

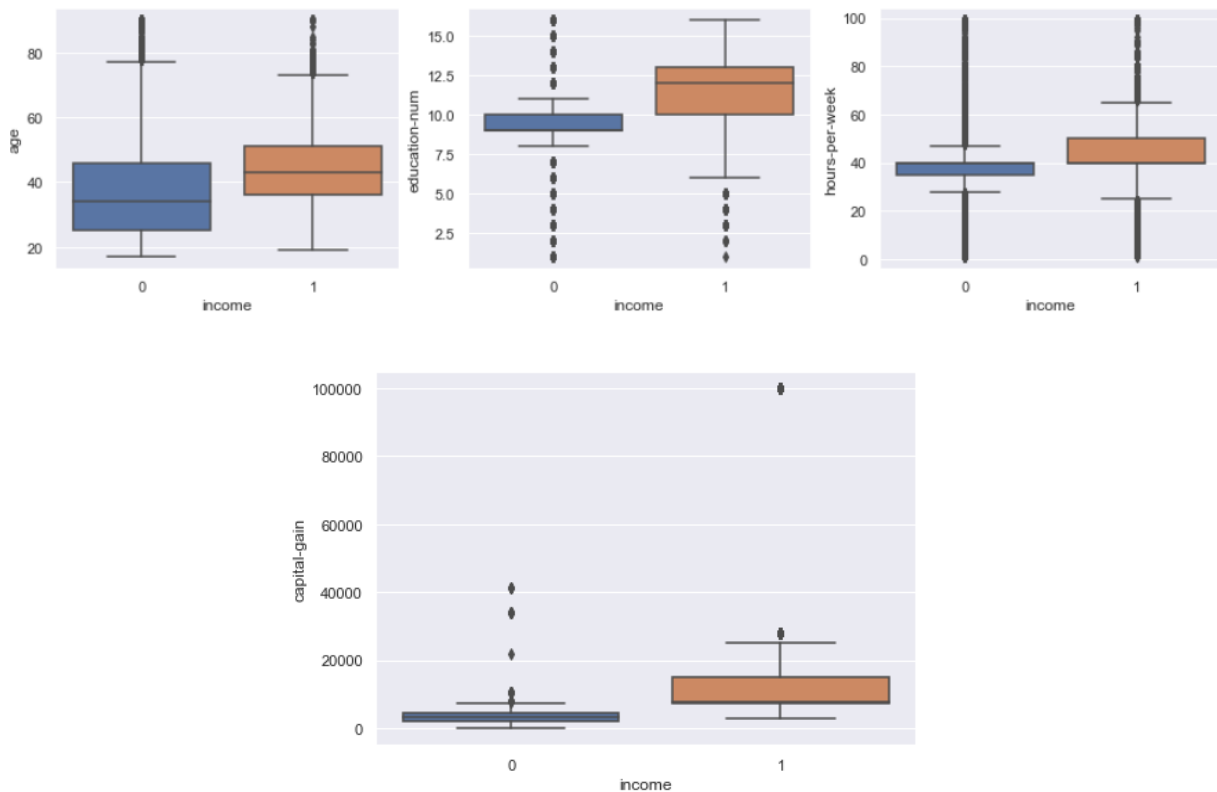
Capstone 2

Exploratory Data Analysis

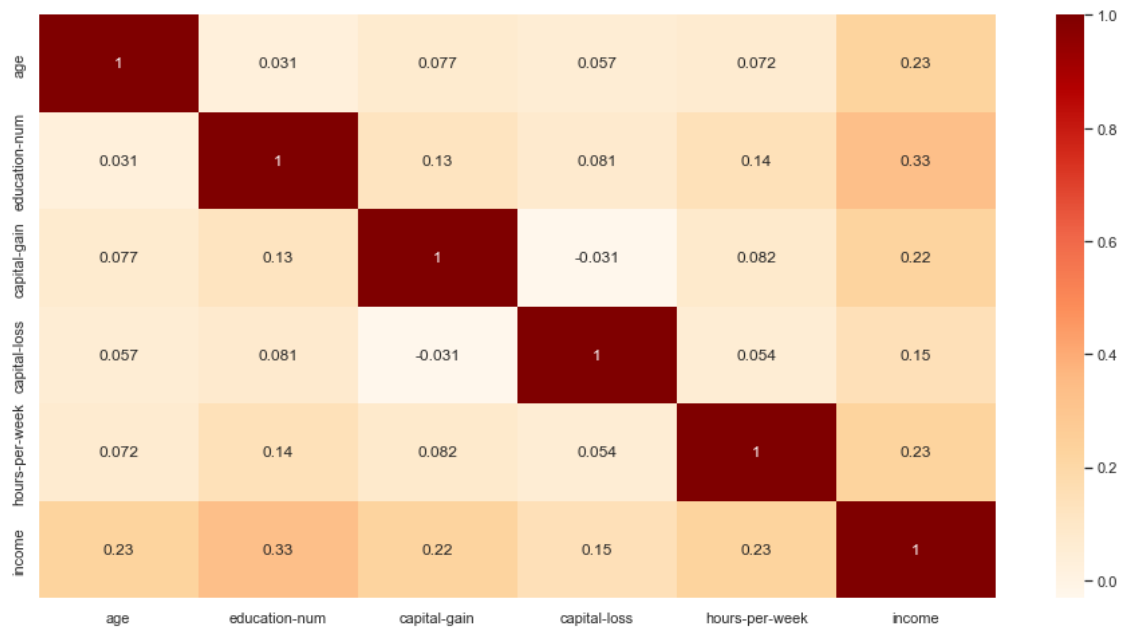
We loaded the dataset from our previous data wrangling exercise. The distribution of our target variable “income” is shown below. About 75% are people with income less or equal than \$50k, the remaining 25% have income greater than \$50k



We then proceeded to do visual inspection of the numerical features. Education years, hours per week and capital-gain showed some differentiation against our income target variable. Also for age, but less pronounced. Capital loss didn't show much variability.

