# STAT451 FINAL PROJECT

FRIDAY M., JOEL M. NAYDEEN M, EVE L.

# INTRODUCTION- DATASET

- UC Irvine Machine Learning Repository

- Drug consumption patterns
    - Demographics
    - Physiological scoring/Personality-related scoring
    - Self-report drug use data for a range of substances

- Gave a decent foundation to start exploring relationships between individual traits and substance use behaviors

- https://archive.ics.uci.edu/dataset/373/drug+consumption+quantified

# CLEANING PROCESS

- Data was quantified

- Dataset had one fictitious drug
  - 7 participants reported use
  - Removed the 7 participants from the dataset

- Chocolate was an included drug, removed from dataset

- Categorized drugs into hard drugs vs. non-hard drugs based on international classifications
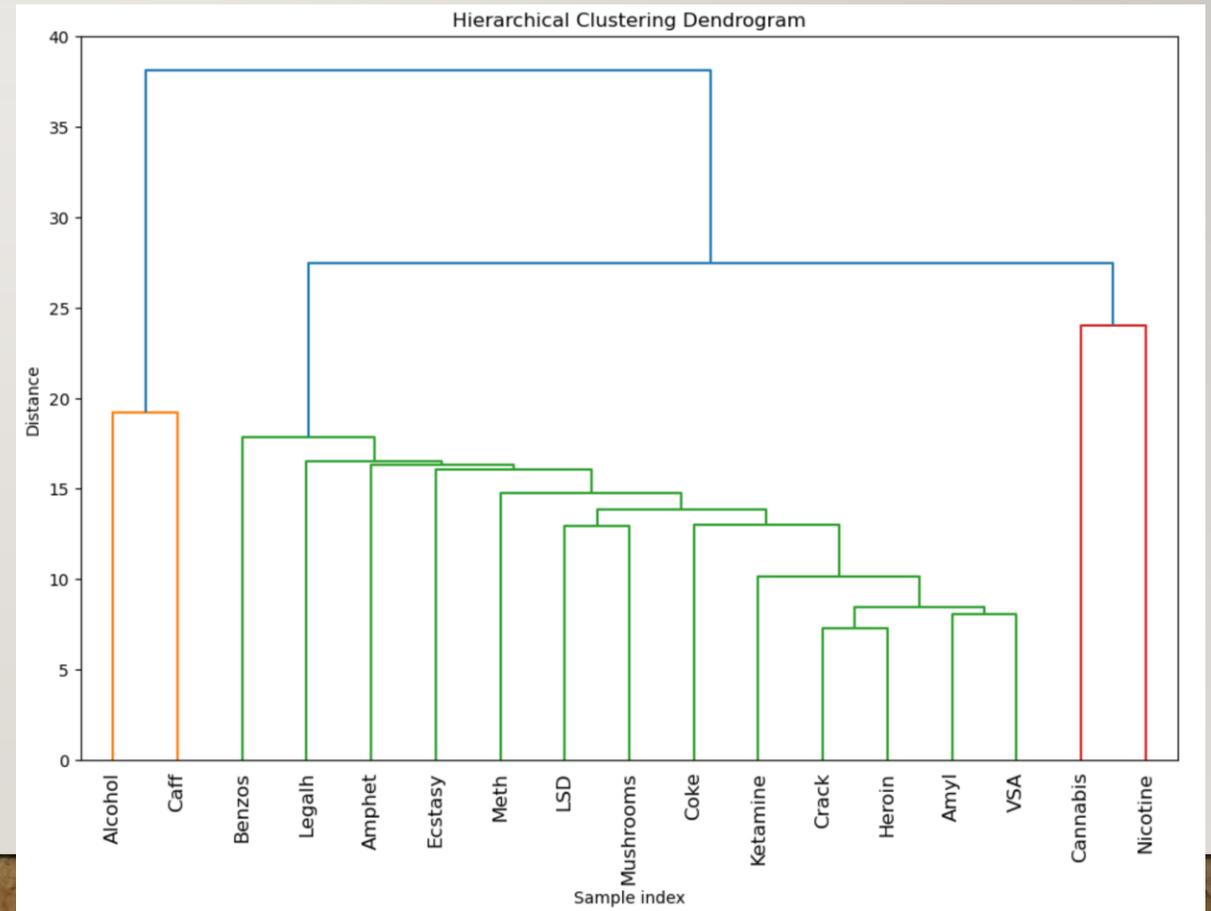
## QUESTIONS WE AIMED TO ANSWER

Question 1: What drugs are often used with other drugs?

Question 2: Which factors most strongly predict whether an individual uses cannabis/ consumes alcohol?

