Before you turn in the homework, make sure everything runs as expected. To do so, select **Kernel**→**Restart & Run All** in the toolbar above. Remember to submit both on **DataHub** and **Gradescope**.

Please fill in your name and include a list of your collaborators below.

```
In [2]:  NAME = "William Sheu"
         COLLABORATORS = ""
```

# Project 2: NYC Taxi Rides

# Part 3: NYC Accidents Data

In the real world, data isn't always nicely bundled in one file; data can be sourced from many places with many formats. Now we will use NYC accident data to try to improve our set of features.

In this part of the project, you'll do some EDA over the combined data set. We'll do a lot of the coding work for you, but there will be a few coding subtasks for you to complete on your own, as well as many results to interpret.

## Note

If your kernel dies unexpectedly, make sure you have shutdown all other notebooks. Each notebook uses valuable memory which we will need for this part of the project.

# Imports

Let us start by loading the Python libraries and custom tools we will use in this part.

```
In [3]:  import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
         import zipfile
         import os
         from pathlib import Path

         sns.set(style="whitegrid", palette="muted")

         plt.rcParams['figure.figsize'] = (12, 9)
         plt.rcParams['font.size'] = 12

         %matplotlib inline
```

## Downloading the Data

We will use the `fetch_and_cache` utility to download the dataset.

```
In [4]:  # Download and cache urls and get the file objects.
         from utils import fetch_and_cache
         data_url = 'https://github.com/DS-100/fa18/raw/gh-pages/assets/datasets,
         file_name = 'collisions.zip'
         dest_path = fetch_and_cache(data_url=data_url, file=file_name)

         print(f'Located at {dest_path}')
```

```
Using version already downloaded: Fri Nov 23 03:41:41 2018
MD5 hash of file: a445b925d24f319cb60bd3ace6e4172b
Located at data/collisions.zip
```

We will store the taxi data locally before loading it.

```
In [5]:  collisions_zip = zipfile.ZipFile(dest_path, 'r')

         #Extract zip files
         collisions_dir = Path('data/collisions')
         collisions_zip.extractall(collisions_dir)
```

