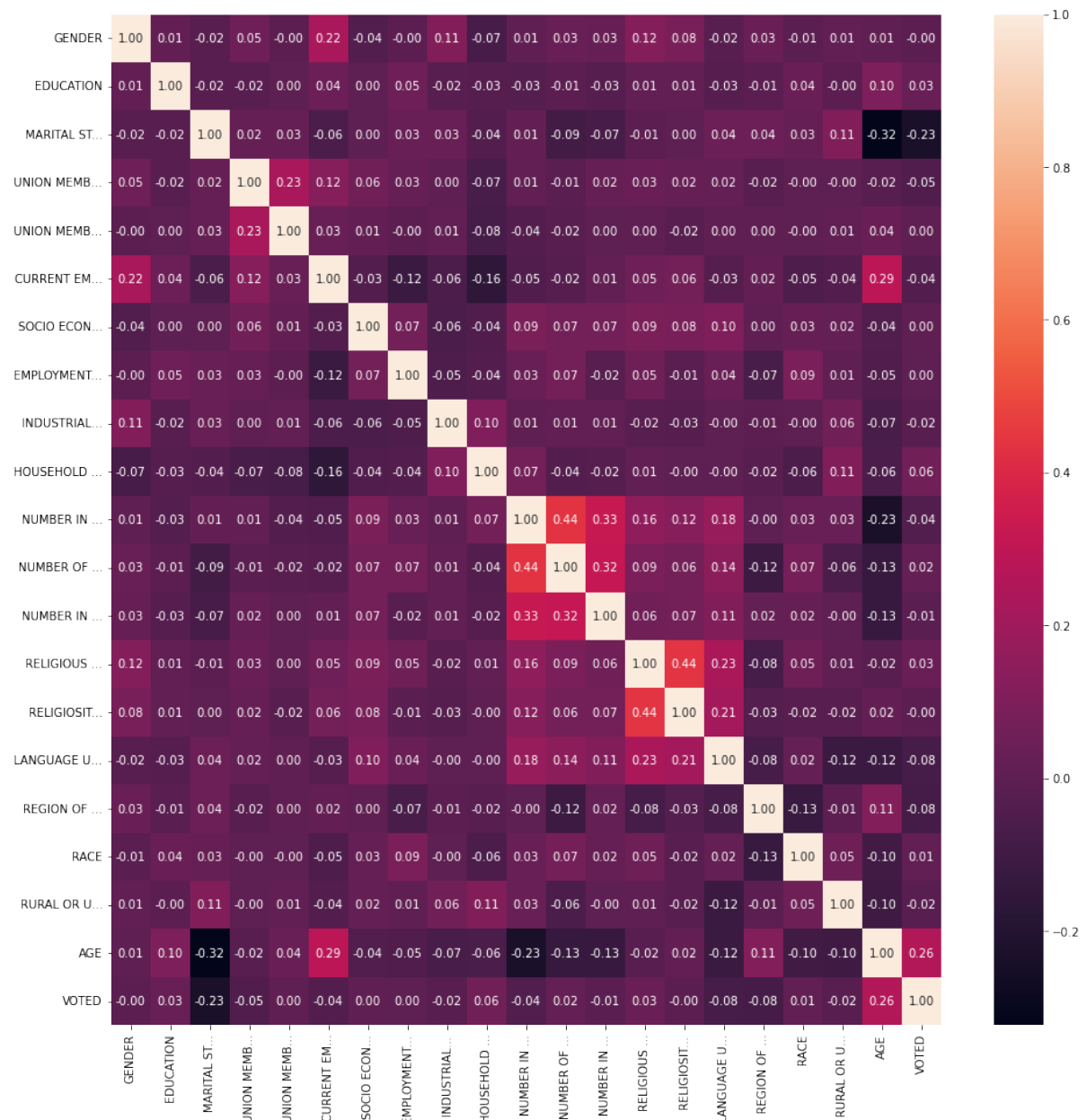












Feature Engineering

- I've only one-hot-encoded the categorical variables, and concatenated all the data with different variable types.
- I decided to use tree-based algorithms. Therefore, I decided not to normalize the data since tree-based algorithms do not require normalization.
- I've splitted the data into train and test sets.

