

# INVESTIGATION OF A MEDICAL PATIENTS DATASET

The goal of this analysis is to investigate the 'noshowappointments dataset' and determine the relationship between the variables.

## Importing Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
matplotlib inline
```

## Importing the Dataset

```
In [2]: #Importing the dataset, highlighting the datetime columns and viewing the first few columns
df = pd.read_csv('C:\Users\Victordano\Desktop\DA\LY LEARN\Udacity\Datasets\noshowappointments-kaggle\
2-May-2016.csv',
                parse_dates=['ScheduledDay', 'AppointmentDay'])
df.head(3)
```

```
Out [2]:
```

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alc
0	2.987250e+13	5642903	F	2016-04-29 18:38:08+00:00	2016-04-29 00:00:00+00:00	62	JARDIM DA PENHA	0	0	1	0
1	5.589978e+14	5642503	M	2016-04-29 16:08:27+00:00	2016-04-29 00:00:00+00:00	56	JARDIM DA PENHA	0	0	0	0
2	4.262962e+12	5642549	F	2016-04-29 16:19:04+00:00	2016-04-29 00:00:00+00:00	62	MATA DA PRAIA	0	0	0	0

## Exploring the dataset

```
In [3]: #checking the information on the dataset and the datatypes
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 110527 entries, 0 to 110526
Data columns (total 14 columns):
0    PatientId      non-null float64
1    AppointmentID  non-null int64
2    Gender         non-null object
3    ScheduledDay    non-null datetime64[ns, UTC]
4    AppointmentDay  non-null datetime64[ns, UTC]
5    Age            non-null int64
6    Neighbourhood   non-null object
7    Scholarship     non-null int64
8    Hypertension    non-null int64
9    Diabetes        non-null int64
10   Alcoholism      non-null int64
11   Handicap        non-null int64
12   SMS_received    non-null int64
13   No-show         non-null object
dtypes: datetime64[ns, UTC](2), float64(1), int64(8), object(3)
memory usage: 10.5+ MB
```

```
In [4]: #checking for datatypes
df.dtypes.count
```

```
Out [4]:
```

```
<bound method Series.count of PatientId      float64
AppointmentID      int64
Gender             object
ScheduledDay       datetime64[ns, UTC]
AppointmentDay     datetime64[ns, UTC]
Age               int64
Neighbourhood      object
Scholarship        int64
Hypertension       int64
Diabetes           int64
Alcoholism         int64
Handicap           int64
SMS_received       int64
No-show            object
dtype: object>
```

The shape of the dataset is 110527 rows and 14 columns.

```
8 integer columns,
3 object(string) columns,
2 datetime columns,
1 float column,
```

```
In [5]: #checking for shape of df
df.shape
```

```
Out [5]: (110527, 14)
```

```
In [6]: #check for null data
df.isnull().sum()
```

```
Out [6]: PatientId      0
AppointmentID  0
Gender         0
ScheduledDay    0
AppointmentDay  0
Age            0
Neighbourhood   0
Scholarship     0
Hypertension    0
Diabetes        0
Alcoholism      0
Handicap        0
SMS_received    0
No-show         0
dtype: int64
```

