

911 Calls Capstone Project

For this capstone project we will be analyzing some 911 call data from Kaggle. The data contains the following fields:

- lat: String variable, Latitude
- lng: String variable, Longitude
- desc: String variable, Description of the Emergency Call
- zip: String variable, Zipcode
- time: String variable, Title
- timeStamp: String variable, YYYY-MM-DD HH:MM:SS
- twp: String variable, Township
- addr: String variable, Address
 - e: String variable, Dummy variable (always 1)

Just go along with this notebook and try to complete the instructions or answer the questions in bold using your Python and Data Science skills!

Data and Setup

```
import numpy and pandas

In [5]: import numpy as np
import pandas as pd

Import visualization libraries and set matplotlib inline.

In [6]: import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

Read in the csv file as a dataframe called df

In [3]: df = pd.read_csv('911.csv')

Check the info() of the df

In [11]: df.info()

<class 'pandas.core.dataframe'>
RangeIndex: 99492 entries, 0 to 99491
Data columns (total 9 columns)
# Column Non-Null Count Dtype
--
0 lat 99492 non-null float64
1 lng 99492 non-null float64
2 desc 99492 non-null object
3 zip 86637 non-null float64
4 title 99492 non-null object
5 timeStamp 99492 non-null object
6 twp 99449 non-null object
7 addr 88973 non-null object
8 e 99492 non-null int64
dtypes: float64(3), int64(1), object(5)
memory usage: 6.8+ MB

Check the head of df

In [12]: df.head()

Out[5]:
```

	lat	lng	desc	zip	title	timeStamp	twp	addr	e
0	40.297876	-75.581284	REINDER CT & DEAD END, NEW HANOVER, Station ...	19525.0	EMS: BACK PAINSINJURY	2015-12-10 17:40:00	NEW HANOVER	REINDER CT & DEAD END	1
1	40.258061	-75.264680	BRIAR PATH & WHITEMARSH LN, HATFIELD TOWNSHIP...	19446.0	EMS: DIABETIC EMERGENCY	2015-12-10 17:40:00	HATFIELD TOWNSHIP	BRIAR PATH & WHITEMARSH LN	1
2	40.121182	-75.351975	HAWS AVE, NORRISTOWN, 2015-12-10 @ 14:39:21 SE...	19401.0	Fire: GAS-ODORLEAK	2015-12-10 17:40:00	NORRISTOWN	HAWS AVE	1
3	40.161513	-75.349513	AIRY ST & SWEDE ST, NORRISTOWN, Station 300A...	19401.0	EMS: CARDIAC EMERGENCY	2015-12-10 17:40:01	NORRISTOWN	AIRY ST & SWEDE ST	1
4	40.251482	-75.603350	CERRYWOOD CT & DEAD END, LOWER POTTSGROVE, S...	NaN	EMS: DIZZINESS	2015-12-10 17:40:01	LOWER POTTSGROVE	CERRYWOOD CT & DEAD END	1

Basic Questions

What are the top 5 zipcodes for 911 calls?

```
In [6]: df['zip'].value_counts().head(5)

Out[6]:
19401.0    6079
19464.0     6043
19403.0     4854
19446.0     4748
19486.0     3174
Name: zip, dtype: int64

What are the top 5 townships (twp) for 911 calls?

In [7]: df['twp'].value_counts().head(5)

Out[7]:
LOWER MERION    8443
ADINKTOWN      5977
NORRISTOWN      5389
UPPER MERION    5227
CHILTEWAW       4575
Name: twp, dtype: int64

Take a look at the 'title' column, how many unique title codes are there?

In [8]: df['title'].nunique()

Out[8]:
118
```

Creating new features

In the titles column there are "Reasons/Departments" specified before the title code. These are EMS, Fire, and Traffic. Use .apply() with a custom lambda expression to create a new column called "Reason" that contains this string value.

