

Analysis of Medical Appointments (Vitória, Brazil in May 2016)

Table of Contents

- Introduction
- Data Wrangling
- Exploratory Data Analysis
- Conclusions

Introduction

In this analysis I will explore a dataset containing 100,000 medical appointments at various hospitals in the city of Vitória in the state of Espírito Santo in Brazil. Special emphasis will be put on investigating if patient age or waiting time from having scheduled an appointment are related to no-show rates for doctors appointments

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

# 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

import pandas as pd
import matplotlib.pyplot as plt

%matplotlib inline
```

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 steps carefully and justify your cleaning decisions.

General Properties

```
In [2]: df = pd.read_csv('noshowappointments-kaggle1ev2-may-2016.csv')

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

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hipertension	Diabetes	A
0	2.987250e+13	5642903	F	2016-04-29T18:38:00Z	2016-04-29T00:00:00Z	62	JARDIM DA PENHA	0	1	0	
1	5.589978e+14	5642503	M	2016-04-29T16:08:27Z	2016-04-29T00:00:00Z	56	JARDIM DA PENHA	0	0	0	
2	4.262962e+12	5642549	F	2016-04-29T16:19:04Z	2016-04-29T00:00:00Z	62	MATA DA PRAIA	0	0	0	
3	8.679512e+11	5642828	F	2016-04-29T17:29:31Z	2016-04-29T00:00:00Z	8	PONTAL DE CAMBURI	0	0	0	
4	8.841986e+12	5642494	F	2016-04-29T16:07:23Z	2016-04-29T00:00:00Z	56	JARDIM DA PENHA	0	1	1	

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 110527 entries, 0 to 110526
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   PatientId              110527 non-null  float64
1   AppointmentID          110527 non-null  object
2   Gender                 110527 non-null  object
3   ScheduledDay           110527 non-null  object
4   AppointmentDay         110527 non-null  object
5   Age                    110527 non-null  int64
6   Neighbourhood         110527 non-null  object
7   Scholarship            110527 non-null  int64
8   Hipertension           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  object
13  No-show                110527 non-null  object
dtypes: float64(1), int64(8), object(5)
memory usage: 11.8+ MB
```

The meaning of the columns in the dataset are as follows:

- PatientId: A unique identifier for a patient
- AppointmentID: A unique identifier for a single appointment
- Gender: The patient's gender
- ScheduledDay: The date and time at which the appointment was scheduled
- AppointmentDay: The actual date of the appointment
- Age: The age of the patient
- Neighbourhood: The neighbourhood (of a localised Hospital) where the appointment is taking place
- Scholarship: whether the patient is part of the Alcosita Family welfare program
- Hipertension: whether the patient suffers from hypertension
- Diabetes: whether the patient suffers from diabetes
- Alcoholism: whether the patient suffers from alcoholism
- Handicap: the degree of disability of the patient
- SMS_received: whether or not the patient received an SMS with the appointment information
- No-show: whether or not the patient showed up for the appointment

There seem to be no values missing, but the `ScheduledDay` and `AppointmentDay` columns are interpreted as strings instead of datetime.

```
In [5]: # convert datetime-string for ScheduledDay and AppointmentDay to proper datetime type
df['ScheduledDay'] = pd.to_datetime(df['ScheduledDay'])
df['AppointmentDay'] = pd.to_datetime(df['AppointmentDay'])

In [6]: # convert No-show column values: Yes -> 1, No -> 0.
# This makes it easier to use summary-statistics like the mean for the No-show column
df['No-show'] = df['No-show'].map({'Yes': 1, 'No': 0})

In [7]: # check if we have duplicated AppointmentIDs and (PatientId, AppointmentID) pairs. These would be problematic
df[['PatientId', 'AppointmentID']].duplicated().value_counts()
```

	PatientId	AppointmentID	Age	Scholarship	Hipertension	Diabetes	Alcoholism	Handicap	SMS
count	110527.0e+05	110527.0e+06	110527.000000	110527.000000	110527.000000	110527.000000	110527.000000	110527.000000	110527.000000
mean	4.172614e+14	5.675305e+06	37.088874	0.098266	0.197246	0.071865	0.030400	0.022248	0.110527
std	2.560949e+14	7.129575e+04	23.110205	0.297675	0.397921	0.258265	0.171686	0.161543	0.110527
min	3.921784e+04	5.030230e+06	-1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	4.172614e+12	5.640286e+06	18.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	3.173184e+13	5.680573e+06	37.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	4.373917e+13	5.725524e+06	55.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	9.999810e+14	5.790484e+06	115.000000	1.000000	1.000000	1.000000	1.000000	4.000000	1.000000

Data Cleaning

Since there are no null/missing values, it might be a good idea to check the `ScheduledDay`, `AppointmentDay` and `Age` columns for plausibility