```
In [7]: #check for duplicates
df.duplicated().sum()
```

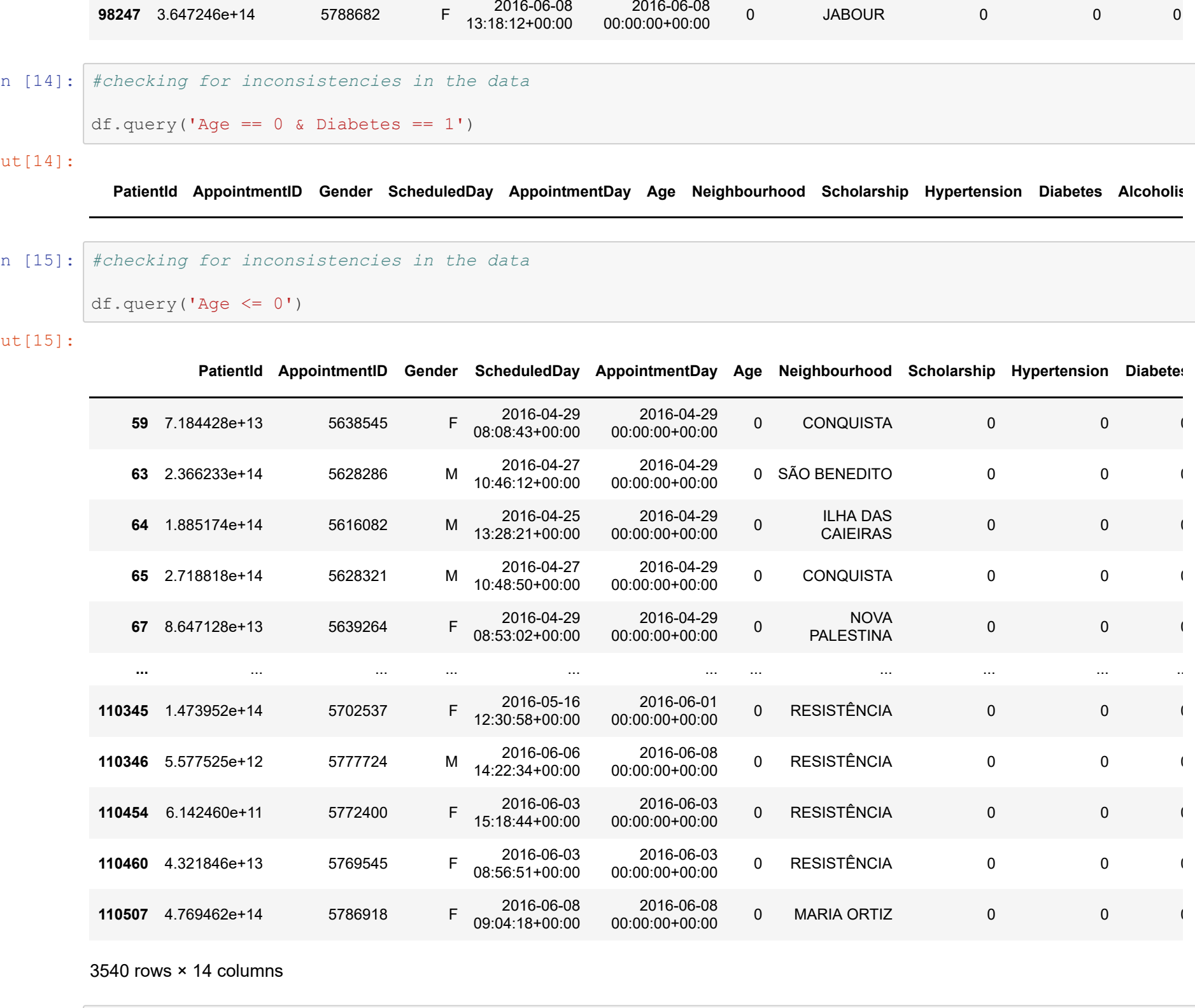
```
Out [7]: 0
```

## Checking for Null Data

As observed above, the dataset has no null data and every column has been filled out.

I also checked for duplicated data and found none in the dataset. Having achieved this, I decided to take a closer look at the data and see if it conforms to the standard data quality dimensions

```
In [8]: #checking the distribution of data in the columns
df.hist(figsize=(10,8))
```



## Some Observations

From the figures above, everything seems to be well distributed but the age and handicap columns have some surprising entries.

The Age column: The entry for this column seems to begin before the 0 mark and this would need to be investigated further

The handicap column: Has values other than 0 or 1 which would have meant that the patient was either handicapped or not. But seeing as entries 0, 1, 2, 3 or 4 were entered, it means the column once categorized their level of physical capabilities and now it has been converted for easy analysis