For example, if the title column value is EMS: BACK PAINSINJURY, the Reason column value would be EMS.

```
In [9]: df['Reason'] = df['title'].apply(lambda title: title.split(':')[0])

In [10]: df['Reason']

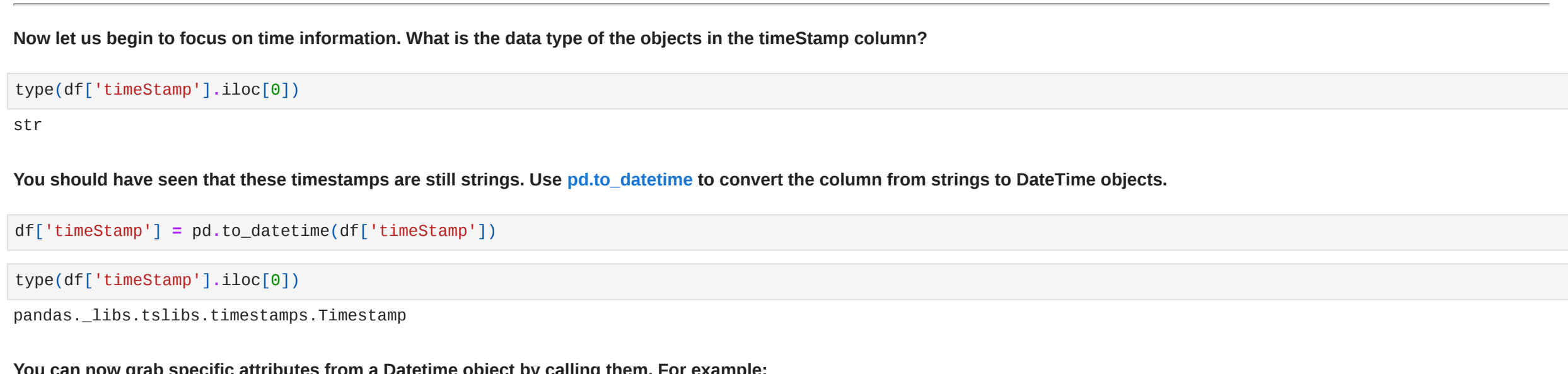
Out[10]:
0      EMS
1      EMS
2      Fire
3      EMS
4      EMS
...
99487   Traffic
99488   Traffic
99489   EMS
99490   EMS
99491   Traffic
Name: Reason, Length: 99492, dtype: object

What is the most common Reason for a 911 call based off of this new column?

In [11]: df['Reason'].value_counts()

Out[11]:
EMS      48877
Traffic  35695
Fire     14928
Name: Reason, dtype: int64

Now use seaborn to create a countplot of 911 calls by Reason.
```



Now let us begin to focus on time information. What is the data type of the objects in the timeStamp column?

```
In [13]: type(df['timeStamp']).iloc[0]
str

You should have seen that these timestamps are still strings. Use pd.to_datetime to convert the column from strings to DateTime objects.

In [14]: df['timeStamp'] = pd.to_datetime(df['timeStamp'])

In [15]: type(df['timeStamp']).iloc[0]
pandas._libs.tslibs.timestamps.Timestamp

You can now grab specific attributes from a Datetime object by calling them. For example:
```

```
time = df['timeStamp'].iloc[0]
time.hour

You can use Jupyter's tab method to explore the various attributes you can call. Now that the timestamp column are actually DateTime objects, use .apply() to create 3 new columns called Hour, Month, and Day of Week. You will create these columns based off of the timeStamp column, reference the solutions if you get stuck on this step.
```

```
In [16]: time = df['timeStamp'].iloc[0]

In [17]: time.hour

Out[17]:
17

In [18]: df['hour'] = df['timeStamp'].apply(lambda time: time.hour)

In [19]: df['month'] = df['timeStamp'].apply(lambda time: time.month)
df['day of week'] = df['timeStamp'].apply(lambda time: time.dayofweek)

In [20]: df.head()

Out[20]:
```

	lat	lng	desc	zip	title	timeStamp	twp	addr	e	Reason	Hour	Month	Day of Week
0	40.297876	-75.581284	REINDER CT & DEAD END, NEW HANOVER, Station ...	19525.0	EMS: BACK PAINSINJURY	2015-12-10 17:40:00	NEW HANOVER	REINDER CT & DEAD END	1	EMS	17	12	3
1	40.258061	-75.264680	BRIAR PATH & WHITEMARSH LN, HATFIELD TOWNSHIP...	19446.0	EMS: DIABETIC EMERGENCY	2015-12-10 17:40:00	HATFIELD TOWNSHIP	BRIAR PATH & WHITEMARSH LN	1	EMS	17	12	3
2	40.121182	-75.351975	HAWS AVE, NORRISTOWN, 2015-12-10 @ 14:39:21 SE...	19401.0	Fire: GAS-ODORLEAK	2015-12-10 17:40:00	NORRISTOWN	HAWS AVE	1	Fire	17	12	3
3	40.161513	-75.349513	AIRY ST & SWEDE ST, NORRISTOWN, Station 300A...	19401.0	EMS: CARDIAC EMERGENCY	2015-12-10 17:40:01	NORRISTOWN	AIRY ST & SWEDE ST	1	EMS	17	12	3
4	40.251482	-75.603350	CERRYWOOD CT & DEAD END, LOWER POTTSGROVE, S...	NaN	EMS: DIZZINESS	2015-12-10 17:40:01	LOWER POTTSGROVE	CERRYWOOD CT & DEAD END	1	EMS	17	12	3

