

Census Data Analytics

Paulo C. Rios, Jr.

Sep 21, 2017

Goals

- Using the “Adult” dataset, predict whether income exceeds \$50K/yr based on census data.
- Choose an accuracy statistic to measure model’s accuracy.
- “Adult” dataset available at <http://archive.ics.uci.edu/ml/datasets/Census+Income>

Actions

- Using a Python Jupyter notebook, I have done the following:
- Training Data Set Importing and Reading
- Data Preprocessing
- Feature Visualization and Transformation
- Model Training and Performance Measurements
- With two models: Random Forest and Stochastic Gradient Descent (SGD)

Data Preprocessing

Steps

- Columns were renamed
- Target feature had its value reset as True/False
- Null values were removed
- Original training data set had 32,560 rows
- Without null values, training data set has 30,161 rows

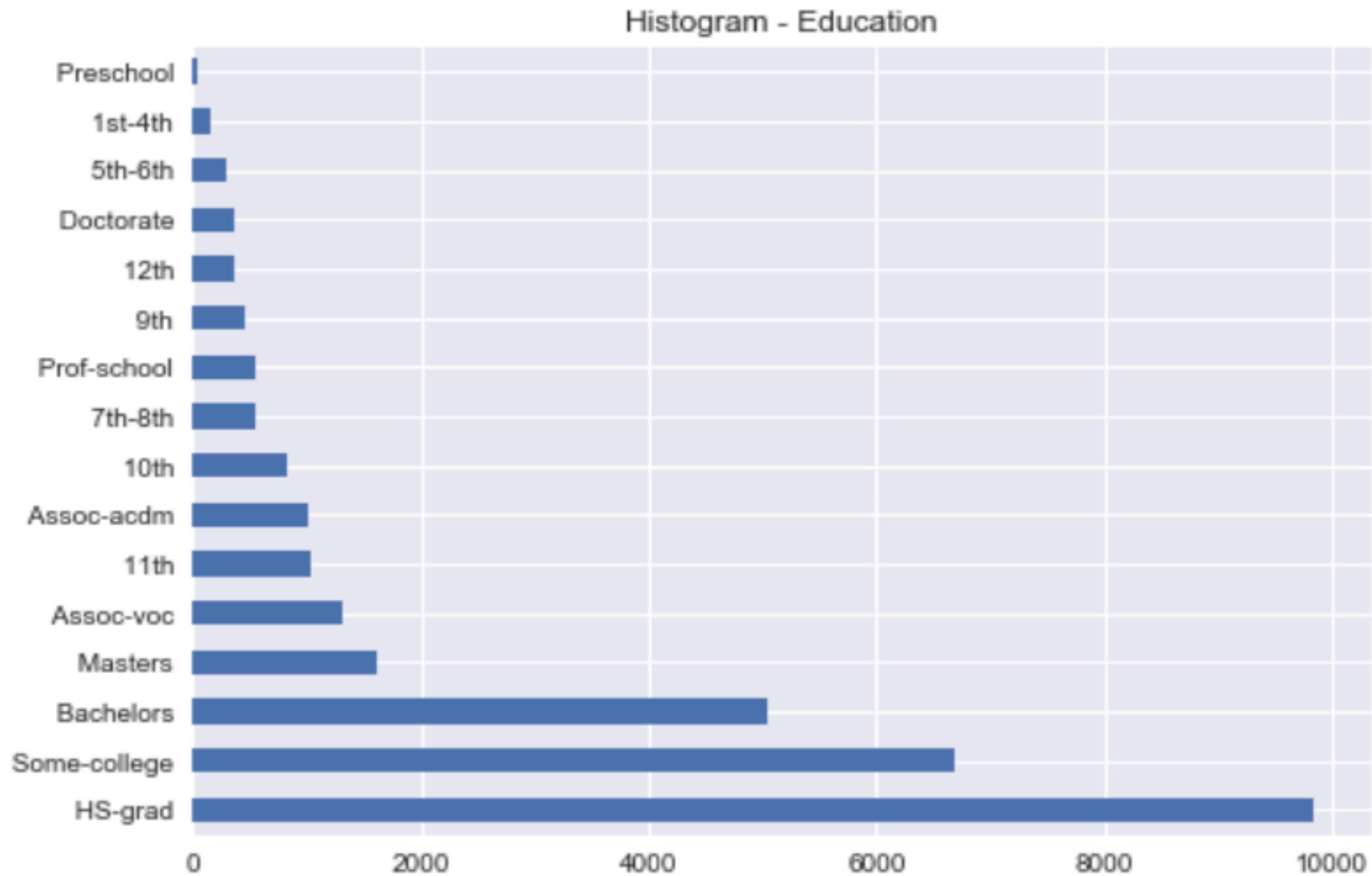
Null Values

```
In [26]: data.isnull().any()
```

```
Out[26]: Age                False  
Work Class              True  
Final Weight           False  
Education              False  
Education Num          False  
Marital Status         False  
Occupation             True  
Relationship           False  
Race                  False  
Sex                   False  
Capital Gain           False  
Capital Loss           False  
Hours per Week         False  
Native Country         True  
<=50K                  False  
dtype: bool
```

