

# STUDY PROJECT REPORT : Drug Consumption

By  
Chloé Coursimault  
Héloïse de Castelnau



# DataSet Exploration





Our dataset from the **UCI Machine Learning** repository is the result of an online survey conducted between 2011 and 2012 among **1,885 respondents** aged 18 and over, from English-speaking countries.

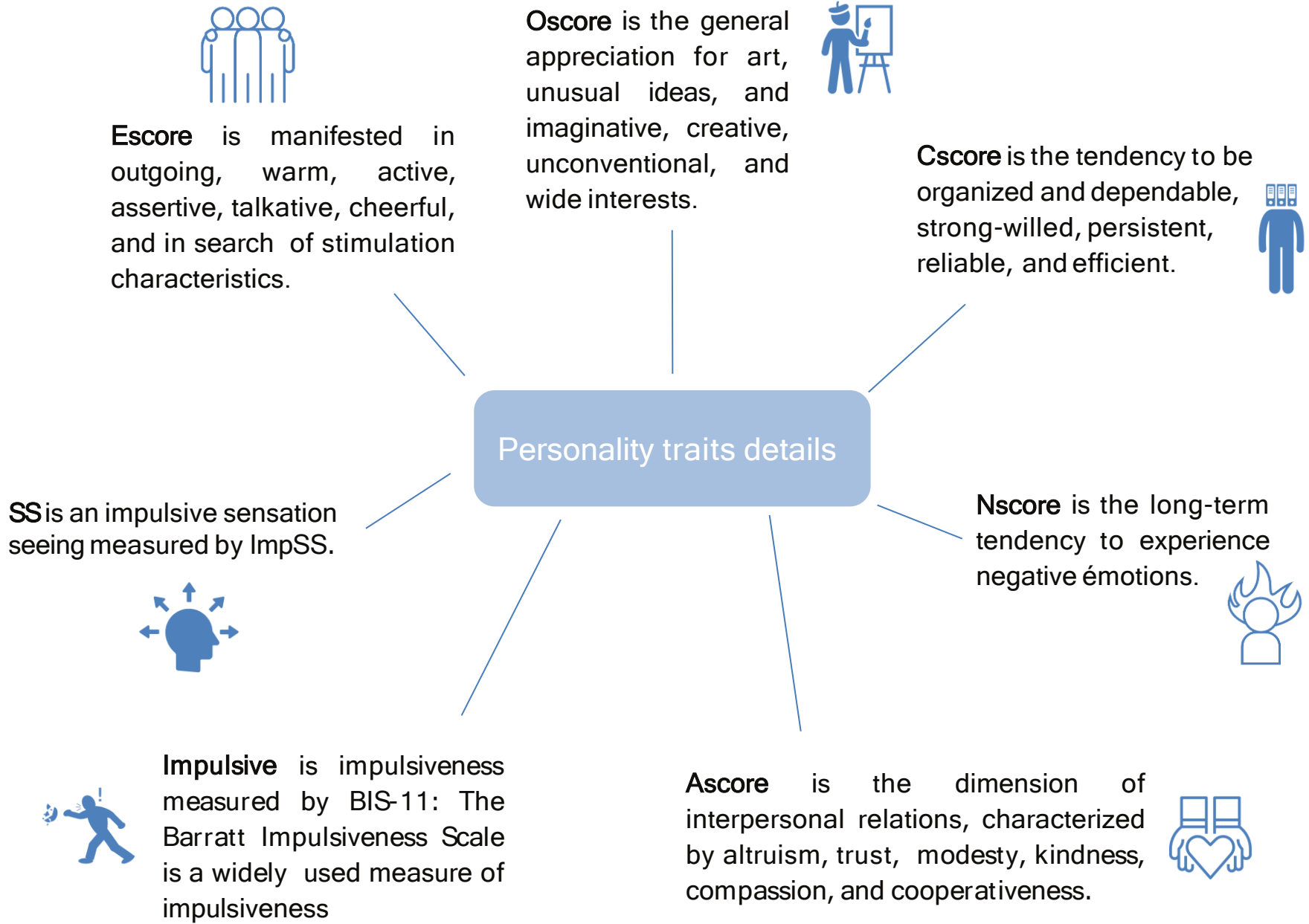
The survey collected data including Big Five personality traits (NEO-FFI-R), impulsivity (BIS-11), sensation seeking (ImpSS), and demographic information.

