

```
import pandas as pd
import sklearn as skl
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
```

```
In [2]: pip install -U notebook==pdf

Requirement already up-to-date: notebook==pdf in /opt/anaconda3/lib/python3.8/site-packages (0.3.1)
Requirement already satisfied, skipping upgrade: ipykernel in /opt/anaconda3/lib/python3.8/site-packages (fr
on notebook==pdf) (0.2.5)
Requirement already satisfied, skipping upgrade: pyppeteer in /opt/anaconda3/lib/python3.8/site-packages (fr
on notebook==pdf) (0.2.5)
Requirement already satisfied, skipping upgrade: PyPDF2 in /opt/anaconda3/lib/python3.8/site-packages (fro
m notebook==pdf) (0.2.5)
Requirement already satisfied, skipping upgrade: jupyter-core in /opt/anaconda3/lib/python3.8/site-packages (
from nbconvert->notebook==pdf) (4.6.3)
Requirement already satisfied, skipping upgrade: nbformat>=4.4 in /opt/anaconda3/lib/python3.8/site-packages (
from nbconvert->notebook==pdf) (5.0.7)
Requirement already satisfied, skipping upgrade: testpath in /opt/anaconda3/lib/python3.8/site-packages (fro
m nbconvert->notebook==pdf) (0.4.4)
Requirement already satisfied, skipping upgrade: mistune<2,>=0.8.1 in /opt/anaconda3/lib/python3.8/site-pack
ages (from nbconvert->notebook==pdf) (0.8.4)
Requirement already satisfied, skipping upgrade: pygments<2.0.0,>=1.4.1 in /opt/anaconda3/lib/python3.8/site-pa
ckages (from nbconvert->notebook==pdf) (2.6.1)
Requirement already satisfied, skipping upgrade: pandocfilters>=1.4.1 in /opt/anaconda3/lib/python3.8/site-p
ackages (from nbconvert->notebook==pdf) (1.4.2)
Requirement already satisfied, skipping upgrade: defusedxml in /opt/anaconda3/lib/python3.8/site-packages (f
rom nbconvert->notebook==pdf) (0.6.0)
Requirement already satisfied, skipping upgrade: mistune<2,>=0.8.1 in /opt/anaconda3/lib/python3.8/site-pack
ages (from nbconvert->notebook==pdf) (0.8.4)
Requirement already satisfied, skipping upgrade: pygments<2.0.0,>=1.4.1 in /opt/anaconda3/lib/python3.8/site-pa
ckages (from nbconvert->notebook==pdf) (2.11.2)
Requirement already satisfied, skipping upgrade: bleach in /opt/anaconda3/lib/python3.8/site-packages (from
nbconvert->notebook==pdf) (3.1.5)
Requirement already satisfied, skipping upgrade: traitlets<4.2,>=0.8.1 in /opt/anaconda3/lib/python3.8/site-packa
ge (from nbconvert->notebook==pdf) (4.3.3)
Requirement already satisfied, skipping upgrade: pyee<9.0.0,>=8.1.0 in /opt/anaconda3/lib/python3.8/site-pac
kages (from pyppeteer->notebook==pdf) (8.1.0)
Requirement already satisfied, skipping upgrade: pytest>0 in /opt/anaconda3/lib/python3.8/site-packages (fr
om nbconvert->notebook==pdf) (0.2.5)
Requirement already satisfied, skipping upgrade: pygments<2.0.0,>=1.4.1 in /opt/anaconda3/lib/python3.8/site-pa
ckages (from nbconvert->notebook==pdf) (2.11.2)
Requirement already satisfied, skipping upgrade: pygments<2.0.0,>=1.4.1 in /opt/anaconda3/lib/python3.8/site-p
ackages (from pyppeteer->notebook==pdf) (2.6.1)
Requirement already satisfied, skipping upgrade: websocket<9.0,>=8.1 in /opt/anaconda3/lib/python3.8/site-p
ackages (from pyppeteer->notebook==pdf) (1.2.5)
Requirement already satisfied, skipping upgrade: jenschema<2.5.0,>=2.4 in /opt/anaconda3/lib/python3.8/sit
-packages (from nbformat>=4.4->nbconvert->notebook==pdf) (3.2.0)
Requirement already satisfied, skipping upgrade: packaging in /opt/anaconda3/lib/python3.8/site-packages (fr
om bleach->nbconvert->notebook==pdf) (1.1.3)
Requirement already satisfied, skipping upgrade: MarkupSafe>=0.23 in /opt/anaconda3/lib/python3.8/site-packa
ges (from Jinja2>=2.0->nbconvert->notebook==pdf) (1.1.3)
Requirement already satisfied, skipping upgrade: webencodings in /opt/anaconda3/lib/python3.8/site-packages (
from bleach->nbconvert->notebook==pdf) (0.5.1)
Requirement already satisfied, skipping upgrade: six>=1.9.0 in /opt/anaconda3/lib/python3.8/site-packages (f
rom traitlets>=4.2->nbconvert->notebook==pdf) (1.15.0)
Requirement already satisfied, skipping upgrade: decorator in /opt/anaconda3/lib/python3.8/site-packages (fr
om traitlets>=4.2->nbconvert->notebook==pdf) (4.4.2)
Requirement already satisfied, skipping upgrade: attr>=17.4.0 in /opt/anaconda3/lib/python3.8/site-packages (
from jenschema<2.5.0,>=2.4->nbformat>=4.4->nbconvert->notebook==pdf) (19.3.0)
Requirement already satisfied, skipping upgrade: pytest>0 in /opt/anaconda3/lib/python3.8/site-packages (fr
om jenschema<2.5.0,>=2.4->nbformat>=4.4->nbconvert->notebook==pdf) (49.2.0.post20200714)
Requirement already satisfied, skipping upgrade: pyrsistent>=0.14.0 in /opt/anaconda3/lib/python3.8/site-pac
kages (from jenschema<2.5.0,>=2.4->nbformat>=4.4->nbconvert->notebook==pdf) (0.16.0)
Requirement already satisfied, skipping upgrade: pyrsistent>=0.2.2 in /opt/anaconda3/lib/python3.8/site-packa
ges (from packaging->bleach->nbconvert->notebook==pdf) (0.4.7)
Note: you may need to restart the kernel to use updated packages.
```

```
In [3]: dataset = pd.read_csv('drug_consumption.csv')
dataset.head()
```

