

# Project: Investigate a Dataset - [No-show appointments]

## Table of Contents

- [Introduction](#)
- [Data Wrangling](#)
- [Exploratory Data Analysis](#)
- [Conclusions](#)

## Introduction

### Dataset Description

This dataset collects information from 100k medical appointments in Brazil and is focused on the question of whether or not patients show up for their appointment. A number of characteristics about the patient are included in each row.

- **ScheduledDay**: tells us on what day the patient set up their appointment.
- **Neighborhood**: indicates the location of the hospital.
- **Scholarship**: indicates whether or not the patient is enrolled in Brazilian welfare program Bolsa Família.
- **No-show**: it says 'No' if the patient showed up to their appointment, and 'Yes' if they did not show up.
- **PatientId**: The Patient Identification number in the hospital
- **AppointmentID**: The appointment identification number in the hospital.
- **Gender**: Tell us about the patient sex.
- **AppointmentDay**: tells us on what day the patient show up their appointment.
- **Age**: tells us about the age of patient.
- **Hypertension**: If the patient has this disease.
- **Diabetes**: If the patient has this disease.
- **Alcoholism**: If the patient has this disease.
- **Handcap\***: If the patient has this disease.
- **SMS\_received**: If the patient received notification for the appointment.

### Question(s) for Analysis

1. Is the number of gender equals through the dataset ?
2. Is the case that a patient get scholarship will help to show up ?
3. What days of week patient show up easily for their appointment ?
4. Is patient going to show up whether they are sick or not ?

```
In [1]: # Use this cell to set up import statements for all of the packages that you
#        plan to use.
import pandas as pd
```

```
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
# Remember to include a 'magic word' so that your visualizations are plotted
# inline with the notebook. See this page for more:
# http://ipython.readthedocs.io/en/stable/interactive/magics.html
```

```
In [ ]: # Upgrade pandas to use dataframe.explode() function.
!pip install --upgrade pandas==0.25.0
```

```
Collecting pandas==0.25.0
  Using cached pandas-0.25.0.tar.gz (12.6 MB)
  Preparing metadata (setup.py) ... done
Requirement already satisfied: python-dateutil>=2.6.1 in /home/ayifa/anaconda3/lib/python3.9/site-packages (from pandas==0.25.0) (2.8.2)
Requirement already satisfied: pytz>=2017.2 in /home/ayifa/anaconda3/lib/python3.9/site-packages (from pandas==0.25.0) (2022.1)
Requirement already satisfied: numpy>=1.13.3 in /home/ayifa/anaconda3/lib/python3.9/site-packages (from pandas==0.25.0) (1.21.5)
Requirement already satisfied: six>=1.5 in /home/ayifa/anaconda3/lib/python3.9/site-packages (from python-dateutil>=2.6.1->pandas==0.25.0) (1.16.0)
Building wheels for collected packages: pandas
  Building wheel for pandas (setup.py) ... |
```

## Data Wrangling

**Tip:** In this section of the report, you will load in the data, check for cleanliness, and then trim and clean your dataset for analysis. Make sure that you **document your data cleaning steps in mark-down cells precisely and justify your cleaning decisions.**

## General Properties

1. Here we are reading our dataset to get known about it contents.
2. After that reload dataset with columns renamed in a way to get ease in our process
3. At this we can try show data sample

```
In [23]: # Load your data and print out a few lines. Perform operations to inspect data
renamed_columns = ['patient_id', 'appointment_id', 'gender', 'scheduled_day',
                    'appointment_day', 'age', 'neighbourhood', 'scholarship', 'hypertension',
                    'diabetes', 'alcoholism', 'handicap', 'sms_received', 'no_show']
df = pd.read_csv('noshowappointments-kaggle2-may-2016.csv', header = 0, names = renamed_columns)
df.sample(10)
```

```
Out[23]:
```

	patient_id	appointment_id	gender	scheduled_day	appointment_day	age	neighbourhood	scholarship
61612	8.977431e+14	5660392	M	2016-05-04T16:23:20Z	2016-05-04T00:00:00Z	50	ITARARÉ	
102077	9.478878e+14	5780481	F	2016-06-07T08:06:29Z	2016-06-07T00:00:00Z	18	ILHA DAS CAIEIRAS	
15144	3.761665e+13	5647238	M	2016-05-02T13:14:18Z	2016-05-06T00:00:00Z	50	ANDORINHAS	
28763	2.728618e+13	5744174	F	2016-05-30T09:54:38Z	2016-05-30T00:00:00Z	1	CRUZAMENTO	
101427	7.259780e+11	5738699	F	2016-05-25T11:23:30Z	2016-06-02T00:00:00Z	19	RESISTÊNCIA	
27879	6.342815e+12	5657907	F	2016-05-	2016-05-	0	MARUÍPE	

				04T10:07:45Z	04T00:00:00Z		
26551	2.226924e+13	5723683	M	2016-05-20T08:30:16Z	2016-05-20T00:00:00Z	28	SÃO PEDRO
21948	8.524478e+13	5624711	F	2016-04-26T16:47:44Z	2016-05-05T00:00:00Z	69	DO CABRAL
90207	1.888439e+14	5753159	F	2016-05-31T13:33:58Z	2016-06-01T00:00:00Z	31	BONFIM
14032	5.366354e+14	5660782	F	2016-05-04T18:58:56Z	2016-05-16T00:00:00Z	40	SÃO CRISTÓVÃO

