

Part 1: Loading and Inspecting our data

. Before we begin inspecting our data we will need to load it by using pandas built in function `read_csv`. This will read our CSV file and load our data to a pandas `DataFrame` Object.

Pandas

. Pandas is a flexible python library we will be using to create our `DataFrame` for data manipulation and analysis.

Matplotlib/Pyplot

. Matplotlib is a plotting python 2D library. Pyplot is a command style function that makes Matplotlib work like MATLAB. Simpler version of Seaborn.

Seaborn

. Seaborn is a python library used for visualization, a more enhanced version of Matplotlib

```
In [1]: #Import the proper libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

#If this is used to make plots show inline
%matplotlib inline

In [2]: #To retrieve our data simply read it as a csv file and assign it to a respective variable
df = pd.read_csv("Users/princeamankwah/Desktop/mock_treatment.csv")
```

Health Care Data Analysis

. We will be exploring an oncology treatment data set which contains information about artificial patients.

```
In [3]: #lets take a look at the first 5 rows of our DataFrame
df.head()
```

PatientID	TreatmentStart	Drug	Dosage
0	PT1	1/14/16	Cisplatin 200
1	PT20	1/2/16	Cisplatin 140
2	PT2	1/10/16	Cisplatin 180
3	PT3	1/24/16	Cisplatin 140
4	PT4	2/14/16	Cisplatin 200

```
Out[3]:
```

```
In [4]: #lets obtain some descriptive statistics on our DataFrame
df.describe()
```

	Dosage
count	34.000000
mean	264.583333
std	528.963843
min	140.000000
25%	190.000000
50%	195.000000
75%	240.000000
max	1800.000000

```
Out[4]:
```

Note: The first thing i noticed after observing basic statistics is the high number of max Dosage. Compared to the other values, max Dosage is significantly higher. This may indicate a typo.

```
In [5]: #Observe the data types of our columns
df.dtypes
```

PatientID	object
TreatmentStart	object
Drug	object
Dosage	int64
dtype:	object

```
Out[5]:
```

```
In [6]: #Observe the number of rows of our data set
len(df)
```

```
Out[6]: 34
```

Accessing columns in our dataframe

```
In [7]: df.columns
```

```
Out[7]: Index(['PatientID', 'TreatmentStart', 'Drug', 'Dosage'], dtype='object')
```

```
In [8]: #Obtain each row in our PatientID column in tabular form. Notice the two brackets which gives it a DataFrame format
df[['PatientID']].head()
```

PatientID
0
1
2
3
4

```
Out[8]:
```

```
In [9]: df[['TreatmentStart']].head()
```

TreatmentStart
0
1
2
3
4

```
Out[9]:
```

```
In [10]: df[['Drug']].head()
```

Drug
0
1
2
3
4

```
Out[10]:
```

```
In [11]: df[['Dosage']].head()
```

Dosage
0
1
2
3
4

```
Out[11]:
```

```
In [12]: #Obtain the unique values for pateintID in Series form.
df['PatientID'].unique()
```

```
Out[12]: array(['PT1', 'PT20', 'PT2', 'PT3', 'PT4', 'PT19', 'PT5', 'PT6', 'PT7', 'PT8', 'PT9', 'PT10', 'PT11', 'PT16', 'PT17', 'PT18', 'PT13', 'PT14', 'PT15', 'PT17', 'PT18'], dtype=object)
```

```
In [13]: #Obtain the unique drugs in our DataFrame in Series form
df['Drug'].unique()
```

```
Out[13]: array(['Cisplatin', 'Nivolumab'], dtype=object)
```

```
In [14]: #Obtain the first row in our DataFrame
df.iloc[0]
```

PatientID	TreatmentStart	Drug	Dosage
0	PT1	1/14/16	Cisplatin 200

```
Out[14]:
```

```
In [15]: #Obtain 5 rows in our DataFrame excluding the first row
df.iloc[1:6]
```

PatientID	TreatmentStart	Drug	Dosage
1	PT20	1/2/16	Cisplatin 140
2	PT2	1/10/16	Cisplatin 180
3	PT3	1/24/16	Cisplatin 140
4	PT4	2/14/16	Cisplatin 200
5	PT19	2/10/16	Cisplatin 180

```
Out[15]:
```

```
In [16]: #Obtain the row(s) where PatientID is PT19
df.loc[df['PatientID'] == 'PT19']
```