## Loading and Formatting Data

The following code loads the collisions data into a Pandas DataFrame.

```
In [6]:  # Run this cell to load the collisions data.
         skiprows = None
         collisions = pd.read_csv(collisions_dir/'collisions_2016.csv', index_co
                                  parse_dates={'DATETIME':["DATE","TIME"]}, skip
         collisions['TIME'] = pd.to_datetime(collisions['DATETIME']).dt.hour
         collisions['DATE'] = pd.to_datetime(collisions['DATETIME']).dt.date
         collisions = collisions.dropna(subset=['LATITUDE', 'LONGITUDE'])
         collisions = collisions[collisions['LATITUDE'] <= 40.85]
         collisions = collisions[collisions['LATITUDE'] >= 40.63]
         collisions = collisions[collisions['LONGITUDE'] <= -73.65]
         collisions = collisions[collisions['LONGITUDE'] >= -74.03]
         collisions.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 116691 entries, 3589202 to 3363795
Data columns (total 30 columns):
DATETIME                          116691 non-null datetime64[ns]
Unnamed: 0                        116691 non-null int64
BOROUGH                           100532 non-null object
ZIP CODE                          100513 non-null float64
LATITUDE                          116691 non-null float64
LONGITUDE                         116691 non-null float64
LOCATION                          116691 non-null object
ON STREET NAME                    95914 non-null object
CROSS STREET NAME                 95757 non-null object
OFF STREET NAME                   61545 non-null object
NUMBER OF PERSONS INJURED         116691 non-null int64
NUMBER OF PERSONS KILLED          116691 non-null int64
NUMBER OF PEDESTRIANS INJURED     116691 non-null int64
NUMBER OF PEDESTRIANS KILLED      116691 non-null int64
NUMBER OF CYCLIST INJURED         116691 non-null int64
NUMBER OF CYCLIST KILLED          116691 non-null int64
NUMBER OF MOTORIST INJURED        116691 non-null int64
NUMBER OF MOTORIST KILLED         116691 non-null int64
CONTRIBUTING FACTOR VEHICLE 1     115162 non-null object
CONTRIBUTING FACTOR VEHICLE 2     101016 non-null object
CONTRIBUTING FACTOR VEHICLE 3     7772 non-null object
CONTRIBUTING FACTOR VEHICLE 4     1829 non-null object
CONTRIBUTING FACTOR VEHICLE 5     434 non-null object
VEHICLE TYPE CODE 1               115181 non-null object
VEHICLE TYPE CODE 2               92815 non-null object
VEHICLE TYPE CODE 3               7260 non-null object
VEHICLE TYPE CODE 4               1692 non-null object
VEHICLE TYPE CODE 5               403 non-null object
TIME                              116691 non-null int64
DATE                              116691 non-null object
dtypes: datetime64[ns](1), float64(3), int64(10), object(16)
memory usage: 27.6+ MB
```
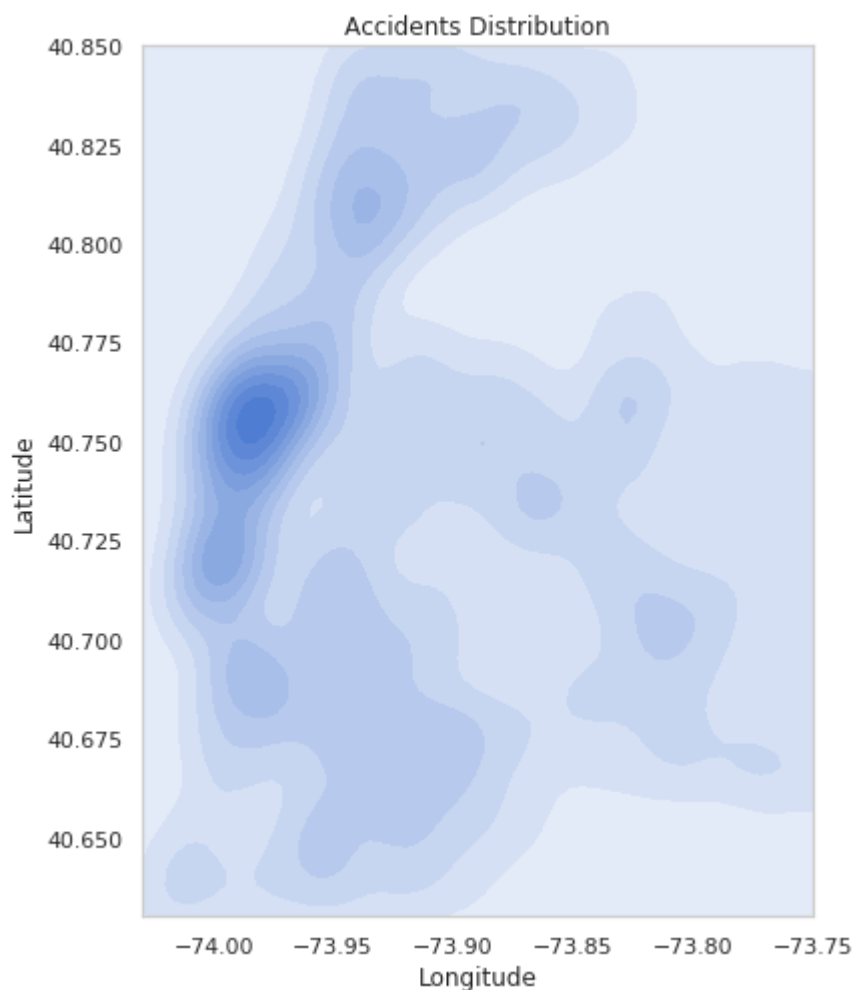
## 1: EDA of Accidents

Let's start by plotting the latitude and longitude where accidents occur. This may give us some insight on taxi ride durations. We sample N times (given) from the collisions dataset and create a 2D KDE plot of the longitude and latitude. We make sure to set the x and y limits according to the boundaries of New York, given below.

Here is a [map of Manhattan](https://www.google.com/maps/place/Manhattan,+New+York,+NY/@40.7590402,-74.0394431,12z73.9712488) for your convenience.

```
In [7]:  # Plot lat/lon of accidents, will take a few seconds
         N = 20000
         city_long_border = (-74.03, -73.75)
         city_lat_border = (40.63, 40.85)

         sample = collisions.sample(N)
         plt.figure(figsize=(6,8))
         sns.kdeplot(sample["LONGITUDE"], sample["LATITUDE"], shade=True)
         plt.xlim(city_long_border)
         plt.ylim(city_lat_border)
         plt.xlabel("Longitude")
         plt.ylabel("Latitude")
         plt.title("Accidents Distribution")
         plt.show();
```



## Question 1a

What can you say about the location density of NYC collisions based on the plot above?

**Hint: Here is a page
(https://www.google.com/maps/place/Manhattan,+New+York,+NY/@40.7590402,-74.0394431
73.9712488) that may be useful, and** another page (https://www.6sqft.com/what-nycs-
population-looks-like-day-vs-night/) **that may be useful.**

In [8]:
```python
q1a_answer = r"""

Most of the accidents are centrallized about midtown Manhattan. This may
"""

# YOUR CODE HERE
#raise NotImplementedError()

print(q1a_answer)
```

```
Most of the accidents are centrallized about midtown Manhattan. This m
ay be because midtown has the most fluxuation of people throughout the
day, leading to more possibilities for collisions between 2 vechicles.
```

We see that an entry in accidents contains information on number of people injured/killed.
Instead of using each of these columns separately, let's combine them into one column called
`'SEVERITY'`. Let's also make columns `FATALITY` and `INJURY`, each aggregating the
fatalities and injuries respectively.

