# Drug Consumption Analysis

Manasa Kandimalla
Carlos Antonio McNulty
Shyam Patel

# Problem Statement

Which personality traits (e.g., neuroticism, extraversion, openness to experience, agreeableness, conscientiousness, impulsiveness, sensation) and other factors (e.g., age, gender, education) make one susceptible to the usage of various illegal drugs?

# Data Sources

E. Fehrman, V. Egan and E. M. Mirkes (2016). UCI Machine Learning Repository [https://archive.ics.uci.edu/ml/datasets/Drug+consumption+%28quantified%29]. Leicester, UK: University of Leicester, Department of Mathematics.

# Data Sources

- 1,885 respondents
- NEO-FFI-R (neuroticism, extraversion, openness to experience, agreeableness, and conscientiousness)
- BIS-11 (impulsivity)
- ImpSS (sensation seeking)
- Level of education
- Age
- Gender
- Country of residence
- Ethnicity

# Illegal Drugs

Amphetamines

Amyl nitrite

Benzodiazepine
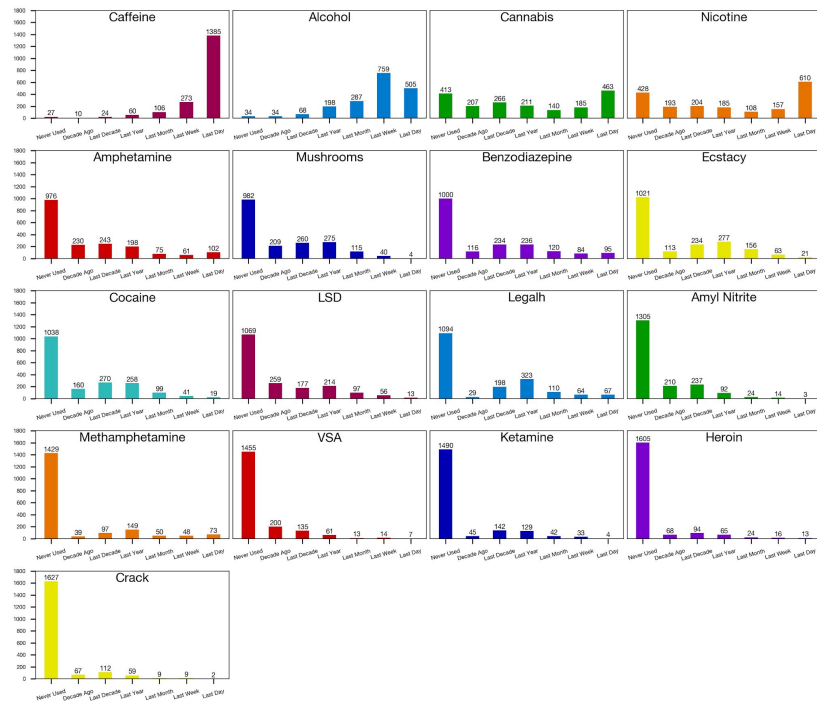
Cannabis

Cocaine

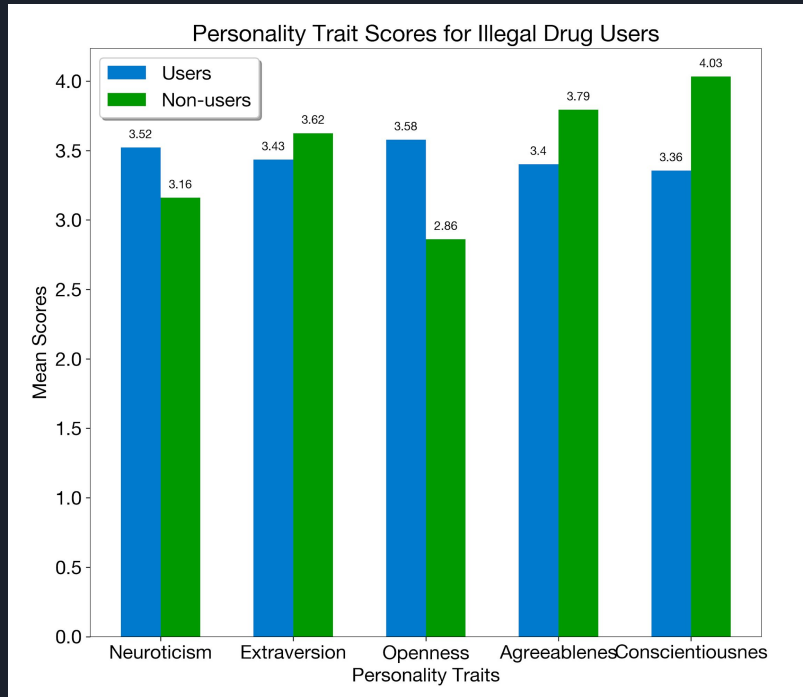Crack

Ecstasy

Heroin

Ketamine

LSD
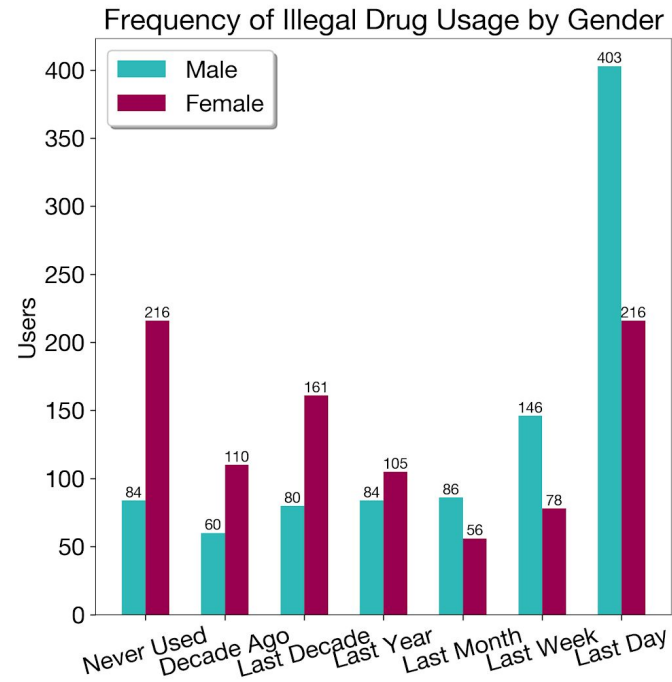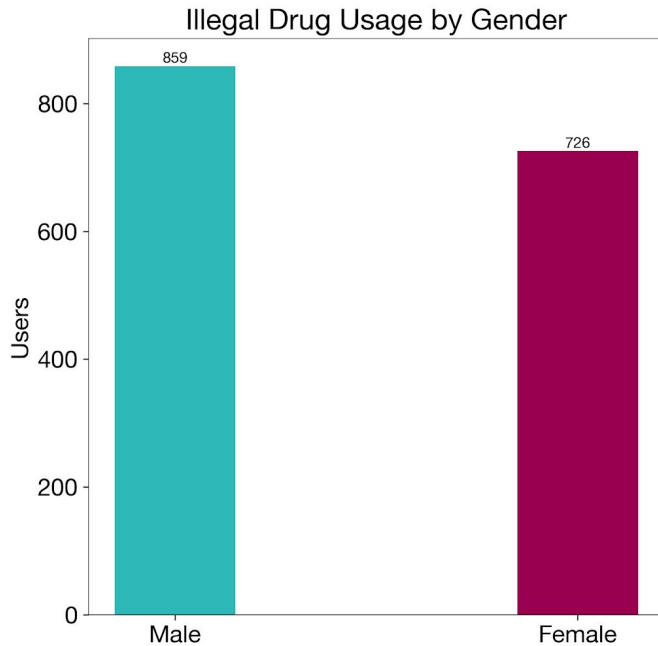
Methadone