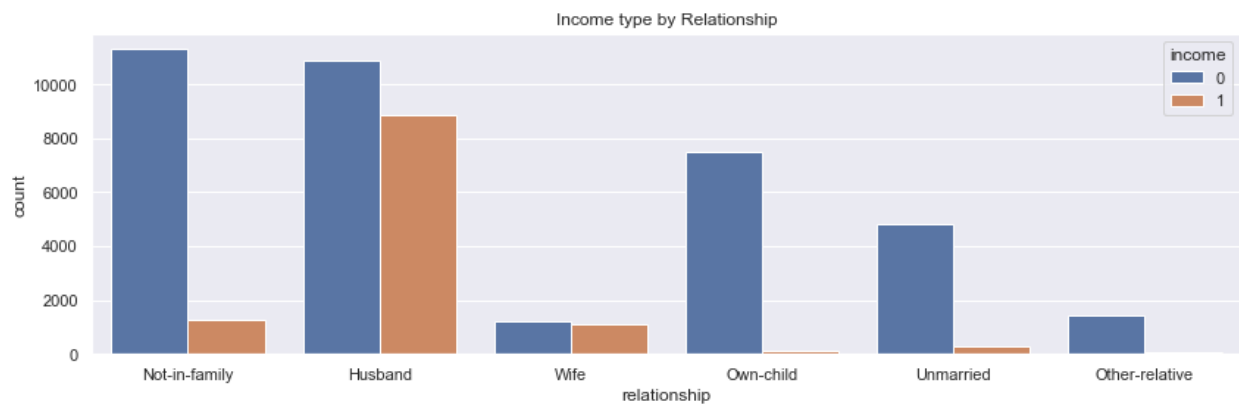


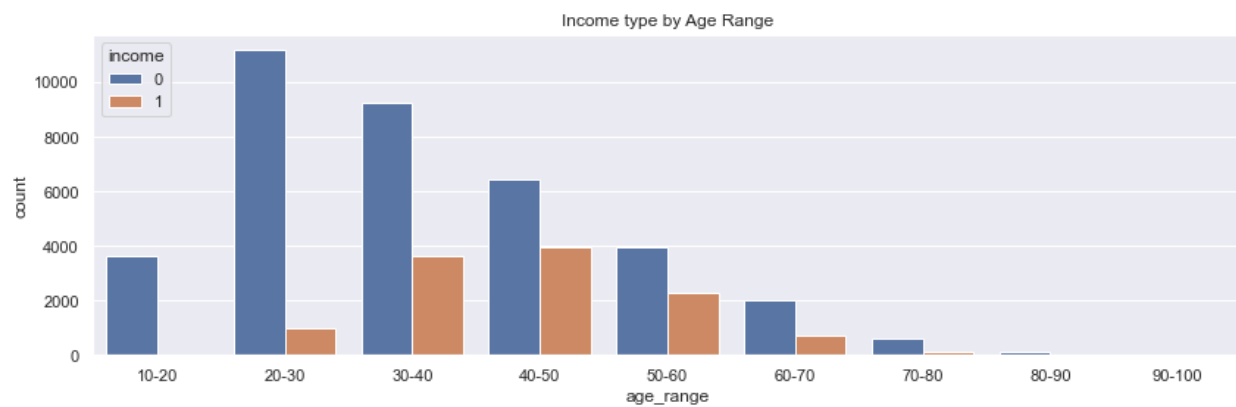
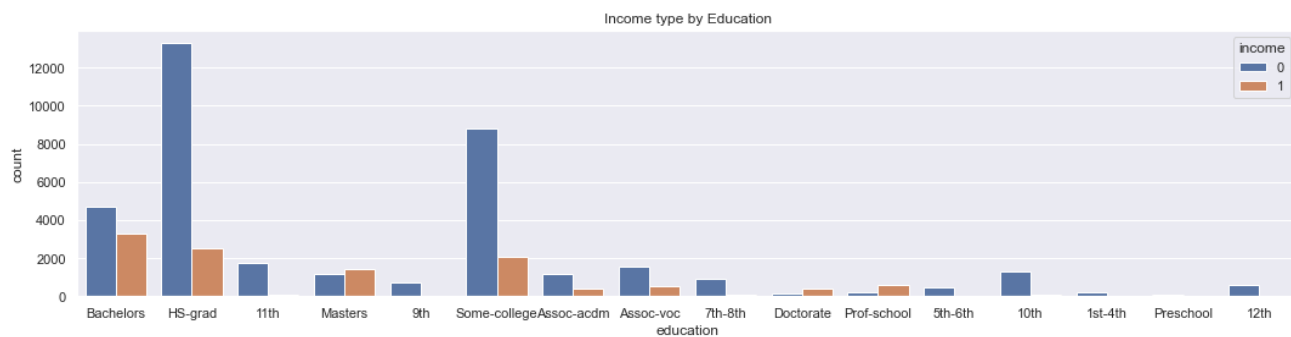
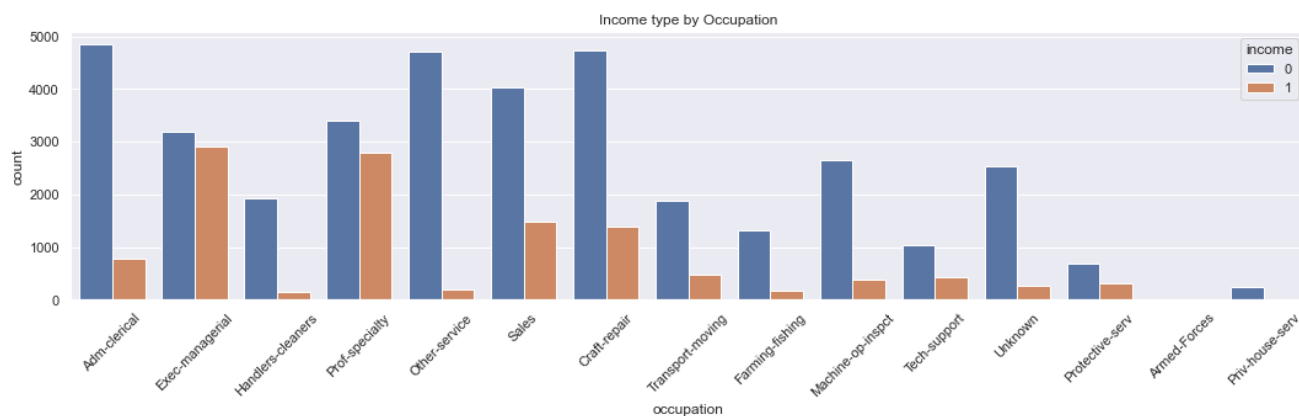