Notice how the Day of Week is an integer 0-6. Use the .map() with this dictionary to map the actual string names to the day of the week:

```
dnmap = {0:'Mon',1:'Tue',2:'Wed',3:'Thu',4:'Fri',5:'Sat',6:'Sun'}

In [21]: dnmap = {0:'Mon',1:'Tue',2:'Wed',3:'Thu',4:'Fri',5:'Sat',6:'Sun'}

In [22]: df['day of week'] = df['day of week'].map(dnmap)

In [23]: df.head()

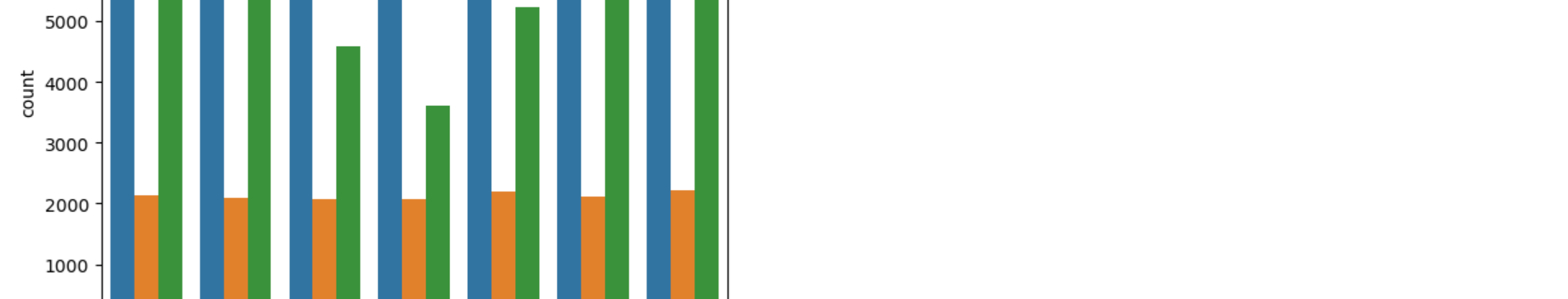
Out[23]:
```

	lat	lng	desc	zip	title	timeStamp	twp	addr	e	Reason	Hour	Month	Day of Week
0	40.297876	-75.581284	REINDER CT & DEAD END, NEW HANOVER, Station ...	19525.0	EMS: BACK PAINSINJURY	2015-12-10 17:40:00	NEW HANOVER	REINDER CT & DEAD END	1	EMS	17	12	Thu
1	40.258061	-75.264680	BRIAR PATH & WHITEMARSH LN, HATFIELD TOWNSHIP...	19446.0	EMS: DIABETIC EMERGENCY	2015-12-10 17:40:00	HATFIELD TOWNSHIP	BRIAR PATH & WHITEMARSH LN	1	EMS	17	12	Thu
2	40.121182	-75.351975	HAWS AVE, NORRISTOWN, 2015-12-10 @ 14:39:21 SE...	19401.0	Fire: GAS-ODORLEAK	2015-12-10 17:40:00	NORRISTOWN	HAWS AVE	1	Fire	17	12	Thu
3	40.161513	-75.349513	AIRY ST & SWEDE ST, NORRISTOWN, Station 300A...	19401.0	EMS: CARDIAC EMERGENCY	2015-12-10 17:40:01	NORRISTOWN	AIRY ST & SWEDE ST	1	EMS	17	12	Thu
4	40.251482	-75.603350	CERRYWOOD CT & DEAD END, LOWER POTTSGROVE, S...	NaN	EMS: DIZZINESS	2015-12-10 17:40:01	LOWER POTTSGROVE	CERRYWOOD CT & DEAD END	1	EMS	17	12	Thu

Now use seaborn to create a countplot of the Day of Week column with the hue based off of the Reason column.