Describing the dataset

```
In [4]: #first we want to have a preview of our dataset
dataset
```

	ID	Age	Gender	Education	Country	Ethnicity	Nscore	Escore	Oscore	Ascore	...	Ecstasy	Heroin	Ketamine
0	1	0.49788	0.48246	-0.05921	0.96082	0.12600	0.31287	-0.57545	-0.58331	-0.91699	...	CL0	CL0	CL0
1	2	-0.07854	-0.48246	-0.05921	0.96082	-0.31685	-0.67825	1.93886	1.43533	0.76096	...	CL4	CL0	CL2
2	3	0.49788	-0.48246	-0.05921	0.96082	-0.31685	-0.46725	0.80523	-0.84732	-1.62090	...	CL0	CL0	CL0
3	4	-0.95197	0.48246	1.16365	0.96082	-0.31685	-0.14882	-0.80015	-0.09282	0.59042	...	CL0	CL0	CL2
4	5	-0.49788	0.48246	1.98437	0.96082	-0.31685	0.73545	-0.63400	-0.45174	-0.30172	...	CL1	CL0	CL0
...														
1880	1884	-0.95197	0.48246	-0.06113	-0.57009	-0.31685	-1.19430	1.74091	1.88511	0.76096	...	CL0	CL0	CL0
1881	1885	-0.95197	-0.48246	-0.06113	-0.57009	-0.31685	-0.24649	1.74091	0.58331	0.76096	...	CL2	CL0	CL0
1882	1886	-0.07854	0.48246	0.45468	-0.57009	-0.31685	1.13281	-1.37639	-1.27553	-1.77200	...	CL0	CL0	CL2
1883	1887	-0.95197	0.48246	-0.06113	-0.57009	-0.31685	0.91093	-1.92173	0.29338	-1.62090	...	CL3	CL0	CL0
1884	1888	-0.95197	-0.48246	-0.06113	0.21128	-0.31685	-0.46725	2.12700	1.65653	1.11408	...	CL3	CL0	CL0

```
1885 rows x 32 columns

In [5]: dataset.describe()
```

	ID	Age	Gender	Education	Country	Ethnicity	Nscore	Escore	Oscore	As
count	1885.000000	1885.000000	1885.000000	1885.000000	1885.000000	1885.000000	1885.000000	1885.000000	1885.000000	1885.00
mean	945.294960	0.03461	-0.000256	-0.003806	0.355542	-0.309577	0.000047	-0.000163	-0.000534	-0.00
std	545.167641	0.07836	0.482588	0.950078	0.700325	0.166226	0.998105	0.997448	0.996229	0.99
min	1.000000	-0.95197	-0.482460	-2.435910	-0.570090	-1.107020	-3.64360	-3.273930	-3.273930	-3.46
25%	474.000000	-0.95197	-0.482460	-0.011130	-0.570090	-0.316850	-0.678250	-0.695090	-0.717720	-0.60
50%	945.000000	-0.07854	0.482460	-0.059210	0.960820	-0.316850	0.042670	0.003320	-0.019280	-0.01
75%	147.000000	0.49788	0.482460	0.454680	0.960820	-0.316850	0.629760	0.637790	0.723300	0.76
max	1888.000000	2.59171	0.482460	1.984370	0.960820	1.907250	3.273930	3.273930	2.901610	3.46

```
In [6]: #shape of the dataset
dataset.shape
```

