

# ***ADDICTIVE PERSONALITIES***

*Lora Johns*



# **DRUG USE & THE BIG FIVE**

**WHAT'S  
YOUR TYPE?**

## DOES PERSONALITY TYPE PREDICT DRUG USE?

- The original study sought to prove that it does (*The Five Factor Model of personality and evaluation of drug consumption risk*, Fehrman et al.)
- They claimed that personality profiles are *strongly associated with* belonging to groups of drug users

THE BIG  
5

- 
- **NEUROTICISM**
  - **EXTRAVERSION**
  - **OPENNESS TO EXPERIENCE**
  - **AGREEABLENESS**
  - **CONSCIENTIOUSNESS**

# ***ADDICTIVE PERSONALITIES***

*Lora Johns*





# **WHAT REALLY "PREDICTS" DRUG USE?**



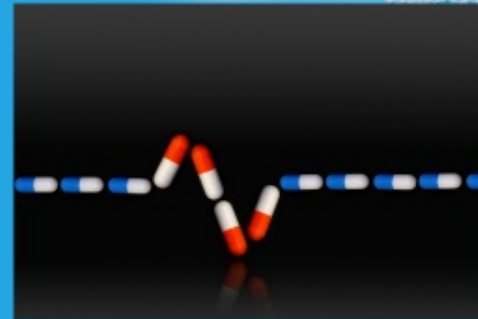
**WHO IS A  
"USER"?**



**THE  
DATA**



- For this analysis, an "active user":
  - Used Heroin, Methadone, Crack, or Cocaine
  - Within the last year
- This definition underlies the target variable. Why?
  - Highest risk population (real-world)
  - Create a reasonable target variable without introducing confusion (pragmatic)



## **FEATURE SELECTION AND ENGINEERING**

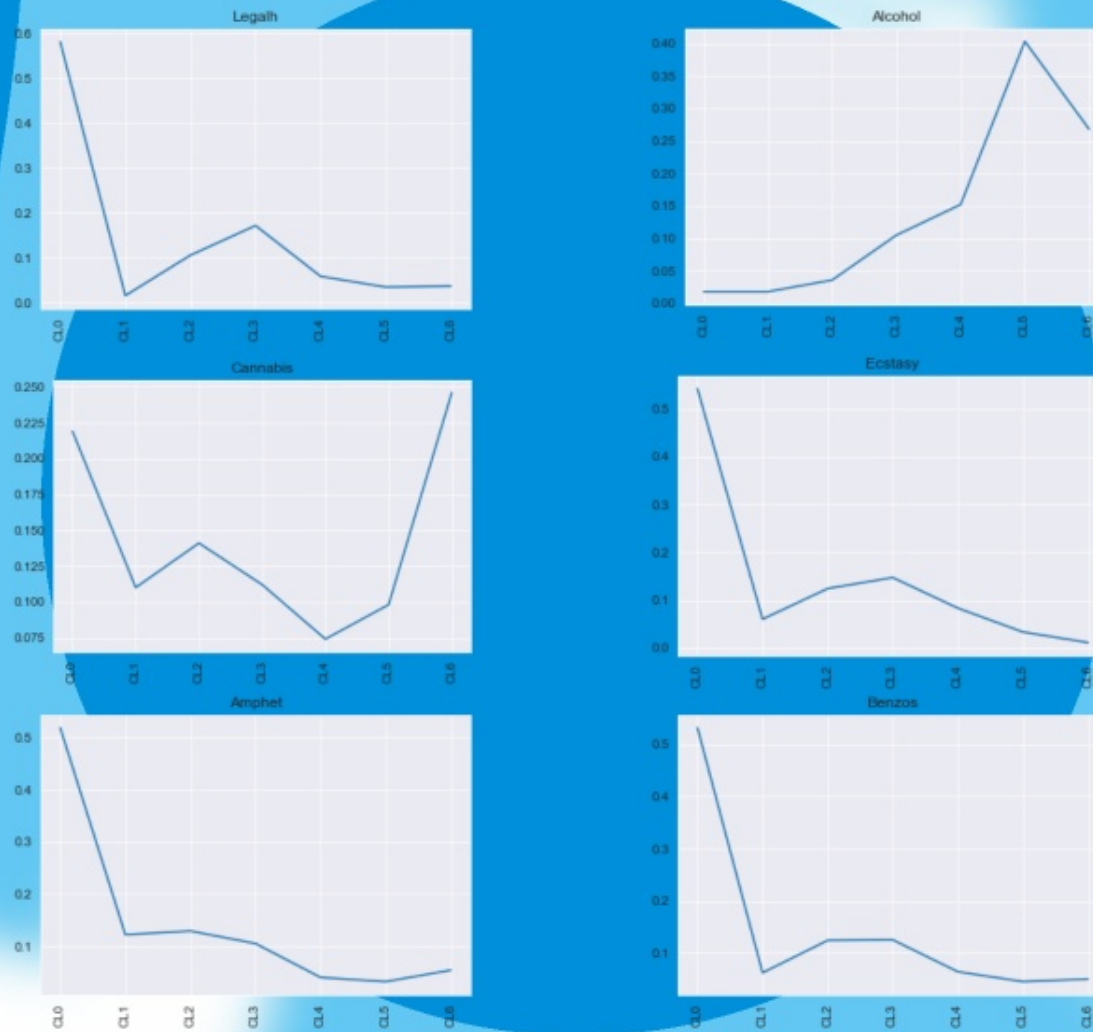
- Engineered:
  - College grad vs. non-grad
  - White vs. Non-white
  - Age and geographic groupings
  - How recently subjects used non-target substances, if ever
- Transformed:
  - Min-max scaling personality