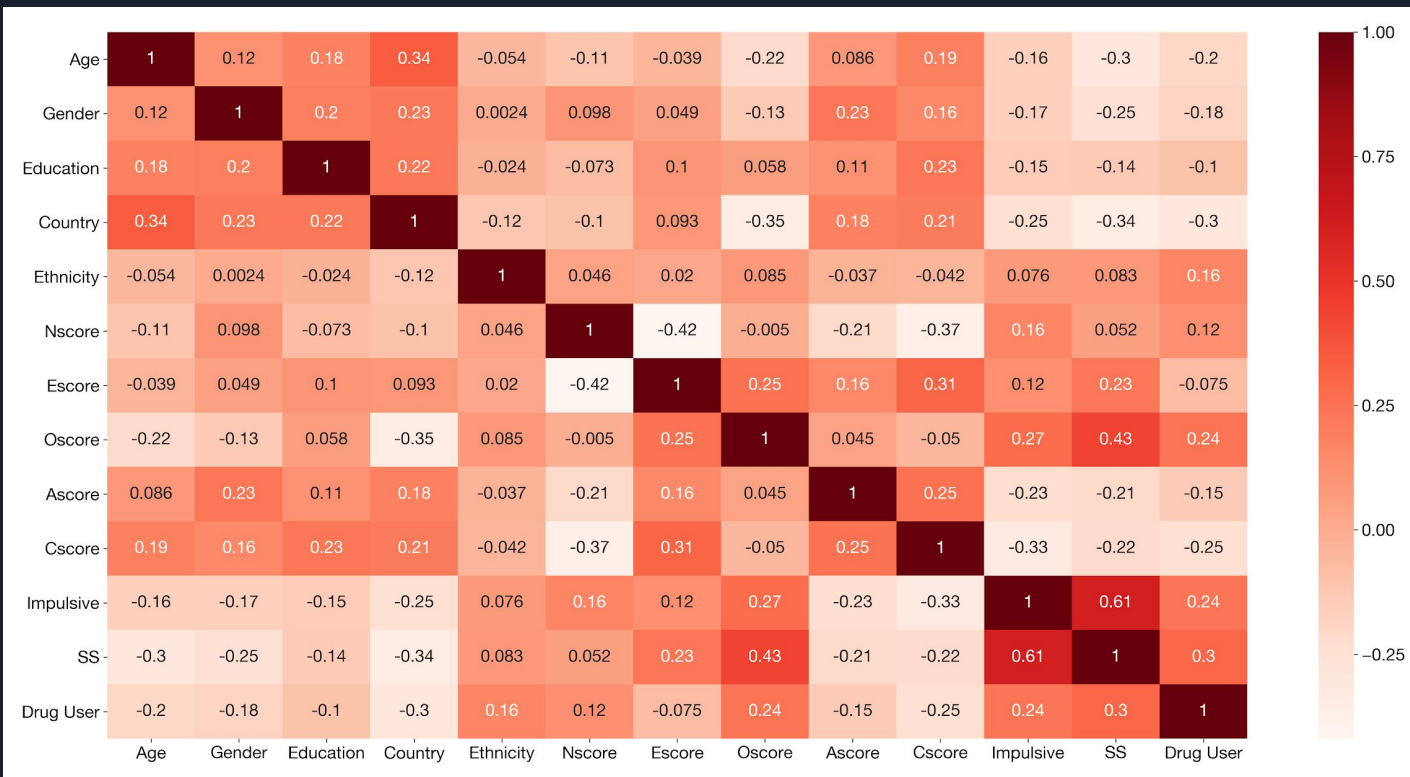Mushrooms

# Drug Usage

# Illegal Drug Usage by Personality Trait



Personality Trait Scores for Illegal Drug Users

|  | t-statistic | p-value |
|---|---|---|
| **Nscore** | 6.3434 | 2.6924e-10 |
| **Escore** | -3.4179 | 0.0003 |
| **Oscore** | 13.2261 | 1.8035e-34 |
| **Ascore** | -6.3017 | 1.8286e-10 |
| **Cscore** | -12.7530 | 1.2412e-32 |

# Illegal Drug Usage and Frequency by Gender

# Heatmap of Correlation Matrix
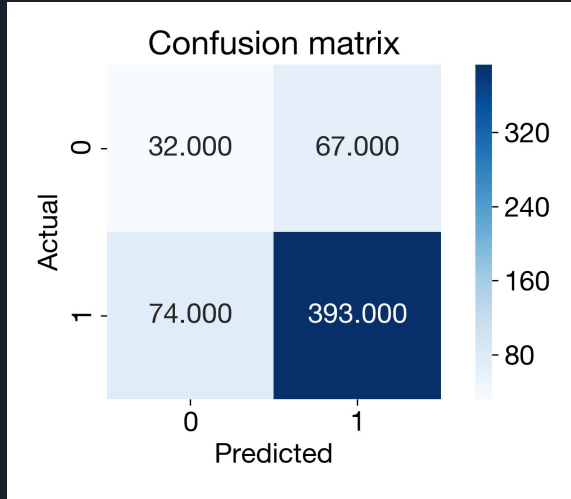
# Data Modeling:
## Decision Tree Classifier

- Used Pearson correlation heatmap to select variables
  - *Age*, *Country*, *Oscore*, *Cscore*, *Impulsive* and *SS* were found to be highly correlated with the output variable *Drug User*
- 475 nodes, may be evidence of overfitting
- 10–fold cross validation → 79.0% mean **accuracy**
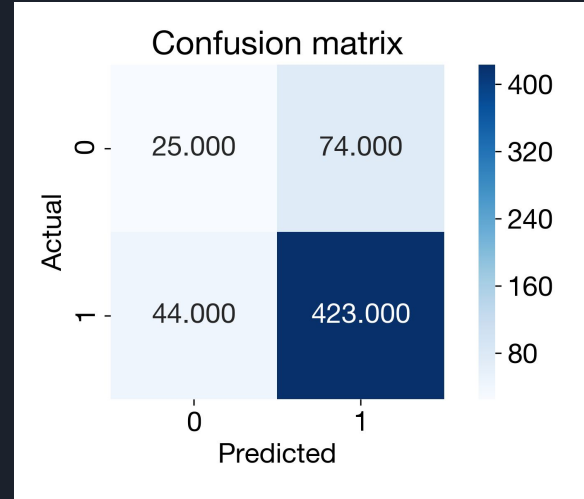  87.4% mean **F1 score**