PatientID	TreatmentStart	Drug	Dosage
5	PT19	2/10/16	Cisplatin 180
22	PT19	6/2/16	Nivolumab 240

```
Out[16]:
```

```
In [17]: #Obtain the row(s) where PatientID is PT20 & the Drug is Cisplatin
df.loc[(df['PatientID'] == 'PT20') & (df['Drug'] == 'Cisplatin')]
```

PatientID	TreatmentStart	Drug	Dosage
1	PT20	1/2/16	Cisplatin 140

```
Out[17]:
```

```
In [18]: #Obtain rows where Dosage is greater than 180
df.loc[df['Dosage'] > 180]
```

PatientID	TreatmentStart	Drug	Dosage
0	PT1	1/14/16	Cisplatin 200
4	PT4	2/14/16	Cisplatin 200
6	PT5	2/6/16	Cisplatin 190
8	PT7	3/1/16	Cisplatin 210
10	PT9	3/27/16	Nivolumab 240
11	PT10	4/7/16	Nivolumab 240
13	PT11	4/17/16	Cisplatin 190
15	PT12	5/15/16	Cisplatin 1800
17	PT14	5/3/16	Nivolumab 240
18	PT15	5/7/16	Nivolumab 240
19	PT1	6/17/16	Nivolumab 240
21	PT18	6/3/16	Nivolumab 240
22	PT19	6/2/16	Nivolumab 240
23	PT20	6/2/16	Nivolumab 240

```
Out[18]:
```

Part 2: Data Cleaning

Lets clean up our Dosage Column by updating a row

```
In [19]: #As you can see we have a typo in our Dosage column. We have a max value of 1800 which is way above our normal value
#lets update our row with 1800 to a value of 180
df['Dosage'].iloc[15] = 180

/Applications/anaconda3/lib/python3.8/site-packages/pandas/core/indexing.py:678: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df.iloc[15].setitem_with_inplace(index=15, value=180)
```

```
In [20]: #Observe our new table with updated Dosage column
df
```

PatientID	TreatmentStart	Drug	Dosage
0	PT1	1/14/16	Cisplatin 200
1	PT20	1/2/16	Cisplatin 140
2	PT2	1/10/16	Cisplatin 180
3	PT3	1/24/16	Cisplatin 140
4	PT4	2/14/16	Cisplatin 200
5	PT19	2/10/16	Cisplatin 180
6	PT5	2/6/16	Cisplatin 190
7	PT6	3/1/16	Cisplatin 180
8	PT7	3/1/16	Cisplatin 210
9	PT8	3/19/16	Cisplatin 180
10	PT9	3/27/16	Nivolumab 240
11	PT10	4/7/16	Nivolumab 240
12	PT10	4/7/16	Cisplatin 140
13	PT11	4/17/16	Cisplatin 190
14	PT16	4/9/16	Cisplatin 160
15	PT12	5/15/16	Cisplatin 180
16	PT13	5/2/16	Cisplatin 180
17	PT14	5/3/16	Nivolumab 240
18	PT15	5/7/16	Nivolumab 240
19	PT1	6/17/16	Nivolumab 240
20	PT17	6/17/16	Cisplatin 160
21	PT18	6/3/16	Nivolumab 240
22	PT19	6/2/16	Nivolumab 240
23	PT20	6/2/16	Nivolumab 240

```
Out[20]:
```

```
In [21]: #Sort DataFrame by Date column
df.sort_values('TreatmentStart').head()
```

PatientID	TreatmentStart	Drug	Dosage
2	PT2	1/10/16	Cisplatin 180
0	PT1	1/14/16	Cisplatin 200
1	PT20	1/2/16	Cisplatin 140
3	PT3	1/24/16	Cisplatin 140
5	PT19	2/10/16	Cisplatin 180

```
Out[21]:
```

Note: Notice how the date column is not properly sorted

. Lets clean up our `TreatmentStart` column

```
In [22]: #We first have to convert the TreatmentStart column to a datetime object then assign it to our Date column
df['TreatmentStart'] = pd.to_datetime(df['TreatmentStart'])
```

```
In [23]: #Observe our new Date column
df.head()
```

PatientID	TreatmentStart	Drug	Dosage
0	PT1	2016-01-14	Cisplatin 200
1	PT20	2016-01-02	Cisplatin 140
2	PT2	2016-01-10	Cisplatin 180
3	PT3	2016-01-24	Cisplatin 140
4	PT4	2016-02-14	Cisplatin 200

```
Out[23]:
```

```
In [24]: #Sort our DataFrame in ascending order then the TreatmentStart in descending order
df.sort_values(['TreatmentStart'], inplace=True)
```

```
In [25]: df.head()
```