**CORRELATE**

**VISUALIZE**







# ***ADDICTIVE PERSONALITIES***

*Lora Johns*





**ANALYSIS**

**METRICS**

## **PRELIMINARY METRICS AND TESTS**

- Chi squared contingency table for predictors indicated statistical significance (test statistic of 36081.0 and p value 1.0)
- Initial multinomial Bayes classification (without feature scaling) produced accuracy of ~ 0.56

# ***ADDICTIVE PERSONALITIES***

*Lora Johns*





# MODELS

- Multinomial Bayes
- Logistic regression
- K-nearest neighbors
- Support vector classifier
- Random forest

**MULTINOMIAL  
BAYES**

**LOGISTIC  
REGRESSION**

**KNN**

**SVC**

**RANDOM  
FOREST**

- Accuracy: 0.801
- Cross-validation:  
stratified k-fold



## LOGISTIC REGRESSION

- F1: 0.739
- Accuracy decreased on testing set
- Ergo, this model probably overfit the training data
- More false labels than other models



## **K-NEAREST NEIGHBORS**

ACCURACY: 0.8453

F1: 0.7245

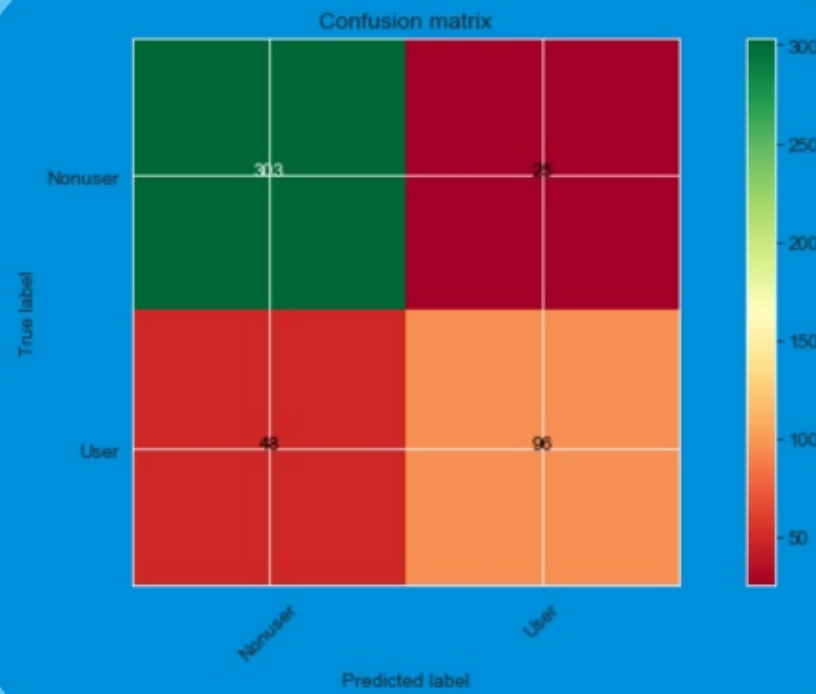
TRAINING RMSE: 0.403

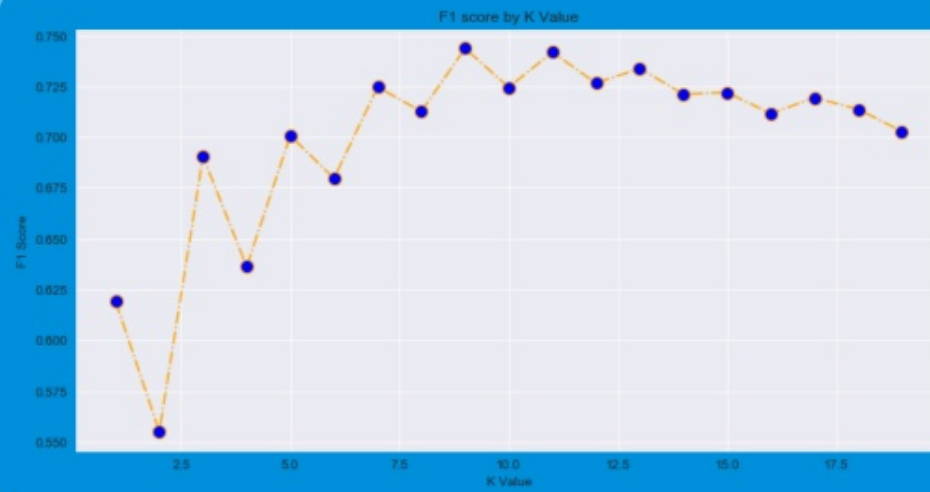
TEST RMSE: 0.393

- Performance improved on test
- Multiple models run to find model with best F1 score

**MATRIX**

**MULTIPLE  
MODELS**







## SUPPORT VECTOR MACHINE

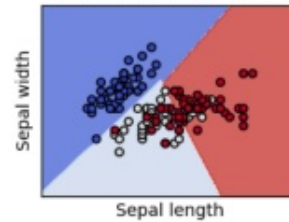
ACCURACY: 0.831

CONFUSION MATRIX:  $\begin{bmatrix} 285 & 43 \\ 37 & 107 \end{bmatrix}$

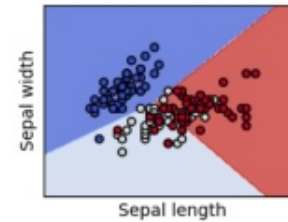
- Tuned hyperparameters with grid search
- Few false negatives (which, if we care about intervention, is good); many false positives (which is problematic in this context)

SVM

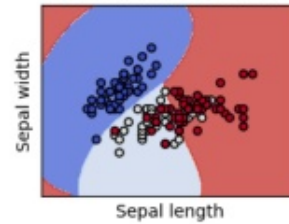
SVC with linear kernel



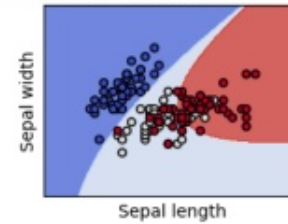
LinearSVC (linear kernel)