```
In [9]: #checking for unique features per column
df.nunique()
```

```
Out [9]: PatientId      62299
AppointmentID  110527
Gender         2
ScheduledDay   103549
AppointmentDay  27
Age           104
Neighbourhood  81
Scholarship    2
Hypertension   2
Diabetes       2
Alcoholism     2
Handicap       5
SMS_received   2
No-show        2
dtype: int64
```

```
In [10]: #why does handicap have 5 unique features
df.Handicap.value_counts()
```

```
Out [10]: 0    108286
1      2042
2       183
3         3
4         3
Name: Handicap, dtype: int64
```

```
In [11]: #renaming columns
df.rename(columns = {'Handicap':'Handicap', 'Hypertension':'Hypertension'}, inplace = True)
```

```
In [12]: #checking for inconsistencies in the data
df.query('Age == 0 & Hypertension == 1')
```

```
Out [12]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alcoholi
-----------	---------------	--------	--------------	----------------	-----	---------------	-------------	--------------	----------	----------

```
In [13]: #checking for inconsistencies in the data
df.query('Age == 0 & Handicap == 1')
```

```
Out [13]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes
98247	3.647246e+14	F	2016-06-08 13:18:12+00:00	2016-06-08 00:00:00+00:00	0	JABOUR	0	0	0

```
In [14]: #checking for inconsistencies in the data
df.query('Age == 0 & Diabetes == 1')
```

```
Out [14]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alcoholi
-----------	---------------	--------	--------------	----------------	-----	---------------	-------------	--------------	----------	----------

```
In [15]: #checking for inconsistencies in the data
df.query('Age <= 0')
```

```
Out [15]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes
69	7.184428e+13	F	2016-05-03 09:15:36+00:00	2016-04-29 00:00:00+00:00	0	CONQUISTA	0	0	0
63	2.366233e+14	M	2016-04-27 10:46:12+00:00	2016-04-29 00:00:00+00:00	0	SAO BENEDITO	0	0	0
64	1.185174e+14	M	2016-04-25 13:28:21+00:00	2016-04-29 00:00:00+00:00	0	ILHA DAS CAIEIRAS	0	0	0
65	2.718818e+14	M	2016-04-27 08:53:02+00:00	2016-04-29 00:00:00+00:00	0	CONQUISTA	0	0	0
67	8.647126e+13	F	2016-05-16 12:30:58+00:00	2016-06-01 00:00:00+00:00	0	NOVA PALESTINA	0	0	0
...	...	...	...	...	...	...	...	...	...
110345	1.473952e+14	F	2016-05-16 12:30:58+00:00	2016-06-01 00:00:00+00:00	0	RESISTÊNCIA	0	0	0
110346	5.577525e+12	M	2016-06-06 14:22:34+00:00	2016-06-08 00:00:00+00:00	0	RESISTÊNCIA	0	0	0
110454	6.142406e+11	F	2016-06-03 15:18:44+00:00	2016-06-03 00:00:00+00:00	0	RESISTÊNCIA	0	0	0
110460	4.321946e+13	F	2016-06-03 08:56:31+00:00	2016-06-03 00:00:00+00:00	0	RESISTÊNCIA	0	0	0
110607	4.769482e+14	F	2016-06-08 09:04:18+00:00	2016-06-08 00:00:00+00:00	0	MARIA ORTIZ	0	0	0

3540 rows × 14 columns

```
In [16]: #checking for inconsistencies in the data
df.query('Age < 0')
```

```
Out [16]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes
99832	4.659432e+14	F	2016-06-06 08:58:13+00:00	2016-06-06 00:00:00+00:00	-1	ROMÃO	0	0	0

## Inconsistencies in the dataset

Upon closer inspection, I noticed some patients were aged 0 and documented, but after abit exploration and confirming that they had no ailments such as Alcoholism, Diabetes or Hypertension. Those patients were included as babies who were all under the age of 1 year old.