Out[6]: (1885, 32)

```
In [7]: #columns of the dataset
dataset.columns
```

Out[7]: Index(['ID', 'Age', 'Gender', 'Education', 'Country', 'Ethnicity', 'Nscore', 'Escore', 'Oscore', 'Ascore', 'Cscore', 'Impulsive', 'SS', 'Alcohol', 'Amphet', 'Amyl', 'Benzos', 'Caff', 'Cannabis', 'Choc', 'Coke', 'Crack', 'Ecstasy', 'Heroin', 'Ketamine', 'Legal', 'LSD', 'Meth', 'Mushrooms', 'Nicotine', 'Semer', 'VSA'], dtype='object')

list of features: ['ID', 'Age', 'Gender', 'Education', 'Country', 'Ethnicity', 'Nscore', 'Escore', 'Oscore', 'Ascore', 'Cscore', 'Impulsive', 'SS', 'Alcohol', 'Amphet', 'Amyl', 'Benzos', 'Caff', 'Cannabis', 'Choc', 'Coke', 'Crack', 'Ecstasy', 'Heroin', 'Ketamine', 'Legal', 'LSD', 'Meth', 'Mushrooms', 'Nicotine', 'Semer', 'VSA']

We need to discuss a bit about the variables :

- ID : number of record in original database
- Age : Age of the of participants
- Gender : Gender of the participant (0.48246 -> Female , -0.48246 -> Male)
- Education : Level of education :
 - 2.43591 Left school before 16
 - 1.73790 Left school at 16 years
 - 1.43719 Left school at 17 years
 - 1.22751 Left school at 18 years
 - 0.61113 Some college or university, no certificate or degree
 - 0.05921 Professional certificate/ diploma
 - 0.45468 University degree
 - 1.16365 Masters degree
 - 1.98437 Doctorate degree
- Country: current residence country :
 - 0.09765 Australia
 - 0.24923 Canada
 - 0.45681 New Zealand
 - 0.28519 Other
 - 0.21128 Republic of Ireland
 - 0.96882 UK
 - 0.57009 USA
- Ethnicity : ethnicity of participants:
 - 0.50212 Asian
 - 1.17082 Black
 - 1.98725 Mixed-Black/Asian
 - 0.12680 Mixed-White/Asian
 - 0.22166 Mixed-White/Black
 - 0.11440 Other
 - 0.31685 White
- Nscore is NEO-FFI-R Neuroticism
- Escore is NEO-FFI-R Extraversion
- Oscore is NEO-FFI-R Openness to experience
- Ascore is NEO-FFI-R Agreeableness
- Cscore is NEO-FFI-R Conscientiousness
- Impulsive is impulsiveness measured by BIS-11
- SS is sensation seeking measured by Imp55

for the 7 previous ones, features are numerical, they describes score. The higher your score, the more you'll be.

- for the last one the features are categorical.
- 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

```
In [8]: #we want also to see the types of the features
dataset.dtypes
```

ID	int64
Age	float64
Gender	float64
Education	float64
Country	float64
Ethnicity	float64
Nscore	float64
Escore	float64
Oscore	float64
Ascore	float64
Cscore	float64
Impulsive	float64
SS	float64
Alcohol	object
Amphet	object
Amyl	object
Benzos	object
Caff	object
Cannabis	object
Choc	object
Coke	object
Crack	object
Ecstasy	object
Heroin	object
Ketamine	object
Legal	object
LSD	object
Meth	object
Mushrooms	object
Nicotine	object
Semer	object
VSA	object
dtype:	object

```
In [9]: dataset.isna().sum()
```

ID	0
Age	0
Gender	0
Education	0
Country	0
Ethnicity	0
Nscore	0
Escore	0
Oscore	0
Ascore	0
Cscore	0
Impulsive	0
SS	0
Alcohol	0
Amphet	0
Amyl	0
Benzos	0
Caff	0
Cannabis	0
Choc	0
Coke	0
Crack	0
Ecstasy	0
Heroin	0
Ketamine	0
Legal	0
LSD	0
Meth	0
Mushrooms	0
Nicotine	0
Semer	0
VSA	0
dtype:	int64

no missing values in columns -> we don't need to delete some rows

Previously we saw we had some categorical variables by studying on the Dataset, let's convert them

```
In [10]: categoricals = ['Age', 'Gender', 'Education', 'Country',
                    , 'Ethnicity', 'Alcohol', 'Amphet', 'Amyl', 'Benzos', 'Caff', 'Cannabis',
                    , 'Choc', 'Coke', 'Crack', 'Ecstasy', 'Heroin',
                    , 'Ketamine', 'Legal', 'LSD', 'Meth', 'Mushrooms', 'Nicotine', 'Semer', 'VSA']
numericals = list(set(dataset.columns)-set(categoricals))
numericals.remove('ID')
```

```
Out[10]: ['Impulsive', 'Escore', 'Cscore', 'Ascore', 'Nscore', 'SS', 'Oscore']
```

```
In [11]: dataset.shape
```