```
In [9]: df.Age.describe()
```

	count	mean	std	min	25%	50%	75%	max	Name: Age, dtype: float64
count	110527.000000								
mean	37.088874								
std	23.110205								
min	-1.000000								
25%	18.000000								
50%	37.000000								
75%	55.000000								
max	115.000000								

```
In [10]: # An age of -1 does not make sense, we seem to have an issue there
df[df.Age < 0]
```

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hipertension	Diabetes
99832	4.659432e+14	5775010	F	2016-06-06 08:58:10+00:00	2016-06-06 00:00:00+00:00	-1	ROMÃO	0	0	0

```
In [11]: missing_age_patient_id = df.at[99832, 'PatientId']

In [12]: df[df.PatientId == missing_age_patient_id]['AppointmentID'].count()
```

	AppointmentID
1	1.4749613e+14

```
In [13]: # Since there are no further appointments for this patient and nothing else out of the ordinary, it should be fine
mean_age = df.Age.mean()
print(mean_age)
37.08887421173107

In [14]: df.at[99832, 'Age'] = mean_age

In [15]: df.at[99832, 'Age']
```

	Age
37	

Now on to the date columns. Since `ScheduledDay` is the day of the Appointment being scheduled, it should have taken place before the `AppointmentDay`.

One thing to notice is the additional precision on the `ScheduledDay` - there is a time of day available for the `ScheduledDay`, while the `AppointmentDay` seems to be lacking this information.

```
In [16]: df['waiting_time_until_appointment'] = df['AppointmentDay'] - df['ScheduledDay']

In [17]: df['waiting_time_until_appointment']
```

	waiting_time_until_appointment
0	-1 days +05:21:52
1	-1 days +07:51:33
2	-1 days +07:40:56
3	-1 days +06:30:29
4	-1 days +07:52:37
...	...
110522	34 days 14:44:25
110523	34 days 16:32:27
110524	40 days 07:56:08
110525	40 days 08:50:37
110526	40 days 10:29:04

Name: waiting_time_until_appointment, Length: 110527, dtype: timedelta64[ns]

There seem to be a couple of instances where the schedule day/time is after the appointment-date. This needs to be investigated further

```
In [18]: import datetime

In [19]: scheduled_after_appointment = df[df['waiting_time_until_appointment'] < datetime.timedelta(0)]

In [20]: scheduled_after_appointment
```

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hipertension	Diabetes
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	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
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
3	8.679512e+11	5642828	F	2016-04-29 17:29:31+00:00	2016-04-29 00:00:00+00:00	8	PONTAL DE CAMBURI	0	0	0
4	8.841986e+12	5642494	F	2016-04-29 16:07:23+00:00	2016-04-29 00:00:00+00:00	56	JARDIM DA PENHA	0	1	1
...
110511	8.235986e+11	5786742	F	2016-06-06 08:50:20+00:00	2016-06-06 00:00:00+00:00	14	MARIA ORTIZ	0	0	0
110512	9.876246e+13	5786368	F	2016-06-08 08:20:11+00:00	2016-06-08 00:00:00+00:00	41	MARIA ORTIZ	0	0	0
110513	8.67478e+13	5785964	M	2016-06-08 07:52:55+00:00	2016-06-08 00:00:00+00:00	2	HONÓRIO	0	0	0
110514	2.695685e+12	5785667	F	2016-06-08 08:35:31+00:00	2016-06-08 00:00:00+00:00	58	MARIA ORTIZ	0	0	0
110517	5.574942e+12	5780122	F	2016-06-07 07:38:34+00:00	2016-06-07 00:00:00+00:00	19	MARIA ORTIZ	0	0	0

38568 rows x 11 columns

```
In [21]: scheduled_after_appointment.describe()['waiting_time_until_appointment']
```

	count	mean	std	min	25%	50%	75%	max	Name: waiting_time_until_appointment, dtype: object
count	38568								
mean	-1 days +13:18:03.704807094								
std	0 days 01:07:36.67164507								
min	-7 days +10:10:40								
25%	-1 days +10:44:12.750000								
50%	-1 days +14:14:25.250000								
75%	-1 days +15:52:27.250000								
max	-1 days +17:50:24								

Name: waiting_time_until_appointment, dtype: object

There seems to be a high number of "same-day" appointments, that we categorized as scheduled after appointment due to the lack of the time-of-day for the appointment. Thus, we need to remove those "same-day" appointments from the `scheduled_after_appointment` frame

```
In [22]: scheduled_after_appointment = scheduled_after_appointment[scheduled_after_appointment.ScheduledDay.dt.day != df['AppointmentDay'].dt.day]