In [9]:
```python
collisions['SEVERITY'] = collisions.filter(regex=r'NUMBER OF *').sum(ax
collisions['FATALITY'] = collisions.filter(regex=r'KILLED').sum(axis=1)
collisions['INJURY'] = collisions.filter(regex=r'INJURED').sum(axis=1)
```

Now let's group by time and compare two aggregations: count vs mean. Below we plot the
number of collisions and the mean severity of collisions by the hour, i.e. the `TIME` column. We
visualize them side by side and set the start of our day to be 6 a.m.

Let's also take a look at the mean number of casualties per hour and the mean number of
injuries per hour, plotted below.

```
In [10]:  fig, axes = plt.subplots(2, 2, figsize=(16,16))
          order = np.roll(np.arange(24), -6)
          ax1 = axes[0,0]
          ax2 = axes[0,1]
          ax3 = axes[1,0]
          ax4 = axes[1,1]

          collisions_count = collisions.groupby('TIME').count()
          collisions_count = collisions_count.reset_index()
          sns.barplot(x='TIME', y='SEVERITY', data=collisions_count, order=order,
          ax1.set_title("Accidents per Hour")
          ax1.set_xlabel("HOUR")
          ax1.set_ylabel('COUNT')


          collisions_mean = collisions.groupby('TIME').mean()
          collisions_mean = collisions_mean.reset_index()
          sns.barplot(x='TIME', y='SEVERITY', data=collisions_mean, order=order, a
          ax2.set_title("Severity of Accidents per Hour")
          ax2.set_xlabel("HOUR")
          ax2.set_ylabel('MEAN SEVERITY')

          fatality_count = collisions.groupby('TIME').mean()
          fatality_count = fatality_count.reset_index()
          sns.barplot(x='TIME', y='FATALITY', data=fatality_count, order=order, a
          ax3.set_title("Fatality per Hour")
          ax3.set_xlabel("HOUR")
          ax3.set_ylabel('MEAN KILLED')

          injury_count = collisions.groupby('TIME').mean()
          injury_count = injury_count.reset_index()
          sns.barplot(x='TIME', y='INJURY', data=injury_count, order=order, ax=ax
          ax4.set_title("Injury per Hour")
          ax4.set_xlabel("HOUR")
          ax4.set_ylabel('MEAN INJURED')

          plt.show();
```
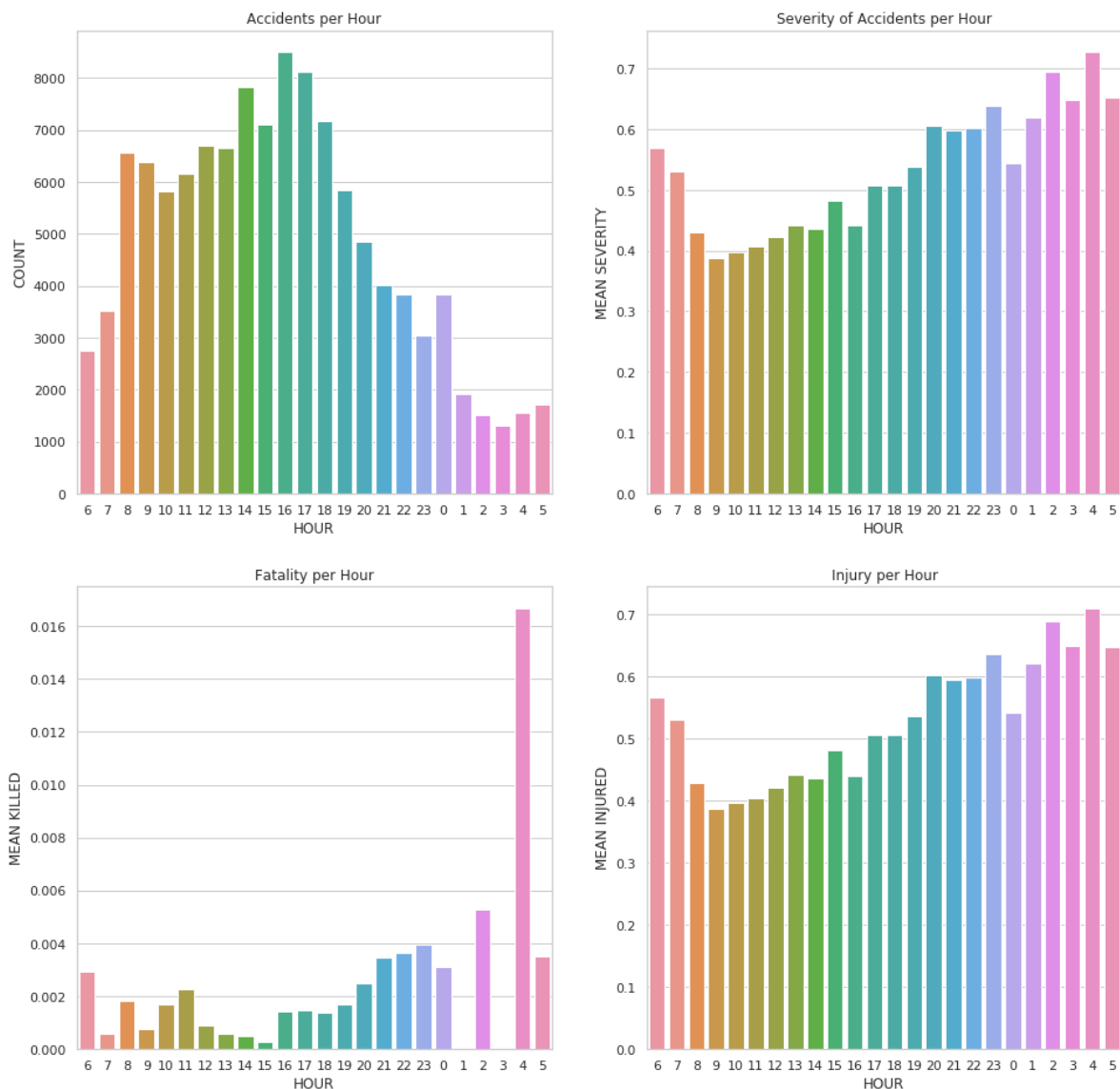
## Question 1b