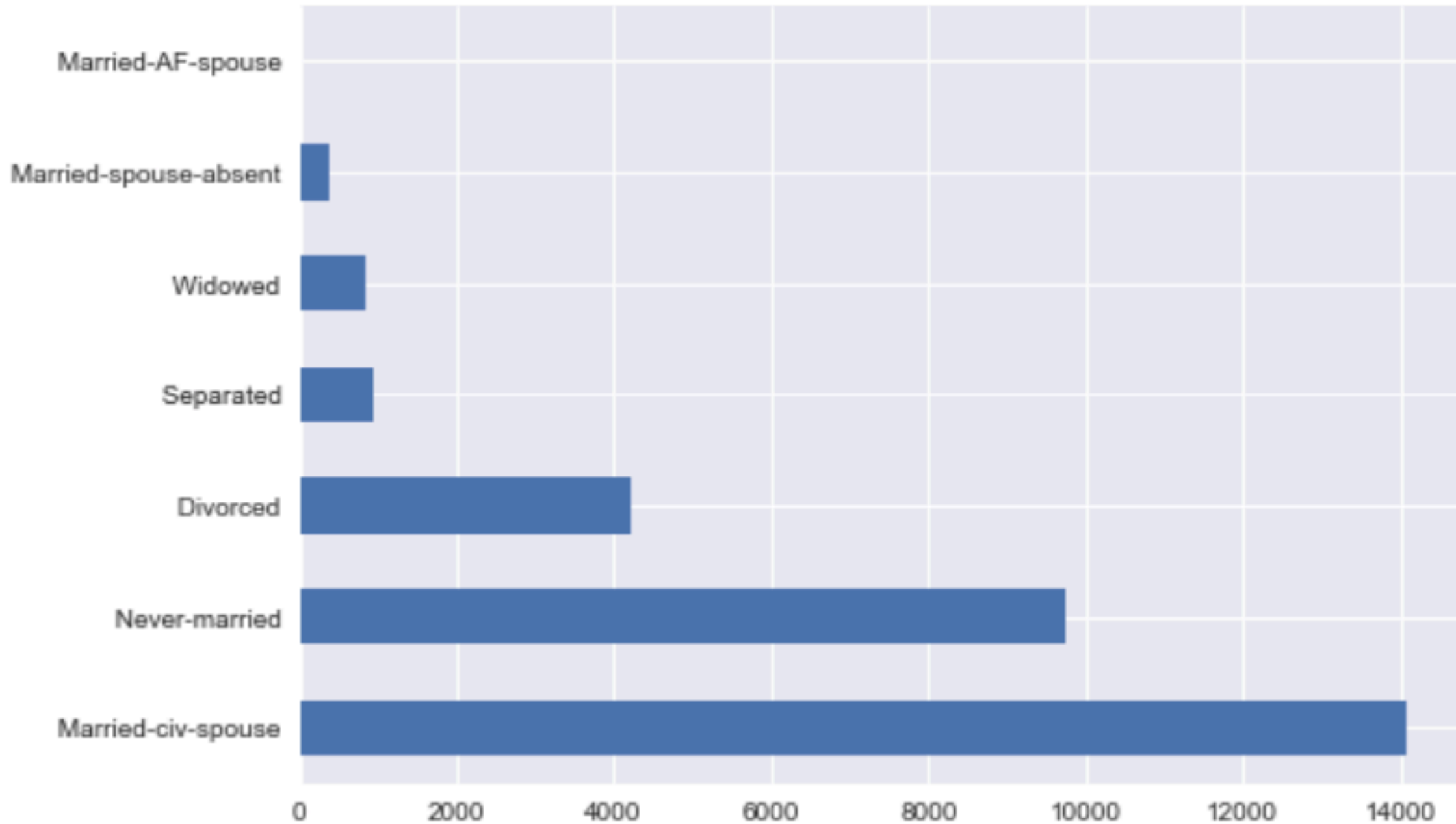
Feature Visualization

Education

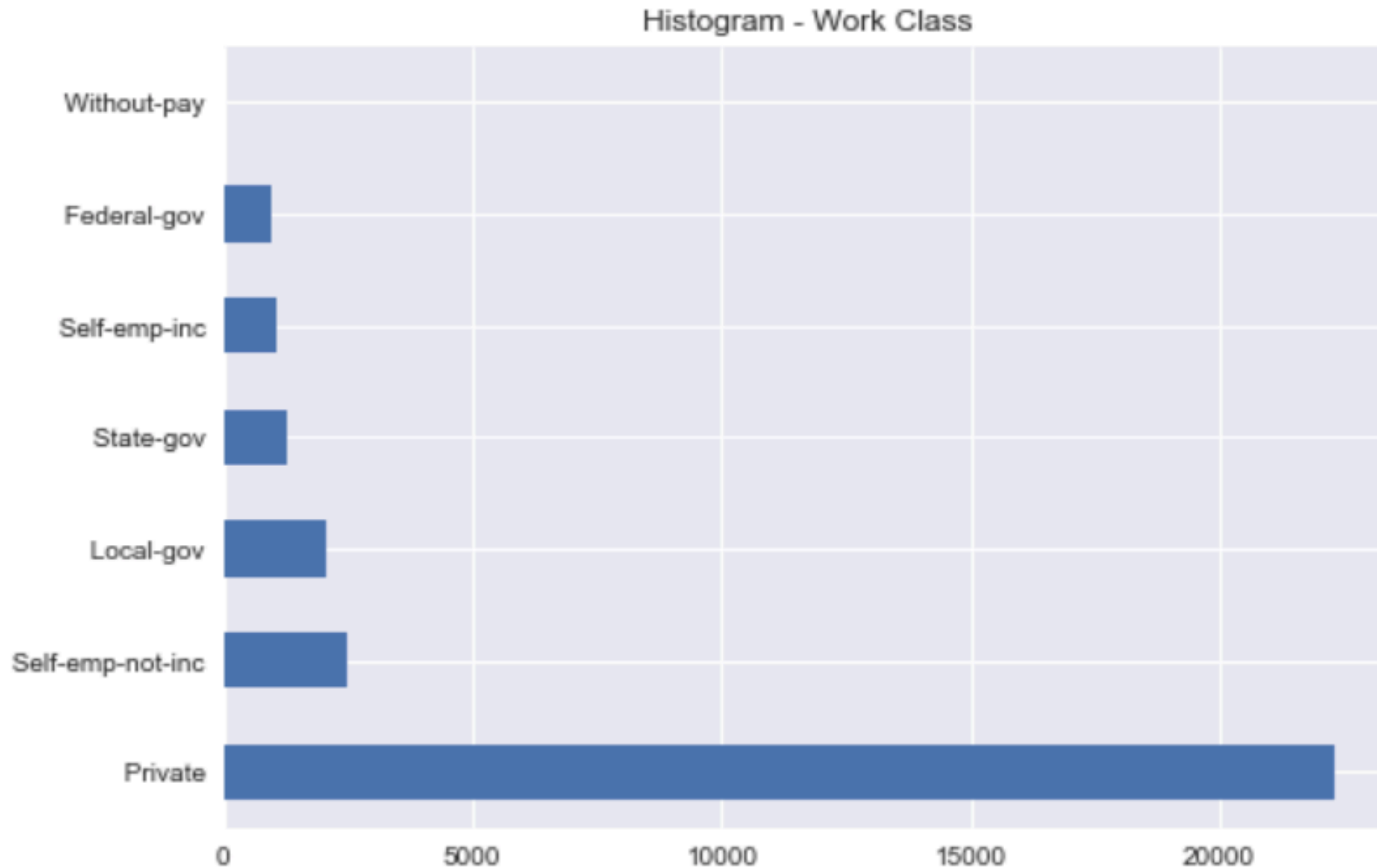


Marital Status

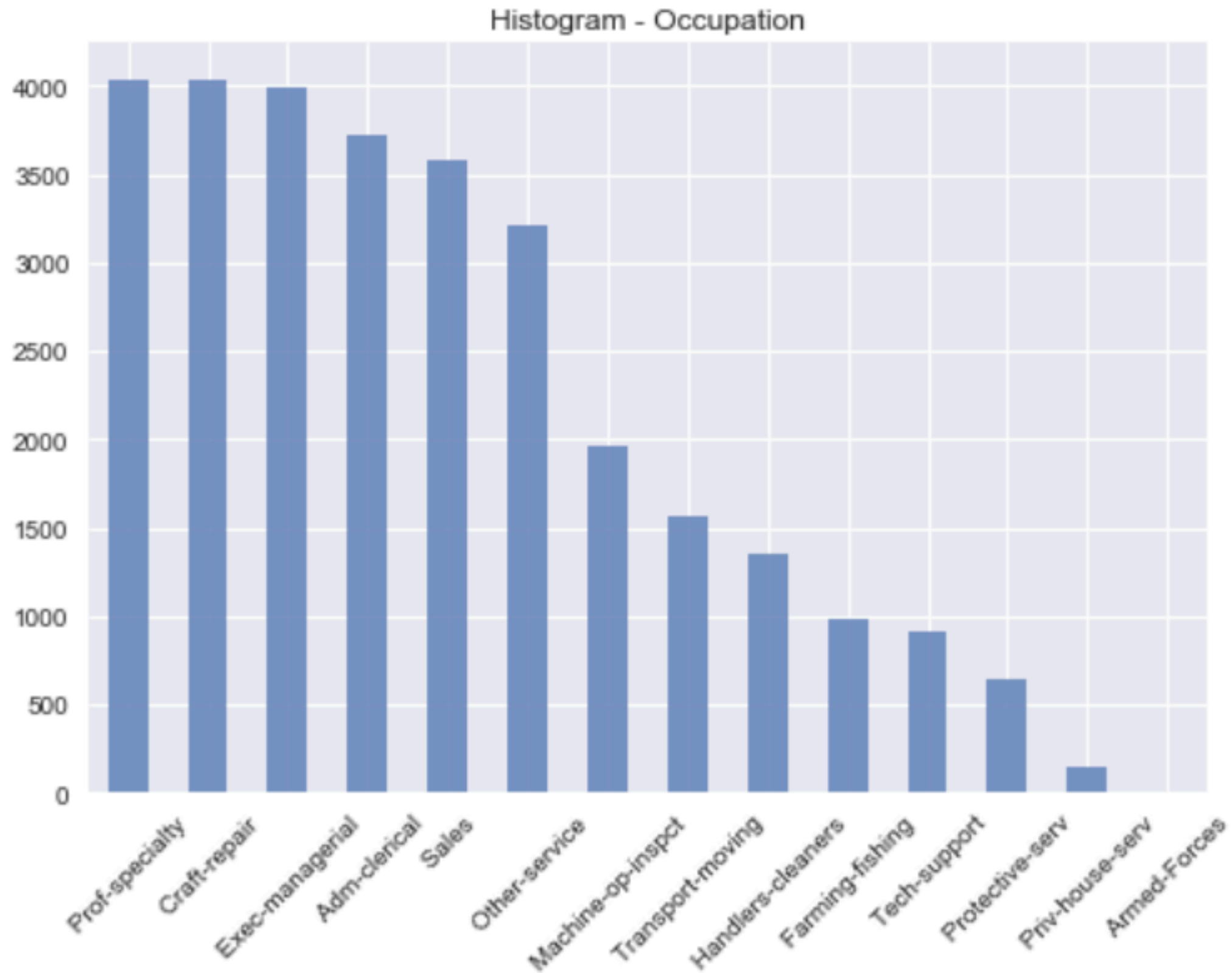
Histogram - Marital Status



