Using decision trees and random forests to predict cannabis consumption

Jorge Rodríguez Peña – INM431 Machine Learning – MSc Data Science – City University of London

Description and motivation

- The original dataset was used and created by Fehrman E., et al. (2017) [1]. The purpose of their study was to predict drug consumption from individual and personality traits.
- It uses the revised **NEO Five Factor Inventory** [2], the reviewed Baratt Impulsiveness Scale [3] and the **Sensation Seeking scale** [4].
- The aim of this project is to predict possible cannabis users on a yearly basis and the influence personality traits might have on it.
- Individual traits such as Age, Gender, Education Country and Ethnicity are also considered.
- The results will be **compared** with those obtained by Fehrman E., et al. (2017) [1].

Users

dnoub 30

each

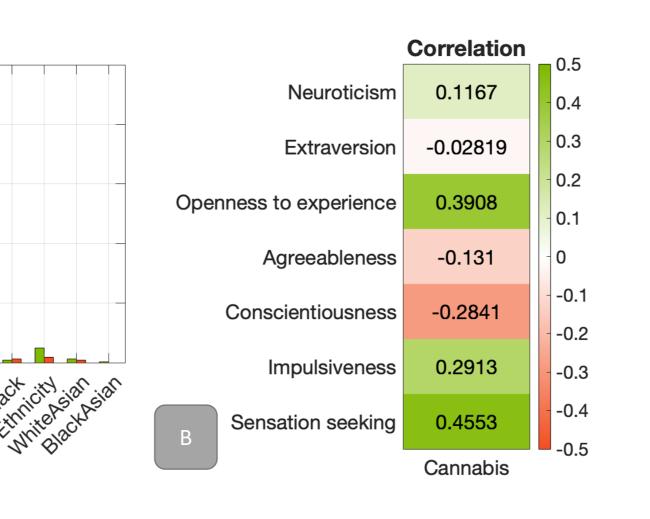
Non-users

Distribution of users and non-users per group

Groups

Analysis of the dataset

- Dataset containing **1885 samples** with 12 predictors:
- **5 categorical**: Age (binned), Gender, Education, Country and Ethnicity.
- 7 numerical (personality tests): Neuroticism, Extraversion, Openness to experience, Agreeableness, Conscientiousness, Impulsiveness and Sensation seeking.
- **1 binary target:** A yearly basis definition of cannabis user.
- The data was collected using the snowball method.
- Original dataset was cleaned to make our study simpler (the original dataset used CatPCA [9]).
- Data imbalance: Country: ~60% UK ~30% USA and Ethnicity: ~95% White (see figure A). Probably due to the snowball method.
- Correlation with cannabis consumption was higher for personality traits: Openness to Experience, Conscientiousness, Impulsiveness and Sensation **Seeking** (see figure B).



Decision trees vs random forest

DECISION TREES:

- **Split** the data by making a decision.
- Each node uses a predictor to split on.
- Use **information gain** to make the decision.
- Generate rules to predict new samples.

RANDOM FORESTS:

- Are built from **decision trees**.
- Train each decision tree **bootstrapping** from the data (bag method) and randomly selecting predictors [8].
- Use majority vote from each tree and average out the results.

Averaging out the results from each decision tree allow random forests to overcome most of decision tree's cons. Adding more trees to the model not only implies getting rid of noise but it also increments the training times and make the model more difficult to understand [6]. See the following table for a summary:

	Decision tree	Random forest
Interpretability	√Readable rules [6]	X Random rules
Noise handling	X [7]	✓Averages out the predictions
Do not overfit	X[7] √Results are averaged out	
Best performance	Small data [8]	Large data
Time	Shorter [9]	Longer

Hypothesis

- The original paper [1] concluded that the decision tree was the **best method** for classifying cannabis users for a decade basis user definition.
- Similar behaviour expected for the yearly basis definition of user used in this project.
- Random forest should take longer time to train than decision trees.
- Experience, to Conscientiousness, Openness Impulsiveness and Sensation Seeking should have a higher impact in the model due to their correlation values with cannabis consumption.
- The biased model should generalize worse the results (see Methodology).

Methodology

Due to the imablance of the data with Country and Ethnicity predictors (see figure A) 4 models will be trained: one decision tree and one random forest per one of the following cases:

All predictors included. Data imbalance or bias are not considered, Country and Ethnicity removed. Considering data imbalance and bias. This will allow us to compared how data imbalance affect the model.