# Data Modeling:
## k-Nearest Neighbors Classifier

- Used same 6 independent variables
  (*Age*, *Country*, *Oscore*, *Cscore*, *Impulsive* and *SS*)
- 10–fold cross validation → 80.9% mean **accuracy**

                                        88.8% mean **F1 score**

# Data Modeling:
## Random Forest Classifier

- Used Grid Search
  - 30,240 base estimators
  - evaluated using 10-fold cross validation
  - 100 iterations
- 10–fold cross validation
  - 91.3% mean **accuracy**
  - 84.4% mean **F1 score**
- 82.7% **accuracy** on the test set
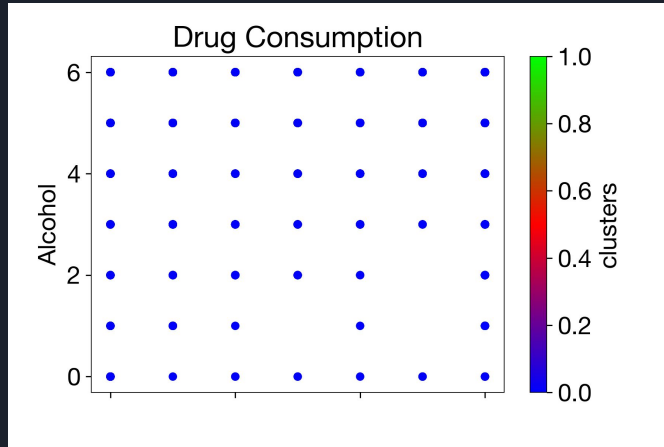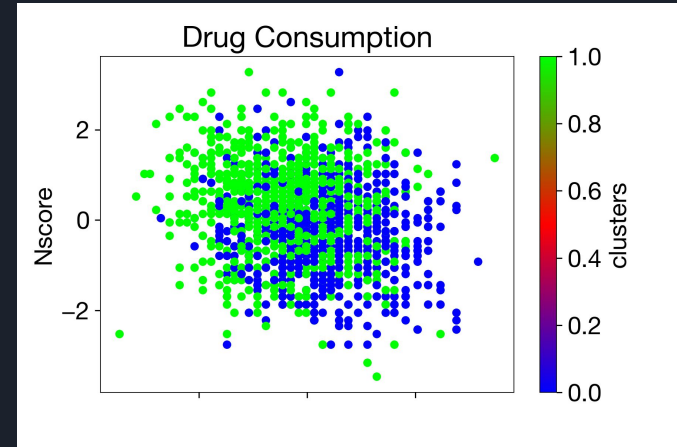- 90.2% **F1 score** on the set

# Data Modeling
## Clustering

- Silhouette Coefficient as the Evaluation Metric
- Clusters with 8 subsets of features
- Hierarchical Clustering with Single Linkage, Hierarchical Clustering with Complete Linkage, K-Means Clustering, DBSCAN applied to each cluster.
- Hierarchical Clustering performed the best and DBSCAN performed the worst.

# Data Modeling
## Clustering

Cluster with only Non-Illegal Drugs(alcohol VS Nicotine)

Cluster with only Personality Traits(Ascore VS Nscore)

# Data Modeling
## Clustering

Silhouette Coefficients

| | Type | Single Linkage | Complete Linkage | K-Means | DBSCAN |
|---|---|---|---|---|---|
| 0 | All Features | 0.507951 | 0.662999 | 0.183011 | -0.189872 |
| 1 | Non_Illegal | 0.823376 | 0.793181 | 0.225031 | 0.094018 |
| 2 | Illegal | 0.608208 | 0.541589 | 0.448770 | 0.516004 |
| 3 | Only Drugs | 0.779997 | 0.743257 | 0.268466 | -0.066339 |
| 4 | Not Using Drugs | -0.211147 | -0.211147 | 0.158041 | -0.115110 |
| 5 | Personality Traits | -0.322901 | -0.322901 | 0.189828 | 0.312260 |
| 6 | Personality Traits and Non-Illegal drugs | 0.220068 | 0.220068 | 0.166378 | -0.211882 |
| 7 | Personality traits and Illegal Drugs | -0.267331 | -0.267331 | 0.240112 | -0.112099 |

# Conclusions

- We observed that the clusters that are made up of Only Illegal, Only Non-Illegal and Only Drugs have the highest Silhouette Coefficient which means they are good clusters compared to the others.
- The best performing binary classifier was the random forest classifier with a mean accuracy of 91.2% when tested with 10-fold cross validation. On the test set, it had an accuracy of 82.7% and an F1 score of 90.2%.
- We are satisfied with our model's ability to predict drug usage by using personality traits and other characteristics.