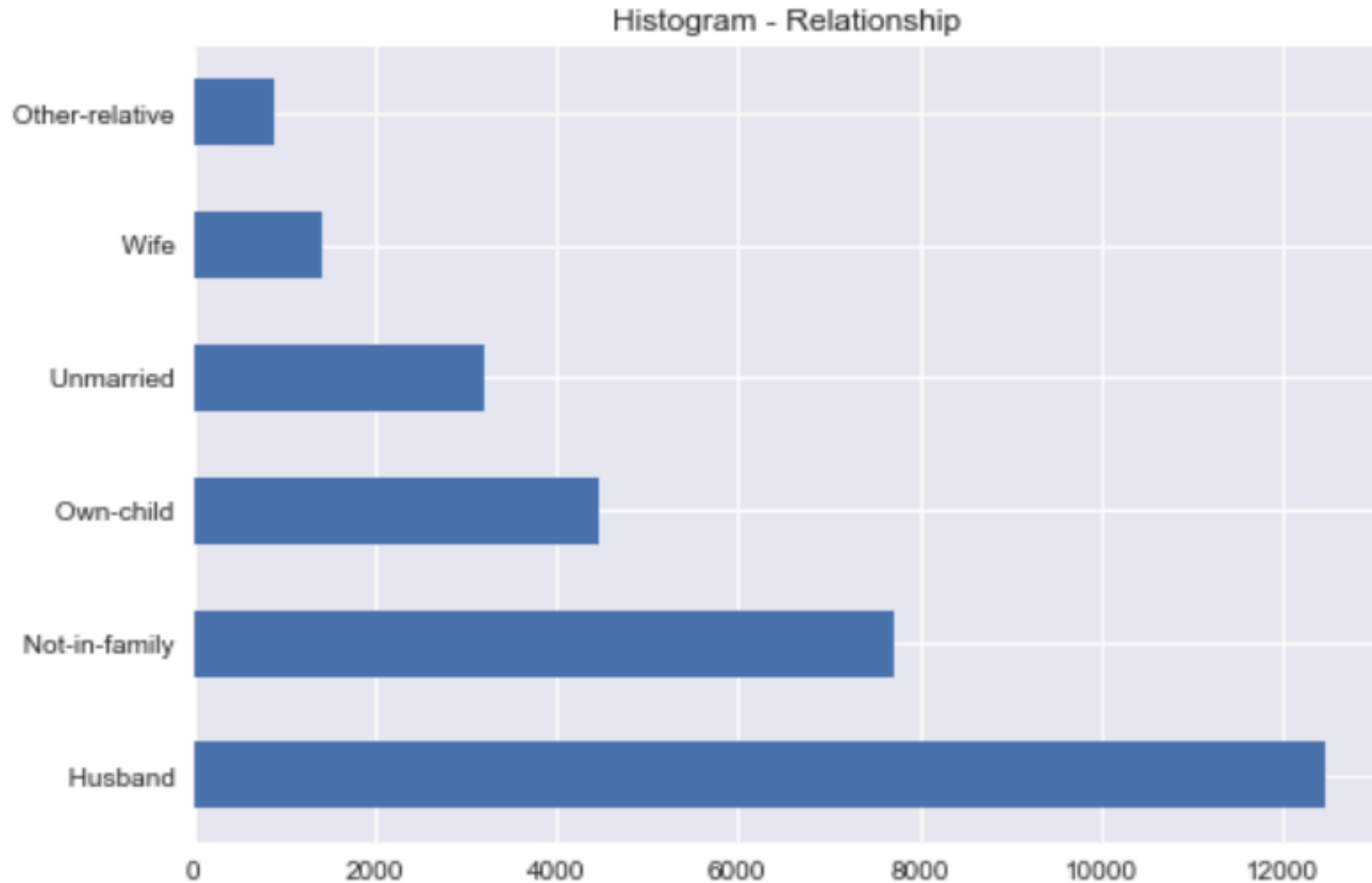
Work Class



Occupation

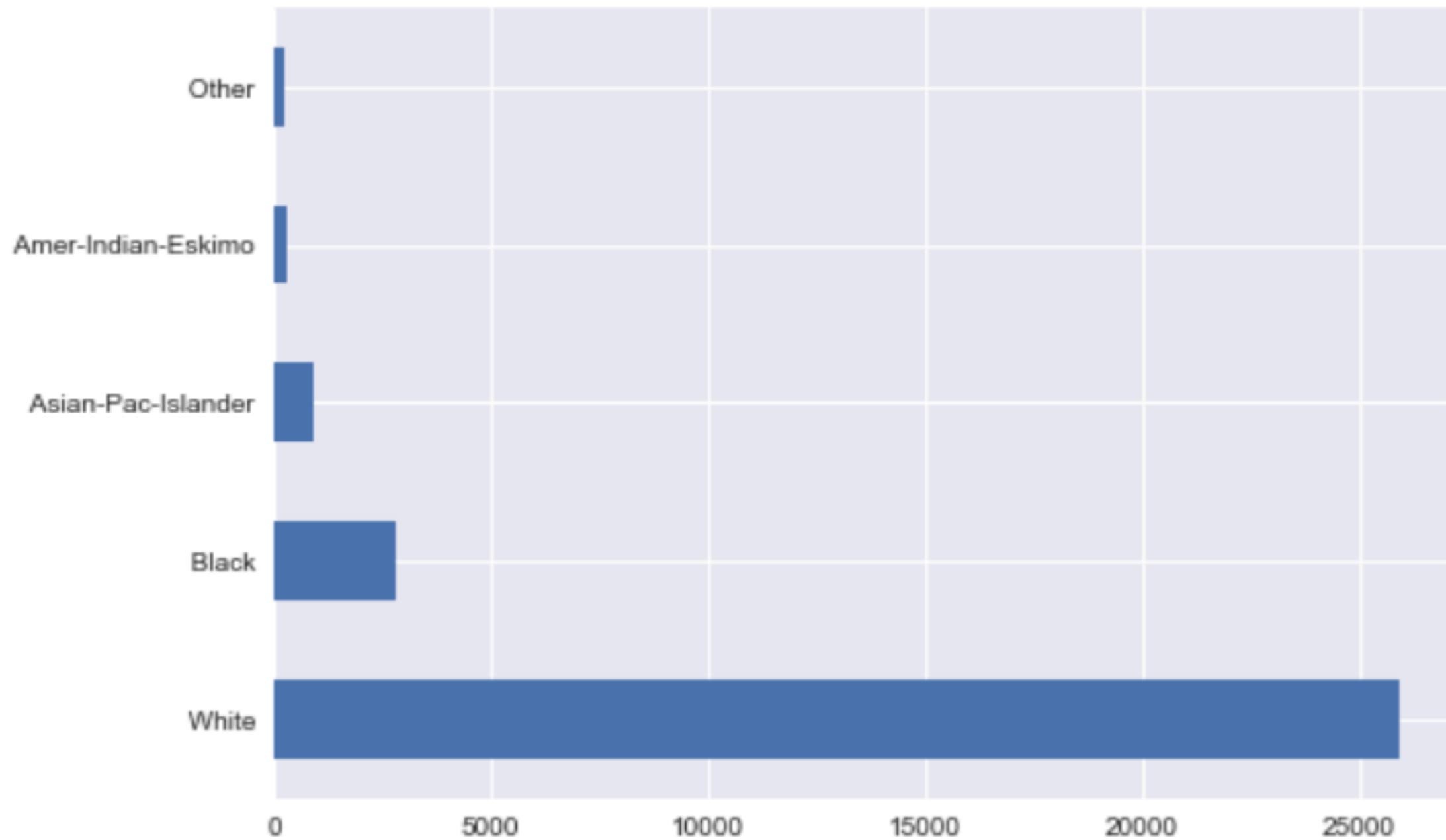


Relationship

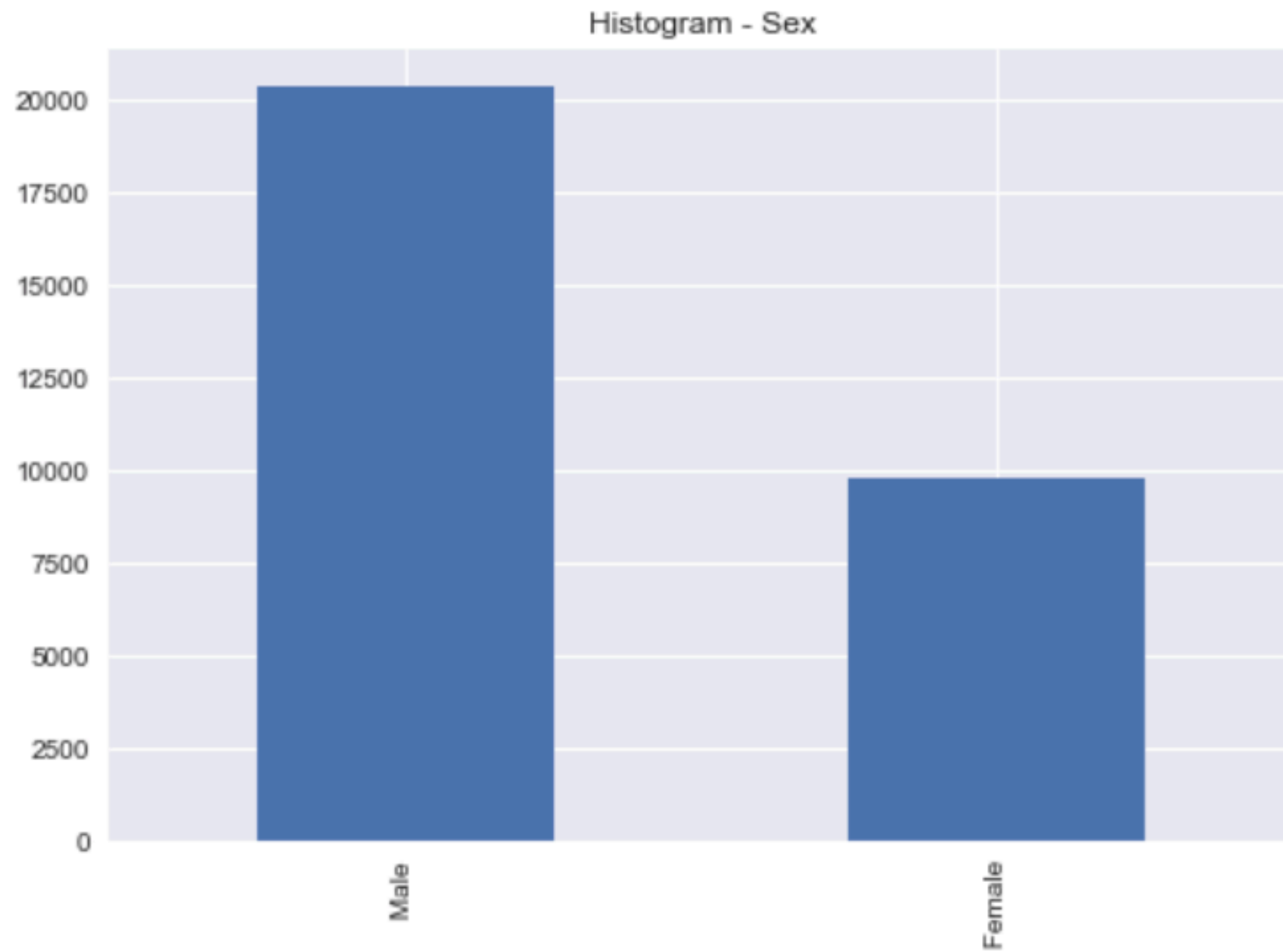


Race

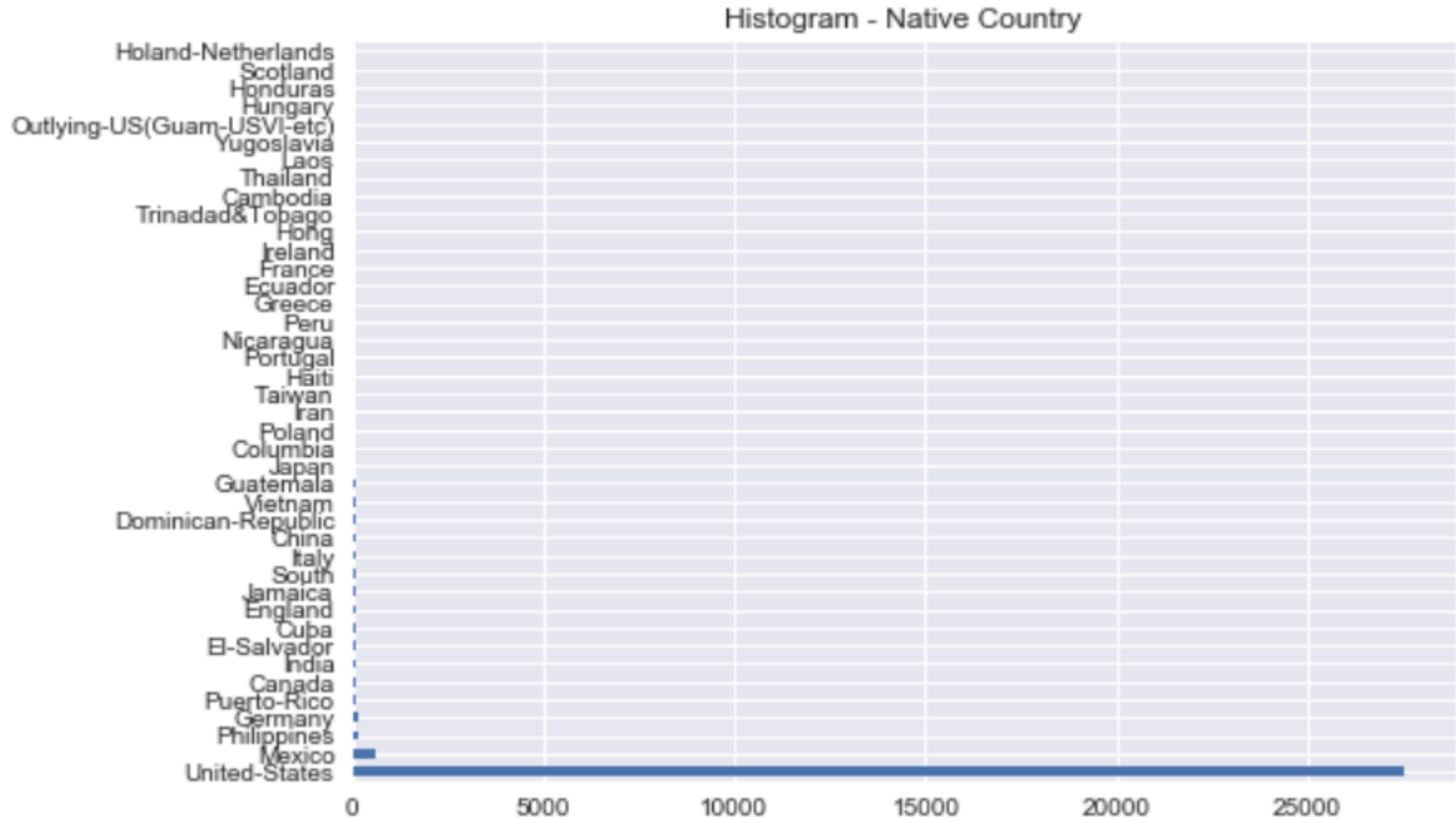
Histogram - Race