## Data info and some quick stats about all information in the dataset

Using the following cell we can describe all our dataset With the information about the dataset we can get all columns data types, numbers of rows and columns and also if there some missing data, duplicated or incorrect data.

```
In [4]: # types and look for instances of missing or possibly errant data.
df.info()
df.describe()
df.duplicated().sum()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 110527 entries, 0 to 110526
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   patient_id            110527 non-null  float64
1   appointment_id        110527 non-null  int64
2   gender                110527 non-null  object
3   scheduled_day         110527 non-null  object
4   appointment_day       110527 non-null  object
5   age                  110527 non-null  int64
6   neighbourhood         110527 non-null  object
7   scholarship          110527 non-null  int64
8   hypertension         110527 non-null  int64
9   diabetes              110527 non-null  int64
10  alcoholism            110527 non-null  int64
11  handicap              110527 non-null  int64
12  sms_received          110527 non-null  int64
13  no_show               110527 non-null  object
dtypes: float64(1), int64(8), object(5)
memory usage: 11.8+ MB
```

Out[4]: 0

```
In [5]: {i : df[i].unique() for i in df.columns}
```

```
Out[5]: {'patient_id': array([2.98724998e+13, 5.58997777e+14, 4.26296230e+12, ...,
7.26331493e+13, 9.96997666e+14, 1.55766317e+13]),
'appointment_id': array([5642903, 5642503, 5642549, ..., 5630692, 5630323, 5629448]),
'gender': array(['F', 'M'], dtype=object),
'scheduled_day': array(['2016-04-29T18:38:08Z', '2016-04-29T16:08:27Z',
'2016-04-29T16:19:04Z', ..., '2016-04-27T16:03:52Z',
'2016-04-27T15:09:23Z', '2016-04-27T13:30:56Z'], dtype=object),
'appointment_day': array(['2016-04-29T00:00:00Z', '2016-05-03T00:00:00Z',
'2016-05-10T00:00:00Z', '2016-05-17T00:00:00Z',
'2016-05-24T00:00:00Z', '2016-05-31T00:00:00Z',
'2016-05-02T00:00:00Z', '2016-05-30T00:00:00Z',
'2016-05-16T00:00:00Z', '2016-05-04T00:00:00Z',
'2016-05-19T00:00:00Z', '2016-05-12T00:00:00Z',
'2016-05-06T00:00:00Z', '2016-05-20T00:00:00Z',
```

```

'2016-05-05T00:00:00Z', '2016-05-13T00:00:00Z',
'2016-05-09T00:00:00Z', '2016-05-25T00:00:00Z',
'2016-05-11T00:00:00Z', '2016-05-18T00:00:00Z',
'2016-05-14T00:00:00Z', '2016-06-02T00:00:00Z',
'2016-06-03T00:00:00Z', '2016-06-06T00:00:00Z',
'2016-06-07T00:00:00Z', '2016-06-01T00:00:00Z',
'2016-06-08T00:00:00Z'], dtype=object),
'age': array([ 62,  56,   8,  76,  23,  39,  21,  19,  30,  29,  22,  28,  54,
   15,  50,  40,  46,   4,  13,  65,  45,  51,  32,  12,  61,  38,
   79,  18,  63,  64,  85,  59,  55,  71,  49,  78,  31,  58,  27,
    6,   2,  11,   7,   0,   3,   1,  69,  68,  60,  67,  36,  10,
   35,  20,  26,  34,  33,  16,  42,   5,  47,  17,  41,  44,  37,
   24,  66,  77,  81,  70,  53,  75,  73,  52,  74,  43,  89,  57,
   14,   9,  48,  83,  72,  25,  80,  87,  88,  84,  82,  90,  94,
   86,  91,  98,  92,  96,  93,  95,  97, 102, 115, 100,  99, -1]),
'neighbourhood': array(['JARDIM DA PENHA', 'MATA DA PRAIA', 'PONTAL DE CAMBURI',
  'REPÚBLICA', 'GOIABEIRAS', 'ANDORINHAS', 'CONQUISTA',
  'NOVA PALESTINA', 'DA PENHA', 'TABUAZEIRO', 'BENTO FERREIRA',
  'SÃO PEDRO', 'SANTA MARTHA', 'SÃO CRISTÓVÃO', 'MARUÍPE',
  'GRANDE VITÓRIA', 'SÃO BENEDITO', 'ILHA DAS CAIEIRAS',
  'SANTO ANDRÉ', 'SOLON BORGES', 'BONFIM', 'JARDIM CAMBURI',
  'MARIA ORTIZ', 'JABOUR', 'ANTÔNIO HONÓRIO', 'RESISTÊNCIA',
  'ILHA DE SANTA MARIA', 'JUCUTUQUARA', 'MONTE BELO',
  'MÁRIO CYPRESTE', 'SANTO ANTÔNIO', 'BELA VISTA', 'PRAIA DO SUÁ',
  'SANTA HELENA', 'ITARARÉ', 'INHANGUETÁ', 'UNIVERSITÁRIO',
  'SÃO JOSÉ', 'REDEÇÃO', 'SANTA CLARA', 'CENTRO', 'PARQUE MOSCOSO',
  'DO MOSCOSO', 'SANTOS DUMONT', 'CARATOÍRA', 'ARIOVALDO FAVALESSA',
  'ILHA DO FRADE', 'GURIGICA', 'JOANA D´ARC', 'CONSOLAÇÃO',
  'PRAIA DO CANTO', 'BOA VISTA', 'MORADA DE CAMBURI', 'SANTA LUÍZA',
  'SANTA LÚCIA', 'BARRO VERMELHO', 'ESTRELINHA', 'FORTE SÃO JOÃO',
  'FONTE GRANDE', 'ENSEADA DO SUÁ', 'SANTOS REIS', 'PIEIDADE',
  'JESUS DE NAZARETH', 'SANTA TEREZA', 'CRUZAMENTO',
  'ILHA DO PRÍNCIPE', 'ROMÃO', 'COMDUSA', 'SANTA CECÍLIA',
  'VILA RUBIM', 'DE LOURDES', 'DO QUADRO', 'DO CABRAL', 'HORTO',
  'SEGURANÇA DO LAR', 'ILHA DO BOI', 'FRADINHOS', 'NAZARETH',
  'AEROPORTO', 'ILHAS OCEÂNICAS DE TRINDADE', 'PARQUE INDUSTRIAL'],
  dtype=object),
'scholarship': array([0, 1]),
'hypertension': array([1, 0]),
'diabetes': array([0, 1]),
'alcoholism': array([0, 1]),
'handicap': array([0, 1, 2, 3, 4]),
'sms_received': array([0, 1]),
'no_show': array(['No', 'Yes'], dtype=object)}