Based on the visualizations above, what can you say about each? Make a comparison between the accidents per hour vs the mean severity per hour. What about the number of fatalities per hour vs the number of injuries per hour? Why do we chose to have our hours start at 6 as opposed to 0?

In [11]:
```python
q1b_answer = r"""

The accidents per hour plot shows about what one expects: there are very

The severty of accidents per hour shows that during the accidents very

The fatality per hour plot shows that there are less fatalities during

The injury per hour plot basically mirrors the severity of accidents per

There seems to be a inverse correlation between the accidents/hour plot

The fatality per hour plot seems to be a scaled version of the injury p

We chose to start our hour at 6 rather than 0 because 6 is usually seen

"""

# YOUR CODE HERE
#raise NotImplementedError()

print(q1b_answer)
```

The accidents per hour plot shows about what one expects: there are ve
ry few accidents/hour early in the morning and late in the night, when
there are less cars about, and much more accidents/hour during the day
time.

The severty of accidents per hour shows that during the accidents very
early in the morning or very late at night, the severety of the accide
nts are on average much higher than during the daytime. This may be du
e to the fact that there is no/less sunlight that could otherwise help
the driver lessen the damage done in an accident.This may also be beca
use people drive more recklessly when nobody else is on the road.

The fatality per hour plot shows that there are less fatalities during
the daytime, and more during the early morning/late night. This may al
so be explained by the points above. There is also mysteriously no fat
alities during hour 1 and hour 3, and a very sharp increase of fatalie
s/hour during hour 4.

The injury per hour plot basically mirrors the severty of accidents pe
r hour plot.

There seems to be a inverse correlation between the accidents/hour plo
t and the severety of accidents/hour. This may be because when there i
s more traffic, there will be more accidents, but less severe; whereas
if there is less traffic, there are less people and thus less accident
s, but as a result, people drive more recklessly, increasing the sever
ity of an accident.

The fatality per hour plot seems to be a scaled version of the injury
per hour plot, with less fatailies and injuries on average during dayt
ime hours, then increase during nighttime/early morning hours. Howeve

r, there is a dramatic spike in the number of fatalities during hour
4, and there is no such spike (to that degree) during that hour in the
injury per hour plot.

We chose to start our hour at 6 rather than 0 because 6 is usually see
n as he start of the day, it being approximately the time of sunrise.

Let's also check the relationship between location and severity. We provide code to visualize a
heat map of collisions, where the x and y coordinate are the location of the collision and the
heat color is the severity of the collision. Again, we sample N points to speed up visualization.