Now do the same for Month:



Did you notice something strange about the Plot?

You should have noticed it was missing some Months, let's see if we can maybe fill in this information by plotting the information in another way, possibly a simple line plot that fills in the missing months, in order to do this, we'll need to do some work with pandas...

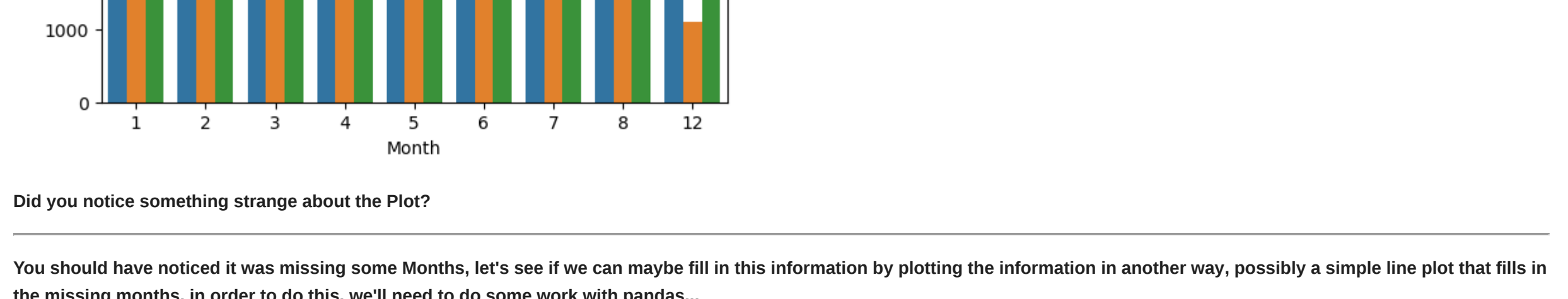
Now create a groupby object called byMonth, where you group the DataFrame by the month column and use the count() method for aggregation. Use the head() method on this returned DataFrame.

```
In [26]: byMonth = df.groupby('Month').count()
byMonth.head()

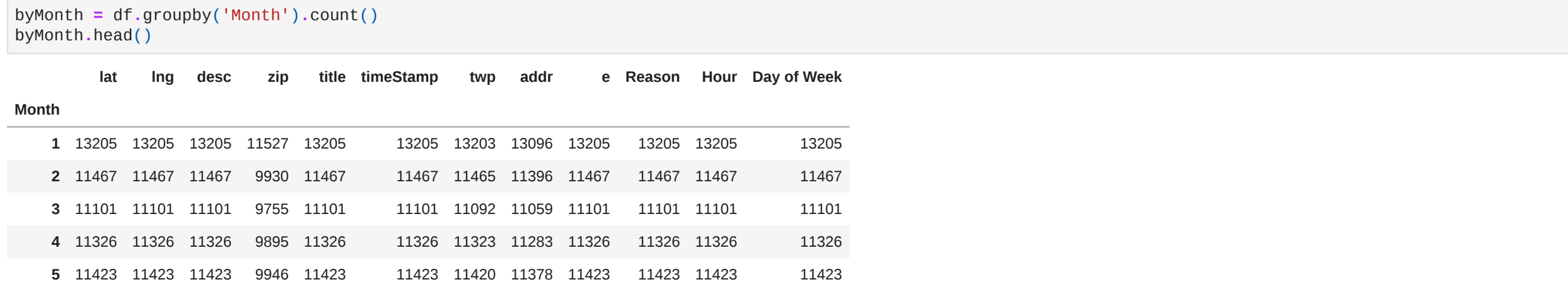
Out[26]:
```

Month	lat	lng	desc	zip	title	timeStamp	twp	addr	e	Reason	Hour	Day of Week
1	13205	13205	13205	11527	13205	13205	13203	13096	13205	13205	13205	
2	11467	11467	11467	9930	11467	11467	11465	11396	11467	11467	11467	
3	11101	11101	11101	9756	11101	11101	11092	11059	11101	11101	11101	
4	11326	11326	11326	8995	11326	11326	11323	11283	11326	11326	11326	
5	11423	11423	11423	9946	11423	11423	11420	11378	11423	11423	11423	