The data set contained information on the consumption of **18 central nervous system psychoactive drugs legal and illegal**. For each of these drugs, each individual had to choose his level of consumption :

- CL0 Never Used
- CL1 Used over a Decade Ago
- CL2 Used in Last Decade
- CL3 Used in Last Year
- CL4 Used in Last Month
- CL5 Used in Last Week
- CL6 Used in Last Day

Demographic features	Personality traits	Legal Drugs	Illegal Drugs
<ul style="list-style-type: none"><li>• Age</li><li>• Gender</li><li>• Education</li><li>• Country</li><li>• Ethnicity</li></ul>	<ul style="list-style-type: none"><li>• Nscore</li><li>• Escore</li><li>• Oscore</li><li>• Ascore</li><li>• Cscore</li><li>• Impulsive</li><li>• SS</li></ul>	<ul style="list-style-type: none"><li>• Alcohol</li><li>• Caff</li><li>• Chocolate</li><li>• Nicotine</li></ul>	<ul style="list-style-type: none"><li>• Amphet</li><li>• Amyl</li><li>• Benzos</li><li>• Cannabis</li><li>• Crack</li><li>• Ecstasy</li><li>• Heroin</li><li>• Ketamine</li><li>• Legalh</li><li>• LSD</li><li>• Meth</li><li>• Mushrooms</li><li>• VSA</li></ul>

Furthermore, to detect Fraud, the survey organizers added a fake drug "sumer".