In [23]: scheduled_after_appointment
```

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hipertension	Diabetes
27033	7.839273e+12	5679978	M	2016-05-10 10:51:53+00:00	2016-05-09 00:00:00+00:00	38	RESISTÊNCIA	0	0	0
55226	7.896294e+12	5718660	F	2016-05-18 14:50:41+00:00	2016-05-17 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
64175	2.425226e+13	5664962	F	2016-05-05 13:43:58+00:00	2016-05-04 00:00:00+00:00	22	CONSOLAÇÃO	0	0	0
71533	9.982316e+14	5686628	F	2016-05-11 13:49:20+00:00	2016-05-05 00:00:00+00:00	81	SANTO ANTONIO	0	0	0
72362	3.787482e+12	5655637	M	2016-05-04 06:50:57+00:00	2016-05-03 00:00:00+00:00	7	TABUAZEIRO	0	0	0

Now we have a number of appointments left which do seem to have faulty (possibly backdated?) `ScheduledDay` or `AppointmentDay` values. What stands out is that all of these appointments are indeed No-shows. This this could possibly and error in the scheduling-software of the hospitals

One final point of interest would be, if the patients associated with these appointments had further appointments

```
In [24]: df[df.PatientId.isin(scheduled_after_appointment['PatientId'])].groupby(['PatientId'])['AppointmentID'].count()
```

	PatientId	AppointmentID
3.787482e+12	5	
7.839273e+12	5	
7.896294e+12	6	
2.425226e+13	2	
9.982316e+14	3	

Name: AppointmentID, dtype: int64

```
In [25]: df[df.PatientId.isin(scheduled_after_appointment['PatientId'])]
```

	PatientId	AppointmentID	Gender	ScheduledDay	AppointmentDay	Age	Neighbourhood	Scholarship	Hipertension	Diabetes
3370	7.839273e+12	5730318	M	2016-05-10 08:27:43+00:00	2016-05-24 00:00:00+00:00	38	RESISTÊNCIA	0	0	0
27033	7.839273e+12	5679978	M	2016-05-10 10:51:53+00:00	2016-05-09 00:00:00+00:00	38	RESISTÊNCIA	0	0	0
35377	7.896294e+12	5653953	F	2016-05-03 13:34:37+00:00	2016-05-05 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
53235	7.896294e+12	5717383	F	2016-05-19 07:54:32+00:00	2016-05-19 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
53250	7.896294e+12	5692824	F	2016-05-12 16:24:16+00:00	2016-05-12 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
53310	7.896294e+12	5693085	F	2016-05-12 17:36:32+00:00	2016-05-12 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
54692	2.425226e+13	5657354	F	2016-05-04 09:15:43+00:00	2016-05-19 00:00:00+00:00	22	CONSOLAÇÃO	0	0	0
55226	7.896294e+12	5715660	F	2016-05-18 14:50:41+00:00	2016-05-17 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
55289	7.896294e+12	5715663	F	2016-05-18 14:51:10+00:00	2016-05-18 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
55687	7.896294e+12	5692824	F	2016-05-12 16:23:09+00:00	2016-05-12 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
64175	2.425226e+13	5664962	F	2016-05-05 13:43:58+00:00	2016-05-04 00:00:00+00:00	22	CONSOLAÇÃO	0	0	0
71526	9.982316e+14	5630880	F	2016-04-27 16:40:44+00:00	2016-05-05 00:00:00+00:00	81	SANTO ANTONIO	0	0	0
71533	9.982316e+14	5686628	F	2016-05-11 13:49:20+00:00	2016-05-05 00:00:00+00:00	81	SANTO ANTONIO	0	0	0
72362	3.787482e+12	5655637	M	2016-05-04 06:50:57+00:00	2016-05-03 00:00:00+00:00	7	TABUAZEIRO	0	0	0
72363	3.787482e+12	5655638	M	2016-05-04 06:51:15+00:00	2016-05-10 00:00:00+00:00	7	TABUAZEIRO	0	0	0
72364	3.787482e+12	5655639	M	2016-05-04 06:51:29+00:00	2016-05-17 00:00:00+00:00	7	TABUAZEIRO	0	0	0
72365	3.787482e+12	5655642	M	2016-05-04 06:51:48+00:00	2016-05-24 00:00:00+00:00	7	TABUAZEIRO	0	0	0
72366	3.787482e+12	5655646	M	2016-05-04 06:52:08+00:00	2016-05-31 00:00:00+00:00	7	TABUAZEIRO	0	0	0
90382	2.425226e+13	5718600	F	2016-05-18 15:38:15+00:00	2016-06-01 00:00:00+00:00	81	SANTO ANTONIO	0	0	0
95335	7.896294e+12	5767269	F	2016-06-02 16:59:44+00:00	2016-06-02 00:00:00+00:00	19	SANTO ANTONIO	0	0	0
100002	7.839273e+12	5787285	M	2016-06-08 09:40:13+00:00	2016-06-08 00:00:00+00:00	38	RESISTÊNCIA	0	0	0
100003	7.839273e+12	5752857	M	2016-05-31 12:56:44+00:00	2016-06-01 00:00:00+00:00	38	RESISTÊNCIA	0	0	0
101919	7.839273e+12	5777702	M	2016-06-08 14:19:28+00:00	2016-06-08 00:00:00+00:00	38	RESISTÊNCIA	0	0	0

All of patients that

