# **Predicting Depression**

Ting Ting | Data Science Part-Time Capstone Project

# Problem Statement

According to the World Health Organization, "**Depression is one of the leading causes of disability. Suicide is the second leading cause of death among 15-29-year-olds.** People with severe mental health conditions die prematurely – as much as two decades early – due to preventable physical conditions.

# Project Goal

This project seeks to use information from an **online assessment** which includes test scores on the **Depression, Anxiety and Stress Scale, Ten Personality Item Scale** and demographic information to **predict depression** in an individual and **explore possible interactions between personality, demographic and depression.**

# Approach and Process

- Computing scores on the Depression, Anxiety and Stress Scale and Ten-Personality Item Scale
- Cleaning Data
- Exploratory Data Analysis
- Modelling and model evaluation
- Interpreting the results

Table 2. Severity levels.

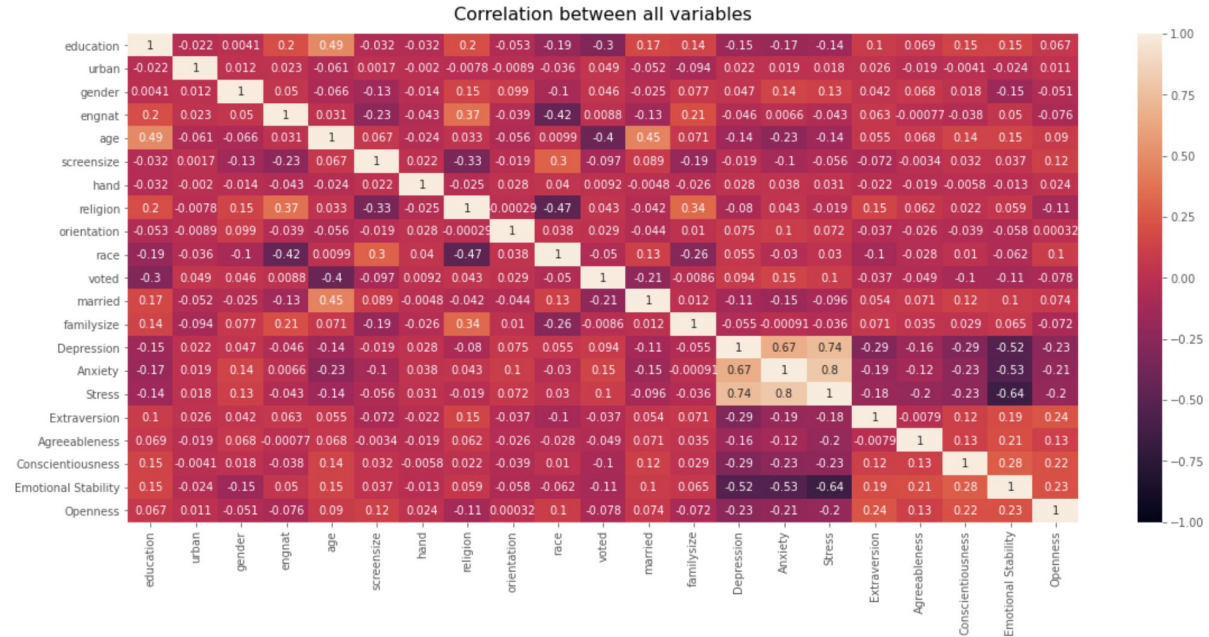|  | Anxiety | Depression | Stress |
|---|---|---|---|
| Normal | 0-7 | 0-9 | 0-14 |
| Mild | 8-9 | 10-13 | 15-18 |
| Moderate | 10-14 | 14-20 | 19-25 |
| Severe | 15-19 | 21-27 | 26-33 |
| Extremely severe | 20+ | 28+ | 33+ |

# Cleaning Data

- Data source was from an **online assessment** with 3 different components

- Only included data where **test takers marked that the results can be used for research**