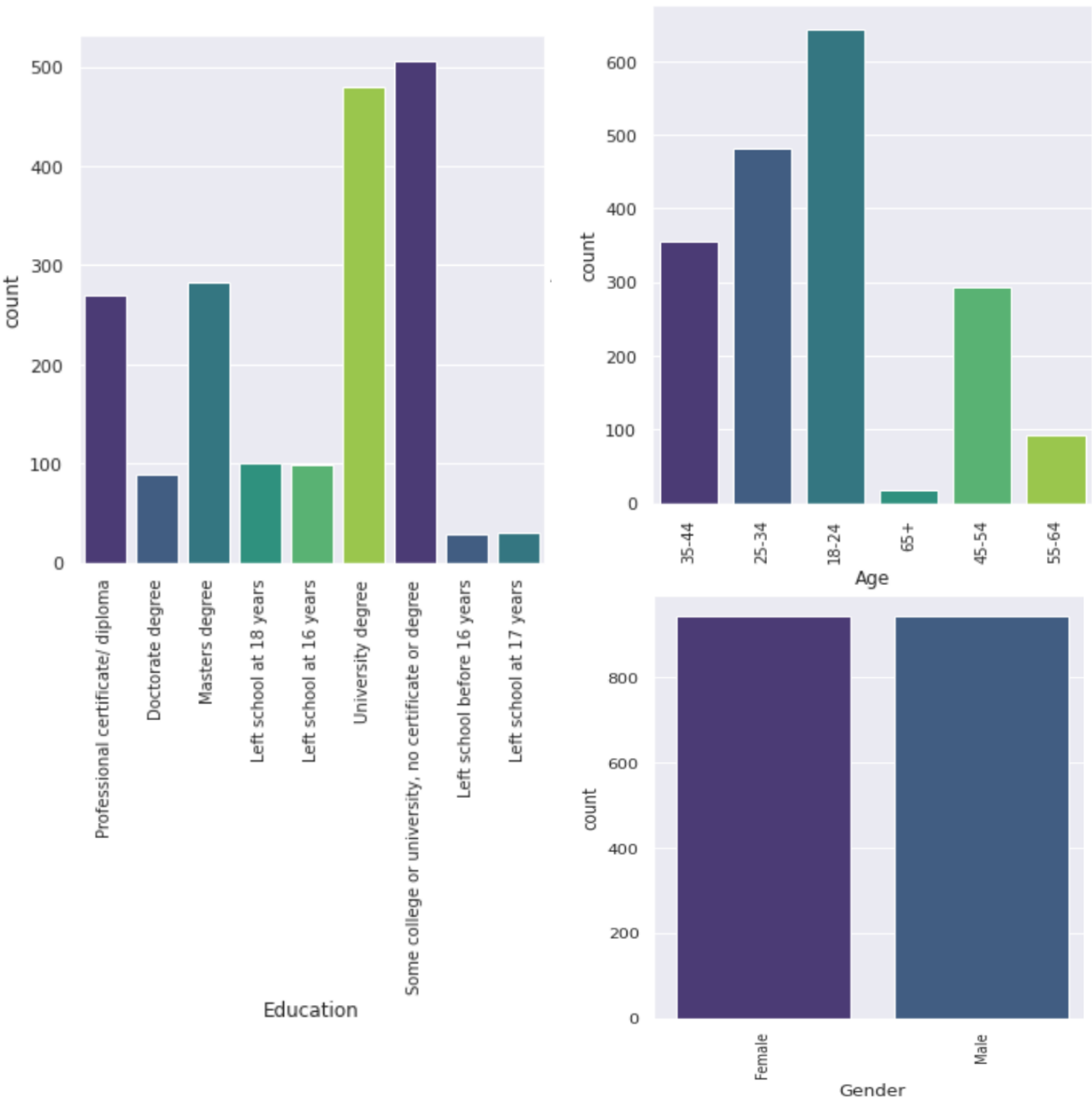
# DataSet Analysis





# DEMOGRAPHIC DATA

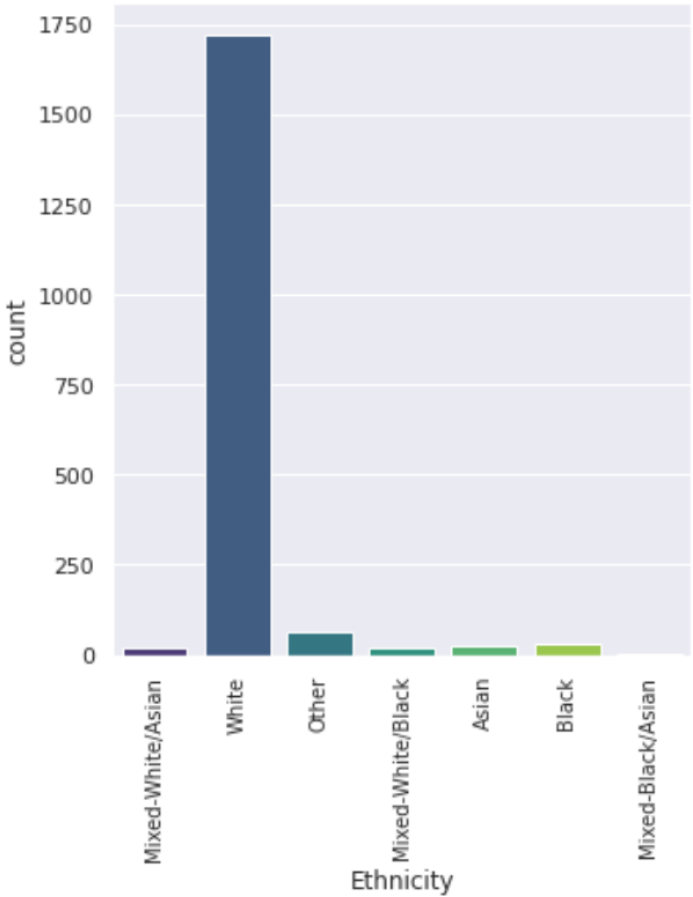
Some categories are almost uniformaly disparate such as **Education** and **Age** or even equally represented wich **Gender**.



We can observe some unevenly represented categories within the columns of the Dataframe, such as **Ethnicity** that count up to 90% of white or **Country** with mostly UK participants.



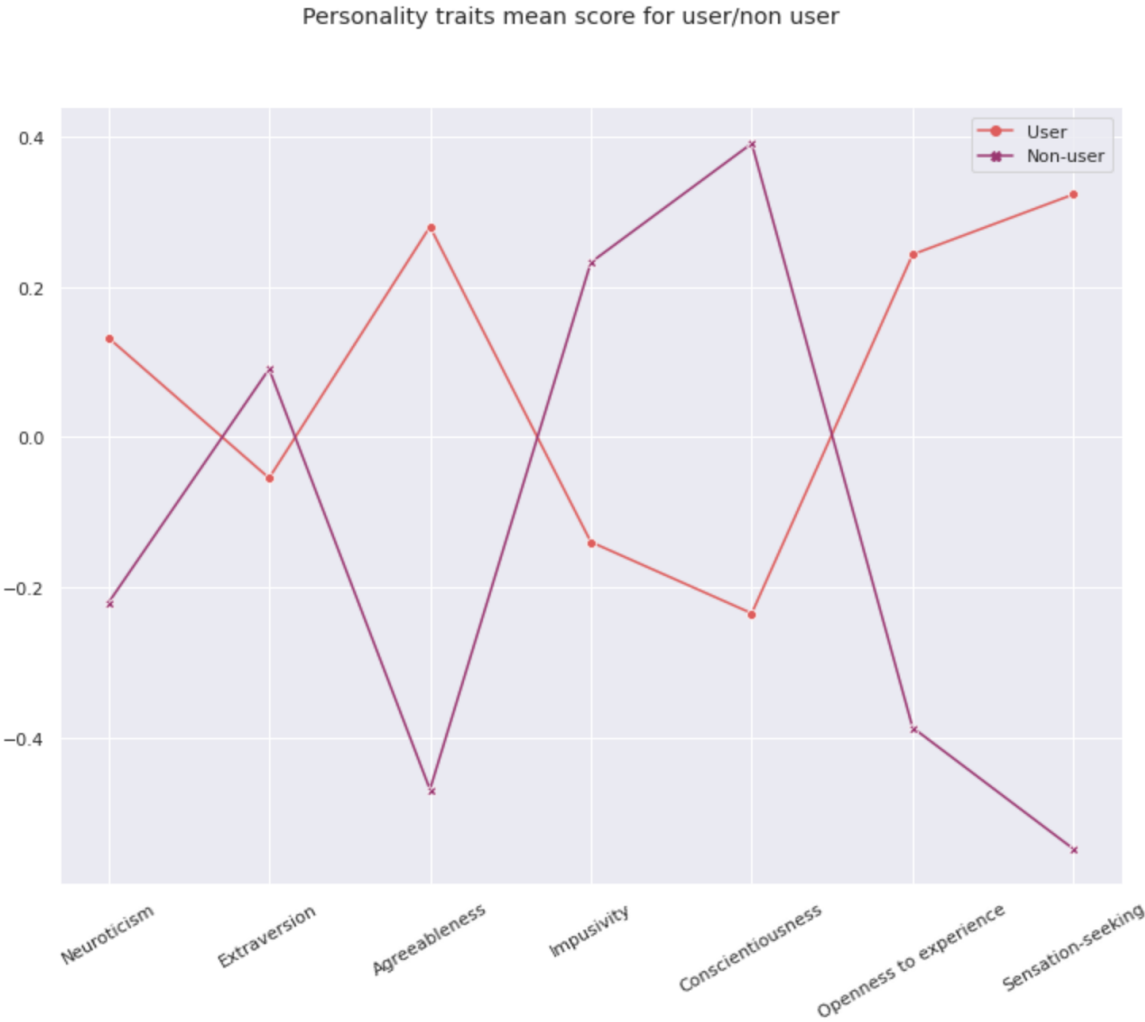
Country	%
UK	55.38
USA	29.55
Other	6.26
Canada	4.61
Australia	2.86
Republic of Ireland	1.06
New Zealand	0.26