```

```
In [6]: df.columns
```

```
Out[6]: Index(['patient_id', 'appointment_id', 'gender', 'scheduled_day',
  'appointment_day', 'age', 'neighbourhood', 'scholarship',
  'hypertension', 'diabetes', 'alcoholism', 'handicap', 'sms_received',
  'no_show'],
  dtype='object')
```

## Data Cleaning

With the last cell we saw that there is **110527** rows and **14** columns. It also mention that we to update some date type especially first for scheduled and appointment days

More important we are almost ready to start our exploration but we first need to transform a little bit our dataframe

```
In [7]: # After discussing the structure of the data and any problems that need to be
# cleaned, perform those cleaning steps in the second part of this section.
```

```
df.drop(['patient_id', 'appointment_id'], axis = 1, inplace = True)
df['appointment_day'] = pd.to_datetime(df['appointment_day'])
df['scheduled_day'] = pd.to_datetime(df['scheduled_day'])
df.head()
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 110527 entries, 0 to 110526
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   gender                 110527 non-null object
1   scheduled_day          110527 non-null datetime64[ns, UTC]
2   appointment_day        110527 non-null datetime64[ns, UTC]
3   age                   110527 non-null int64
4   neighbourhood          110527 non-null object
5   scholarship            110527 non-null int64
6   hypertension           110527 non-null int64
7   diabetes               110527 non-null int64
8   alcoholism             110527 non-null int64
9   handicap               110527 non-null int64
10  sms_received           110527 non-null int64
11  no_show                110527 non-null object
dtypes: datetime64[ns, UTC](2), int64(7), object(3)
memory usage: 10.1+ MB
```

```
In [8]: # Let's transform the no_show columns but not required
df['no_show'].replace({'Yes': 1, 'No': 0}, inplace = True)
```

```
In [9]: df.sample(10)
```

```
Out[9]:
```

	gender	scheduled_day	appointment_day	age	neighbourhood	scholarship	hypertension	diabetes	alc
23725	M	2016-04-18 09:03:53+00:00	2016-05-12 00:00:00+00:00	57	ARIOVALDO FAVALESSA	0	0	0	
98873	F	2016-06-01 11:22:35+00:00	2016-06-01 00:00:00+00:00	0	JOANA D'ARC	0	0	0	
22192	M	2016-05-20 15:22:28+00:00	2016-05-31 00:00:00+00:00	0	ILHA DAS CAIEIRAS	0	0	0	
91199	F	2016-06-01 13:57:29+00:00	2016-06-03 00:00:00+00:00	61	JARDIM CAMBURI	0	0	0	
41564	F	2016-05-17 12:28:33+00:00	2016-05-19 00:00:00+00:00	26	SANTO ANTÔNIO	0	0	0	
6312	F	2016-04-20 16:31:35+00:00	2016-05-12 00:00:00+00:00	56	NOVA PALESTINA	0	1	0	
32440	F	2016-05-13 07:58:56+00:00	2016-05-20 00:00:00+00:00	20	JABOUR	1	0	0	
49415	M	2016-05-17 07:39:18+00:00	2016-05-24 00:00:00+00:00	7	SANTOS DUMONT	0	0	0	
86634	M	2016-06-06 10:06:56+00:00	2016-06-06 00:00:00+00:00	78	MARIA ORTIZ	0	0	0	
92838	F	2016-06-02 09:42:59+00:00	2016-06-06 00:00:00+00:00	64	JARDIM CAMBURI	0	0	0	