- Data set was relatively clean and required **minimal data cleaning** (removing outliers and null values were in features that are inconsequential to our target e.g. major)
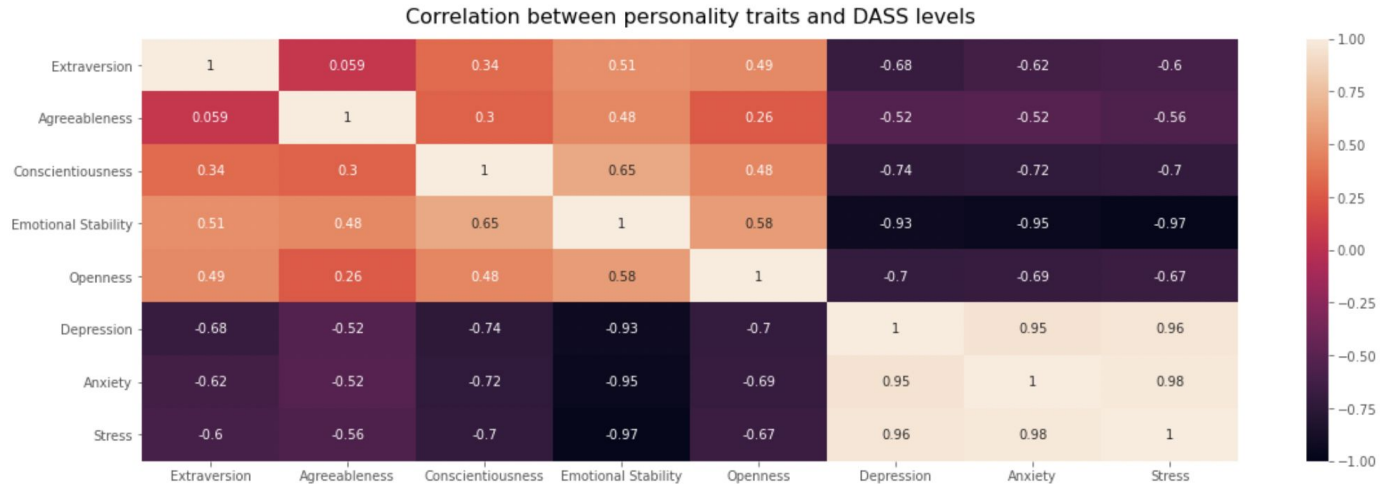
# Exploratory Data Analysis

● None of the independent variables saw particularly high correlation with another independent variable
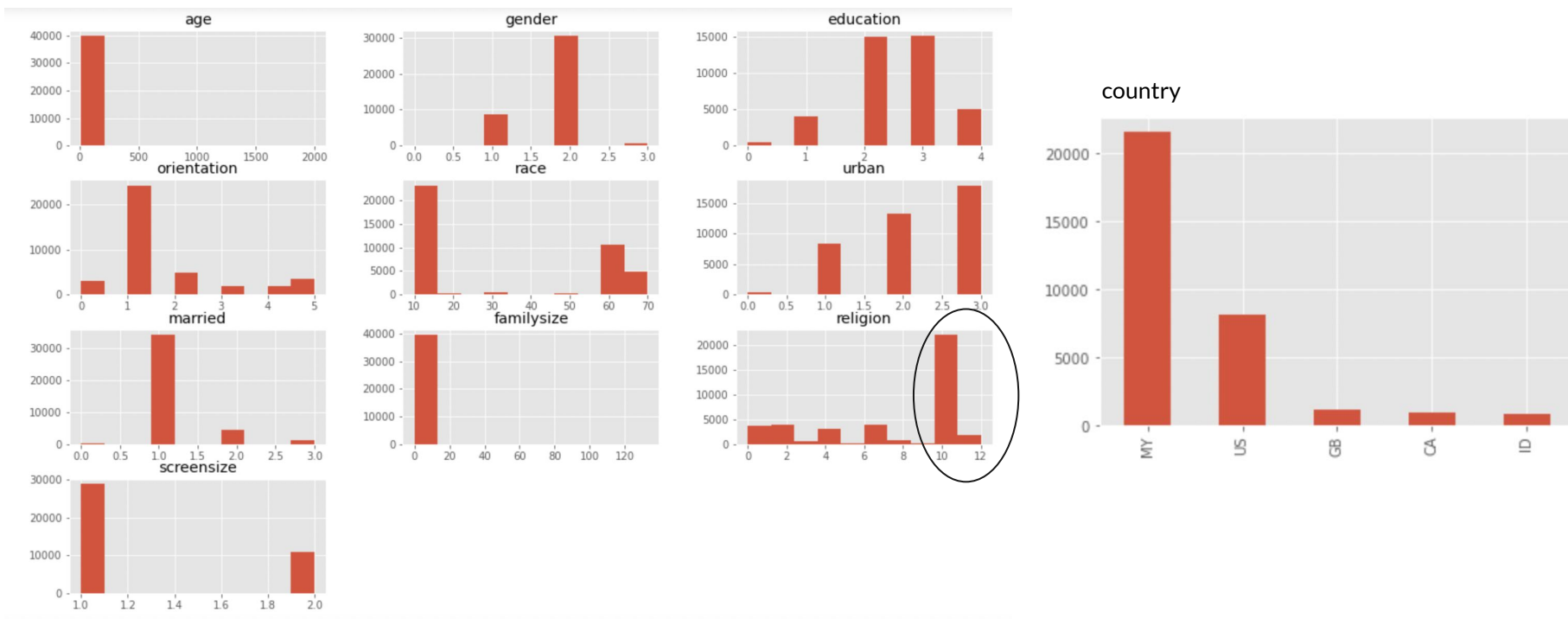