I also noticed a female patient with age of -1 present in the dataset and I went ahead to drop this column.

```
In [17]: #drop the column with -1 age
df.tail()
```

```
Out [17]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes
110522	2.572134e+12	F	2016-05-03 09:15:36+00:00	2016-06-07 00:00:00+00:00	56	MARIA ORTIZ	0	0	0
110523	3.596206e+12	F	2016-05-03 07:27:33+00:00	2016-06-07 00:00:00+00:00	51	MARIA ORTIZ	0	0	0
110524	1.557963e+13	F	2016-04-27 16:03:52+00:00	2016-06-07 00:00:00+00:00	21	MARIA ORTIZ	0	0	0
110525	9.213493e+13	F	2016-04-27 15:09:23+00:00	2016-06-07 00:00:00+00:00	38	MARIA ORTIZ	0	0	0
110526	3.757151e+14	F	2016-04-27 13:30:56+00:00	2016-06-07 00:00:00+00:00	54	MARIA ORTIZ	0	0	0

```
In [19]: #checking the measures of spread and central tendency of the dataset
df.describe()
```

```
Out [19]:
```

	PatientId	AppointmentID	Age	Scholarship	Hypertension	Diabetes	Alcoholism	Handicap	SMS_rev
count	1.105270e+05	1.105270e+05	110527.000000	110527.000000	110527.000000	110527.000000	110527.000000	110527.000000	110527.0
mean	1.474963e+14	5.675305e+06	37.088874	0.098266	0.197246	0.071865	0.030400	0.022248	0.3
std	2.560949e+14	7.129575e+04	23.110205	0.297675	0.397921	0.258265	0.171686	0.161543	0.4
min	1.312178e+04	5.030230e+06	-1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
25%	4.172614e+12	5.640286e+06	18.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
50%	3.173184e+13	5.680573e+06	37.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0
50%	9.439172e+13	5.725554e+06	55.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.0
max	9.999816e+14	5.780484e+06	115.000000	1.000000	1.000000	1.000000	1.000000	4.000000	1.0

```
In [20]: df.head(3)
```

```
Out [20]:
```

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alc
0	2.987250e+13	5642903	F	2016-04-29 18:38:08+00:00	2016-04-29 00:00:00+00:00	62	JARDIM DA PENHA	0	0	1	0
1	5.589978e+14	5642503	M	2016-04-29 16:08:27+00:00	2016-04-29 00:00:00+00:00	56	JARDIM DA PENHA	0	0	0	0
2	4.262962e+12	5642549	F	2016-04-29 16:19:04+00:00	2016-04-29 00:00:00+00:00	62	MATA DA PRAIA	0	0	0	0

## Creating an age group column to categorize the ages of all patients

```
cut_age = [-2, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 120]
df['AgeGroup'] = pd.cut(df['Age'], cut_age)
```

```
In [22]: df.head(3)
```

```
Out [22]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alc
0	2.987250e+13	F	2016-04-29 18:38:08+00:00	2016-04-29 00:00:00+00:00	62	JARDIM DA PENHA	0	0	1	0
1	5.589978e+14	M	2016-04-29 16:08:27+00:00	2016-04-29 00:00:00+00:00	56	JARDIM DA PENHA	0	0	0	0
2	4.262962e+12	F	2016-04-29 16:19:04+00:00	2016-04-29 00:00:00+00:00	62	MATA DA PRAIA	0	0	0	0

## Categorizing the Age Column

A new column 'AgeGroup' was created to convert the data in the 'Age' column from continuous data into a category so it would be easier to get an insight into the ages of the patients easily. The ages were spread 10 years apart.

It should also be noted that for this analysis, the legal age for adulthood is 21 years old.

```
In [ ]:
```

```
In [23]: #checking for null data in the new column to ensure all columns are filled
df[df.AgeGroup.isnull()]
```

```
Out [23]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alcoholi
-----------	---------------	--------	--------------	----------------	-----	---------------	-------------	--------------	----------	----------

```
In [24]: #gender distribution in the dataset
df['Gender'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [24]:
```

Gender Distribution of Patients'

```
In [25]: #What Percentage of patients showed up for their appointment?
df['No-show'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [25]:
```

Percentage of Appointments

No: Represents the 79.8% that showed up for their appointment Yes: Represents the 20.2% that missed the appointment