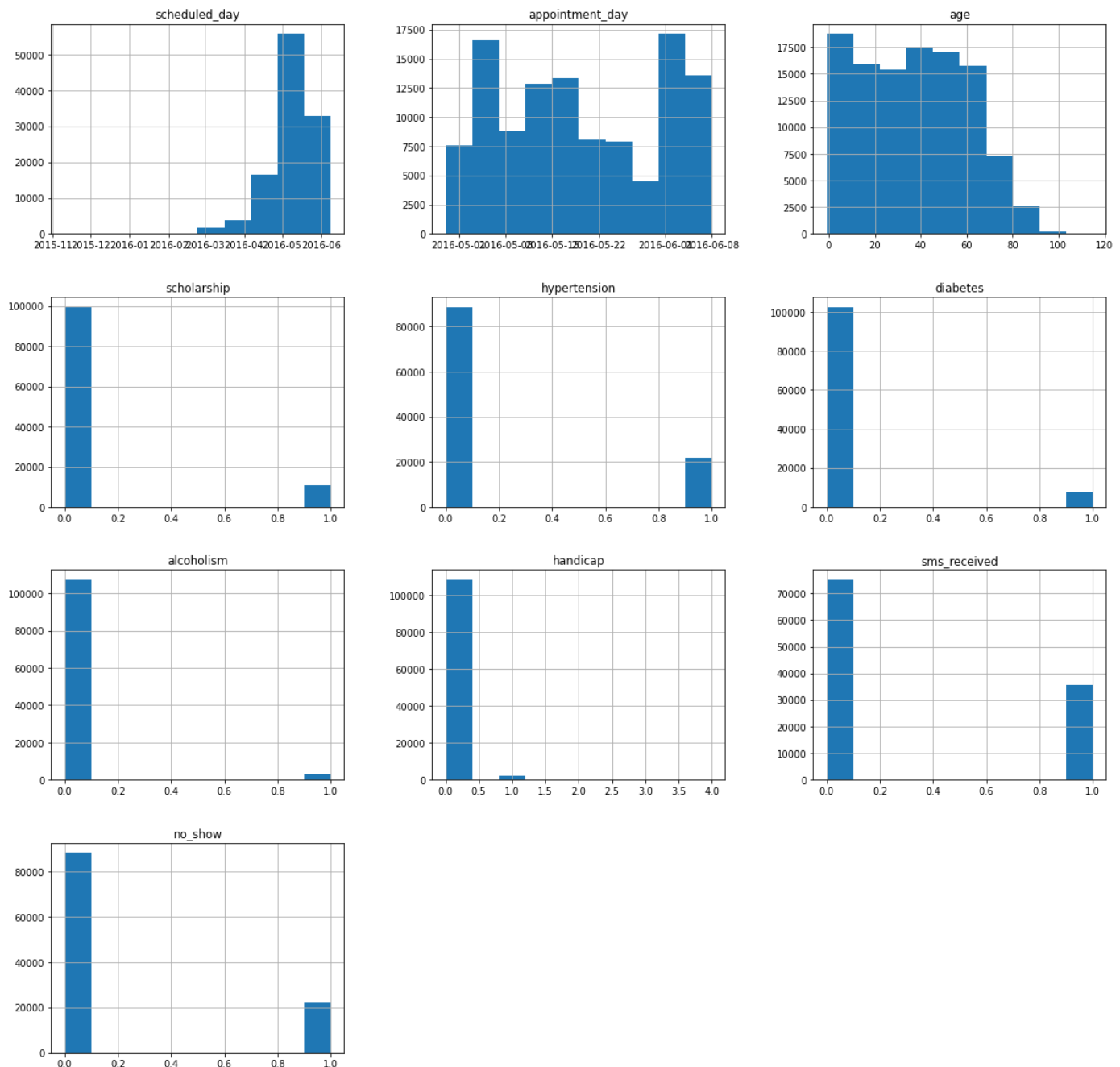
## Exploratory Data Analysis

Let's use some quick graphics to get answer for our early asked questions

# Research Question 1 (Is the number of gender equals through the dataset ?)

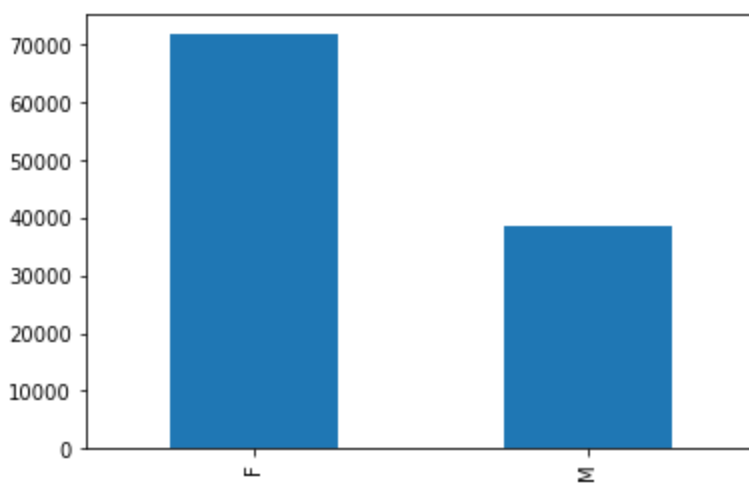
```
In [10]: # Use this, and more code cells, to explore your data. Don't forget to add
# Markdown cells to document your observations and findings.
df.hist(figsize = [20, 20])
```

```
Out[10]: array([[<AxesSubplot:title={'center':'scheduled_day'}>,
      <AxesSubplot:title={'center':'appointment_day'}>,
      <AxesSubplot:title={'center':'age'}>],
      [<AxesSubplot:title={'center':'scholarship'}>,
      <AxesSubplot:title={'center':'hypertension'}>,
      <AxesSubplot:title={'center':'diabetes'}>],
      [<AxesSubplot:title={'center':'alcoholism'}>,
      <AxesSubplot:title={'center':'handicap'}>,
      <AxesSubplot:title={'center':'sms_received'}>],
      [<AxesSubplot:title={'center':'no_show'}>, <AxesSubplot:>,
      <AxesSubplot:>]], dtype=object)
```



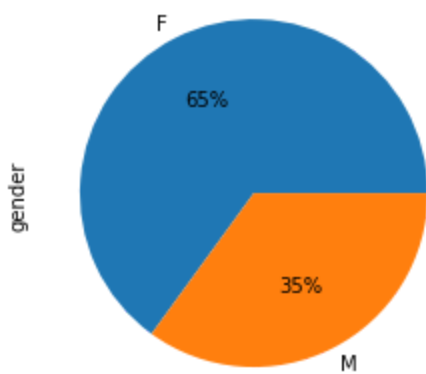
```
In [11]: df['gender'].value_counts().plot(kind = 'bar')
```

```
Out[11]: <AxesSubplot:>
```



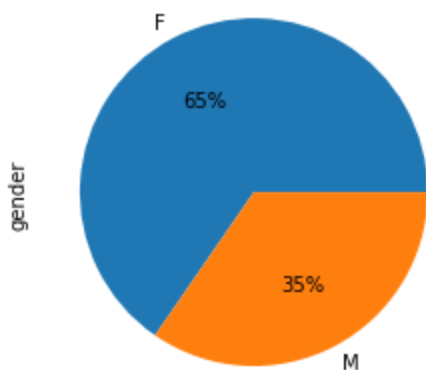
```
In [12]: df.gender.value_counts().plot(kind='pie', autopct='%.0f%%')
```

```
Out[12]: <AxesSubplot:ylabel='gender'>
```