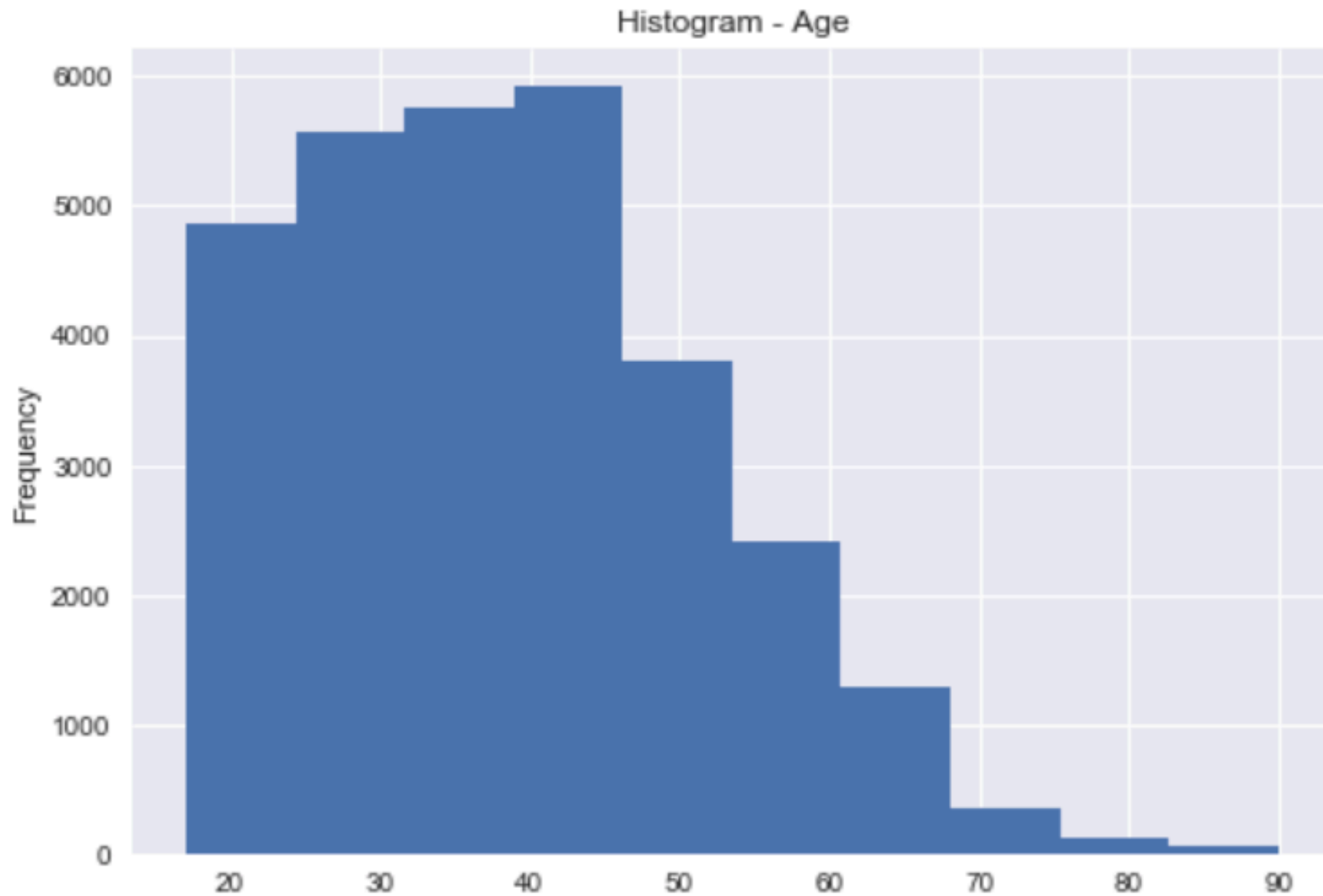
Sex



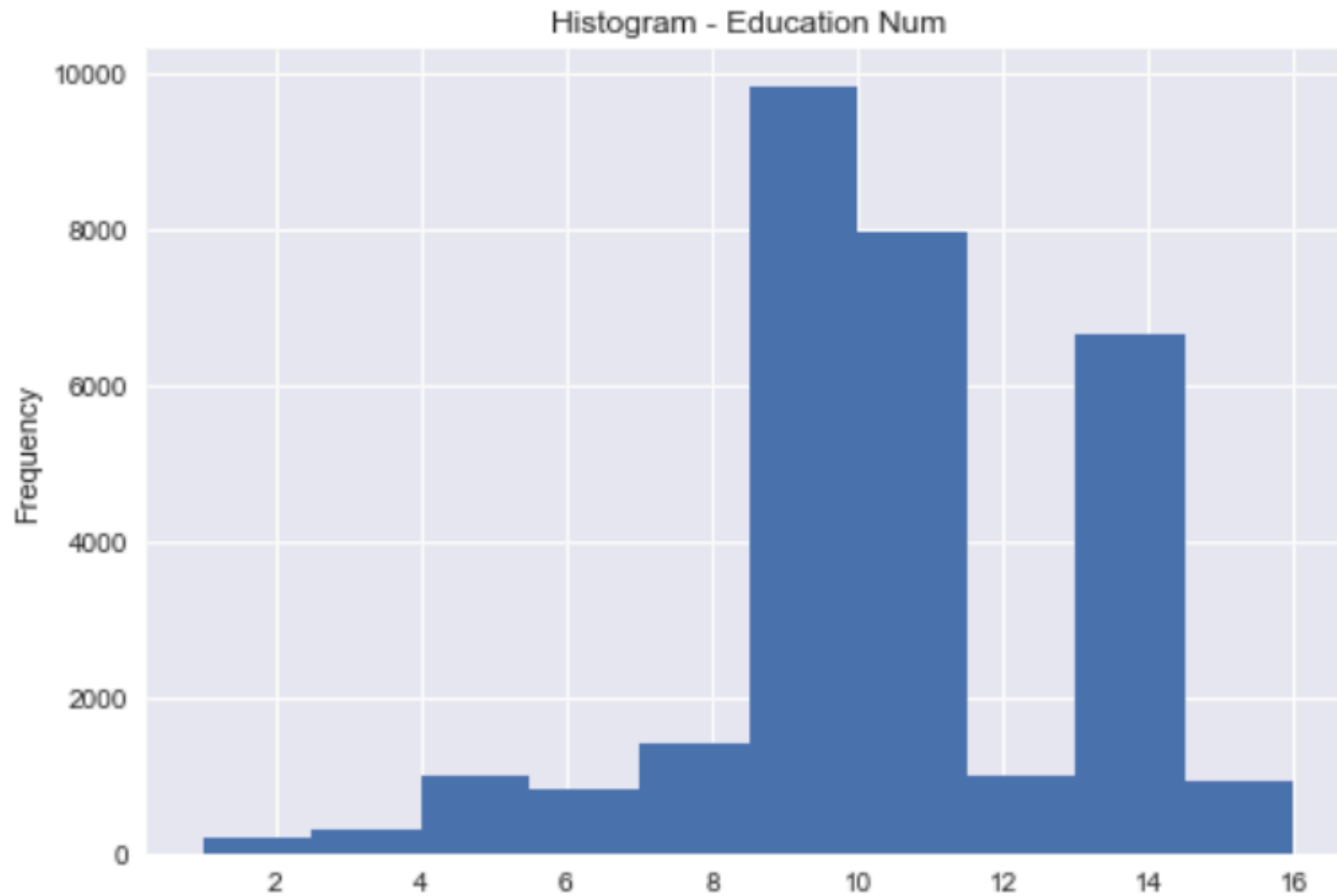
Native Country



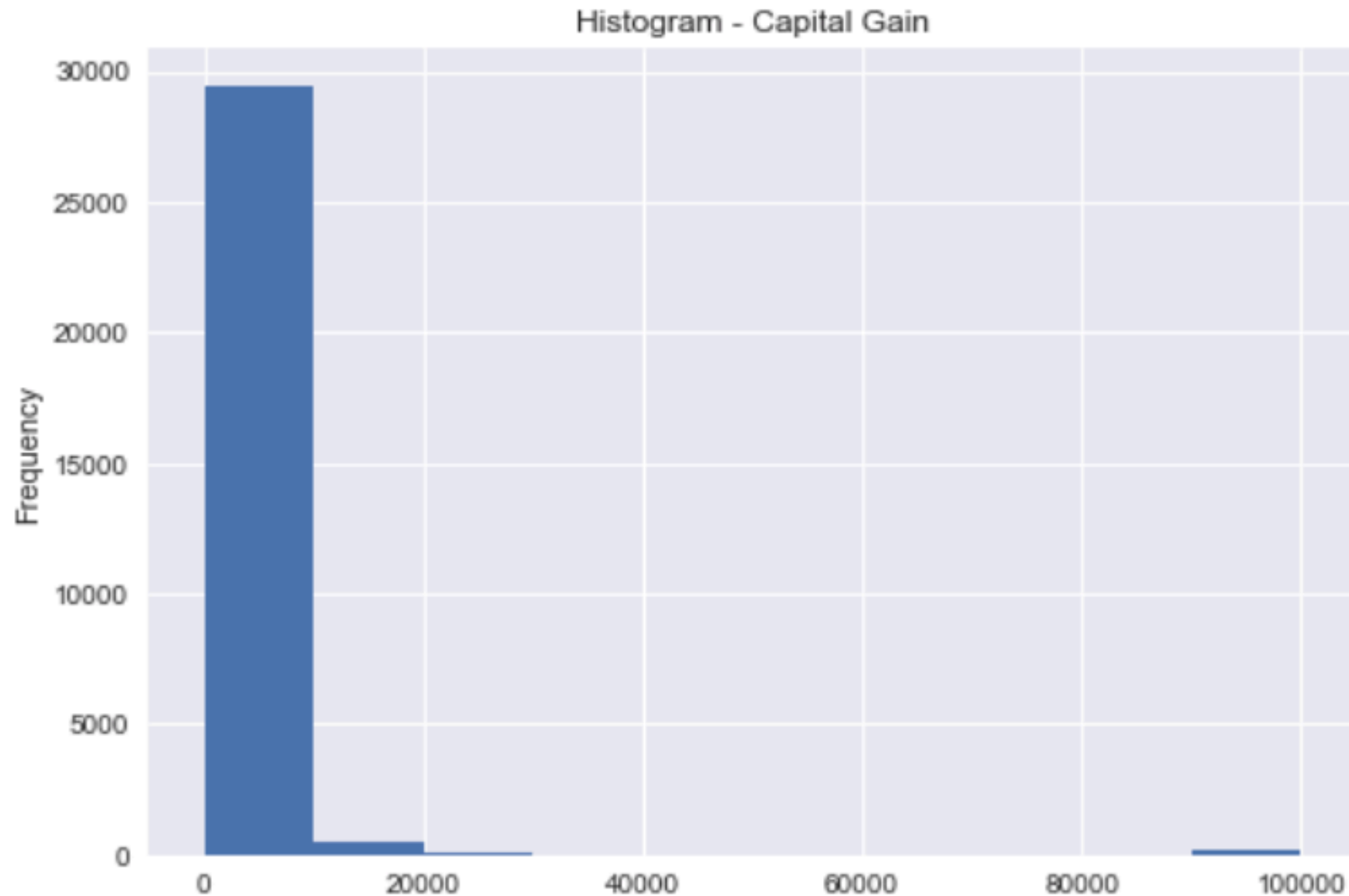
Age



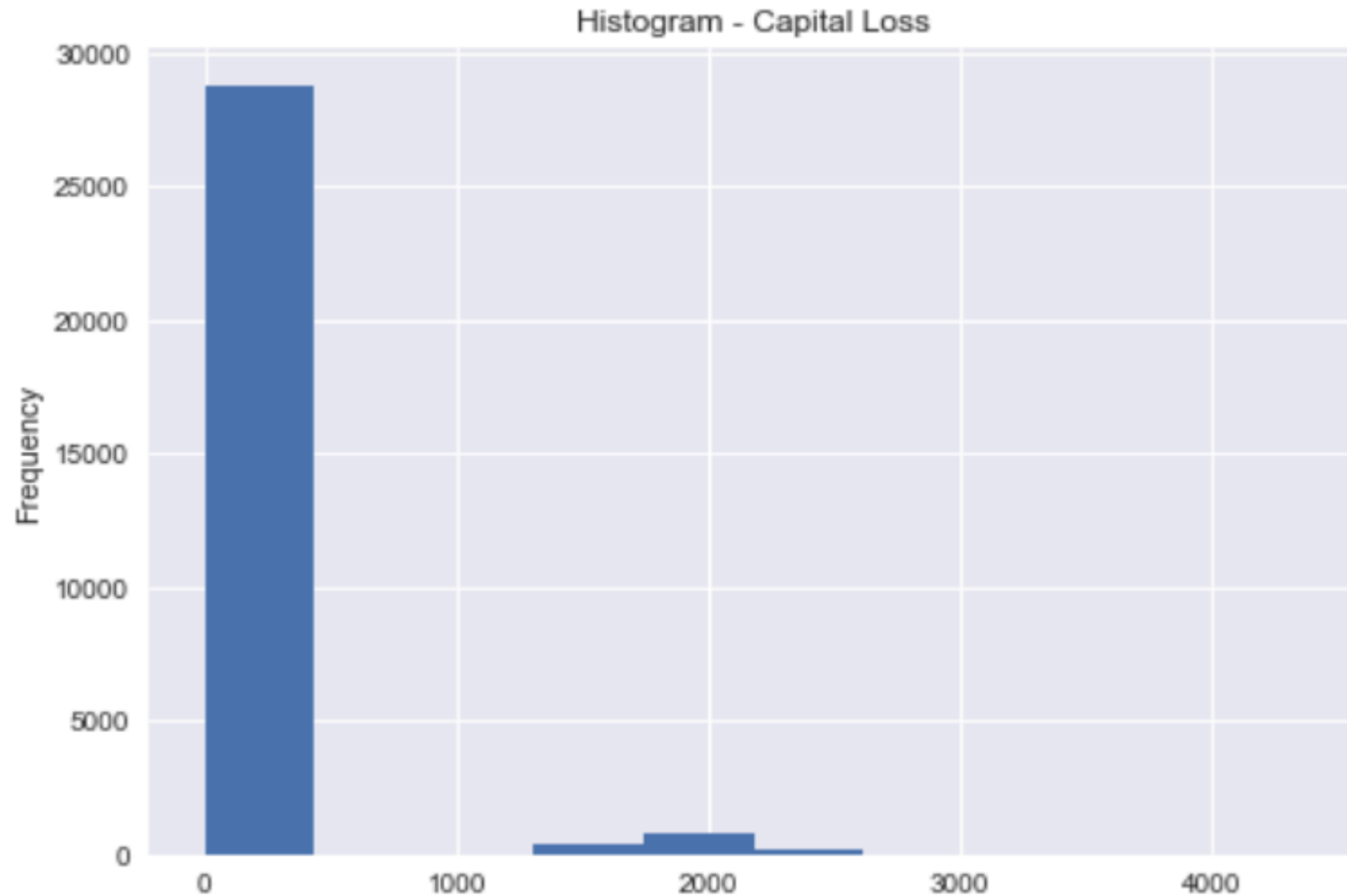
Education - Num