Correlation between all variables

# Exploratory Data Analysis

- Isolating personality traits and depression, anxiety and stress levels, anxiety and stress correlates strongly with depression and emotional stability negatively correlatives with all 3

Correlation between personality traits and DASS levels

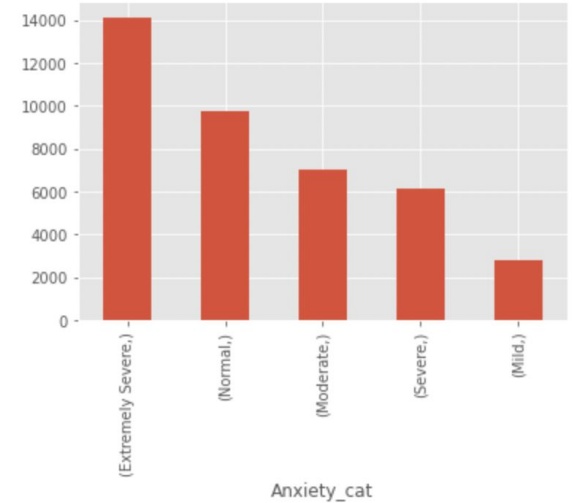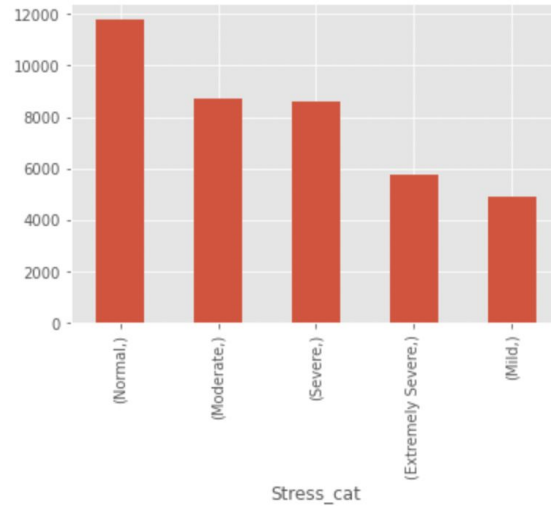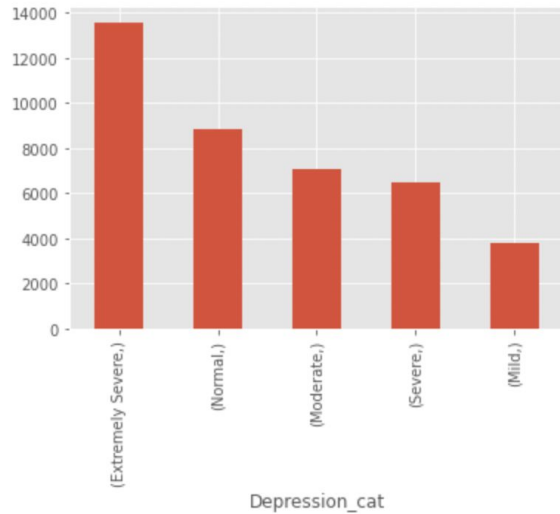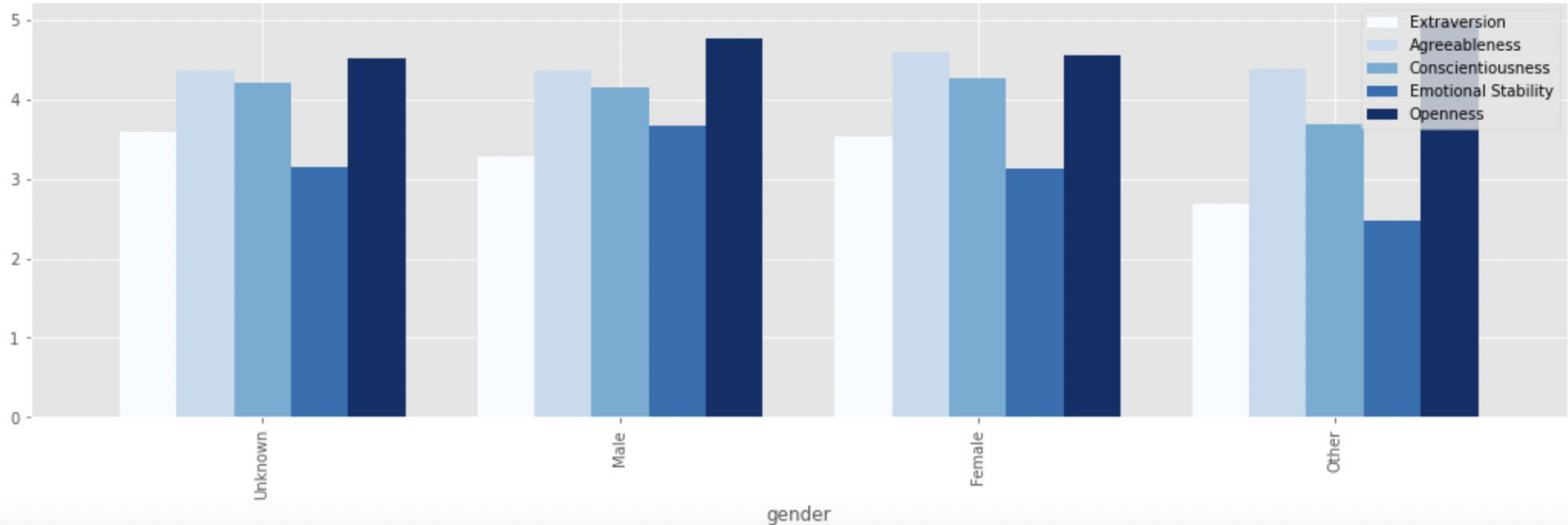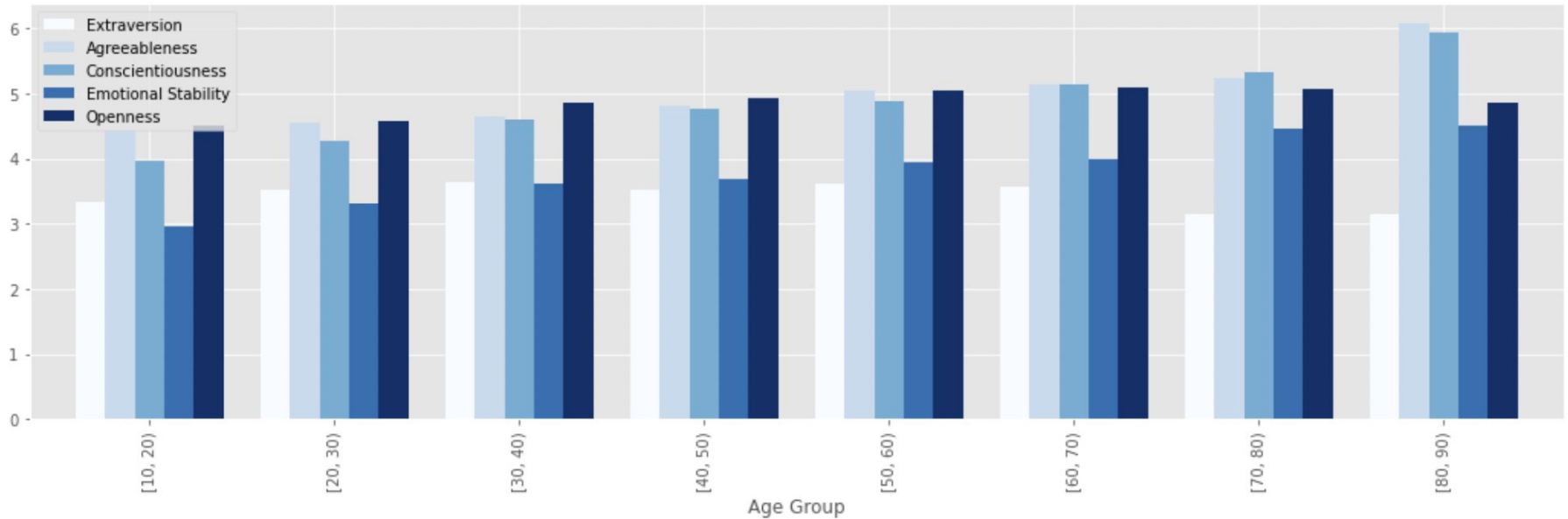|                     | Extraversion | Agreeableness | Conscientiousness | Emotional Stability | Openness | Depression | Anxiety | Stress |
|---------------------|--------------|---------------|-------------------|---------------------|----------|------------|---------|--------|
| Extraversion        | 1            | 0.059         | 0.34              | 0.51                | 0.49     | -0.68      | -0.62   | -0.6   |
| Agreeableness       | 0.059        | 1             | 0.3               | 0.48                | 0.26     | -0.52      | -0.52   | -0.56  |
| Conscientiousness   | 0.34         | 0.3           | 1                 | 0.65                | 0.48     | -0.74      | -0.72   | -0.7   |
| Emotional Stability | 0.51         | 0.48          | 0.65              | 1                   | 0.58     | -0.93      | -0.95   | -0.97  |
| Openness            | 0.49         | 0.26          | 0.48              | 0.58                | 1        | -0.7       | -0.69   | -0.67  |
| Depression          | -0.68        | -0.52         | -0.74             | -0.93               | -0.7     | 1          | 0.95    | 0.96   |
| Anxiety             | -0.62        | -0.52         | -0.72             | -0.95               | -0.69    | 0.95       | 1       | 0.98   |
| Stress              | -0.6         | -0.56         | -0.7              | -0.97               | -0.67    | 0.96       | 0.98    | 1      |