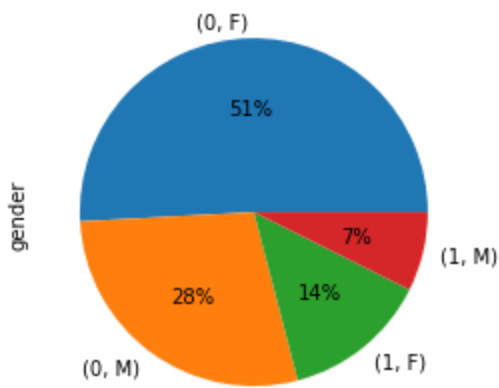
```
In [13]: df_noshow = df.query('no_show == 1')
df_noshow.gender.value_counts().plot(kind='pie', autopct='%.0f%%')
```

```
Out[13]: <AxesSubplot:ylabel='gender'>
```



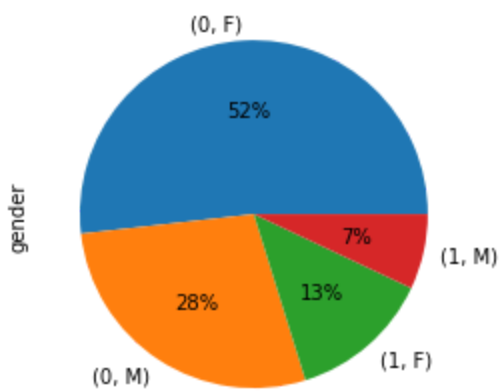
```
In [14]: df_no_show_not_seek = df.query('alcoholism == 0 and hypertension == 0 and diabetes == 0')
df_no_show_not_seek.gender.value_counts().plot(kind='pie', autopct='%.0f%%')
```

```
Out[14]: <AxesSubplot:ylabel='gender'>
```



```
In [15]: df_no_show_seek = df.query('alcoholism == 0 or hypertension == 0 or diabetes == 0 or h
df_no_show_seek.gender.value_counts().plot(kind='pie', autopct='%0f%%')
```

```
Out[15]: <AxesSubplot:ylabel='gender'>
```

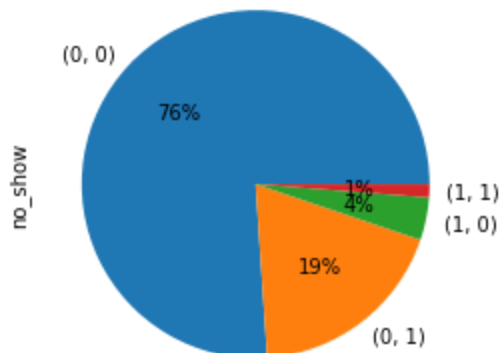


Research Question 2 and 4 (Is the case that a patient get scholarship will help to show up ? / Is patient going to show up wether they are sick or not ?)

```
In [16]: # Explore the dataset based on scholarship
df_scholar_men = df.query('gender == "M"').groupby(['scholarship'])

df_scholar_men.no_show.value_counts().plot(kind='pie', autopct='%0f%%')
```

```
Out[16]: <AxesSubplot:ylabel='no_show'>
```

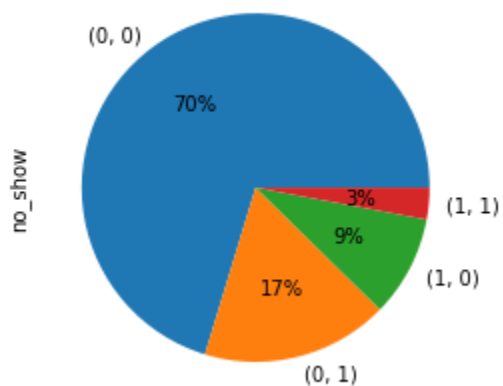


```
In [17]: df_scholar_women = df.query('gender == "F"').groupby(['scholarship'])
```



```
df_scholar_women.no_show.value_counts().plot(kind='pie', autopct='%.0f%%')
```

```
Out[17]: <AxesSubplot:ylabel='no_show'>
```



### Research Question 3 (What days of week patient sow up easily for they appointment ?)

```
In [18]: # Let's first transform day to get what days is suitable for patients.
df_day_week = df.copy()
df_day_week['appointment_day'] = df['appointment_day'].dt.dayofweek
df_day_week['scheduled_day'] = df['scheduled_day'].dt.dayofweek
df_day_week.head()
```