# INTERESTINGS OBSERVATIONS

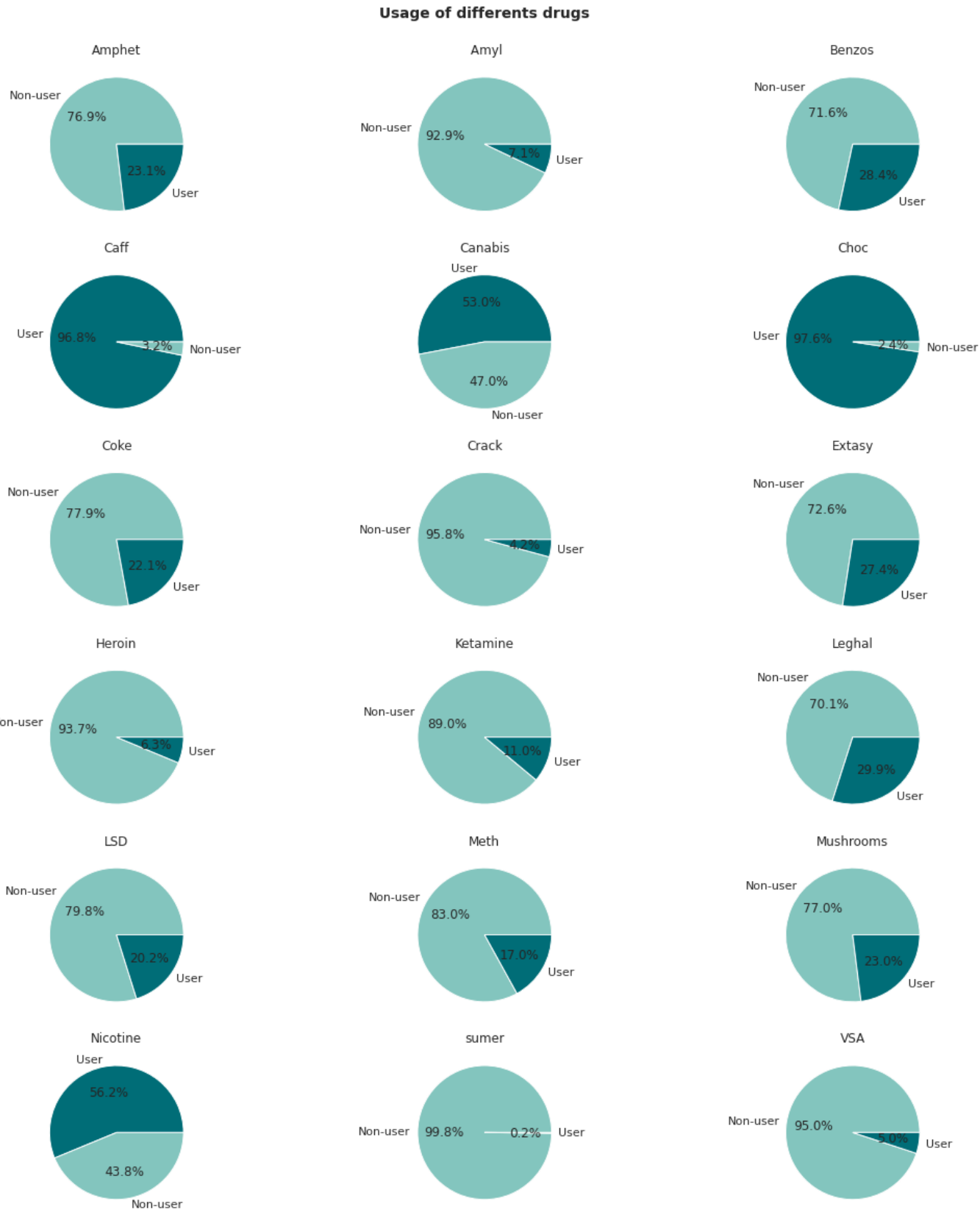
By creating a new column of **user/non-user** based on the annual consumption of illegal drugs, we plotted the mean score for each **personality traits** for cosumer or not consumers.We can see a strong correlation between the consumption and the changing behaviour, the personality traits are **diametrically opposed**. It's interesting because we may argue that it's drugs use that may change the behaviour of the participants. Because the question of how personality change would affect how we will develop our ML model later.