# Who were the test takers?

# How participants scored on the DASS

# How participants scored on the Personality Scale by gender

# How participants scored on the Personality Scale by age

# Model and solution

List of models I applied before parameter hypertuning and their corresponding accuracy and f1 scores

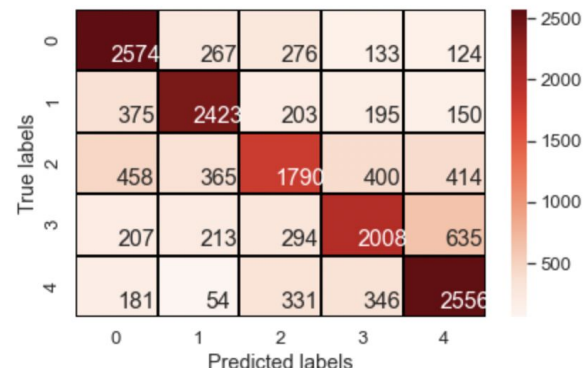| Model | Before SMOTE | After SMOTE |
|---|---|---|
| Random Forest Classifier | 0.51, 0.35 | 0.67, 0.67 |
| Support Vector Classification | 0.52, 0.33 | - |
| Logistic Regression | 0.52, 0.33 | 0.49, 0.47 |
| K Nearest Neighbors Classifier | 0.45, 0.36 | 0.65, 0.64 |
| XGBoost | 0.51, 0.37 | 0.55, 0.54 |
| Dumb Classifier | 0.34 Zero Rate, 0.23 Random Weight Guessing | |