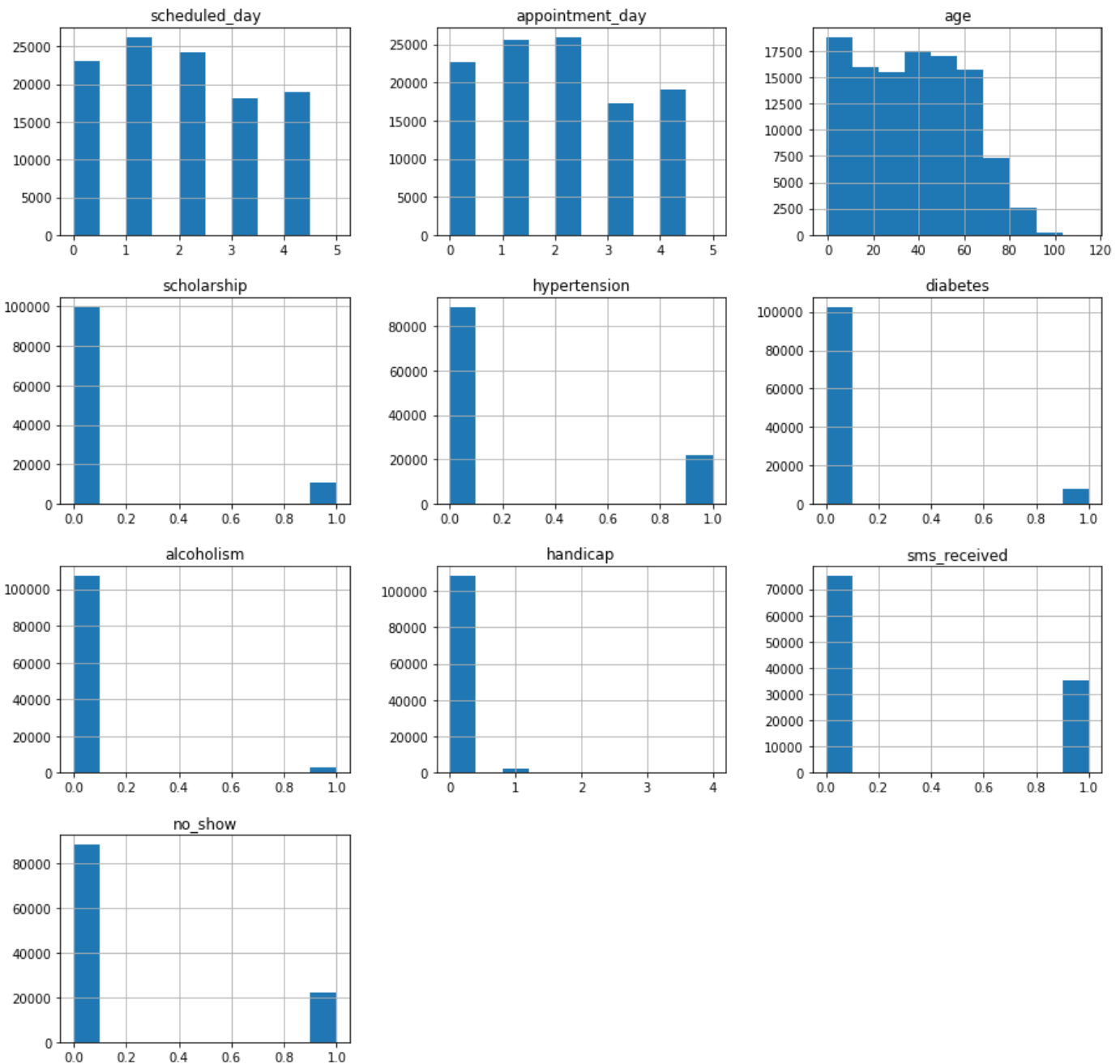
Out[18]:	gender	scheduled_day	appointment_day	age	neighbourhood	scholarship	hypertension	diabetes	alcohol
0	F	4	4	62	JARDIM DA PENHA	0	1	0	
1	M	4	4	56	JARDIM DA PENHA	0	0	0	
2	F	4	4	62	MATA DA PRAIA	0	0	0	
3	F	4	4	8	PONTAL DE CAMBURI	0	0	0	
4	F	4	4	56	JARDIM DA PENHA	0	1	1	

```
In [19]: df_day_week.sample(10)
```