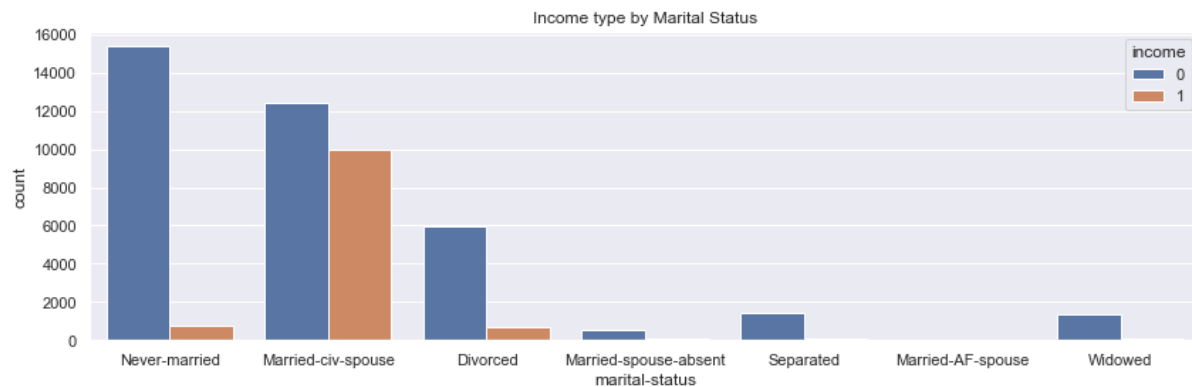
Correlation heatmap for numerical variables also confirmed these 4 features having some correlation with target variable income.