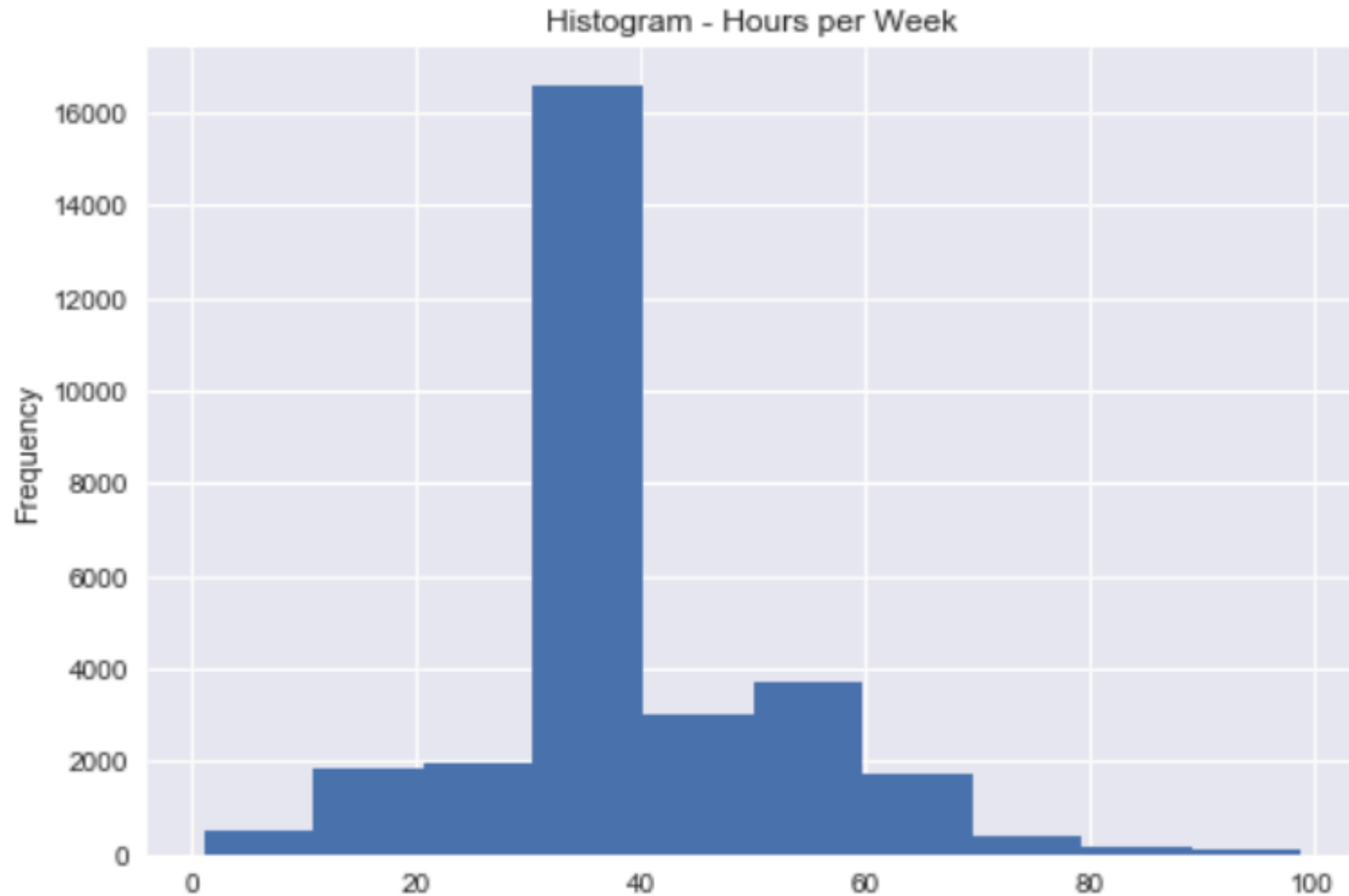
Capital Gain



Capital Loss



Hours Per Week



Model Training and Performance Measurements

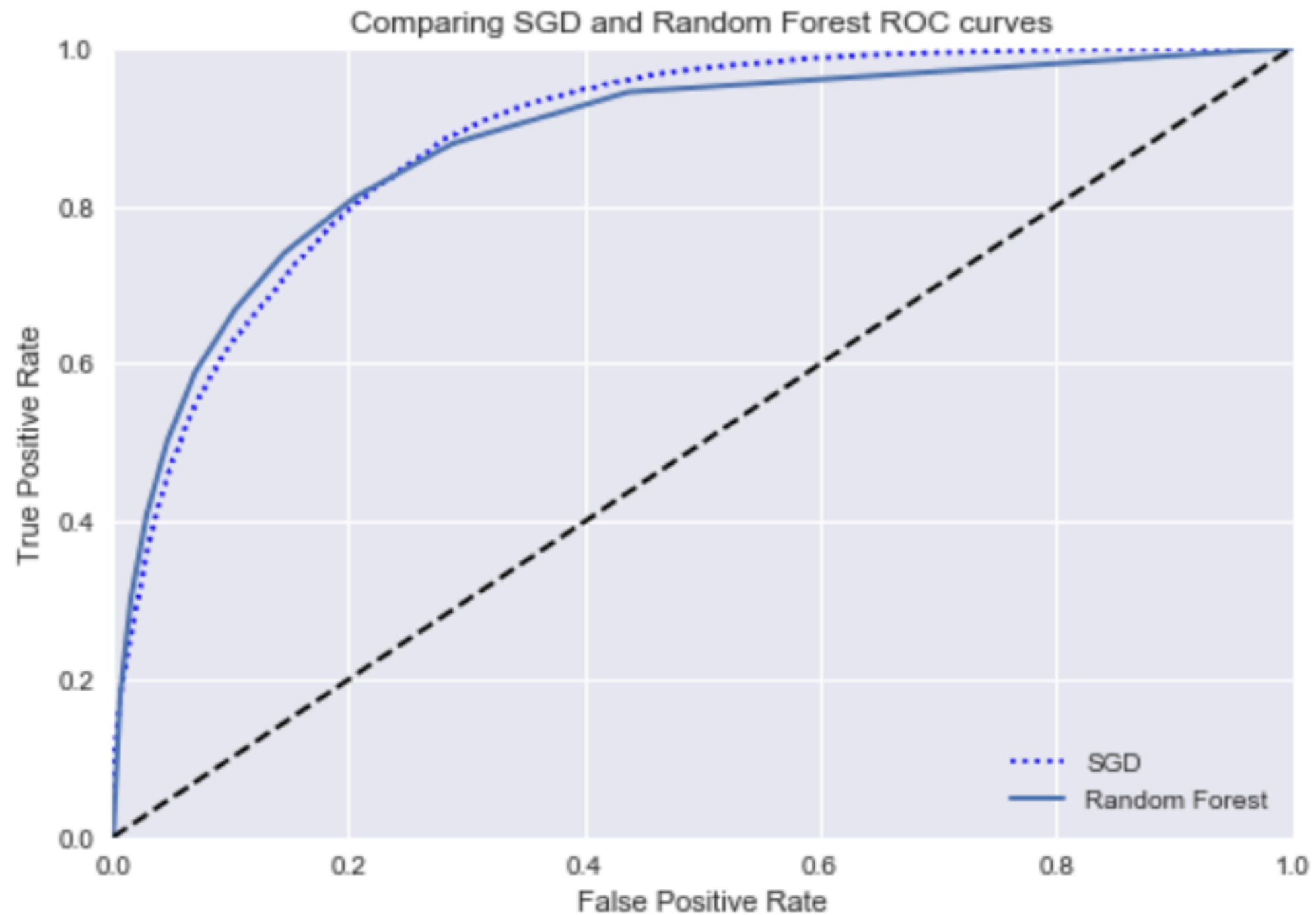
Cross-Validation Accuracy Scores

	Random Forest	SGD
1st CV run	0.845236	0.816590
2nd CV run	0.841058	0.838074
3rd CV run	0.849597	0.844425

Confusion Matrix Scores

	Random Forest	SGD
Precision Score	0.736991	0.684533
Recall Score	0.588572	0.610682
F1 Score	0.654473	0.645502