Question 3: Based on demographic and personality traits, can we classify individuals into low, moderate, or high-risk groups for substance abuse?

# QUESTION 1: WHAT DRUGS ARE COMMONLY USED TOGETHER?

- Dendrogram
  - Groups drugs based on similarities in usage patterns, using hierarchical clustering
- Goal: identify clusters of drugs that are often used together or have similar usage behaviors
- Key findings
  - Alcohol and Caffeine are closely related, suggesting they are commonly used together
  - Cannabis and Nicotine form another distinct pair
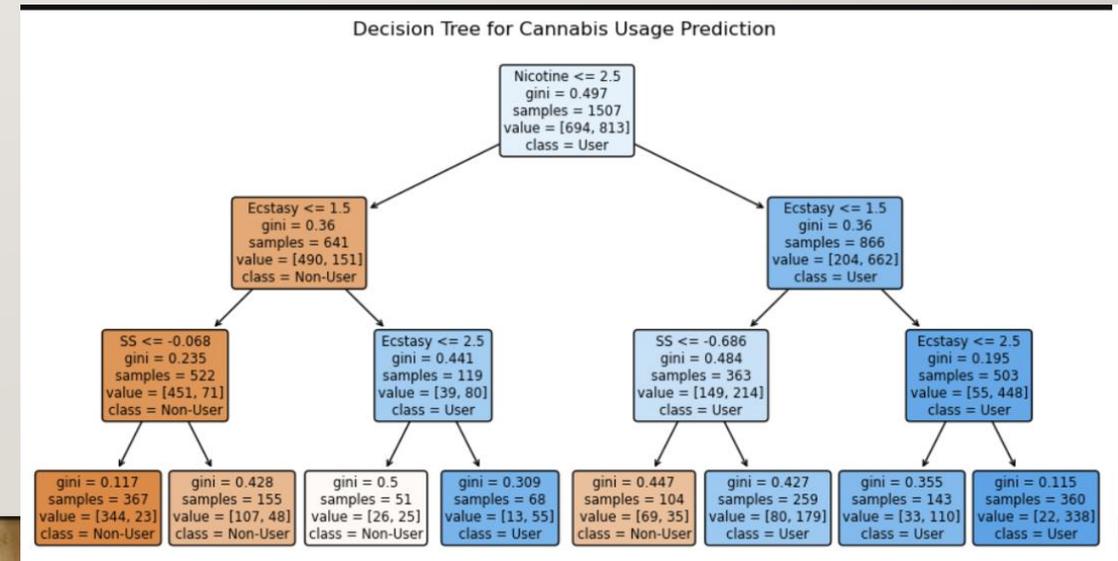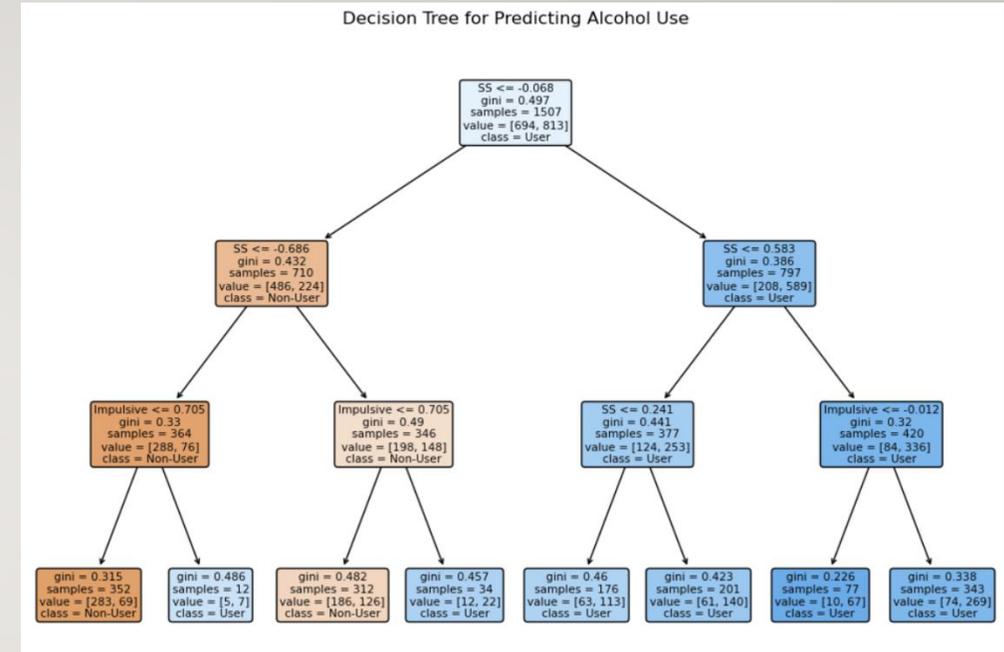  - All other substances fall into broader, less tightly related groups