Out[11]: (1885, 32)

let's try to replace some values by their meanings

```
In [12]: dataset['Gender'].unique()
d = dict(zip(dataset['Age'].unique(), ['Female', 'Male']))
dataset['Gender'].replace(d, inplace=True)
```

```
In [13]: #replacing age categories
temp=dict(zip(dataset['Age'].unique(), ['35-44', '25-34', '18-24', '65+', '45-54', '55-64']))
dataset['Age'].replace(temp, inplace=True)
```

```
In [14]: data['Education'].unique()
temp=dict(zip(dataset['Education'].unique(), ['Professional Deg', 'Doctorate Deg',
'Master Deg', '18 Left', '16 Left', 'Uni Deg',
'Uni No Deg', '16 Left', '17 Left']))
dataset['Education'].replace(temp, inplace=True)
```

```
In [16]: #replacing countries
temp=dict(zip(dataset['Country'].unique(), ['UK', 'Canada', 'USA', 'Other', 'Australia', 'Ireland',
, 'NZ']))
dataset['Country'].replace(temp, inplace=True)
```

```
In [17]: dataset['Country'].value_counts()
# We can see we have the same stats -> successful replacement
```

UK	1044
USA	557
Other	118
Canada	87
Australia	54
Ireland	20
NZ	5
Name: Country, dtype: int64	

```
In [18]: #Replacing ethnicity
temp=dict(zip(dataset['Ethnicity'].unique(),
['White/Asian', 'White', 'Other', 'White/Black', 'Asian', 'Black', 'Black/Asian']))
dataset['Ethnicity'].replace(temp, inplace=True)
```

```
In [19]: dataset[categoricals] = dataset[categoricals].astype("category")
dataset[numericals] = dataset[numericals].astype("float")
dataset.dtypes
```

ID	int64
Age	category
Gender	category
Education	category
Country	category
Ethnicity	category
Nscore	float64
Escore	float64
Oscore	float64
Cscore	float64
Ascore	float64
Impulsive	float64
SS	float64
Alcohol	category
Amphet	category
Amyl	category
Benzos	category
Caff	category
Cannabis	category
Choc	category
Coke	category
Crack	category
Ecstasy	category
Heroin	category
Ketamine	category
Legal	category
LSD	category
Meth	category
Mushrooms	category
Nicotine	category
Semer	category
VSA	category
dtype:	object

```
In [20]: dataset['Ethnicity'].value_counts()
#successful replacement too
```

White	1720
Other	63
Black	33
Asian	26
White/Black	20
White/Asian	20
Black/Asian	8
Name: Ethnicity, dtype: int64	

```
In [21]: dataset.unique()
```

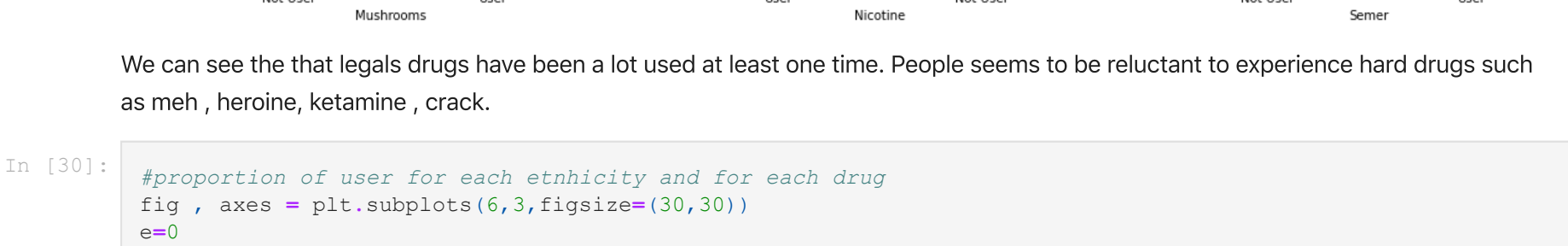
Age	1885
Gender	2
Education	8
Country	7
Ethnicity	7
Nscore	49
Escore	42
Oscore	35
Ascore	41
Cscore	41
Impulsive	10
SS	11
Alcohol	7
Amphet	7
Amyl	7
Benzos	7
Caff	7
Cannabis	7
Choc	7
Coke	7
Crack	7
Ecstasy	7
Heroin	7
Ketamine	7
Legal	7
LSD	7
Meth	7
Mushrooms	7
Nicotine	5
Semer	5
VSA	7
dtype:	int64

```
In [22]: dataset['Semer'].unique()
```

Out[22]: [CL0, CL2, CL3, CL4, CL1]
Categories (5, object): [CL0, CL2, CL3, CL4, CL1]

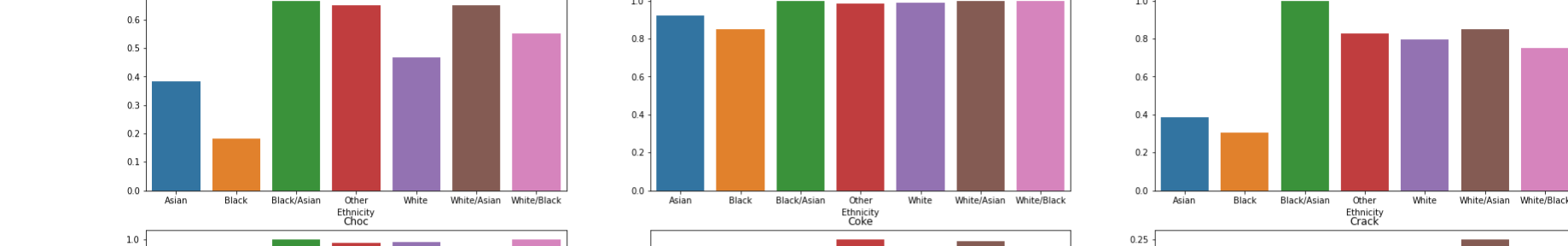
Ok, so now let's try to visualise the datas

```
In [23]: # see proportions of the first 5 features
sns.countplot(dataset['Gender'])
# we have almost the same number of male and female in our dataset
```