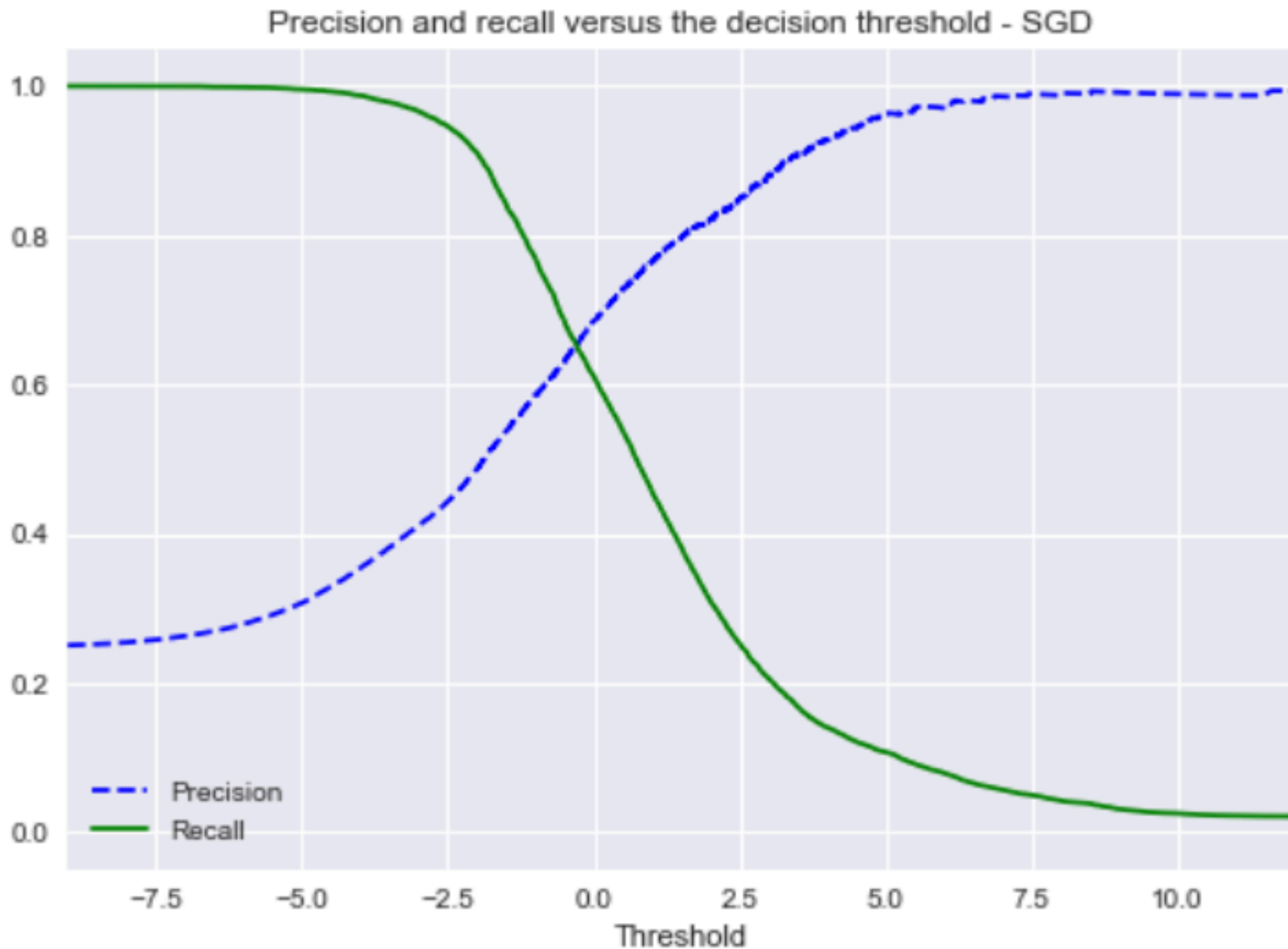
ROC Curves



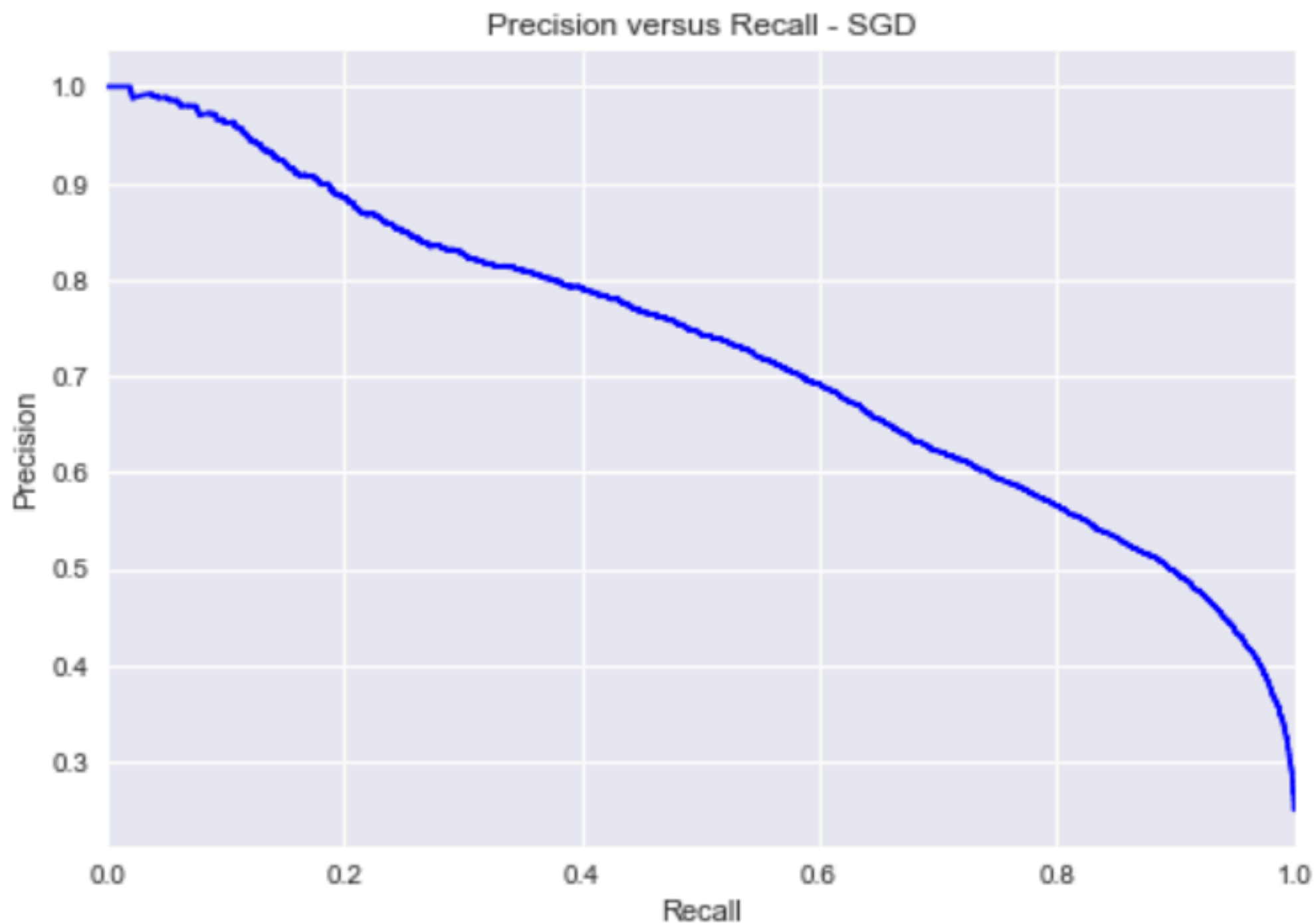
ROC AUC Scores

	Random Forest	SGD
ROC AUC Scores	0.87842	0.885639

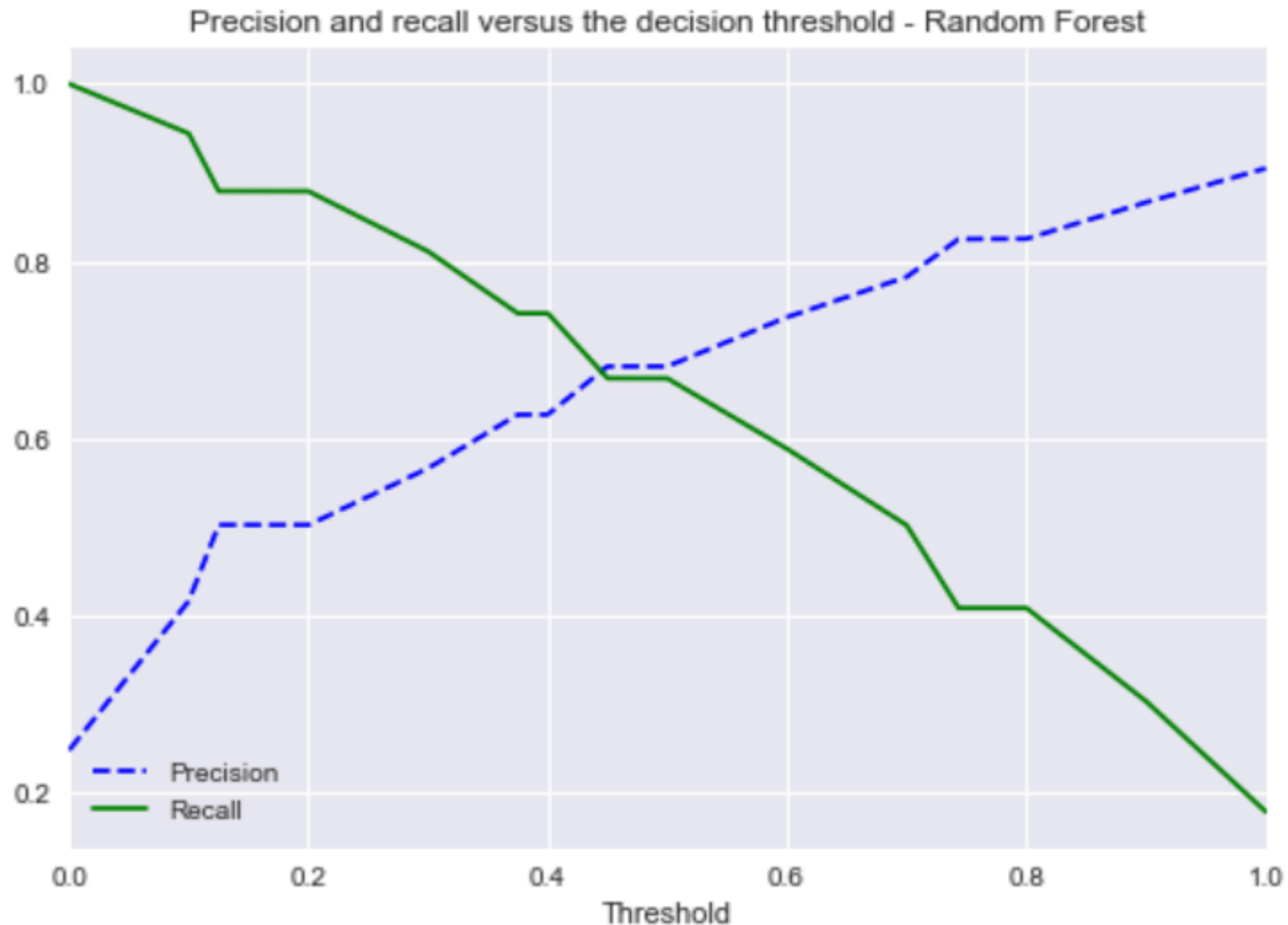
Precision and Recall - SGD



Precision and Recall - SGD

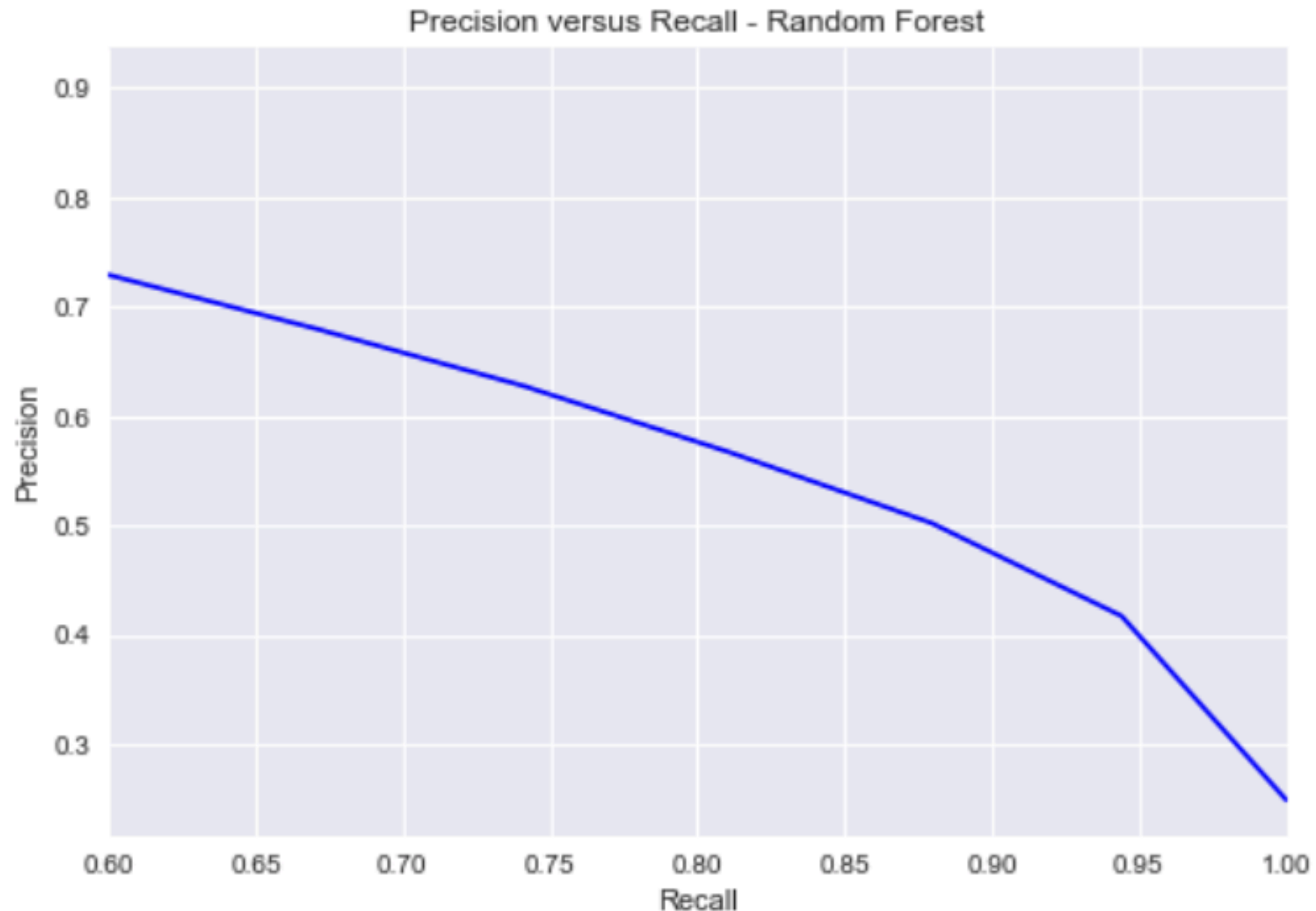


Precision and Recall Random Forest



Precision and Recall

Random Forest



Next Steps

- Remove outliers and check effect on performance measurements
- Choose best precision/recall tradeoff for best classifier
- Fine tune best classifier with Grid Search
- Perform data preprocessing and transformation on the test set
- Evaluate best performing classifier on the test set