Hierarchical Clustering Dendrogram
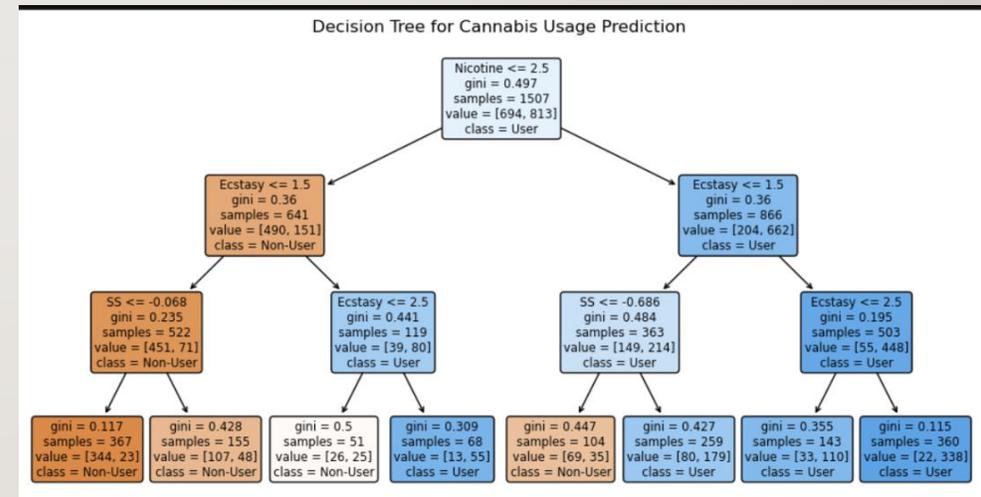
# K-MEANS CLUSTERING (QUESTION 1)

## QUESTION 2: WHICH FACTORS MOST STRONGLY PREDICT WHETHER AN INDIVIDUAL USES CANNABIS/ CONSUMES ALCOHOL?

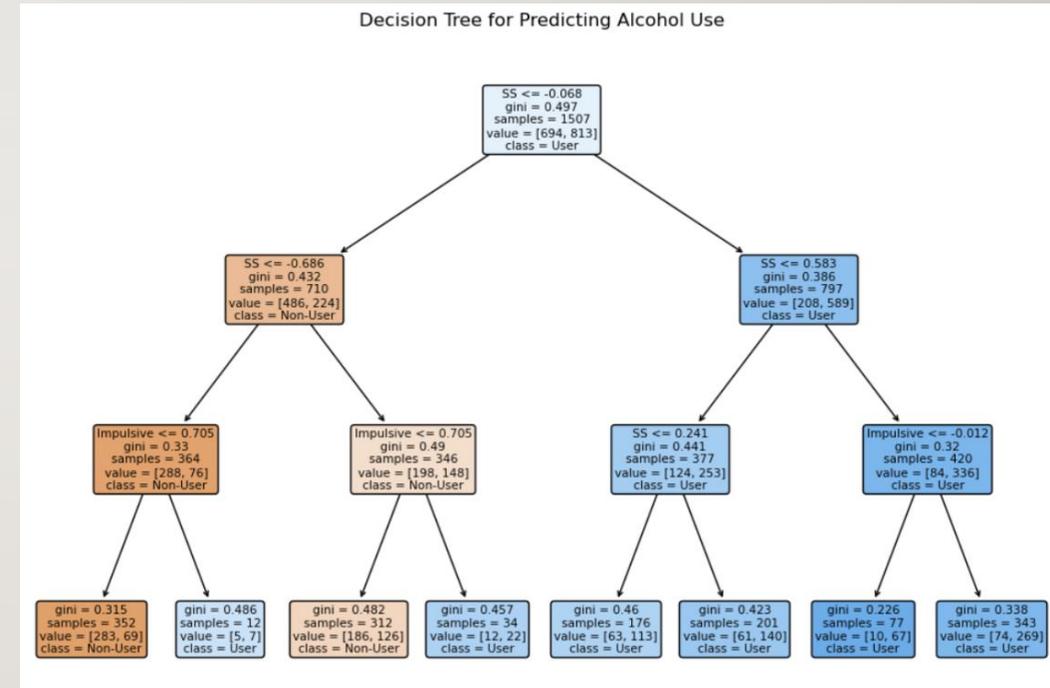- Decision Tree for Alcohol

- Decision Tree for Cannabis



Decision Tree for Predicting Alcohol Use



Decision Tree for Cannabis Usage Prediction

# CANNABIS PREDICTION

- Starting point: root node starts with Nicotine, split at <= 0.25.