SVC with RBF kernel



SVC with polynomial (degree 3) kernel



## **RANDOM FOREST DECISION TREES**

*ACCURACY: 0.970*

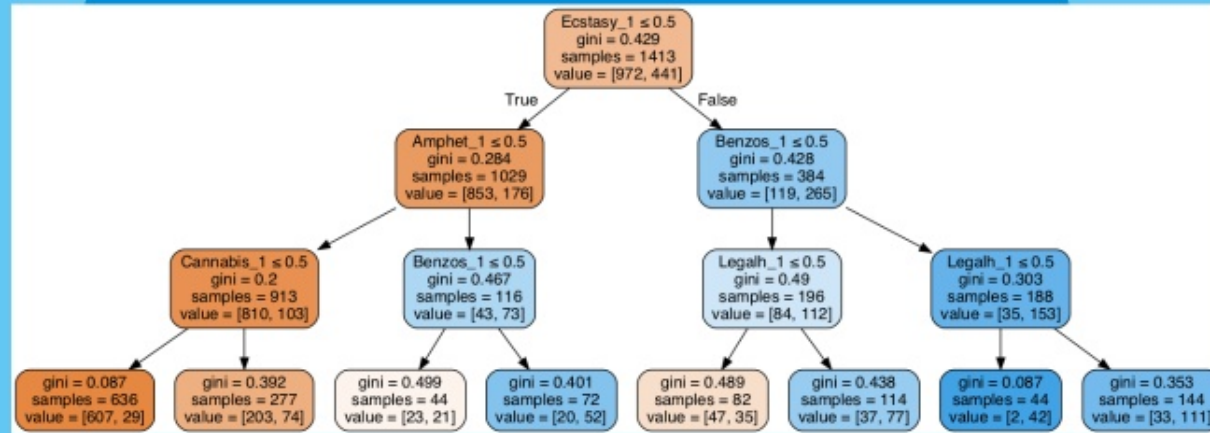
*AVERAGE ERROR: 0.03*

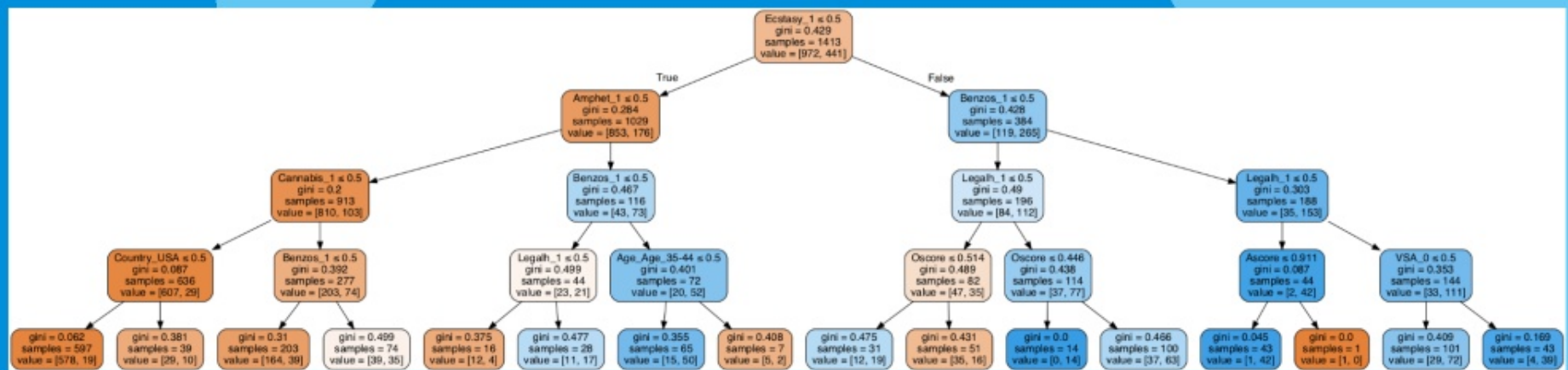
- In all models, ecstasy was the top discriminator (followed by benzos and legal highs)
- Tuned hyperparameters with randomized search
- Performance improved from training to testing (RMSE dropped from 0.190 to 0.03)