PatientID	TreatmentStart	Drug	Dosage
1	PT20	2016-01-02	Cisplatin 140
2	PT2	2016-01-10	Cisplatin 180
0	PT1	2016-01-14	Cisplatin 200
3	PT3	2016-01-24	Cisplatin 140
6	PT5	2016-02-06	Cisplatin 190

```
Out[25]:
```

```
In [26]: #Sort multiple columns. First sort the Patient column in ascending order then the TreatmentStart in descending order
df.sort_values(['PatientID', 'TreatmentStart'], ascending = [True, False]).head()
```

PatientID	TreatmentStart	Drug	Dosage
19	PT1	2016-06-17	Nivolumab 240
0	PT1	2016-01-14	Cisplatin 200
12	PT10	2016-04-07	Cisplatin 140
13	PT10	2016-04-07	Nivolumab 240
13	PT11	2016-04-17	Cisplatin 190

```
Out[26]:
```

```
In [27]: #Sort the DataFrame by Dosage in descending order
df.sort_values('dosage', ascending = False)
```

PatientID	TreatmentStart	Drug	Dosage
11	PT10	2016-04-07	Nivolumab 240
19	PT1	2016-06-17	Nivolumab 240
18	PT15	2016-05-07	Nivolumab 240
22	PT19	2016-06-02	Nivolumab 240
10	PT19	2016-03-27	Nivolumab 240
21	PT18	2016-06-03	Nivolumab 240
17	PT14	2016-05-03	Nivolumab 240
23	PT20	2016-06-02	Nivolumab 240
8	PT7	2016-03-01	Cisplatin 210
0	PT1	2016-01-14	Cisplatin 200
4	PT4	2016-02-14	Cisplatin 200
13	PT11	2016-04-17	Cisplatin 190
6	PT5	2016-02-06	Cisplatin 190
15	PT12	2016-05-15	Cisplatin 180
9	PT8	2016-03-19	Cisplatin 180
2	PT2	2016-01-10	Cisplatin 180
7	PT6	2016-03-01	Cisplatin 180
5	PT19	2016-02-10	Cisplatin 180
16	PT13	2016-05-21	Cisplatin 180
20	PT17	2016-06-17	Cisplatin 160
14	PT16	2016-04-09	Cisplatin 160
12	PT10	2016-04-07	Cisplatin 140
3	PT3	2016-01-24	Cisplatin 140
1	PT20	2016-01-02	Cisplatin 140

```
Out[27]:
```

Data Analysis

```
In [28]: #Observe the number of patients
len(df['PatientID'])
```

```
Out[28]: 34
```

```
In [29]: #Observe the number of unique patients
len(df['PatientID'].unique())
```

```
Out[29]: 20
```

Note: As you can see there are 20 unique patients and the length of patients is 24, which means there are 4 Patients on two different medications

Patients treated at practice

How many Patients does the practice treat?

```
In [30]: print("The practice treats:", len(df['PatientID'].unique()), 'patients')

The practice treats: 20 patients
```

Drug usage at practice

What drugs are used at practice?

```
In [31]: #Implement the groupby() function to group the drugs in a DataFrame
df.groupby('Drug').count()
```

Drug	PatientID	TreatmentStart	Dosage
Cisplatin	16	16	16
Nivolumab	8	8	8

```
Out[31]:
```

```
In [32]: #Format the above code correctly for better reading by using reset_index() function
df.groupby('Drug').count()[['PatientID', 'TreatmentStart', 'Dosage']].reset_index()
```

Drug	PatientID	TreatmentStart	Dosage
0	Cisplatin	16	16
1	Nivolumab	8	8

```
Out[32]:
```

```
In [33]: #Our DataFrame is still in the format we want. The above code did not permanently modify our DataFrame.
#To permanently modify our DataFrame we could've used inplace=True
df.head()
```

PatientID	TreatmentStart	Drug	Dosage
1	PT20	2016-01-02	Cisplatin 140
2	PT2	2016-01-10	Cisplatin 180
0	PT1	2016-01-14	Cisplatin 200
3	PT3	2016-01-24	Cisplatin 140
6	PT5	2016-02-06	Cisplatin 190

```
Out[33]:
```