We can observe a widely disparate representation of user/non-user for illegal and legal drugs.

Fisrt, we can spot a small proportion of cheater 0.2%, but because it may add noise we will get rid of this column in our futur prediction.

Then we can see that the illegal drug the most used is the Cannabis with up to 47% of users, this drug is weel-known as a Getaway drug. *A gateway drug is a habit-forming drug that can lead to the use of other, more addictive drugs. They include alcohol, marijuana, nicotine, and prescription drugs.*



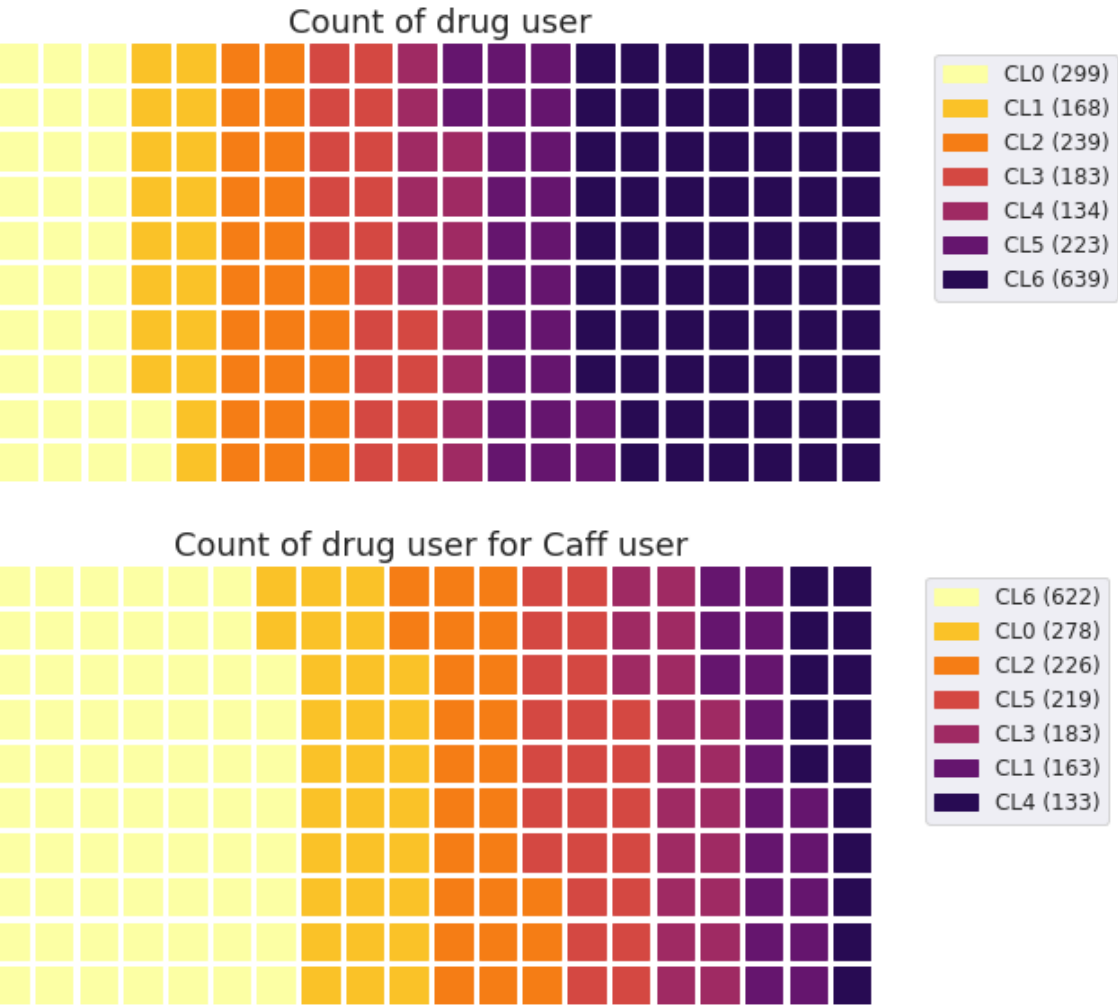




# INTERESTINGS OBSERVATIONS (2)

## Normal use

We plotted the proportion of different categories of drugs usage within the dataset. We can see a widely disparate representation of those classes, with a majority of users up to a year (rf Machine Learning part). Before studying the Getaway drugs we also poted a waffle chart for Caffeine as a referent.