In [12]:
```python
N = 10000
sample = collisions.sample(N)

# Round / bin the latitude and longitudes
sample['lat_bin'] = np.round(sample['LATITUDE'], 3)
sample['lng_bin'] = np.round(sample['LONGITUDE'], 3)

# Average severity for regions
gby_cols = ['lat_bin', 'lng_bin']

coord_stats = (sample.groupby(gby_cols)
               .agg({'SEVERITY': 'mean'})
               .reset_index())

# Visualize the average severity per region
city_long_border = (-74.03, -73.75)
city_lat_border = (40.63, 40.85)
fig, ax = plt.subplots(ncols=1, nrows=1, figsize=(14, 10))

scatter_trips = ax.scatter(sample['LONGITUDE'].values,
                           sample['LATITUDE'].values,
                           color='grey', s=1, alpha=0.5)

scatter_cmap = ax.scatter(coord_stats['lng_bin'].values,
                          coord_stats['lat_bin'].values,
                          c=coord_stats['SEVERITY'].values,
                          cmap='viridis', s=10, alpha=0.9)

cbar = fig.colorbar(scatter_cmap)
cbar.set_label("Manhattan average severity")
ax.set_xlim(city_long_border)
ax.set_ylim(city_lat_border)
ax.set_xlabel('Longitude')
ax.set_ylabel('Latitude')
plt.title('Heatmap of Manhattan average severity')
plt.axis('off');
```

Heatmap of Manhattan average severity



## Question 1c

Do you think the location of the accident has a significant impact on the severity based on the visualization above? Additionally, identify something that could be improved in the plot above and describe how we could improve it.

```
In [13]: q1c_answer = r"""

         No, it seems that accident location does not significantly impact the av

         The dots on the plot seem to overlap and cover other dots, which can obs

         """

         # YOUR CODE HERE
         #raise NotImplementedError()

         print(q1c_answer)
```

No, it seems that accident location does not significantly impact the average severity of an accident, as the map seems to keep a relatively constant average severity.

The dots on the plot seem to overlap and cover other dots, which can o bscure data. Perhaps adding a transparency to each point will resolve this issue.

## Question 1d

Create a plot to visualize one or more features of the `collisions` table.

In [14]:
```
# YOUR CODE HERE
collisions1=collisions.copy()
collisions1['MONTH'] = [int(str(i)[5:7]) for i in collisions['DATETIME'

collisions1_mean = collisions1.groupby('MONTH').mean()
collisions1_mean = collisions1.reset_index()
ax=sns.barplot(x='MONTH', y='SEVERITY', data=collisions1_mean)
ax.set_title("Severity of Accidents per Month")
ax.set_xlabel("MONTH")
ax.set_ylabel('MEAN SEVERITY')

plt.show();
#raise NotImplementedError()
```



## Question 1e

Answer the following questions regarding your plot in 1d.

1. What feature you're visualization
2. Why you chose this feature
3. Why you chose this visualization method

```
In [15]: q1e_answer = r"""

I am visualizing the average severity of accidents per month. I decided

"""
# YOUR CODE HERE
#raise NotImplementedError()
print(q1e_answer)
```

I am visualizing the average severity of accidents per month. I decided to chose this feature, thinking that some months will have more severe accidents than others, since snow and other weather conditions would greatly impact severity. I used a bar plot since it best highlights the differences between mean severity among the different months by displaying a difference in length, and since months are not a continious distribution, I opted for a bar chart rather than a histogram.

## 2: Combining External Datasets

It seems like accident timing and location may influence the duration of a taxi ride. Let's start to join our NYC Taxi data with our collisions data.

Let's assume that an accident will influence traffic in the surrounding area for around 1 hour. Below, we create two columns, START and END :

- START : contains the recorded time of the accident
- END : 1 hours after START

**Note:** We chose 1 hour somewhat arbitrarily, feel free to experiment with other time intervals outside this notebook.

```
In [16]: collisions['START'] = collisions['DATETIME']
         collisions['END'] = collisions['START'] + pd.Timedelta(hours=1)
```

### Question 2a

Drop all of the columns besides the following: DATETIME , TIME , START , END , DATE , LATITUDE , LONGITUDE , SEVERITY . Feel free to experiment with other subsets outside of this notebook.