ML Algorithm

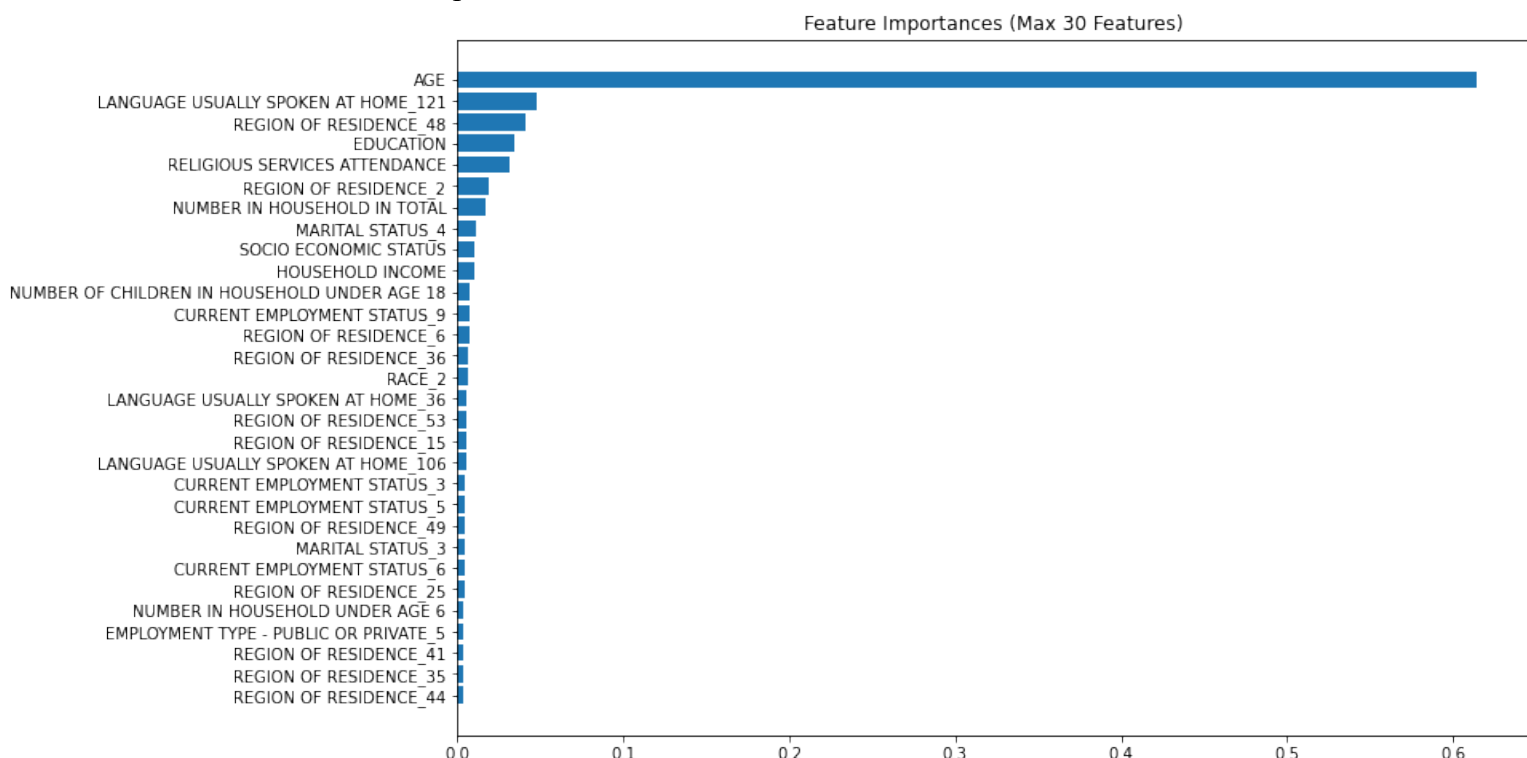
- Initially, I've tried the following four algorithms from Scikit-Learn: Decision Tree Classifier, Extra Trees Classifier, Random Forest Classifier, Gradient Boosting Classifier.

	DecisionTreeClassifier	ExtraTreesClassifier	RandomForestClassifier	GradientBoostingClassifier
Train Accuracy	1.000000	1.000000	1.000000	0.871014
Test Accuracy	0.792650	0.846435	0.859333	0.861037
Precision	0.874041	0.978171	0.981711	0.974926
Recall	0.874557	0.856184	0.865765	0.871802
F1	0.874299	0.913121	0.920100	0.920485
ROC AUC	0.641162	0.771312	0.815901	0.802253

- All models except the Gradient Boosting Classifier resulted in overfitting.

Feature Selection

- Using Gradient Boosting Classifier, I've found feature importances. The following plot shows the most important 30 features.



- Then, I have run the Gradient Boosting Classifier on the data set using [5, 10, 15, ... 105] features each time. I got the best result with **75 features**. After this point, I have continued with those 75 features.

Hyperparameter Tuning

- I've used Scikit-Learn's GridSearchCV with the following parameters and 3-fold cross-validation:

```
params = {  
    'learning_rate': [0.1, 0.3, 0.5, 0.7, 0.9],  
    'max_depth': [1,3,5,7,9],  
    'min_samples_split': [0.1, 0.3, 0.5, 0.7, 0.9],  
    'min_samples_leaf': [0.1, 0.3, 0.5, 1],  
}
```

- Finally, I have compared the initial results of Gradient Boosting Classifier, the results after reducing to 45 variables and the results of it with hyperparameters tuned.

	GBC_45_features	Initial_GBC	GBC_tuned
Train Accuracy	0.866579	0.867538	0.865020
Test Accuracy	0.859577	0.859577	0.859577
Precision	0.976991	0.976106	0.869451
Recall	0.869063	0.869645	0.976401
F1	0.919872	0.919805	0.919828
ROC AUC	0.803659	0.801599	0.642582