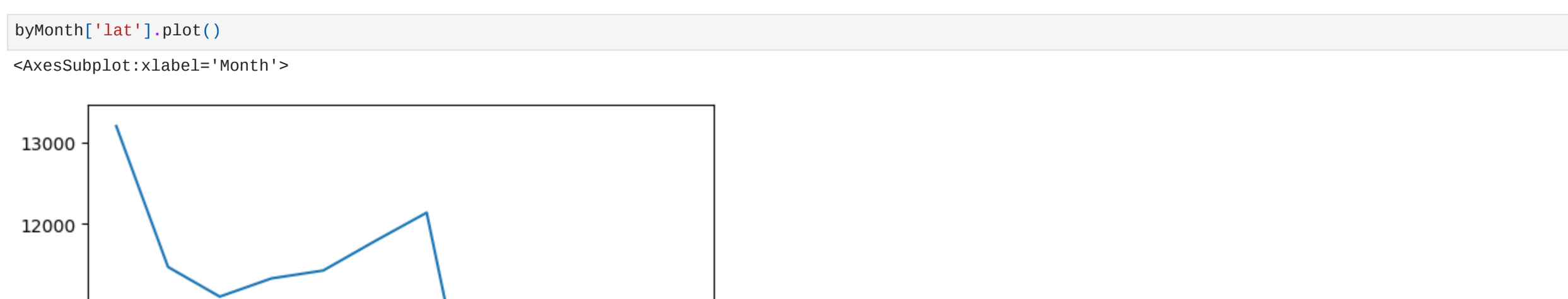
Now create a simple plot of df of the dataframe indicating the count of calls per month.



No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



Now see if you can use seaborn's lmplo() to create a linear fit on the number of calls per month. Keep in mind you may need to reset the index to a column.



Create a new column called 'date' that contains the date from the timeStamp column. You'll need to use apply along with the .date() method.

```
In [30]: t = df['timeStamp'].iloc[0]
t.date()

Out[30]:
datetime.date(2015, 12, 10)

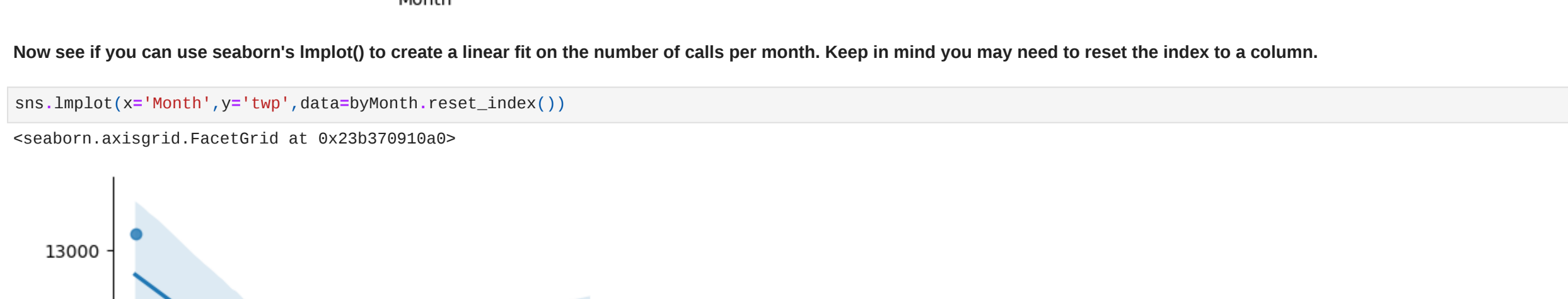
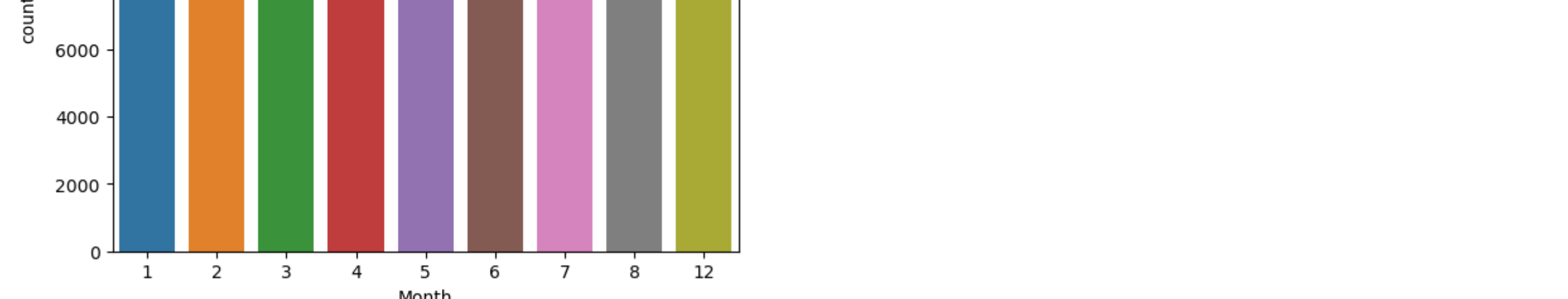
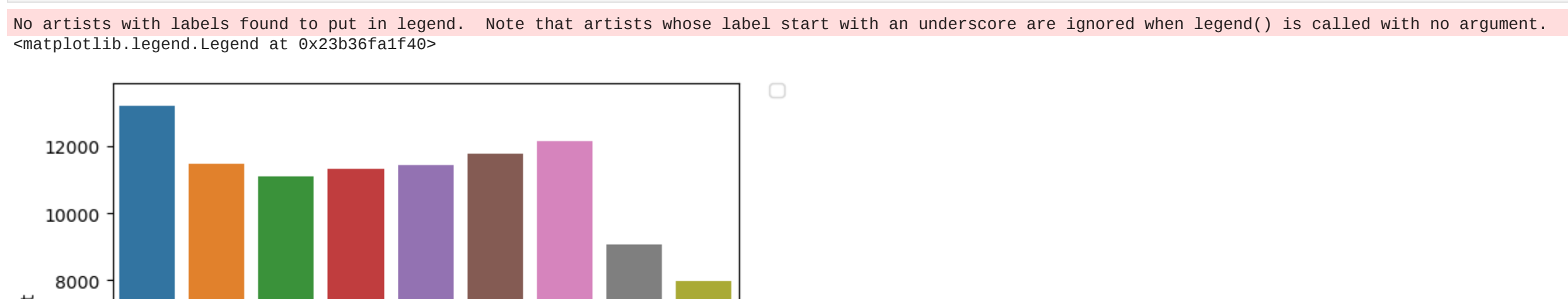
In [31]: df['date'] = df['timeStamp'].apply(lambda t: t.date())

In [32]: df.head()

Out[32]:
```