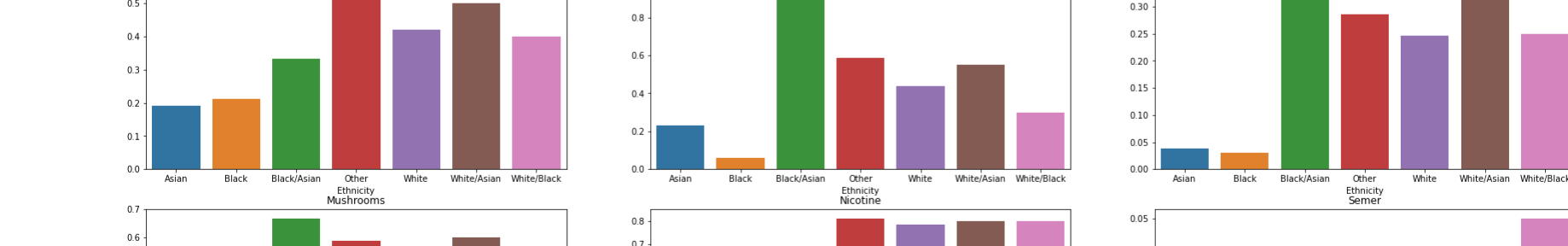
Out[23]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb76e66a760>

```
In [24]: # visualize the proportions of current residence country
sns.countplot(dataset['Country'])
#UK and USA are the most represented here
```



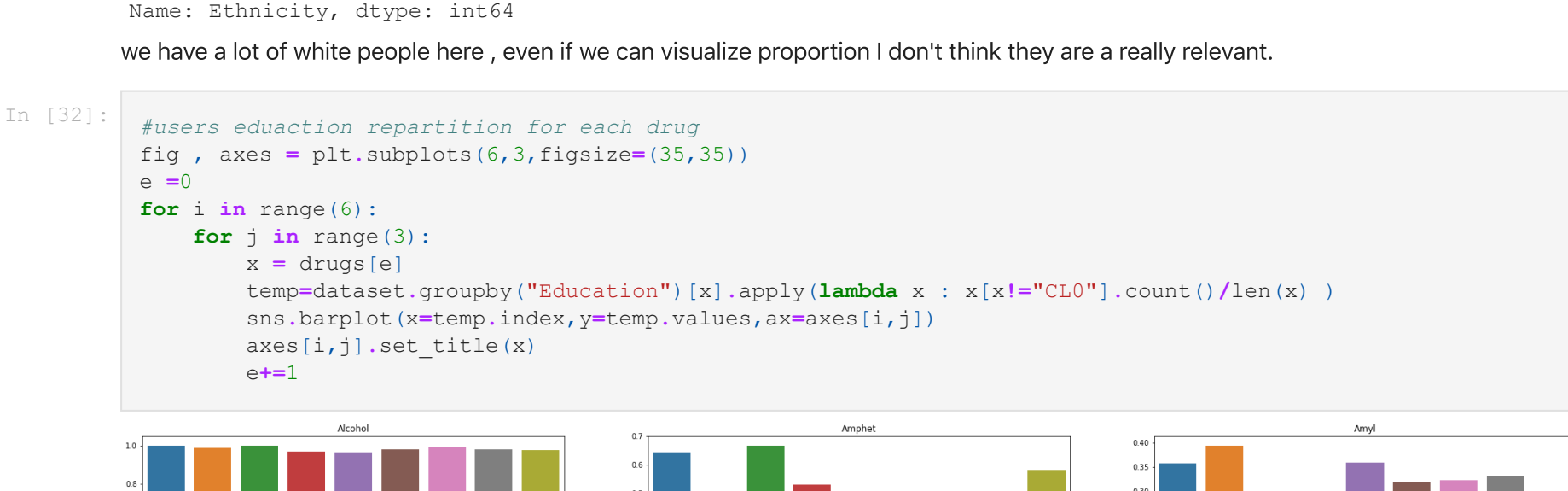
Out[24]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb7687825e0>

```
In [25]: #Same thing for ethnicity
sns.countplot(dataset['Ethnicity'])
#huge proportion of participants are white
```