The approach will be the same in the four models:

- Split the data in a training (80%) and test set (20%) at random.
- Perform bayesian optimization aiming to maximize the AUC with:
 - Parent size and leaf size for decision trees to control the deph of them and prevent overfitting.
 - **Leaf size**, **number of predictors** and **number of trees** for *random forest*. Run 25 times and used the median as the optimal parameter due to the
- randomness of the process [5]. Evaluate the models using 10-fold cross validation and bootstrap validation
- (10 bootstraps for decision trees and out of bag prediction for random forests).
- **Train** the model with the *training set* and its *optimal parameters*.
- **Predict** with the *test set* and compare the results.

Main results

The AUC obtained for each of the models is shown in the table bellow.

- 1st model includes all predictors
- 2nd model removes Country and Ethnicity.
- Bootstrap means:
 - Bootstrap validation for decision trees
- Out of bag validation for random forests

We see lower AUC scores for the 2nd model.

	AUC estimation	Decision tree	Random forest
1 st model	CV	0.86 ± 0.03	0.89 ± 0.02
	Bootstrap	0.84 ± 0.02	0.89
	Test	0.87	0.89
2 nd model	CV	0.83 ± 0.04	0.87 ± 0.03
	Bootstrap	0.80 ± 0.03	0.87
	Test	0.82	0.85