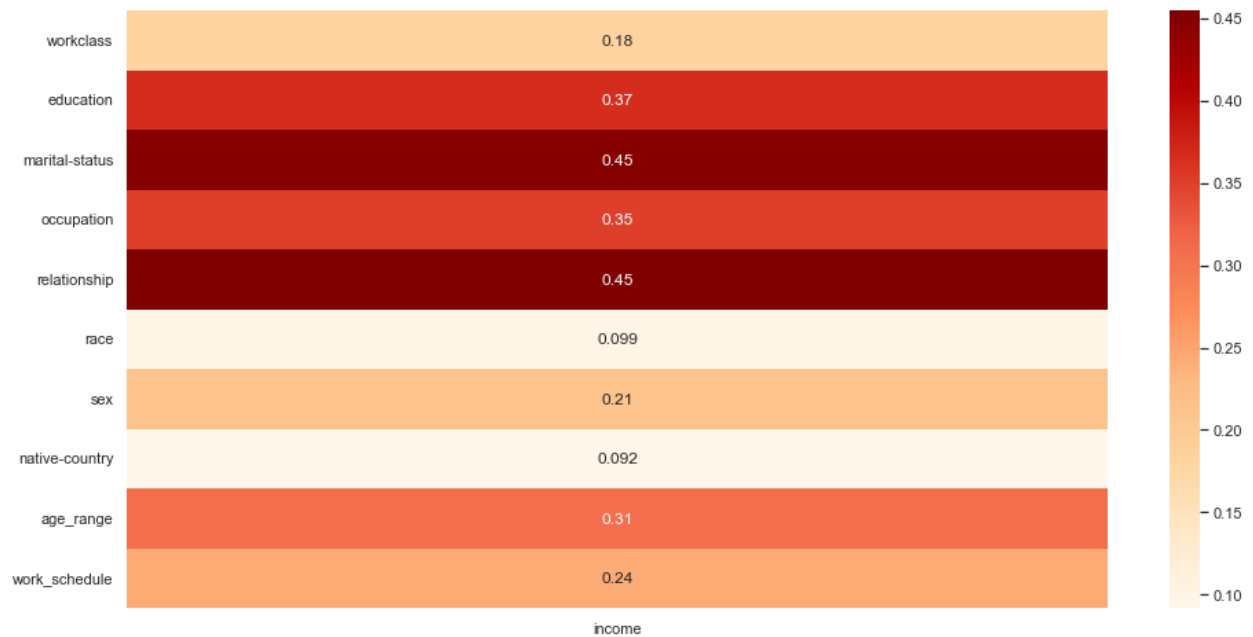
Before proceeding to visually inspect our categorical features, we added 2 more features: age range and work schedule based on hours worked per week.