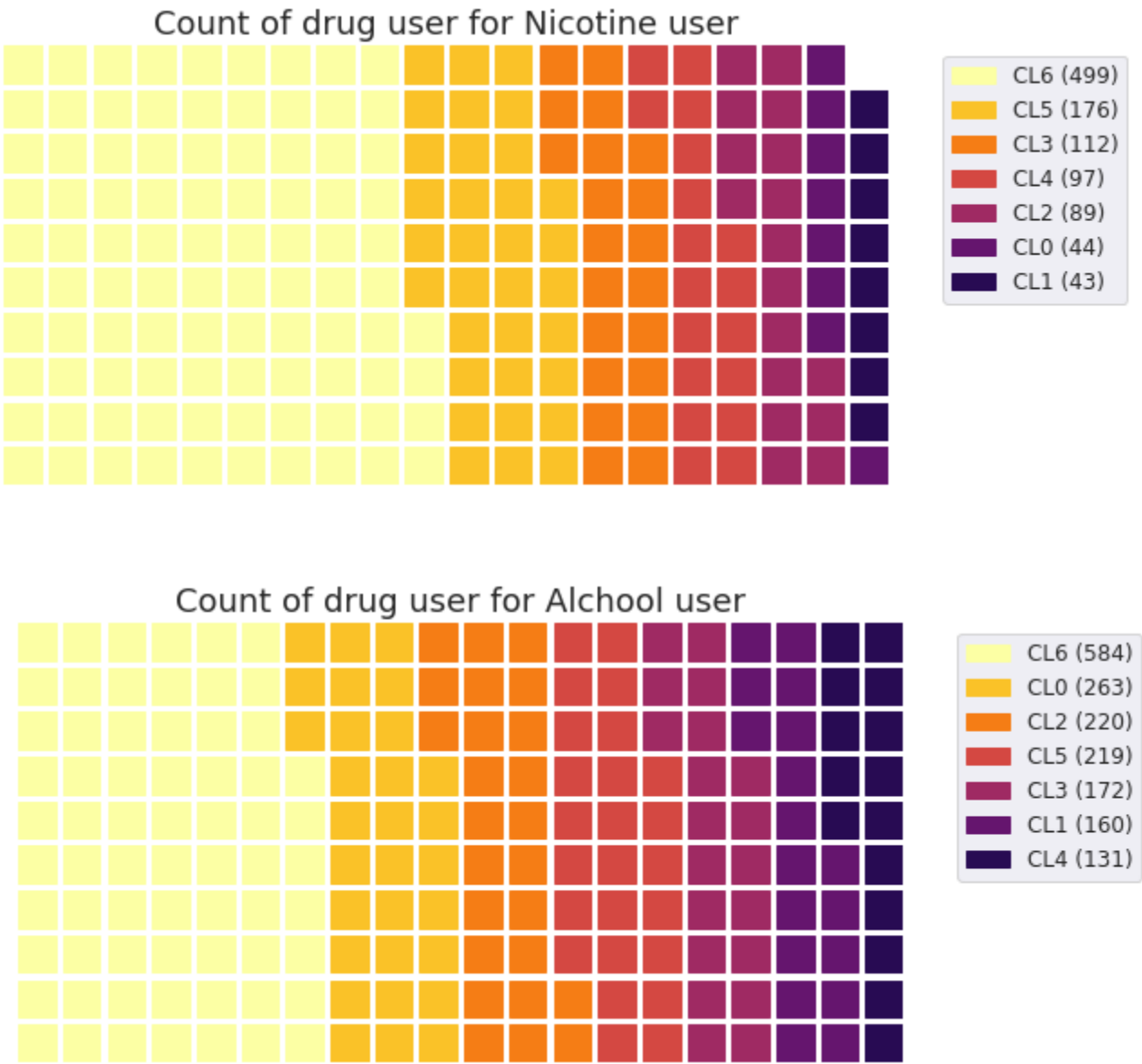


## Getaway drugs

We can observe here that the *Getaway drugs* **Nicotine** & **Alchool** seems to have a more highly count of illegal drug user.

However, we didn't get the timeline for each participants of their consumption, so it's complicated to say if the use of alchool could impact the use of Amphets for example.

We will see thoses correlations more detaillled in the Machine Learning part, but it do seems that the consumption of Nicotine or Alchool is *strongly correlated* with the consumption of others illegal drugs.



# Machine Learning







## Chosen Output

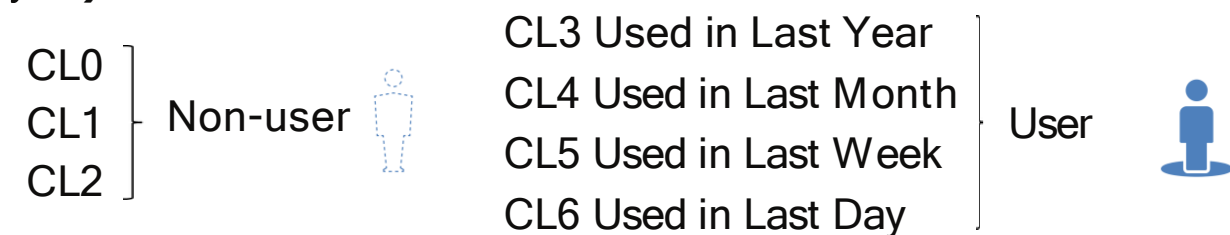
### Why not multiclass - State of Art

This dataset is well-known and a lot of work has already been done on it: from a seven-classes classification, to binary classification for each drug, going through a multiclass classification clarifying the use state of each user (*Never Used, Used over a Decade Ago, etc.*).