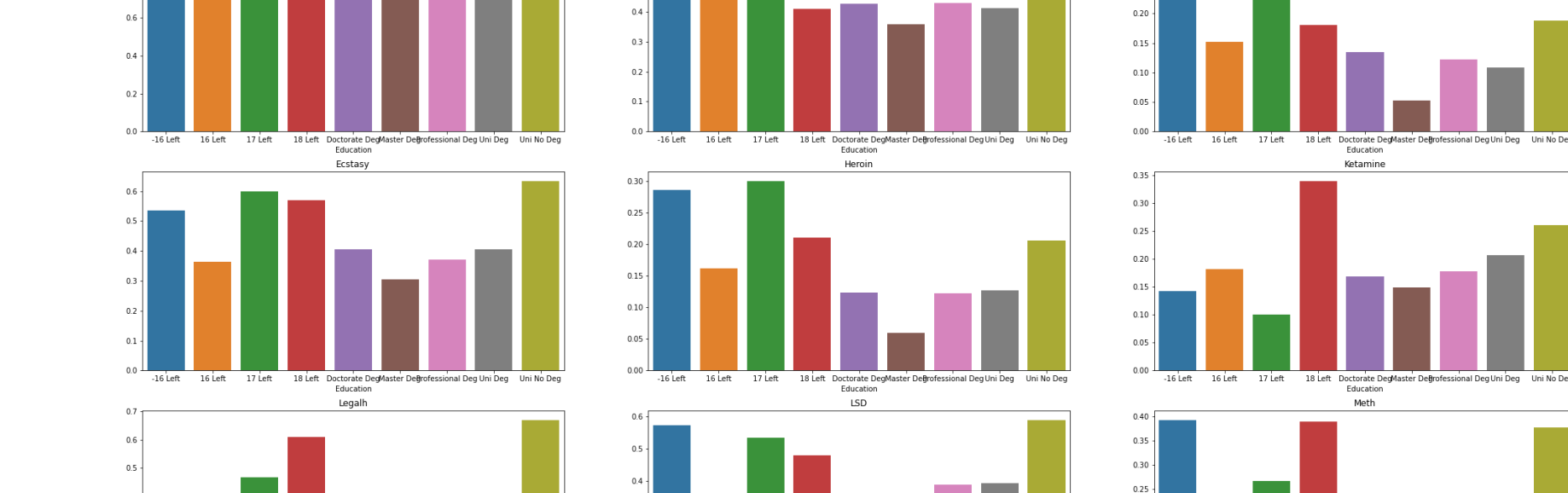
Out[25]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb7687825e0>

```
In [26]: # Now, let's try to visualise the distribution of dimensions of NEO Personality Inventory on our sample.
NEO = list(set(numericals)-set(['SS', 'Impulsive']))
n = len(NEO)
i=0
fig, axes = plt.subplots(1, n, figsize=(15, 4), sharey=True)
for element in NEO:
    sns.distplot(dataset[element], ax=axes[i])
    i+=1
```