- Other important predictors: Ecstasy usage and the Sensation Seeking (SS) variable.

- Outcomes
  - Final classifications: User or Non-User
  - If nicotine >2.5 and Ecstasy <= 1.5, individual predicted to be Cannabis User
  - Inverse, likely to not be a User



Decision Tree for Cannabis Usage Prediction
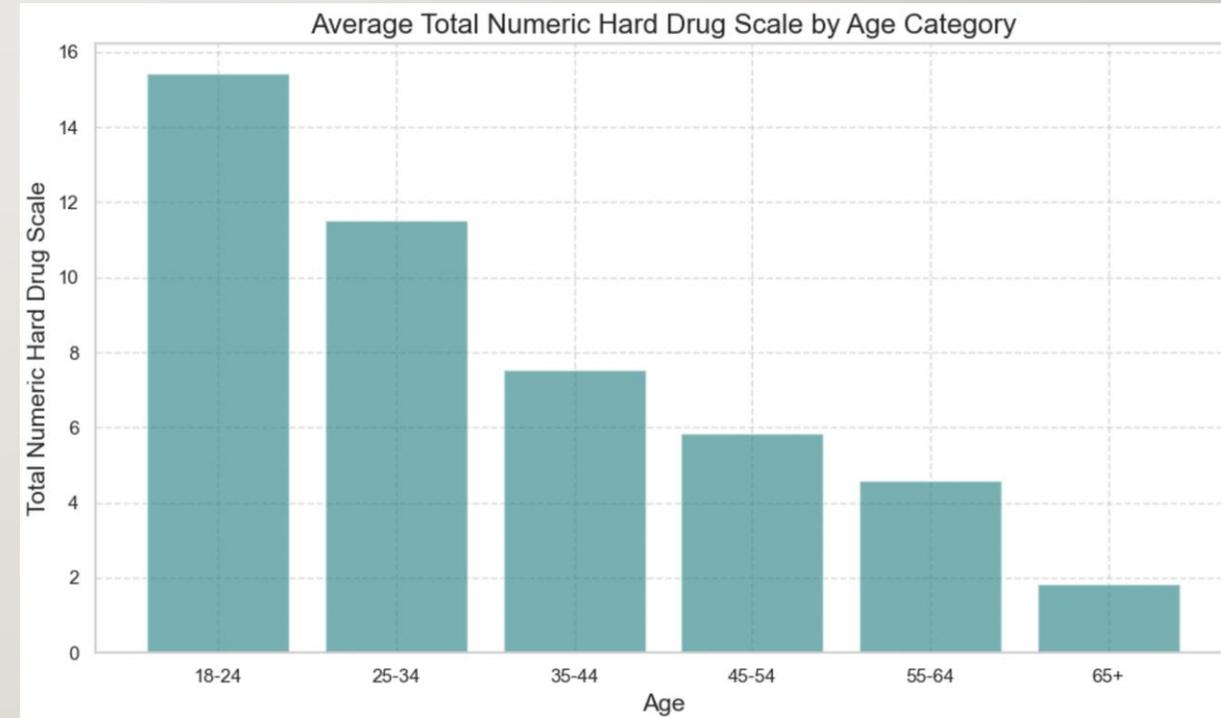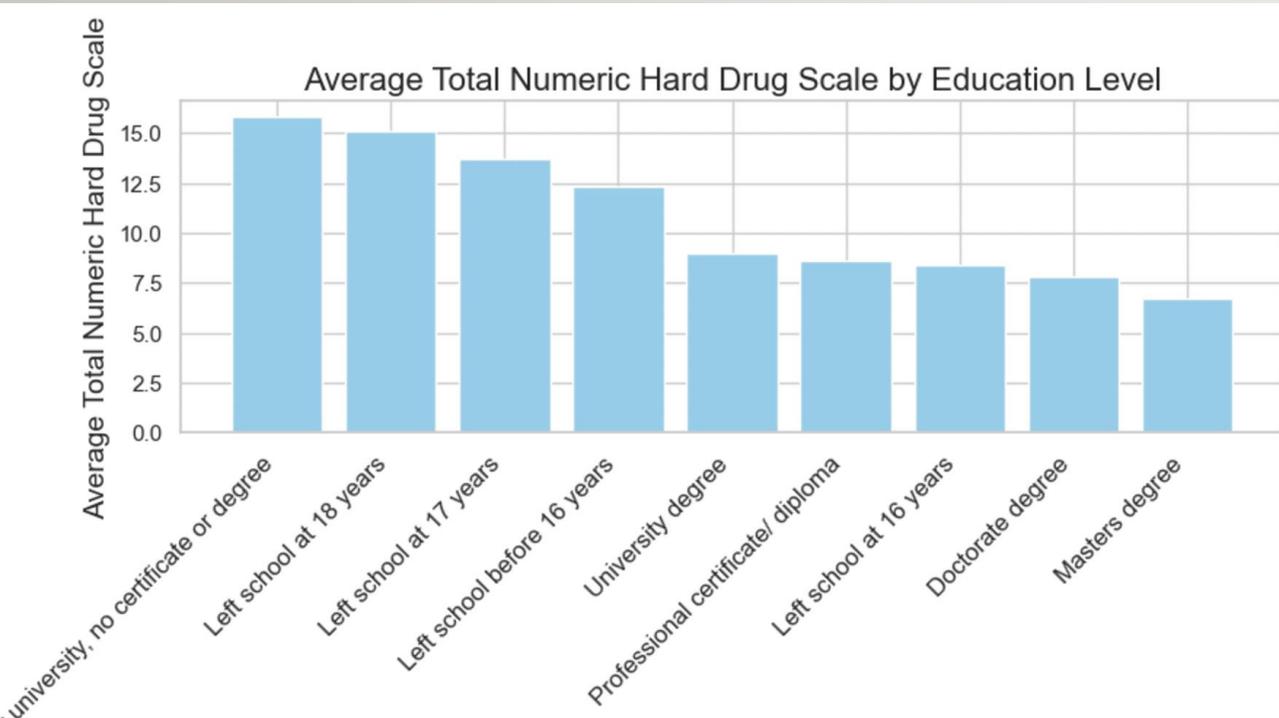
# ALCOHOL PREDICTION