Thus, we wanted to bring something new to all the studies already done. We first thought about multiclass classification predicting the consumption level of the most used drug for a specific user. However, this turned out to be more complicated than expected and we decided to switch sides. Looking back at the research papers we noticed that no one had actually predict whether an individual was a user or not. This led to our choice of binary classification.

### Binary Classification

Has the individual with the characteristics X used illegal drugs lately (up to one year)?



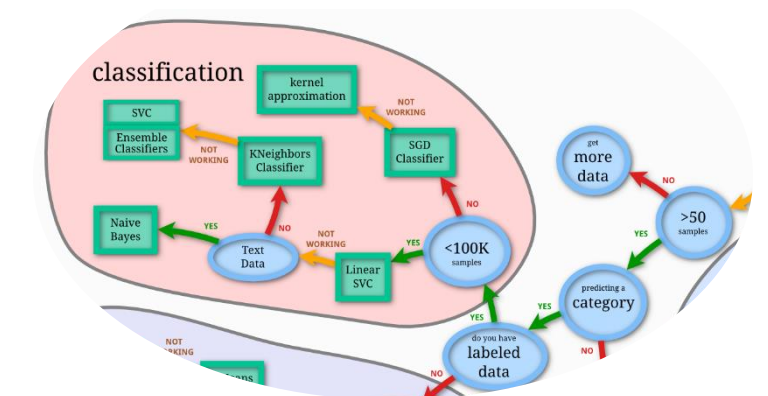
	ID	Age	Gender	Education	Nscore	Escore	Oscore	Ascore	Cscore	Impulsive	SS	Alcohol	Caff	Choc	Nicotine	Illegal drug use	Illegal drug user/non-user
0	1	0.49788	0.48246	-0.05921	0.31287	-0.57545	-0.58331	-0.91699	-0.00665	-0.21712	-1.18084	0	0	0	0	CL2	Non-user
1	2	-0.07854	-0.48246	1.98437	-0.67825	1.93886	1.43533	0.76096	-0.14277	-0.71126	-0.21575	0	0	0	0	CL4	User
2	3	0.49788	-0.48246	-0.05921	-0.46725	0.80523	-0.84732	-1.62090	-1.01450	-1.37983	0.40148	0	0	0	0	CL3	User
3	4	-0.95197	0.48246	1.16365	-0.14882	-0.80615	-0.01928	0.59042	0.58489	-1.37983	-1.18084	0	0	0	0	CL3	User
4	5	0.49788	0.48246	1.98437	0.73545	-1.63340	-0.45174	-0.30172	1.30612	-0.21712	-0.21575	0	0	0	0	CL3	User
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...



### What Algorithms did we implemented?

We tried multiple Machine Learning Algorithms from the scikit-learn library with the « Illegal drug user/non-user » as an output, by checking for each individual if he used an illegal drug up to a year before the survey. Here is a little overview of the algorithms we tested :

- Logistic Regression
- Linear SVC
- Gaussian NB
- Random Forest
- Decision Tree
- KNN
- XGBOOST
- SVM
- Neural Network



[Here](#) is a little cheat sheet of the scikit-learn library.





- ✓ XGBOOST, Linear Regression & Neural network

We employ XGBoost, a fast and strong gradient boosting framework that is also based on multiple decision tree learning, in this part. Gradient boosting is a technique for transforming weak learners into strong ones. When combined with prior models, it is based on the premise that predicting the next tree inside a model can minimize prediction errors. XGBoost, like Random Forest, provides a number of benefits, including: trains quicker and more effectively while consuming less memory. In general, very efficient, especially on huge datasets, with more complicated trees using a level-wise level-wise growth strategy.