Out[19]:	gender	scheduled_day	appointment_day	age	neighbourhood	scholarship	hypertension	diabetes	alc
7575	F	3	3	10	MÁRIO CYPRESTE	0	0	0	
40408	M	1	2	47	SÃO CRISTÓVÃO	0	0	0	
95395	F	2	3	11	MARIA ORTIZ	0	0	0	
39105	M	1	3	29	JOANA D'ARC	0	0	0	
62035	F	0	0	31	ILHA DO BOI	0	0	0	
27175	F	4	4	77	RESISTÊNCIA	0	1	0	
44947	F	4	4	22	SÃO JOSÉ	0	0	0	
80225	F	4	4	11	FONTE GRANDE	0	0	0	
31930	M	2	4	35	RESISTÊNCIA	0	0	0	

```
In [20]: df_day_week.hist(figsize= [15,15])
```

```
Out[20]: array([[<AxesSubplot:title={'center':'scheduled_day'}>,
      <AxesSubplot:title={'center':'appointment_day'}>,
      <AxesSubplot:title={'center':'age'}>],
      [<AxesSubplot:title={'center':'scholarship'}>,
      <AxesSubplot:title={'center':'hypertension'}>,
      <AxesSubplot:title={'center':'diabetes'}>],
      [<AxesSubplot:title={'center':'alcoholism'}>,
      <AxesSubplot:title={'center':'handicap'}>,
      <AxesSubplot:title={'center':'sms_received'}>],
      [<AxesSubplot:title={'center':'no_show'}>, <AxesSubplot:>,
      <AxesSubplot:>]], dtype=object)
```



```
In [21]: # Quick Explore with neighbourhood
df_day_week_group_age = df_day_week.query('no_show == 1').groupby(['neighbourhood'])

pd.plotting.scatter_matrix(df_day_week, figsize = (25,25))
```

```
Out[21]: array([[<AxesSubplot:xlabel='scheduled_day', ylabel='scheduled_day'>,
      <AxesSubplot:xlabel='appointment_day', ylabel='scheduled_day'>],
```

```

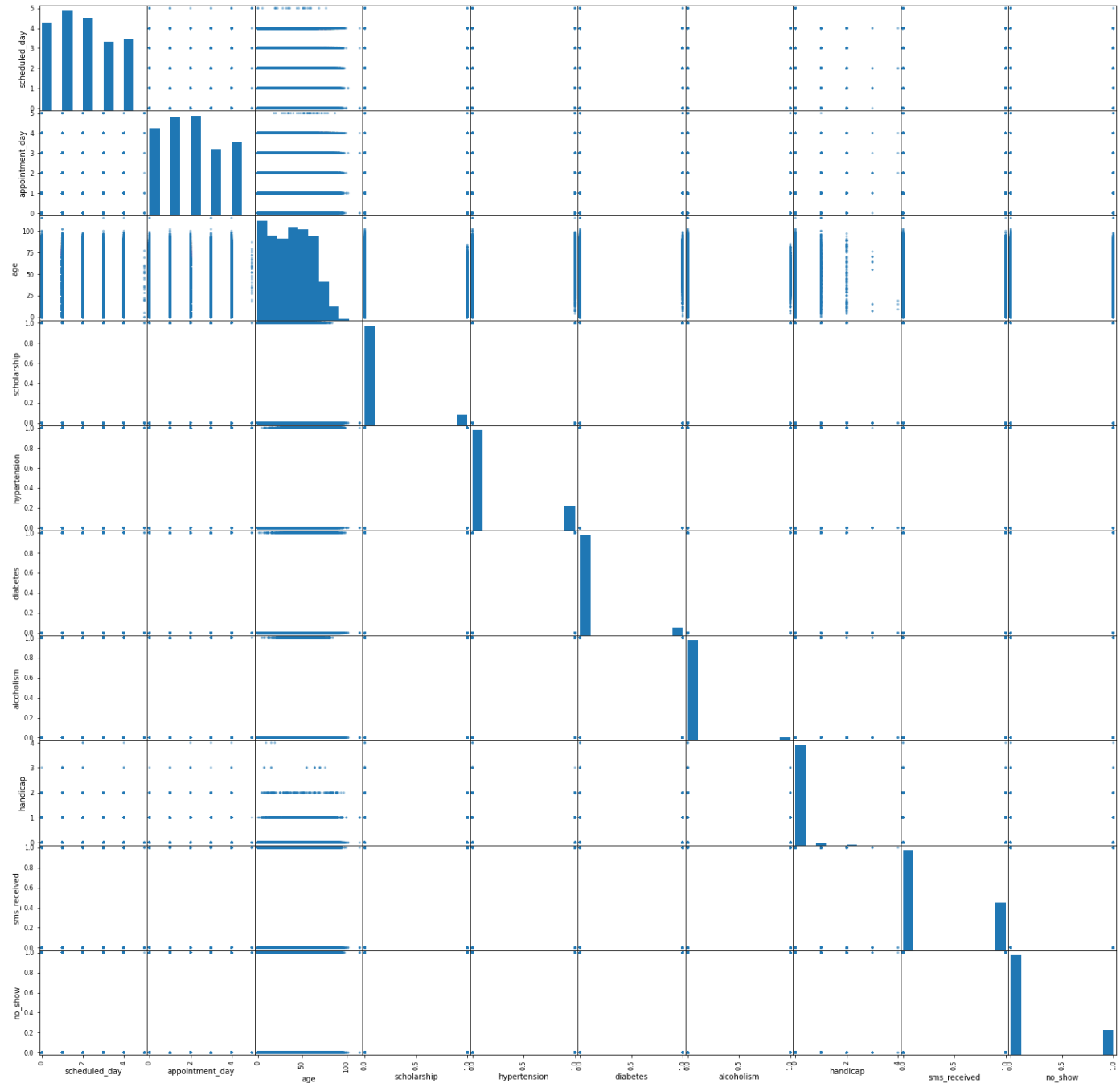
<AxesSubplot:xlabel='age', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='scholarship', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='hypertension', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='diabetes', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='alcoholism', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='handicap', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='sms_received', ylabel='scheduled_day'>,
<AxesSubplot:xlabel='no_show', ylabel='scheduled_day'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='appointment_day'>,
<AxesSubplot:xlabel='appointment_day', ylabel='appointment_day'>,
<AxesSubplot:xlabel='age', ylabel='appointment_day'>,
<AxesSubplot:xlabel='scholarship', ylabel='appointment_day'>,
<AxesSubplot:xlabel='hypertension', ylabel='appointment_day'>,
<AxesSubplot:xlabel='diabetes', ylabel='appointment_day'>,
<AxesSubplot:xlabel='alcoholism', ylabel='appointment_day'>,
<AxesSubplot:xlabel='handicap', ylabel='appointment_day'>,
<AxesSubplot:xlabel='sms_received', ylabel='appointment_day'>,
<AxesSubplot:xlabel='no_show', ylabel='appointment_day'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='age'>,
<AxesSubplot:xlabel='appointment_day', ylabel='age'>,
<AxesSubplot:xlabel='age', ylabel='age'>,
<AxesSubplot:xlabel='scholarship', ylabel='age'>,
<AxesSubplot:xlabel='hypertension', ylabel='age'>,
<AxesSubplot:xlabel='diabetes', ylabel='age'>,
<AxesSubplot:xlabel='alcoholism', ylabel='age'>,
<AxesSubplot:xlabel='handicap', ylabel='age'>,
<AxesSubplot:xlabel='sms_received', ylabel='age'>,
<AxesSubplot:xlabel='no_show', ylabel='age'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='scholarship'>,
<AxesSubplot:xlabel='appointment_day', ylabel='scholarship'>,
<AxesSubplot:xlabel='age', ylabel='scholarship'>,
<AxesSubplot:xlabel='scholarship', ylabel='scholarship'>,
<AxesSubplot:xlabel='hypertension', ylabel='scholarship'>,
<AxesSubplot:xlabel='diabetes', ylabel='scholarship'>,
<AxesSubplot:xlabel='alcoholism', ylabel='scholarship'>,
<AxesSubplot:xlabel='handicap', ylabel='scholarship'>,
<AxesSubplot:xlabel='sms_received', ylabel='scholarship'>,
<AxesSubplot:xlabel='no_show', ylabel='scholarship'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='hypertension'>,
<AxesSubplot:xlabel='appointment_day', ylabel='hypertension'>,
<AxesSubplot:xlabel='age', ylabel='hypertension'>,
<AxesSubplot:xlabel='scholarship', ylabel='hypertension'>,
<AxesSubplot:xlabel='hypertension', ylabel='hypertension'>,
<AxesSubplot:xlabel='diabetes', ylabel='hypertension'>,
<AxesSubplot:xlabel='alcoholism', ylabel='hypertension'>,
<AxesSubplot:xlabel='handicap', ylabel='hypertension'>,
<AxesSubplot:xlabel='sms_received', ylabel='hypertension'>,
<AxesSubplot:xlabel='no_show', ylabel='hypertension'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='diabetes'>,
<AxesSubplot:xlabel='appointment_day', ylabel='diabetes'>,
<AxesSubplot:xlabel='age', ylabel='diabetes'>,
<AxesSubplot:xlabel='scholarship', ylabel='diabetes'>,
<AxesSubplot:xlabel='hypertension', ylabel='diabetes'>,
<AxesSubplot:xlabel='diabetes', ylabel='diabetes'>,
<AxesSubplot:xlabel='alcoholism', ylabel='diabetes'>,
<AxesSubplot:xlabel='handicap', ylabel='diabetes'>,
<AxesSubplot:xlabel='sms_received', ylabel='diabetes'>,
<AxesSubplot:xlabel='no_show', ylabel='diabetes'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='alcoholism'>,
<AxesSubplot:xlabel='appointment_day', ylabel='alcoholism'>,
<AxesSubplot:xlabel='age', ylabel='alcoholism'>,
<AxesSubplot:xlabel='scholarship', ylabel='alcoholism'>,
<AxesSubplot:xlabel='hypertension', ylabel='alcoholism'>,
<AxesSubplot:xlabel='diabetes', ylabel='alcoholism'>,
<AxesSubplot:xlabel='alcoholism', ylabel='alcoholism'>,
<AxesSubplot:xlabel='handicap', ylabel='alcoholism'>,