- Starting point: root node starts with the Sensation Seeking score, split at 0.068.
  - Indicates individuals with a low SS are more likely to not use alcohol
  - Other important factors was the Impulsiveness personality score
- Outcomes
  - SS > 0.068 and Impulsive <= 0.705, individual is likely to be a user
  - Inverse outcome, likely to not use alcohol
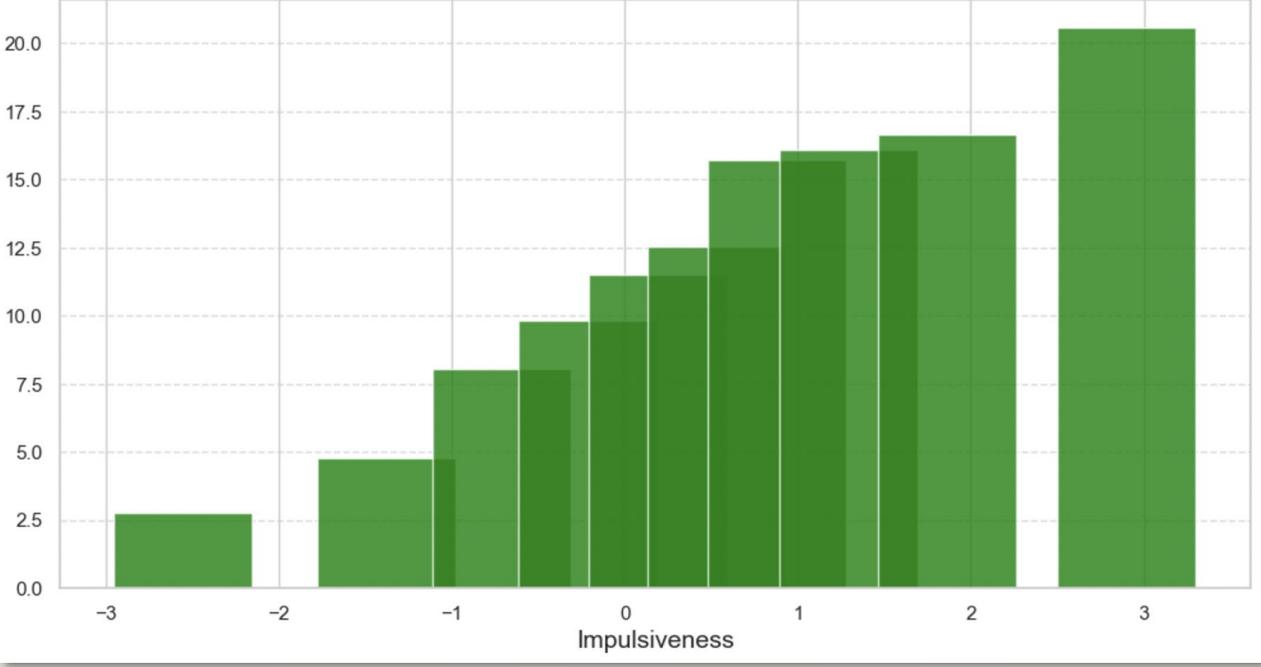