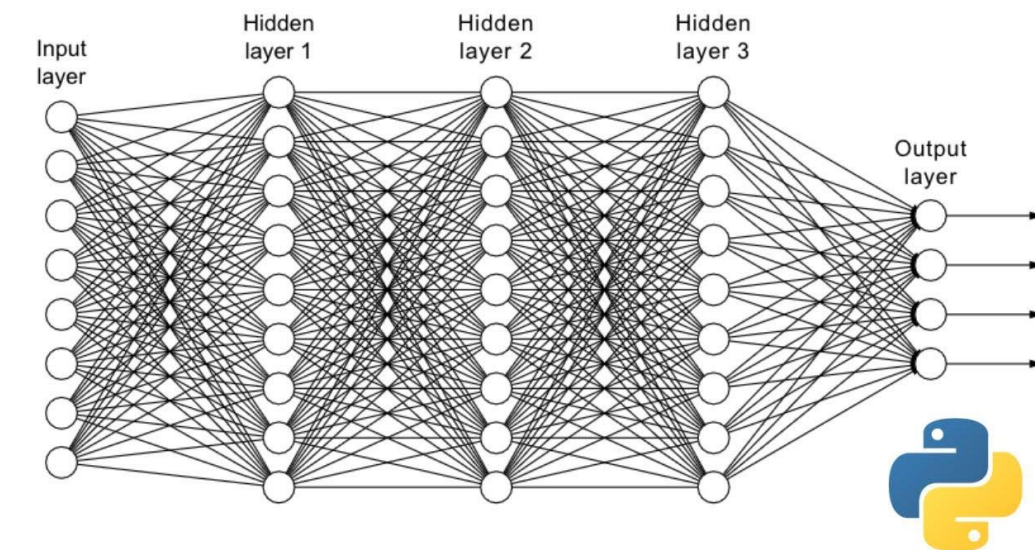
```

graph TD
    Root(f10 < -0.0662712529) -- yes --> N1(f1 < -0.60139823)
    Root -- no --> N2(f1 < 0.215821832)
    
    N1 -- yes, missing --> N3(f6 < -0.254743606)
    N1 -- no --> N4(f8 < -0.201739043)
    
    N3 -- yes, missing --> N5(f5 < -0.199988365)
    N3 -- no --> L1(leaf = 0.064000003)
    
    N5 -- yes, missing --> L2(leaf = 0.0474074073)
    N5 -- no --> L3(leaf = 0.0208695643)
    
    N4 -- yes, missing --> N6(f6 < -0.522618473)
    N4 -- no --> N7(f6 < -0.522618473)
    
    N6 -- yes, missing --> N8(f8 < -0.705397546)
    N6 -- no --> N9(f10 < -0.381085575)
    
    N8 -- yes, missing --> L4(leaf = 0.00761904754)
    N8 -- no --> L5(leaf = -0.0399999991)
    
    N9 -- yes, missing --> L6(leaf = 0.036923077)
    N9 -- no --> L7(leaf = -0)
    
    N7 -- yes, missing --> N10(f5 < -0.835874081)
    N7 -- no --> N11(f1 < 0.877009988)
    
    N10 -- yes, missing --> L8(leaf = -0.0347826071)
    N10 -- no --> L9(leaf = -0.0667961165)
    
    N11 -- yes, missing --> N12(f6 < 0.207202718)
    N11 -- no --> N13(f8 < 1.03473556)
    
    N12 -- yes, missing --> N14(f8 < 0.502556384)
    N12 -- no --> L10(leaf = 0.0242424253)
    
    N14 -- yes, missing --> L11(leaf = -0.0208695643)
    N14 -- no --> L12(leaf = -0.0320000015)
    
    N13 -- yes, missing --> L13(leaf = 0.0340000018)
    N13 -- no --> L14(leaf = -0.0700000003)
    
    N2 -- yes, missing --> N15(f6 < -0.522618473)
    N2 -- no --> N16(f5 < -0.58299711)
    
    N15 -- yes, missing --> N17(f8 < -0.201739043)
    N15 -- no --> N18(f8 < -0.19597277)
    
    N17 -- yes, missing --> L15(leaf = 0.0533333346)
    N17 -- no --> L16(leaf = -0.00551724108)
    
    N18 -- yes, missing --> L17(leaf = 0.0747826099)
    N18 -- no --> N19(f6 < -0.640072584)
    
    N19 -- yes, missing --> N20(f5 < -0.568299711)
    N19 -- no --> L18(leaf = 0.0666666627)
    
    N20 -- yes, missing --> L19(leaf = 0.0533333346)
    N20 -- no --> L20(leaf = -0.0254545435)
    
    N16 -- yes, missing --> N21(f6 < -0.390766442)
    N16 -- no --> N22(f6 < -0.358354628)
    
    N21 -- yes, missing --> L21(leaf = 0.00432432443)
    N21 -- no --> N23(f10 < -0.610730171)
    
    N23 -- yes, missing --> L22(leaf = 0.0320000015)
    N23 -- no --> L23(leaf = 0.0613333322)
    
    N22 -- yes, missing --> L24(leaf = -0.0399999991)
    N22 -- no --> L25(leaf = -0.0166666657)
  
```

For the base-line model, we consider the multiple linear regression model from sklearn.linear\_model library LogisticRegression.

In essence, neural networks are self-optimizing functions that map inputs to the desired outputs. The function may then take a fresh input and predict an output based on the function it developed with the training data.

- An input layer that receives data and pass it on
- A hidden layer
- An output layer
- Weights between the layers
- A deliberate activation function for every hidden layer.



For each of those algorithms we used a GridSearch with cross-validation to select the hyperparameters that are most suited. It entails putting a machine learning model through many permutations of hyperparameters in order to assess its performance and make improvements.

# Conclusion





## Which ones are the best model ?

### What metrics did we used ?

The goal of our projections is to anticipate whether or not someone will take drugs. We particularly want our system to enable specialists such as psychologists and others to detect current or future users. To accomplish so, we want to **decrease the rate of false positives** while **increasing our recall**.

$$\text{Recall} = \frac{TP}{TP + FN}$$

However, we'll keep an eye on other metrics like precision and accuracy, as well as the f1-score (wich is the report between recall and precision)