Lets plot some information

We will be comparing our Drug & PatientID Columns

Before plotting the comparison between Drug & PatientID, lets rename our PatientID column to PatientCount

```
In [34]: df.rename(columns={'PatientID' : 'PatientCount'}, inplace=True)

df.head()
```

PatientCount	TreatmentStart	Drug	Dosage
1	PT20	2016-01-02	Cisplatin 140
2	PT2	2016-01-10	Cisplatin 180
0	PT1	2016-01-14	Cisplatin 200
3	PT3	2016-01-24	Cisplatin 140
6	PT5	2016-02-06	Cisplatin 190

```
Out[35]:
```

Before plotting the comparison between Drug & Patient Count we have to assign that DataFrame to a variable for flexibility.

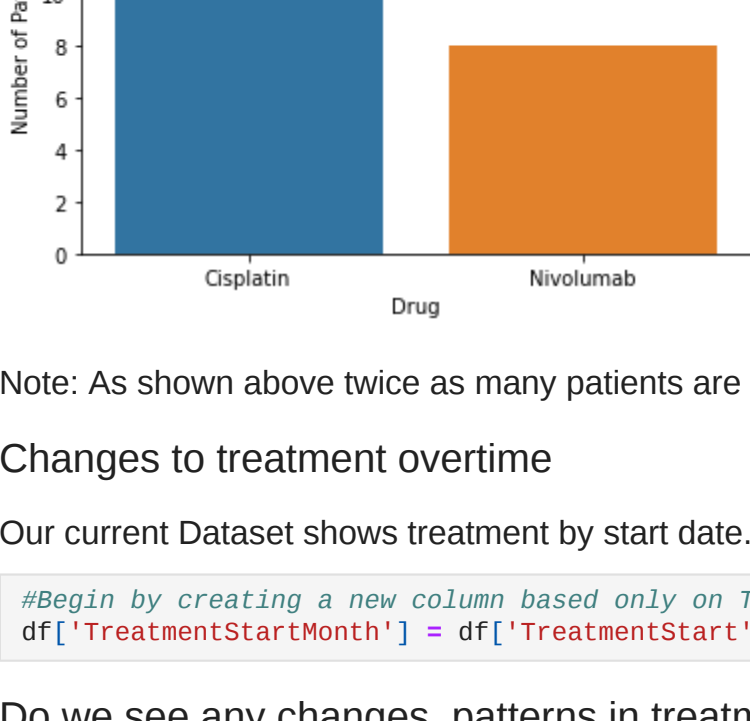
```
In [36]: counts = df.groupby('Drug').count()[['PatientCount']].reset_index()

counts
```

Drug	PatientCount
0	Cisplatin 16
1	Nivolumab 8

```
Out[37]:
```

```
In [38]: #lets begin by using a simple barchart to compare our 2 columns
fig = sns.barplot(data = counts, x = 'Drug', y = 'PatientCount')
plt.title('Drug & Patient Comparison')
plt.xlabel('Drug')
plt.ylabel('Number of Patients')
plt.show(fig)
```



Note: As shown above twice as many patients are taking Cisplatin compared to Nivolumab

Changes to treatment overtime

Our current Dataset shows treatment by start date. But we would like to compare the changes in treatments by the month to see how many patients start a given drug overtime

```
In [39]: #begin by creating a new column based only on TreatmentStart month
df['TreatmentStartMonth'] = df['TreatmentStart'].astype('datetime64[M]')
```

Do we see any changes, patterns in treatment overtime?

```
In [40]: #Observe our new column
df.head()
```

PatientCount	TreatmentStart	Drug	Dosage	TreatmentStartMonth
1	PT20	2016-01-02	Cisplatin 140	2016-01-01
2	PT2	2016-01-10	Cisplatin 180	2016-01-01
0	PT1	2016-01-14	Cisplatin 200	2016-01-01
3	PT3	2016-01-24	Cisplatin 140	2016-01-01
6	PT5	2016-02-06	Cisplatin 190	2016-02-01

```
Out[40]:
```

```
In [41]: #Now implement the groupby() function to compare TreatmentStartMonth, Drug, only counting Patients
drug_by_month = df.groupby(['TreatmentStartMonth', 'Drug']).count()[['PatientCount']].reset_index()
```

```
In [42]: #Show below is the the drug count started based on the month
drug_by_month
```