```
In [26]: #What percentage of patients have a Health Scholarship
df['Scholarship'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [26]:
```

Percentage of Patients with a Health Scholarship

```
In [27]: #What Percentage of patients have alcoholism?
df['Alcoholism'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [27]:
```

Alcoholism

```
In [28]: #What Percentage of patients have diabetes?
df['Diabetes'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [28]:
```

Percentage of Patients dealing with Diabetes

```
In [29]: #What percentage of patients received an SMS reminder about their appointment?
df['SMS_received'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [29]:
```

Percentage of Patients that received an SMS

```
In [30]: #how many of the patients are hypertensive?
df['Hypertension'].value_counts().plot.pie(autopct = '%.1f%%')
```

```
Out [30]:
```

Percentage of Patients dealing with Hypertension

From the above information:

- 19.7% are dealing with Hypertension
- 3% of patients are dealing with alcoholism
- 7.2% of patients are dealing with Diabetes
- 9.8% of patients have a health scholarship
- 65% of patients are females
- 35% of patients are males
- 79.8% of patients showed up for their scheduled appointments
- 20.2% of patients didnt show up for their appointments

```
In [31]: #Age Group Distribution of Patients in the dataset
df.AgeGroup.value_counts().sort_index().plot.bar()
```

```
Out [31]:
```

## Does the dataset have Patients who have the four ailments?

```
In [32]: #df1 is a dataset containing patients with 4 ailments
df1 = df.loc[(df['Alcoholism'] == 1) & (df['Handicap'] >= 1)
            & (df['Diabetes'] == 1) & (df['Hypertension'] == 1)]
```

```
In [33]: df1.shape
```

```
Out [33]: (13, 15)
```

```
In [34]: #which of these patients are on a health scholarship?
df1.loc[(df1['Scholarship']==1)]
```

```
Out [34]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes	Alcoholi
-----------	---------------	--------	--------------	----------------	-----	---------------	-------------	--------------	----------	----------

```
In [35]: #which of these patients received an SMS reminding them about their appointment?
df1.loc[(df1['SMS_received']== 1)]
```

```
Out [35]:
```

PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hypertension	Diabetes
103545	8.588652e+12	F	2016-05-19 07:47:46+00:00	2016-06-02 00:00:00+00:00	54	JESUS DE NAZARETH	0	0	1
106517	8.588652e+12	F	2016-05-12 10:04:19+00:00	2016-06-01 00:00:00+00:00	54	JESUS DE NAZARETH	0	0	1
106544	8.588652e+12	F	2016-06-01 10:09:39+00:00	2016-06-07 00:00:00+00:00	54	JESUS DE NAZARETH	0	0	1

```
In [36]: #Horizontal Bar chart showing Patients with All 4 ailments
plt.figure(figsize=(10,10))
df1['Gender'].value_counts().plot.barh()
```

```
Out [36]:
```

Gender of patients who have all three ailments

```
In [37]: #Age group of patients with 4 ailments
df1['AgeGroup'].value_counts().plot.barh()
```

```
Out [37]:
```

Age Groups of patients who have all three ailments

```
In [38]: #patients who missed and attended the appointment
df1['No-show'].value_counts()
```

```
Out [38]:
```