Decision Tree for Predicting Alcohol Use

# COMPARING THE TREES

- Cannabis:
  - Nicotine and Ecstasy usage are key predictors, highlighting the role of co-usage with other substances

- Alcohol:
  - Behavioral Traits like Sensation Seeking and Impulsiveness are more significant, reflecting a stronger psychological component.
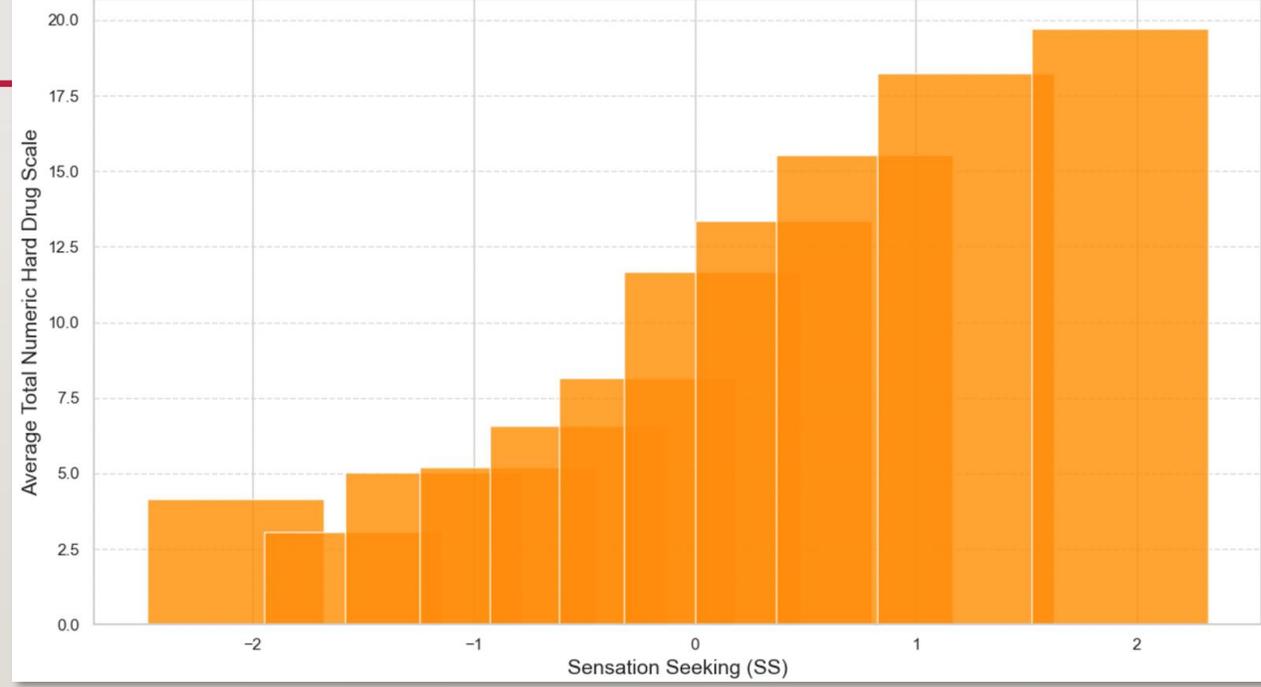
# QUESTION 3: BASED ON DEMOGRAPHIC AND PERSONALITY TRAITS, CAN WE CLASSIFY INDIVIDUALS INTO LOW, MODERATE, OR HIGH-RISK GROUPS FOR SUBSTANCE ABUSE?
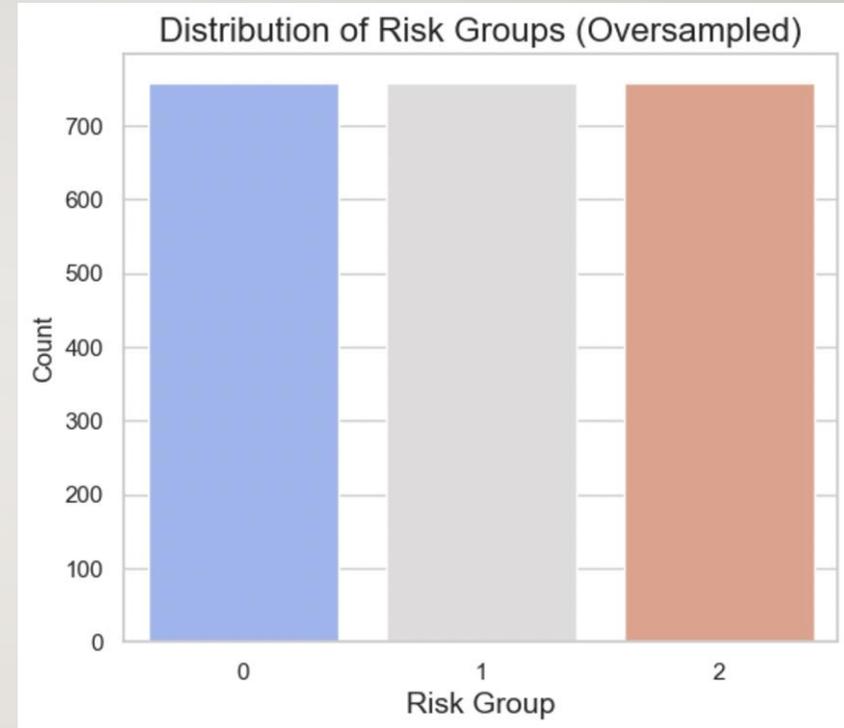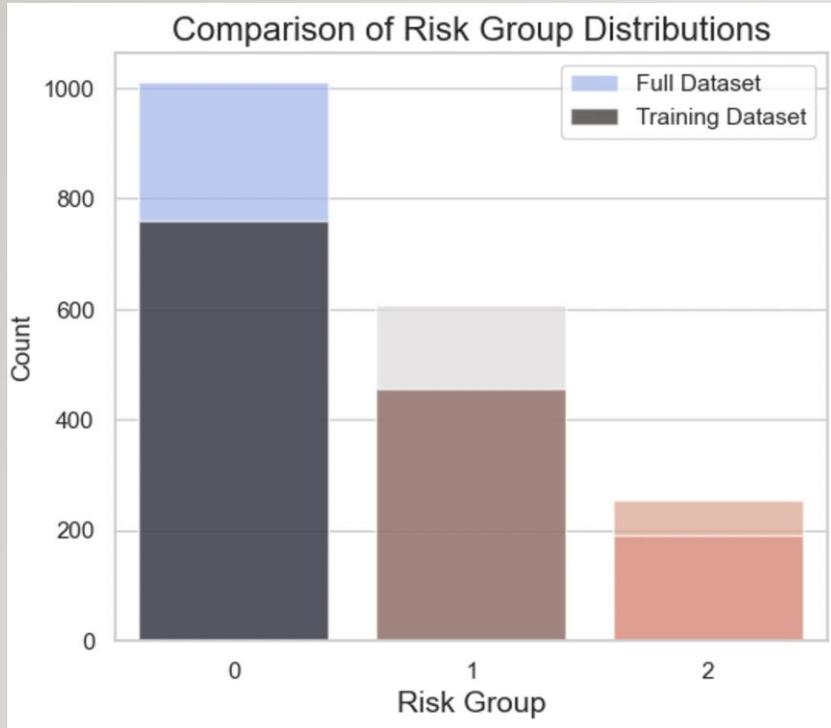
Average Total Numeric Hard Drug Scale by Impulsiveness