In [17]:
```
collisions_subset = collisions.loc[:,['DATETIME', 'TIME', 'START', 'END
# YOUR CODE HERE
#raise NotImplementedError()
collisions_subset.head(5)
```

Out[17]:

| UNIQUE KEY | DATETIME | TIME | START | END | DATE | LATITUDE | LONGITUDE | SEVERITY |
|---|---|---|---|---|---|---|---|---|
| 3589202 | 2016-12-29 00:00:00 | 0 | 2016-12-29 00:00:00 | 2016-12-29 01:00:00 | 2016-12-29 | 40.844107 | -73.897997 | 0 |
| 3587413 | 2016-12-26 14:30:00 | 14 | 2016-12-26 14:30:00 | 2016-12-26 15:30:00 | 2016-12-26 | 40.692347 | -73.881778 | 0 |
| 3578151 | 2016-11-30 22:50:00 | 22 | 2016-11-30 22:50:00 | 2016-11-30 23:50:00 | 2016-11-30 | 40.755480 | -73.741730 | 2 |
| 3567096 | 2016-11-23 20:11:00 | 20 | 2016-11-23 20:11:00 | 2016-11-23 21:11:00 | 2016-11-23 | 40.771122 | -73.869635 | 0 |
| 3565211 | 2016-11-21 14:11:00 | 14 | 2016-11-21 14:11:00 | 2016-11-21 15:11:00 | 2016-11-21 | 40.828918 | -73.838403 | 0 |

In [18]:
```
assert collisions_subset.shape == (116691, 8)
```

## Question 2b

Now, let's merge our `collisions_subset` table with `train_df`. Start by merging with only the date. We will filter by a time window in a later question.

We should be performing a left join, where our `train_df` is the left table. This is because we want to preserve all of the taxi rides in our end result. It happens that an inner join will also work, since both tables contain data on each date.

Note that the resulting `merged` table will have multiple rows for every taxi ride row in the original `train_df` table. For example, `merged` will have 483 rows with `index` equal to 16709, because there were 483 accidents that occurred on the same date as ride #16709.

Because of memory limitation, we will select the third week of 2016 to analyze. Feel free to change to it week 1 or 2 to see if the observation is general.

In [19]:
```
data_file = Path("./", "cleaned_data.hdf")
train_df = pd.read_hdf(data_file, "train")
train_df = train_df.reset_index()
train_df = train_df[['index', 'tpep_pickup_datetime', 'pickup_longitude
train_df['date'] = train_df['tpep_pickup_datetime'].dt.date
```

```
In [20]: collisions_subset = collisions_subset[collisions_subset['DATETIME'].dt.w
         train_df = train_df[train_df['tpep_pickup_datetime'].dt.weekofyear == 3
```

```
In [21]: # merge the dataframe here

         merged = train_df.merge(collisions_subset, left_on='date', right_on='DA

         # YOUR CODE HERE
         #raise NotImplementedError()

         merged.head()
```

Out[21]:

| | index | tpep_pickup_datetime | pickup_longitude | pickup_latitude | duration | date | DATETIME | T |
|---|---|---|---|---|---|---|---|---|
| **0** | 16709 | 2016-01-21 22:28:17 | -73.997986 | 40.741215 | 736.0 | 2016-01-21 | 2016-01-21 10:35:00 | |
| **1** | 16709 | 2016-01-21 22:28:17 | -73.997986 | 40.741215 | 736.0 | 2016-01-21 | 2016-01-21 13:20:00 | |
| **2** | 16709 | 2016-01-21 22:28:17 | -73.997986 | 40.741215 | 736.0 | 2016-01-21 | 2016-01-21 16:00:00 | |
| **3** | 16709 | 2016-01-21 22:28:17 | -73.997986 | 40.741215 | 736.0 | 2016-01-21 | 2016-01-21 18:30:00 | |
| **4** | 16709 | 2016-01-21 22:28:17 | -73.997986 | 40.741215 | 736.0 | 2016-01-21 | 2016-01-21 00:05:00 | |

```
In [22]: assert merged.shape == (1528162, 14)
```

## Question 2c

Now that our tables are merged, let's use temporal and spatial proximity to condition on the duration of the average length of a taxi ride. Let's operate under the following assumptions.

Accidents only influence the duration of a taxi ride if the following are satisfied:

1) The haversine distance between the the pickup location of the taxi ride and location of the recorded accident is within 5 (km). This is roughly 3.1 miles.

2) The start time of a taxi ride is within a 1 hour interval between the start and end of an accident.