1000

800

600

400

200

1000

1000

Bayesian Optimization

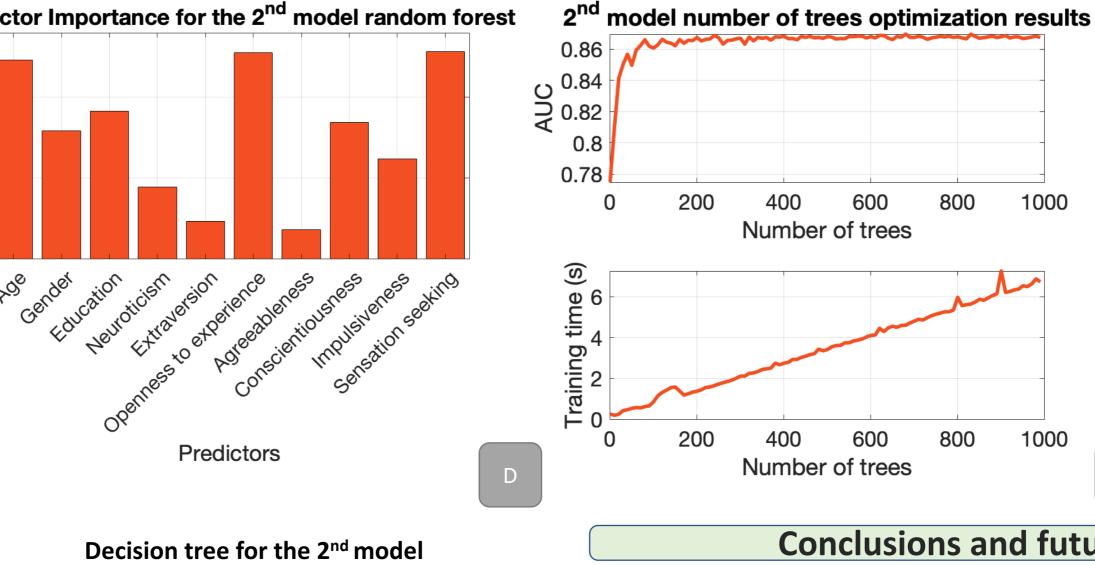
Number of trees

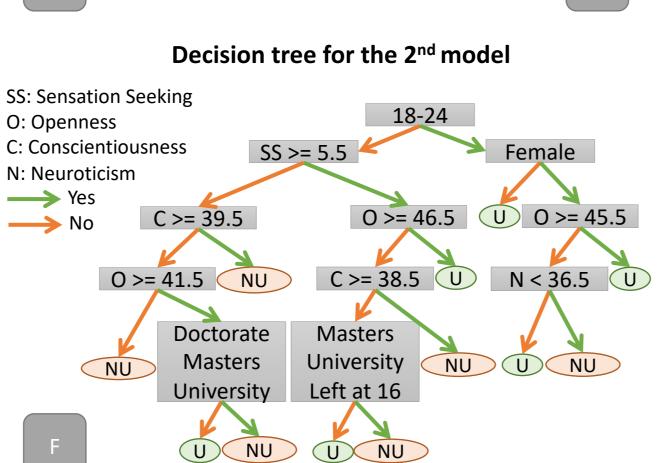
Discussion of results

- **MODEL PERFORMANCE:**
- Random forests performed better than decision trees in all of the models.
- The AUC obtained for decision trees is sometimes inside the error range of the AUC obtained for random forests, suggesting the decision trees may outperform random forests in some cases as shown in the original paper [1].
- **VALIDATION RESULTS:**
- Out of bag validation is consistent with cross validation in random forests.
- Bootstrap validation tends to underestimate the AUC value. Bootstrap validation has been studied as a low variance but higher bias validation method [8].
- **TRAINING TIME:**
- Random forests take around 4 seconds to train (see figure D), decision trees around **0.03 seconds**.
- **DIFFERENCES BETWEEN MODELS:**
- Models including all predictors performed slightly better than those excluding Country and Ethnicity.
- Country had the highest predictor importance in the 1st model. This suggests a very biased dataset and not representative of a real-world situation.
- The 2nd model proves to be very accurate in predicting and generalising results with AUC values over 0.8.
- Opennes to Experience, Conscientiousness and Sensation Seeking were the most relevant personality traits in our 2nd model (see figure C).
- The **decision tree** generated for the 2nd model is **very** readable (see figure F) and shows an easy way to classify cannabis users with a very decent AUC score.
- **HYPERPARAMETER TUNING:**
- **Decision trees** tended to have **bigger leaf sizes** (14, 32) as compared to random forests (3, 3). This is to reduce overfitting avoided in random forest by averaging out the results.
- Random forests performed better sampling a few predictors (1, 2), therefore decreasing the correlation between trees and improving the model [9].
- Number of trees presented very high variance in Bayesian optimization (see figure E), suggesting that the number of trees was **not that important** to optimize.
- Evaluation of **number of trees** against time and against AUC proved that *incrementing* the **number of trees did not improve** by much *the model* from around 600 trees, but it added extra training time (see figure D). Models were trained with 600 trees.
- The best decision tree for the 2nd model had 8 pruning levels (see figure F) and 11 for the 1st model. Decreasing the size of the tree by 3 levels produced very similar AUC scores, proving the strength of the model.

9.

Predictor Importance for the 2nd model random forest Importance 5.0 **Predictors**





Conclusions and future work

- **Decision trees** are **less computationally expensive** than random forests and, given the size of this dataset (1885 samples), perform almost as good.
- For quickier results, decision trees are a good option in this case.
- Imbalanced data meant higher values of the AUC, suggesting that the dataset is not representative of the real world.
 - Data collection used snowball method [1], other approaches could **improve** the accuracy with reality.
- Most relevant personality traits to determine consumption: Sensation Seeking, Opennes to experience and Conscientiousness. Similar to Fehrman, E. et al. (2017) [1].
- The original paper [1] trained its best model (a decision tree) using only 6 parameters (here decision trees use 12 and 10). Feature selection from feature importance results might be useful to improve AUC scores.
- Fehrman, E. et al. (2017). The Five Factor Model of Personality and Evaluation of Drug Consumption Risk. 10.1007/978-3-319-55723-6_18.
- Robert R. McCrae, Paul T. Costa, 'A contemplated revision of the NEO Five-Factor Inventory'. Personality and Individual Differences, Volume 36, Issue 3, 2004, Pages 587-596, ISSN 0191-8869, https://doi.org/10.1016/S0191-8869(03)00118-1.
- Stanford MS, Mathias CW, Dougherty DM, Lake SL, Anderson NE, Patton JH. Fifty years of the Barratt Impulsiveness Scale: An update and review. Personality and Individual Differences. 2009; 47(5):385–395.

Ali, Jehad & Khan, Rehanullah & Ahmad, Nasir & Magsood, Imran. (2012). Random Forests and Decision Trees. International Journal of Computer Science Issues (IJCSI). 9.

- Zuckerman M. Behavioral expressions and biosocial bases of sensation seeking. New York: Cambridge University Press; 1994. Snoek, J., H. Larochelle, R. P. Adams. Practical Bayesian Optimization of Machine Learning Algorithms. https://arxiv.org/abs/1206.2944, 2012. Yang, P., Hwa Yang, Y., Zhou, B., Zomaya, Y., et al.: "A review of ensemble methods in bioinformatics". Current Bioinformatics 5(4), 296–308 (2010).
- 7. Kohavi, Ron. (2001). A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection. 14. 8.
 - Breiman, L. Random Forests. *Machine Learning* 45, 5–32 (2001). https://doi.org/10.1023/A:1010933404324 Linting M, van der Kooij A. Nonlinear Principal Components Analysis with CATPCA: A tutorial. Journal of Personality Assessment. 2012; 94(1):12–25.