```

```

<AxesSubplot:xlabel='sms_received', ylabel='alcoholism'>,
<AxesSubplot:xlabel='no_show', ylabel='alcoholism'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='handicap'>,
<AxesSubplot:xlabel='appointment_day', ylabel='handicap'>,
<AxesSubplot:xlabel='age', ylabel='handicap'>,
<AxesSubplot:xlabel='scholarship', ylabel='handicap'>,
<AxesSubplot:xlabel='hypertension', ylabel='handicap'>,
<AxesSubplot:xlabel='diabetes', ylabel='handicap'>,
<AxesSubplot:xlabel='alcoholism', ylabel='handicap'>,
<AxesSubplot:xlabel='handicap', ylabel='handicap'>,
<AxesSubplot:xlabel='sms_received', ylabel='handicap'>,
<AxesSubplot:xlabel='no_show', ylabel='handicap'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='sms_received'>,
<AxesSubplot:xlabel='appointment_day', ylabel='sms_received'>,
<AxesSubplot:xlabel='age', ylabel='sms_received'>,
<AxesSubplot:xlabel='scholarship', ylabel='sms_received'>,
<AxesSubplot:xlabel='hypertension', ylabel='sms_received'>,
<AxesSubplot:xlabel='diabetes', ylabel='sms_received'>,
<AxesSubplot:xlabel='alcoholism', ylabel='sms_received'>,
<AxesSubplot:xlabel='handicap', ylabel='sms_received'>,
<AxesSubplot:xlabel='sms_received', ylabel='sms_received'>,
<AxesSubplot:xlabel='no_show', ylabel='sms_received'>],
[<AxesSubplot:xlabel='scheduled_day', ylabel='no_show'>,
<AxesSubplot:xlabel='appointment_day', ylabel='no_show'>,
<AxesSubplot:xlabel='age', ylabel='no_show'>,
<AxesSubplot:xlabel='scholarship', ylabel='no_show'>,
<AxesSubplot:xlabel='hypertension', ylabel='no_show'>,
<AxesSubplot:xlabel='diabetes', ylabel='no_show'>,
<AxesSubplot:xlabel='alcoholism', ylabel='no_show'>,
<AxesSubplot:xlabel='handicap', ylabel='no_show'>,
<AxesSubplot:xlabel='sms_received', ylabel='no_show'>,
<AxesSubplot:xlabel='no_show', ylabel='no_show'>]], dtype=object)

```



## Conclusions

1. Is the number of gender equals through the dataset ? **Answer** : We have **65% Women** and **35% Men**. So no the dataset is not equally provided for gender feature
2. Is the case that a patient get scholarship will help to show up ? **Answer** :  

The scholarship is not deterministic for this exploration.
3. What days of week patient show up easily for they appointment ? **Answer** : Quick comparison of appointment days of week and scheduled days of week show that patients come regularly for appointment on (Monday, Wednesday and Friday). **Tuesday and Thursday patients number decrease a little**
4. Is patient going to show up whether they are sick or not ? **Answer** : No matter they are sick or not the pie charts show us that we have around the same behaviour for all patients

So patients sheduled for (Monday, Wednesday and Friday)  
will mostly show\_up

## Submitting your Project

```
In [ ]: from subprocess import call  
call(['python', '-m', 'nbconvert', 'Investigate_a_Dataset.ipynb'])
```