	lat	lng	desc	zip	title	timeStamp	twp	addr	e	Reason	Hour	Month	Day of Week	date
0	40.297876	-75.581284	REINDER CT & DEAD END, NEW HANOVER, Station ...	19525.0	EMS: BACK PAINSINJURY	2015-12-10 17:40:00	NEW HANOVER	REINDER CT & DEAD END	1	EMS	17	12	Thu	2015-12-10
1	40.258061	-75.264680	BRIAR PATH & WHITEMARSH LN, HATFIELD TOWNSHIP...	19446.0	EMS: DIABETIC EMERGENCY	2015-12-10 17:40:00	HATFIELD TOWNSHIP	BRIAR PATH & WHITEMARSH LN	1	EMS	17	12	Thu	2015-12-10
2	40.121182	-75.351975	HAWS AVE, NORRISTOWN, 2015-12-10 @ 14:39:21 SE...	19401.0	Fire: GAS-ODORLEAK	2015-12-10 17:40:00	NORRISTOWN	HAWS AVE	1	Fire	17	12	Thu	2015-12-10
3	40.161513	-75.349513	AIRY ST & SWEDE ST, NORRISTOWN, Station 300A...	19401.0	EMS: CARDIAC EMERGENCY	2015-12-10 17:40:01	NORRISTOWN	AIRY ST & SWEDE ST	1	EMS	17	12	Thu	2015-12-10
4	40.251482	-75.603350	CERRYWOOD CT & DEAD END, LOWER POTTSGROVE, S...	NaN	EMS: DIZZINESS	2015-12-10 17:40:01	LOWER POTTSGROVE	CERRYWOOD CT & DEAD END	1	EMS	17	12	Thu	2015-12-10

Now groupby this Date column with the count() aggregate and create a plot of counts of 911 calls.



Now let's move on to creating heatmaps with seaborn and our data. We'll first need to restructure the dataframe so that the columns become the Hours and the index becomes the Day of the Week. There are lots of ways to do this, but I would recommend trying to combine groupby with an .unstack() method. Reference the solutions if you get stuck on this!