# Performance evaluation

- SVC scored similarly to other models but needed 35X more run time. Dropped after the first round of model evaluation
- Through randomized search and hyperparameter tuning, we found the best model (**Random Forest Classifier**) and **improved f1 score by 5% (from 0.67 to 0.7)**
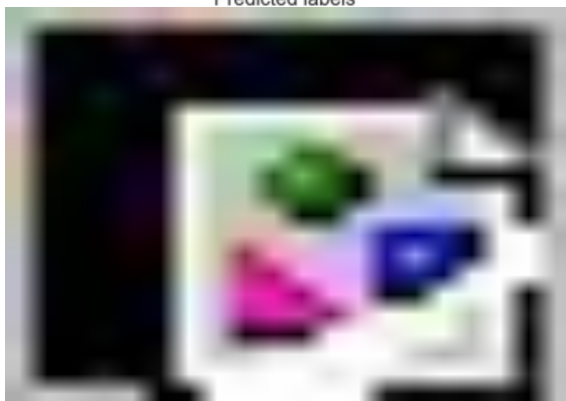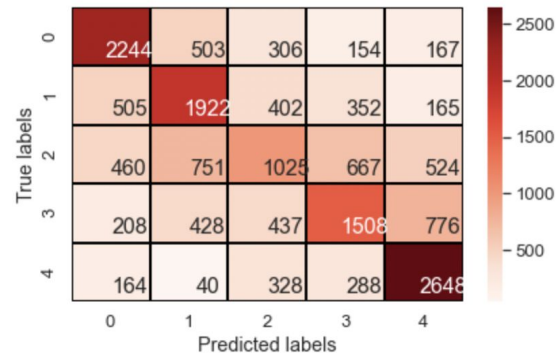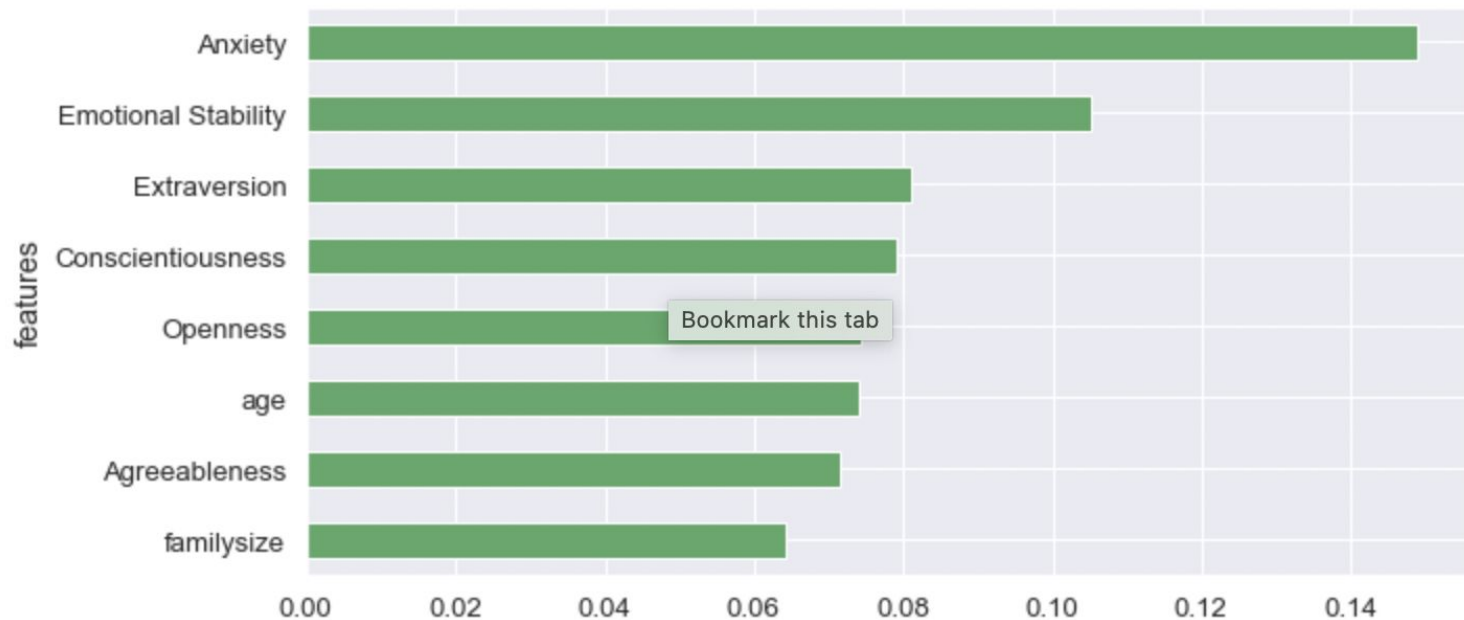
# Confusion Matrix

# Feature Importance

# Findings

- **Personality traits, particularly Emotional Stability**, are strong predictors of depression. They take precedence over test takers' background (e.g. gender, education received etc. )
- **Younger people on average scored lower on emotional stability**, males in general scored higher on emotional stability
- **Stress and anxiety levels had the largest effects on predicting depression**. They were also strongly positive correlated.

# Limitations

- Data collected was skewed to test takers in one country,  gender, age, marriage status and religion

- Classes were imbalanced: online assessments can be biased due to self-selection

# Recommendations

- Collecting more data and model improvement: **diversify dataset** to include test takers in other countries, genders and age  so that the model will be more representative of the general population
- Address class imbalance by administering test to a target group of people through random sampling or representative sampling instead