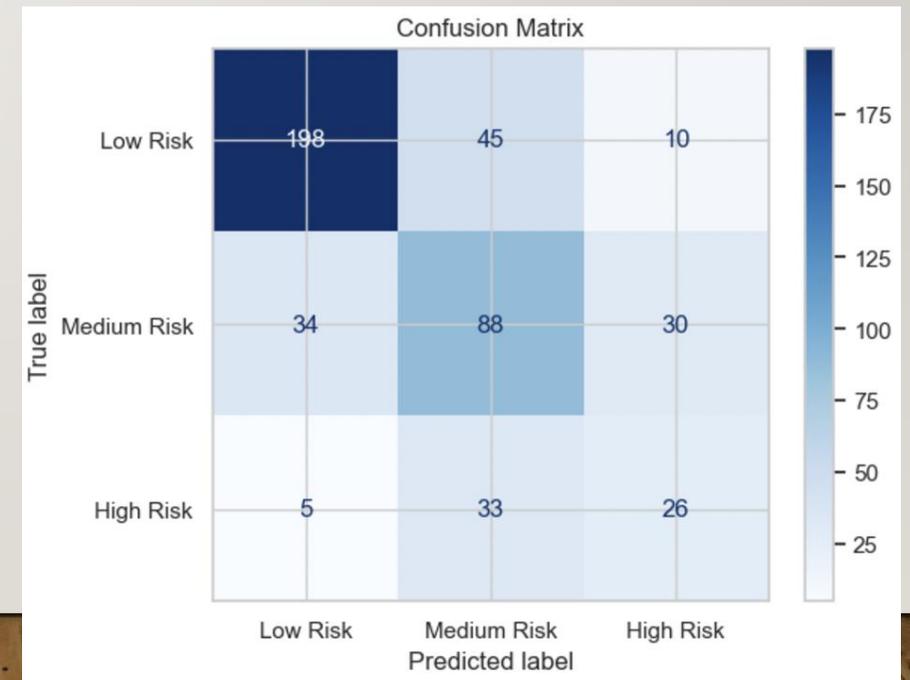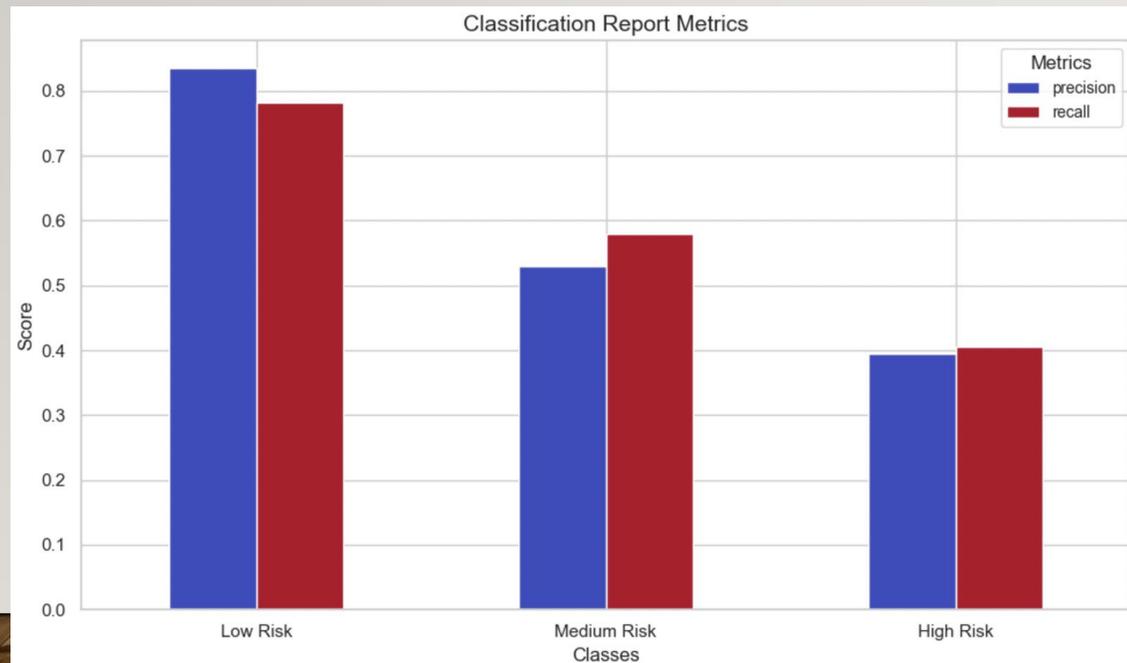Average Total Numeric Hard Drug Scale by Sensation Seeking

# OVERSAMPLING

# RANDOMFORESTCLASSIFIER

- Mean CV Score: 0.8564

- Test Accuracy: 0.6652

# CONSIDERATIONS - BUILDING A MORE ACCURATE MODEL

- Two classes - Low Risk, High Risk

- More features
  - Childhood maltreatment
  - Familial substance abuse

- Individual Risk factors
  - Neurodevelopmental disorders
  - Mental illness

- Social risk factors
  - Bullying

  - Gang affiliation

# CONCLUSION

- **Patterns of Co-Usage:**
Alcohol and caffeine, as well as cannabis and nicotine, show strong co-usage tendencies, reflecting shared social or behavioral contexts.

- **Substance Use Predictors:**

- **Cannabis:** Nicotine and Ecstasy usage, along with Sensation Seeking scores, were strong predictors.

- **Alcohol:** Psychological traits such as Sensation Seeking and Impulsiveness dominated prediction accuracy.

- **Risk Classification:**
While the Random Forest model achieved promising results (Mean CV: 85.6%), further improvements could involve deeper demographic and psychosocial variables.