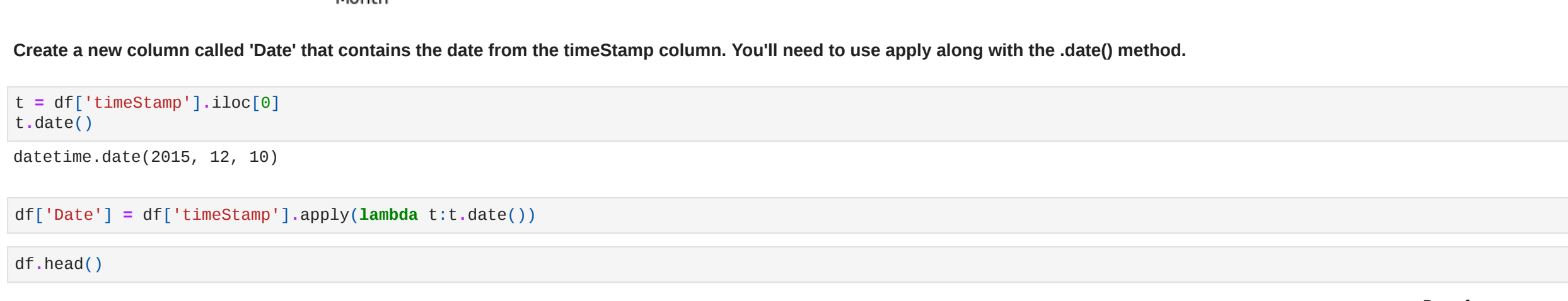
```
In [42]: dayHour = df.groupby(by=['day of week', 'hour']).count()['Reason'].unstack()

In [43]: dayHour

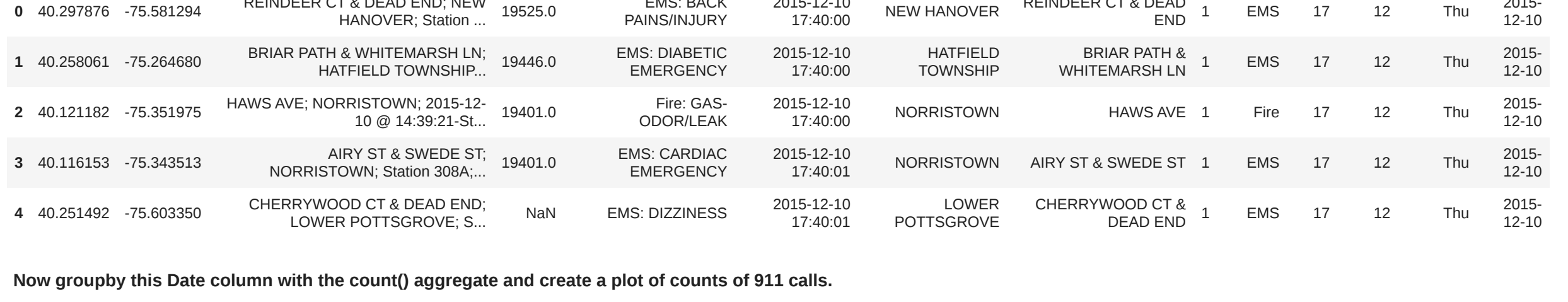
Out[43]:
```

Day of Week	Hour	0	1	2	3	4	5	6	7	8	9	...	14	15	16	17	18	19	20	21	22	23
Fri	275	235	191	175	201	194	372	598	742	752	...	932	980	1039	980	820	696	667	559	514	474	
Mon	282	221	201	194	204	257	397	653	819	785	...	899	813	989	997	885	746	613	497	472	325	
Sat	375	301	263	260	224	231	257	391	459	640	...	789	796	648	757	778	696	628	572	506	467	
Sun	383	306	286	268	242	240	300	402	483	620	...	684	691	663	714	670	655	537	461	415	330	
Thu	278	202	236	158	202	352	570	777	828	...	676	969	925	1013	810	698	617	571	462	274		
Tue	209	240	186	170	209	239	415	655	889	880	...	943	938	1026	1019	905	731	647	571	462	274	
Wed	250	216	189	209	156	255	410	701	875	808	...	904	867	990	1027	894	686	575	490	335		

Now create a HeatMap using this new DataFrame.



Now create a clustermap using this DataFrame.



Now repeat these same plots and operations, for a DataFrame that shows the Month as the column.