Complete the code below to create an `'accident_close'` column in the `merged` table that indicates if an accident was close or not according to the assumptions above.

```
In [23]: def haversine(lat1, lng1, lat2, lng2):
             """
             Compute haversine distance
             """
             lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
             average_earth_radius = 6371
             lat = lat2 - lat1
             lng = lng2 - lng1
             d = np.sin(lat * 0.5) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(l
             h = 2 * average_earth_radius * np.arcsin(np.sqrt(d))
             return h

         def manhattan_distance(lat1, lng1, lat2, lng2):
             """
             Compute Manhattan distance
             """
             a = haversine(lat1, lng1, lat1, lng2)
             b = haversine(lat1, lng1, lat2, lng1)
             return a + b
```

```
In [26]: start_to_accident = haversine(merged['pickup_latitude'].values,
                                        merged['pickup_longitude'].values,
                                        merged['LATITUDE'].values,
                                        merged['LONGITUDE'].values)
         merged['start_to_accident'] = start_to_accident

         # initialze accident_close column to all 0 first
         merged['accident_close'] = 0

         # Boolean pd.Series to select the indices for which accident_close shou
         # (1) record's start_to_accident <= 5
         # (2) pick up time is between start and end
         is_accident_close = (merged['start_to_accident'] <= 5) & (merged['tpep_

         # YOUR CODE HERE
         #raise NotImplementedError()

         merged.loc[is_accident_close, 'accident_close'] = 1
         #for i in np.arange(len(is_accident_close)):
         #    if is_accident_close[i]:
         #        print(merged.iloc[[i]])
         #merged[is_accident_close]
```

```
In [27]: assert merged['accident_close'].sum() > 16000
```

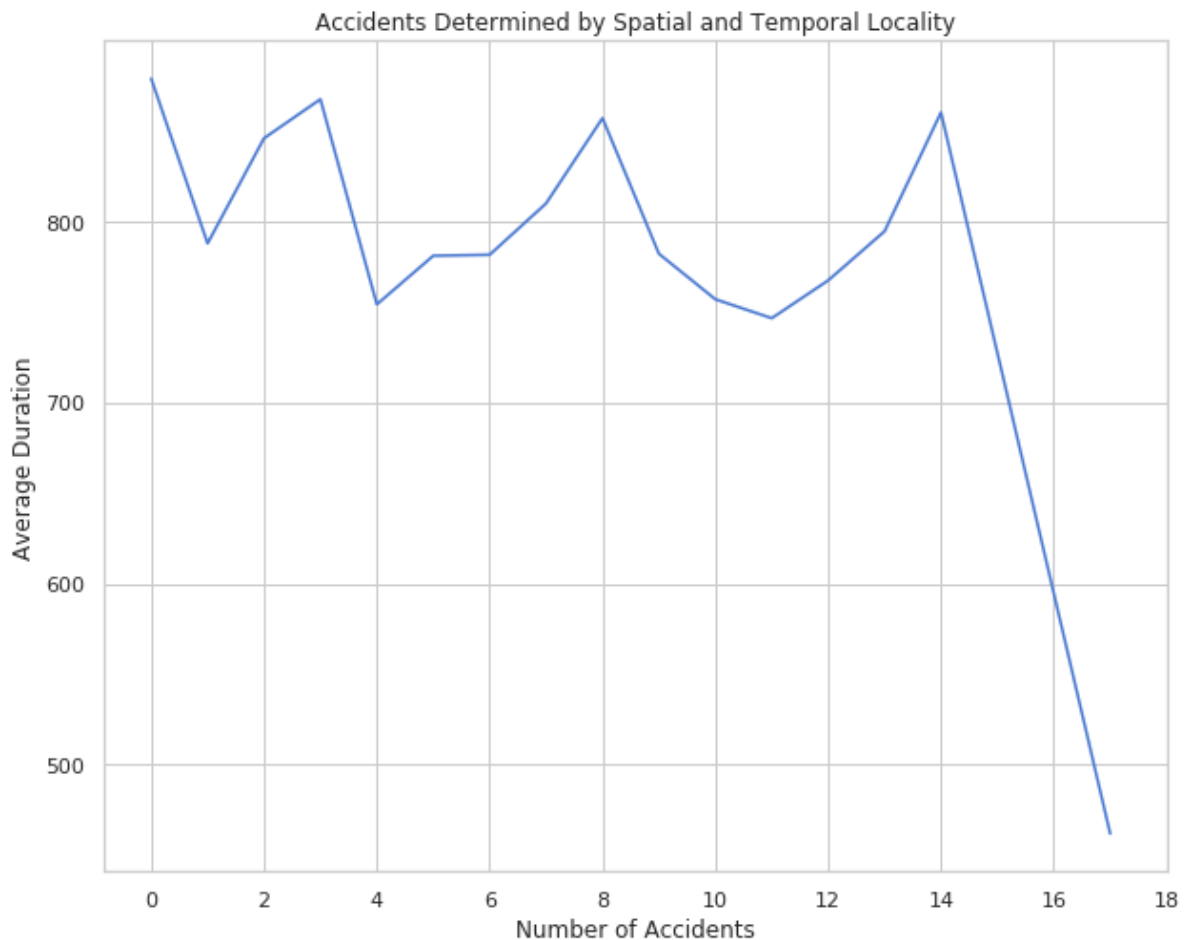The last step is to aggregate the total number of proximal accidents. We want to count the total
number of accidents that were close spatially and temporally and condition on that data.