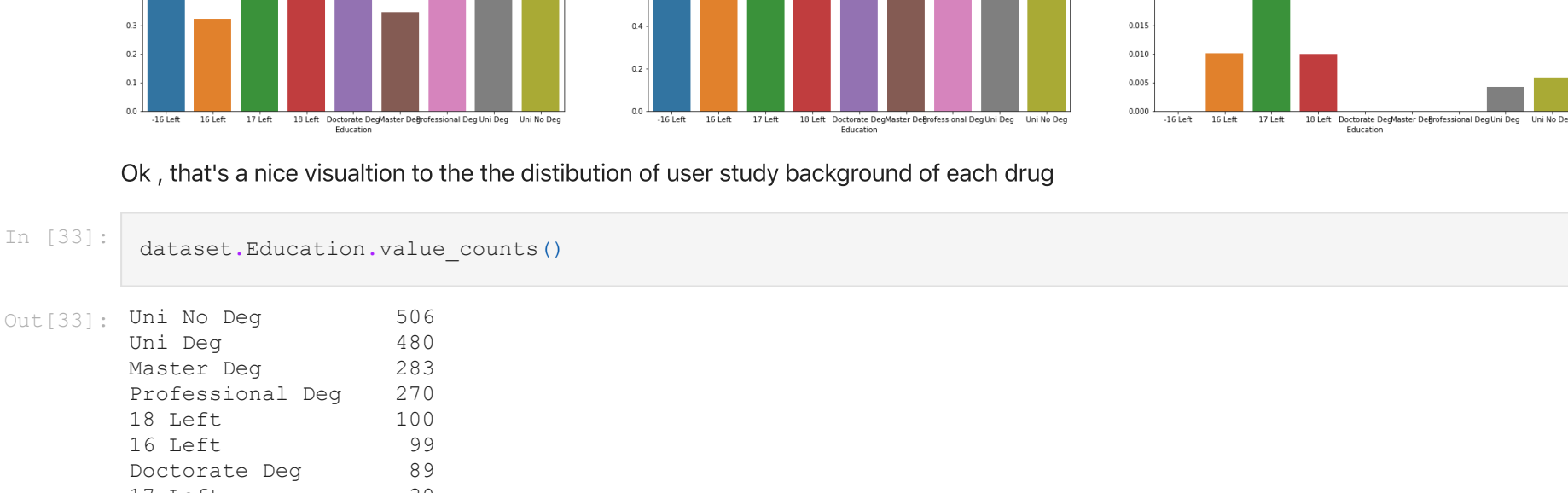


Score seems to be distributed like a Normal distribution -> need to scale or not ?

```
In [27]: # For SS and Impulsive
fig, axes = plt.subplots(1, 2, figsize=(15, 4))
i=0
for element in SS_IMP:
    sns.distplot(dataset[element], ax=axes[i])
    i+=1
```



```
In [28]: # visualization for the drugs
drugs = ['Alcohol',
'Amphet', 'Amyl', 'Benzos', 'Caff', 'Cannabis', 'Choc', 'Coke', 'Crack',
'Ecstasy', 'Heroin', 'Ketamine', 'Legal', 'LSD', 'Meth', 'Mushrooms',
'Nicotine', 'Semer', 'VSA']
c=0
for i in range(6):
    x = drugs[i]
    for j in range(3):
        sns.countplot(dataset[drugs[i]], ax=axes[c+j])
        c+=1
```



We can see that "legal" or "almost-legal" drugs are highly use (caf, choc, nicotine, cannabis, legal, benzos)

Maybe we could try to classify drugs into two subparts such as legal and not legal ?

Maybe it could help to see the influence of the law about drug consumption

Maybe we can visualise the proportions of user and not user for each drug too

```
In [29]: #user vs not user for each drug
fig, axes = plt.subplots(6, 3, figsize=(20, 20))
c=0
for i in range(6):
    for j in range(3):
        element = [ 'Not User' if c=="CL0" else 'User' for c in dataset[x] ]
        x = drugs[i]
        sns.countplot(element, columns = [x])
        axes[i,j].set_title(x+'', total User = [j].format(temp.values[i]+temp.values[0]))
        c+=1
```




It seems that mens are more likely to experiment drugs.

```
In [35]: from sklearn.preprocessing import LabelEncoder
col = list(dataset(drugs))
for c in col:
    le = LabelEncoder()
    data[c] = le.fit_transform( data[c])

In [36]: data.drop("ID",axis=1,inplace=True)

In [37]: data

Out[37]:
```

	Age	Gender	Education	Country	Ethnicity	Nscore	Escore	Oscore	Ascore	Cscore	...	Ecstasy	Heroin	Katam
0	0.49788	0.48246	-0.05921	0.96082	0.12600	0.31287	-0.57545	-0.58331	-0.91699	-0.00665	...	0	0	0
1	-0.07854	-0.48246	1.98437	0.96082	-0.31685	-0.67825	1.93886	1.43531	0.76096	-0.14277	...	4	0	0
2	0.49788	-0.48246	-0.05921	0.96082	-0.31685	-0.46725	0.80623	-0.84732	-1.62090	-0.14160	...	0	0	0
3	-0.95197	0.48246	1.16365	0.96082	-0.31685	-0.14882	-0.80615	-0.01928	0.59042	0.58489	...	0	0	0
4	0.49788	0.48246	1.98437	0.96082	-0.31685	0.73545	-1.63340	-0.45174	-0.30172	1.30612	...	1	0	0
...
1880	-0.95197	0.48246	-0.61113	-0.57009	-0.31685	-1.19430	1.74091	1.88511	0.76096	-1.13788	...	0	0	0
1881	-0.95197	-0.48246	-0.61113	-0.57009	-0.31685	-0.24649	1.74091	0.58331	0.76096	-1.51840	...	2	0	0
1882	-0.07854	0.48246	0.45466	-0.57009	-0.31685	1.13281	-1.37639	-1.27553	-1.77200	-1.38502	...	4	0	0
1883	-0.95197	-0.48246	-0.61113	-0.57009	-0.31685	0.91093	-1.92173	0.29338	-1.62090	-2.57309	...	3	0	0
1884	-0.95197	0.48246	-0.61113	0.21128	-0.31685	-0.46725	2.12700	1.65553	1.11408	0.41594	...	3	0	0

1885 rows x 31 columns

```
In [ ]: corrmat = data.corr()

plt.figure(figsize=(20,20))

sns.heatmap(corrmat,cmap = 'Reds',annot=True,
            yticklabels = data.columns, xticklabels = data.columns)
plt.xticks(rotation=90)

plt.title("Correlation b/w Different Features",fontsize=18)
plt.show()
```

we can see there a strong correlation for some drugs, it could be expected because if you have experienced some hard drugs you are more likely to experienced others.

Now it's to time to build a prediction model to evaluate consumption of Cannabis knowing score and informations on a person.

KNN

```
In [ ]: features = ['Age', 'Gender', 'Education', 'Country', 'Ethnicity', 'Nscore',
                'Escore', 'Oscore', 'Ascore', 'Cscore', 'Impulsive', 'SS']
target = ['Cannabis']

In [ ]: from sklearn.model_selection import train_test_split
X = data[features]
y = data[target].astype("category")

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20)
```

```
In [ ]: from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score
knn = KNeighborsClassifier(n_neighbors=1)
mod.fit(X_train,y_train.values.ravel())
pred = mod.predict(X_test)
accuracy_score(y_test,pred)
```