No: 9  
Yes: 4  
Name: No-show, dtype: int64

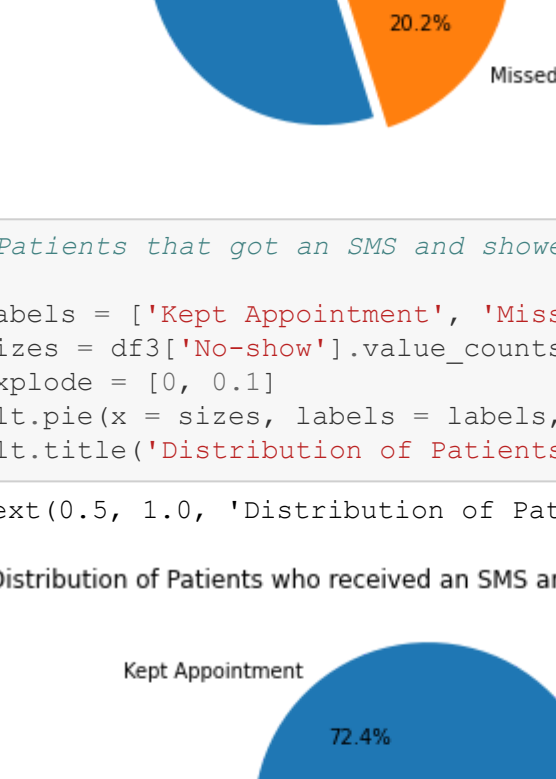
In the entirety of the dataset, Fourteen out of all the patients have alcoholism, diabetes and are also handicap and none of them are have a health scholarship. These fourteen patients are all from the ages



```
In [51]: #Patients in the dataset that showed up for the Appointment
labels = ['Kept Appointment', 'Missed']
sizes = df[['No-show']].value_counts()
explode = [0, 0.1]
plt.pie(x = sizes, labels = labels, autopct = '%.1f%%', explode = explode)
plt.title('Distribution of Patients kept their appointment')
```

Out[51]: Text(0.5, 1.0, 'Distribution of Patients kept their appointment')

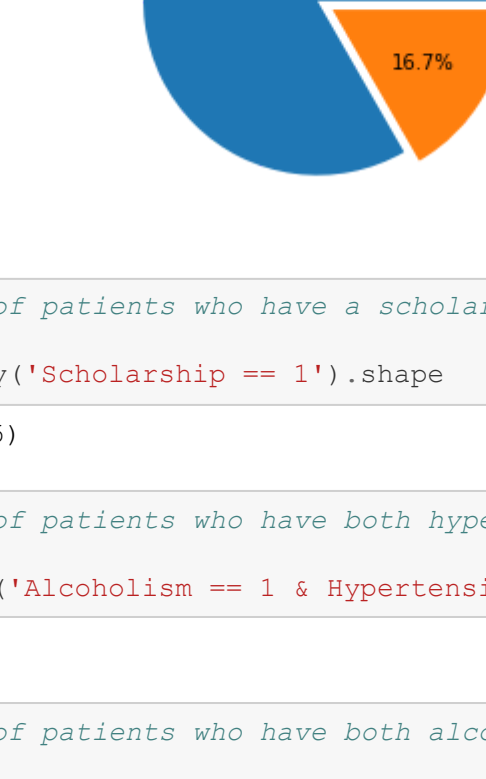
Distribution of Patients kept their appointment



```
In [52]: #Patients that got an SMS and showed up for the appointment
labels = ['Kept Appointment', 'Missed']
sizes = df.loc[(df['SMS_received'] == 0)][['No-show']].value_counts()
explode = [0, 0.1]
plt.pie(x = sizes, labels = labels, autopct = '%.1f%%', explode = explode)
plt.title('Distribution of Patients who received an SMS and kept their appointment')
```

Out[52]: Text(0.5, 1.0, 'Distribution of Patients who received an SMS and kept their appointment')