TreatmentStartMonth	Drug	PatientCount
0	2016-01-01	Cisplatin 4
1	2016-02-01	Cisplatin 3
2	2016-03-01	Cisplatin 3
3	2016-03-01	Nivolumab 1
4	2016-04-01	Cisplatin 3
5	2016-04-01	Nivolumab 1
6	2016-05-01	Cisplatin 2
7	2016-05-01	Nivolumab 2
8	2016-06-01	Cisplatin 1
9	2016-06-01	Nivolumab 4

```
Out[42]:
```

```
In [43]: #Now lets make use of Seaborn once again to visualize how many patients started each drug by month
fig = sns.lineplot(data = drug_by_month, x = 'TreatmentStartMonth', y = 'PatientCount', hue = 'Drug')
plt.title('Monthly Drug Comparison')
plt.show(fig)
```



What is the average dosage of each drug?

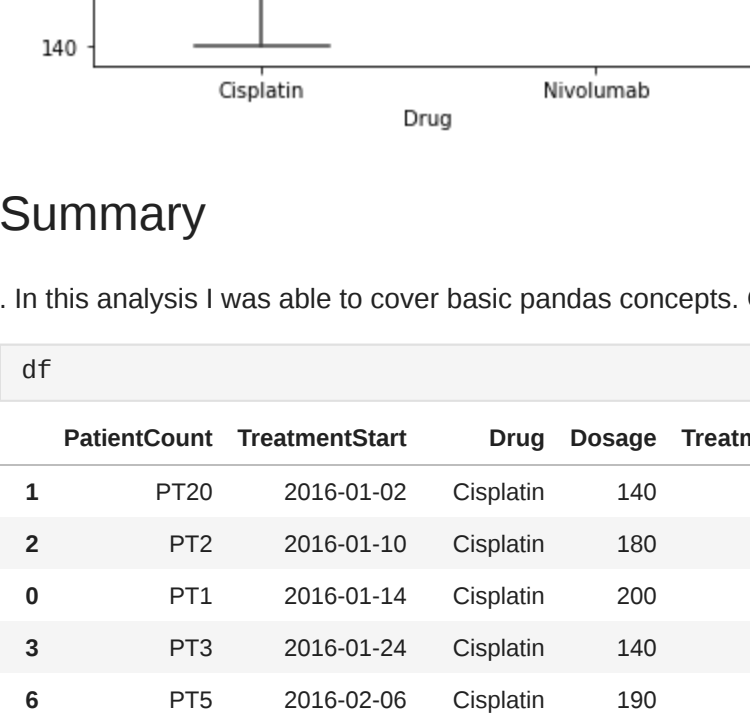
```
In [44]: #Observe some descriptive statistics based on the Drug Column
df.groupby(['Drug']).describe().reset_index()
```

	Drug	count	mean	std	min	25%	50%	75%	max
0	Cisplatin	16.0	175.625	21.899391	140.0	160.0	180.0	190.0	210.0
1	Nivolumab	8.0	240.000	0.000000	240.0	240.0	240.0	240.0	240.0

```
Out[44]:
```

Note: As shown above each value in row 1 besides count/standard deviation is 240 which means each patient who did not receive Nivolumab received a total dosage of 240.

```
In [45]: #Below we're able to visually compare the average Dosage of each drug
fig = sns.boxplot(data = df, x = 'Drug', y = 'Dosage', showFliers=False)
plt.show(fig)
```



Summary

. In this analysis I was able to cover basic pandas concepts. Overall I was able to explore our data set with the help of pandas and seaborn library. Some data cleaning was also done to further analyze our data.

```
In [46]: df
```

PatientCount	TreatmentStart	Drug	Dosage	TreatmentStartMonth
1	PT20	2016-01-02	Cisplatin 140	2016-01-01
2	PT2	2016-01-10	Cisplatin 180	2016-01-01
0	PT1	2016-01-14	Cisplatin 200	2016-01-01
3	PT3	2016-01-24	Cisplatin 140	2016-01-01
6	PT5	2016-02-06	Cisplatin 190	2016-02-01
5	PT19	2016-02-10	Cisplatin 180	2016-02-01
4	PT4	2016-02-14	Cisplatin 200	2016-02-01
8	PT6	2016-03-01	Cisplatin 180	2016-03-01
9	PT7	2016-03-19	Cisplatin 210	2016-03-01
10	PT8	2016-03-27	Nivolumab 240	2016-03-01
12	PT10	2016-04-07	Cisplatin 140	2016-04-01
11	PT10	2016-04-07	Nivolumab 240	2016-0