We then proceeded with visual charts for our categorical variables with the following features being the most relevant: Relationship, marital status, education, occupation and age range.

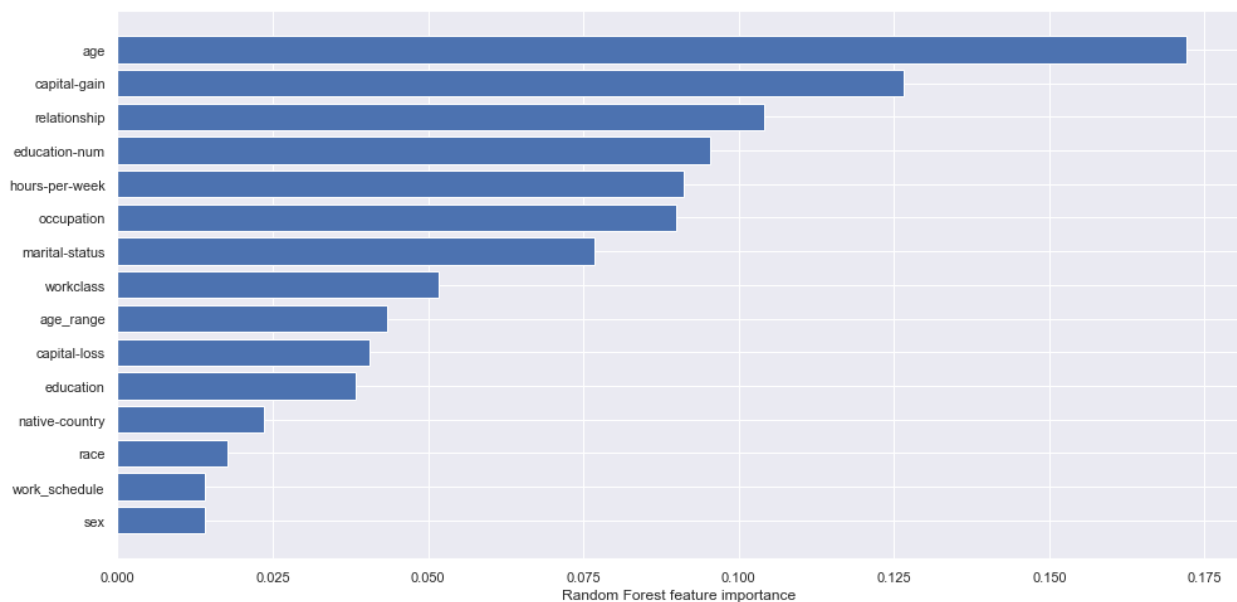




Correlation heatmap for categorical variables confirmed these features having some correlation with target variable income.



Finally we built a basic Random Forest model to obtain the features importances, with the following results:



This gave us a different point of view of some features we had seen before, but not strongly correlated with our target variable income. Age and capital gain had some correlation, but we can see them here as the top 2 features with more importance.

The 2 features we added (*workclass* & *age range*) might be redundant, since their original features have more importance than the bucketed new features.