Trying to know which is the best number of neighbours

```
In [ ]: from sklearn.neighbors import KNeighborsClassifier

graph= dict()
for i in range(1,100):
    temp = KNeighborsClassifier(n_neighbors=i)
    temp.fit(X_train,y_train.values.ravel())
    pred_temp = temp.predict(X_test)
    graph[i]=accuracy_score(y_test,pred_temp)

x=list(graph.keys())
y=list(graph.values())
plt.plot(x,y)
plt.plot(x,y)
sorted_d = sorted(graph.items(), key=lambda x: x[1],reverse=True)
sorted_d[10]
```

```
In [ ]: from sklearn.metrics import confusion_matrix
knn = KNeighborsClassifier(sorted_d[10][0])
knn.fit(X_train,y_train.values.ravel())
final_pred = knn.predict(X_test)
kaccres = precision_score(y_test2,yfinal),accuracy_score(y_test2,yfinal)
kaccres
```

```
In [ ]: sns.heatmap(confusion_matrix(y_test2,yfinal),annot=True)

print(kaccres)
```

prediction is much better

Random Forrest

```
In [ ]: from sklearn.ensemble import RandomForestClassifier

In [ ]: rf = RandomForestClassifier()
logscores = precision_score(y_test2,rf.predict(), accuracy_score(y_test2,rf.predict)
logscores
```

searching the best number of tree to use like we did for knn

rep = dict() for i in range (10,1000, 10): rf = RandomForestClassifier() rf.fit(X_train2,y_train2.values.ravel())
rdpredict=rf.predict(X_test2) rep[i]= accuracy_score(rdpredict,y_test2) plt.plot(list(rep.keys()),list(rep.values()))
max_acc = sorted(rep.items(), key=lambda x: x[1],reverse=True)[0][0] max_acc[0]
rfnal = RandomForestClassifier(n_estimators=max_acc[0]) rf.fit(X_train2,y_train2.values.ravel()) rdpredict=rf.predict(X_test2)
precision_score(rdpredict,y_test2)

Logistic Regression

```
In [ ]: from sklearn.linear_model import LogisticRegression
lr = LogisticRegression(max_iter=100)
lr.fit(X_train2,y_train2)
y_rf = lr.predict(X_test2)
rdscores = precision_score(y_test2,y_rf), accuracy_score(y_test2,y_rf)
rdscores
```

```
In [ ]: from sklearn.ensemble import GradientBoostingClassifier
GBC = GradientBoostingClassifier()
GBC.fit(X_train2,y_train2)
pred_GBC = GBC.predict(X_test2)
GBCscore=precision_score(y_test2,pred_GBC),accuracy_score(y_test2,pred_GBC)
GBCscores
```

```
In [ ]: models = ["KNN",
               "Random Forest",
               "Logistic Regression",
               "Gradient Boosting"]
precisions = [kscores[0],
              rdscores[0],
              logscores[0],
              GBCscores[0]]
accuracies = [kscores[1],
              rdscores[1],
              logscores[1],
              GBCscores[1]]
df = pd.DataFrame(zip(models,precisions,accuracies),columns=["Model","Prec","Acc"])
```

```
In [ ]: fig, axes = plt.subplots(1,2,figsize=(20,5))

sns.barplot(x=df["Model"],
            y=df["Prec"],
            ax=axes[0],palette="Greens_d")

axes[0].set_title("Precision comparison for each model")

for ax in axes:
    for p in ax.patches:
        h = p.get_height()
        ax.text(p.get_x() + (p.get_width()/2), h*2.75/3, '{}'.format(round(h,3)), ha="center", color="black")
sns.barplot(x=df["Model"],
            y=df["Acc"],
            ax=axes[1],
            palette="Greens_d")

axes[1].set_title("Accuracy comparison for each model")

for ax in axes:
    for p in ax.patches:
        h = p.get_height()
        ax.text(p.get_x() + (p.get_width()/2), h*2.75/3, '{}'.format(round(h,3)), ha="center", color="black")
```

Random forest seems to be the best model to deal with precision and accuracy

```
In [ ]: sns.heatmap(confusion_matrix(rf_predict,y_test2),annot=True,cmap="Reds")
```

Let's try to improve our Random Forest regression model

```
In [ ]: from sklearn.model_selection import GridSearchCV

In [ ]: param_grid = {
    'bootstrap': [True],
    'criterion': ['gini','entropy'],
    'min_samples_leaf': [1,5,10],
    'max_depth': range(1,12),
    'max_features': ['sqrt'],
    'n_estimators': [75, 100, 200, 300, 500]
}
```

```
In [ ]: grid_search = GridSearchCV(estimator = rf, scoring = ['accuracy','precision'], param_grid = param_grid,
                             cv = 5, n_jobs = -1, verbose = 2, refit='precision')
```

```
In [ ]: grid_search.fit(X_train2, y_train2)
```

```
In [ ]: best_grid_search.best_estimator_

In [ ]: grid_search.best_score_

In [ ]: params = grid_search.best_params_
params
```

```
In [ ]: best_score(X_test2, y_test2)

In [ ]: precision_score(y_test2,best.predict(X_test2))
```

from joblib import dump best_fit(X_train2,y_train2)

dump(best, 'best_model.joblib')

```
In [ ]:

In [ ]:

In [ ]:
```