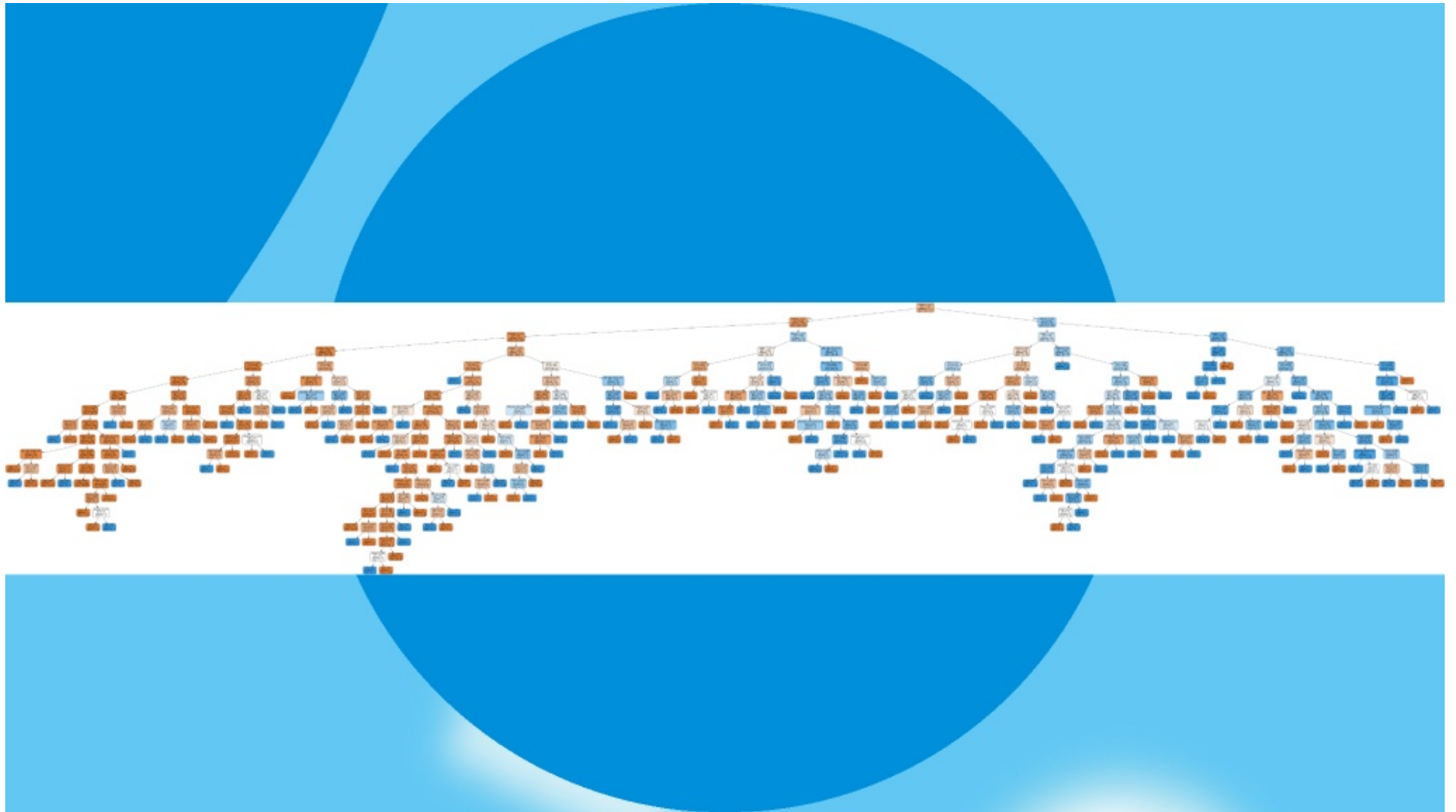
**TREE**

**TREE**

**TREE**









# **ADDICTIVE PERSONALITIES**

*Lora Johns*



# TAKEAWAYS

- Accuracy alone isn't enough
- Interrogate the underlying research
- The purpose to which you put your data matters
- Examine bias (in the data and the people)

# ***ADDICTIVE PERSONALITIES***

*Lora Johns*