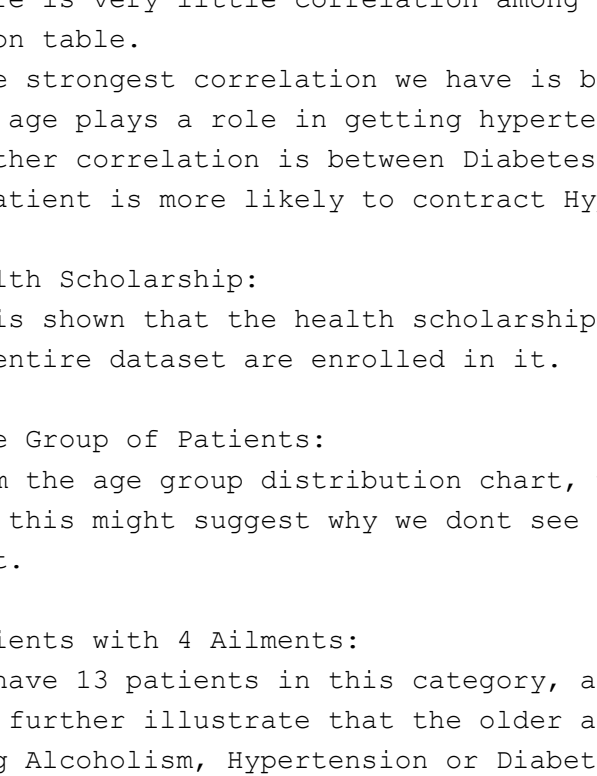
Distribution of Patients who received an SMS and kept their appointment



```
In [53]: #distribution of patients that didnt get an SMS
labels = ['Kept Appointment', 'Missed']
sizes = df.loc[(df['SMS_received'] == 0)][['No-show']].value_counts()
explode = [0.1, 0]
plt.pie(x = sizes, labels = labels, autopct = '%.1f%%', explode = explode)
plt.title('Distribution of Patients who didnt receive an SMS but kept their appointment')
```

Out[53]: Text(0.5, 1.0, 'Distribution of Patients who didnt receive an SMS but kept their appointment')

Distribution of Patients who didnt receive an SMS but kept their appointment



```
In [54]: #number of patients who have a scholarship
df3.query('Scholarship == 1').shape
```

Out[54]: (3505, 15)

```
In [55]: #number of patients who have both hypertension and alcoholism and handicap
df.query('Alcoholism == 1 & Hypertension == 1 & Handicap == 1').shape
```

Out[55]: (38, 15)

```
In [56]: #number of patients who have both alcoholism and hypertension and diabetes
df.query('Alcoholism == 1 & Hypertension == 1 & Diabetes == 1').shape
```

Out[56]: (256, 15)

```
In [57]: df[['Age', 'Alcoholism', 'Hypertension', 'Diabetes', 'Handicap', 'Scholarship']].corr()
```

Out[57]:

	Age	Alcoholism	Hypertension	Diabetes	Handicap	Scholarship
Age	1.000000	0.095811	0.504586	0.292391	0.078033	-0.092457
Alcoholism	0.095811	1.000000	0.087971	0.016474	0.004648	0.035022
Hypertension	0.504586	0.087971	1.000000	0.433086	0.080083	-0.019729
Diabetes	0.292391	0.016474	0.433086	1.000000	0.057530	-0.024894
Handicap	0.078033	0.004648	0.080083	0.057530	1.000000	-0.008586
Scholarship	-0.092457	0.035022	-0.019729	-0.024894	-0.008586	1.000000

### Conclusion

After a thorough investigation of the dataset, I posed some questions and used my analysis of the dataset to arrive at some conclusions.

**Correlation:**  
There is very little correlation among the data provided in the dataset as shown from the correlation table.  
The strongest correlation we have is between age and hypertension which is 0.5 which suggests that age plays a role in getting hypertensive.  
Another correlation is between Diabetes and Hypertension of 0.4 which may suggest that a diabetic patient is more likely to contract Hypertension

**Health Scholarship:**  
It is shown that the health scholarship is not popular among the patients, as less than 10% of the entire dataset are enrolled in it.

**Age Group of Patients:**  
From the age group distribution chart, we see that the data is skewed towards the younger patients, this might suggest why we don't see higher cases of hypertension and diabetes amongst the dataset.

**Patients with 4 Ailments:**  
We have 13 patients in this category, and all patients are from the ages of 40 to 70 which goes to further illustrate that the older a patient gets, the more likely his/her chances of contracting Alcoholism, Hypertension or Diabetes.

**Underage Intake of Alcohol:**  
We can see that some patients consume alcohol illegally and that the male underage patients are more prone to it than the female underage patients.

**SMS Reminder for Appointment**  
The information gathered from the dataset shows that the SMS reminders do not play any significant role in getting a patient to keep their appointment as more turn outs were gotten from patients who didn't get any reminders

```
In [ ]:
```