The code below create a new data frame called `train_accidents`, which is a copy of
`train_df`, but with a new column that counts the number of accidents that were close
(spatially and temporally) to the pickup location/time.

In [28]:
```python
train_df = train_df.set_index('index')
num_accidents = merged.groupby(['index'])['accident_close'].sum().to_fra
train_accidents = train_df.copy()
train_accidents['num_accidents'] = num_accidents
```

Next, for each value of `num_accidents`, we plot the average `duration` of rides with that number of accidents.

In [29]:
```python
plt.figure(figsize=(10,8))
train_accidents.groupby('num_accidents')['duration'].mean().plot(xticks=
plt.title("Accidents Determined by Spatial and Temporal Locality")
plt.xlabel("Number of Accidents")
plt.ylabel("Average Duration")
plt.show();
```



It seems that using both spatial and temporal proximity doesn't give us much insight on if collisions increase taxi ride durations. Let's try conditioning on spatial proximity and temporal proximity separately and see if there are more interesting results there.

In [30]:
```python
# Temporal locality

# Condition on time
index = (((merged['tpep_pickup_datetime'] >= merged['START']) & \
          (merged['tpep_pickup_datetime'] <= merged['END']))))

# Count accidents
merged['accident_close'] = 0
merged.loc[index, 'accident_close'] = 1
num_accidents = merged.groupby(['index'])['accident_close'].sum().to_fra
train_accidents_temporal = train_df.copy()
train_accidents_temporal['num_accidents'] = num_accidents

# Plot
plt.figure(figsize=(10,8))
train_accidents_temporal.groupby('num_accidents')['duration'].mean().pl
plt.title("Accidents Determined by Temporal Locality")
plt.xlabel("Number of Accidents")
plt.ylabel("Average Duration")
plt.show();
```

In [31]:
```python
# Spatial locality

# Condition on space
index = (merged['start_to_accident'] <= 5)

# Count accidents
merged['accident_close'] = 0
merged.loc[index, 'accident_close'] = 1
num_accidents = merged.groupby(['index'])['accident_close'].sum().to_fra
train_accidents_spatial = train_df.copy()
train_accidents_spatial['num_accidents'] = num_accidents

# Plot
plt.figure(figsize=(10,8))
train_accidents_spatial.groupby('num_accidents')['duration'].mean().plot
plt.title("Accidents Determined by Spatial Locality")
plt.xlabel("Number of Accidents")
plt.ylabel("Average Duration")
plt.show();
```



## Question 2d

By conditioning on temporal and spatial proximity separately, we reveal different trends in average ride duration as a function of number of accidents nearby.

What can you say about the temporal and spatial proximity of accidents to taxi rides and the effect on ride duration? Think of a new hypothesis regarding accidents and taxi ride durations and explain how you would test it.

Additionally, comment on some of the assumptions being made when we condition on temporal and spatial proximity separately. What are the implications of only considering one and not the other?

In [32]:
```python
q2d_answer = r"""

For using just a temporal proximity of accidents, we find that that the

For using just a spatial proximity of accidents, we find that there is

Therefore, I predict that if there are many accidents during the time o

An assumption for the temporal proximity is that if there are more acci

An assumption for the spatial proximity is that if there are less than
"""

# YOUR CODE HERE
#raise NotImplementedError()

print(q2d_answer)
```

For using just a temporal proximity of accidents, we find that that th
ere is a slight positive correlation between the number of accidents o
f each ride and the average duration. This may be the case because if
there are more accidents taking place at a certain time, then there ar
e less roads avaliable to drive on, leading to more traffic and thus a
overall longer taxi ride.

For using just a spatial proximity of accidents, we find that there is
a strong spike for the average duration if the number of accidents dur
ing the day in that area less than ~20, then a relative small and cons
tant average duration for an area with 20+ accidents in the area. My t
heory for why this is occuring is that any taxi ride that had 20+ acci
dents within the vicinity of the pickup location was in Manhattan, and
most trips starting at Manhattan also end in Manhattan. However, any t
rips with <20 accidents in the vicinity of the pickup are from the air
ports or the area surrounding Manhattan, where these trips should be o
n average travelling farther and thus longer. This might be the case,
since there are bound to be more accidents within the cluttered street
s of Manhattan than outside of Manhattan.

Therefore, I predict that if there are many accidents during the time
of the pickup, the duration should also be longer. Also, if there are
many accidents (20+) within the vicinity of the pickup during the day,
then it is resonable to assume a significantly lower duration than if
there are less accidents (<20) in the area. I can test this by buildin
g a model based off of this hypothesis and off of the functions above,
test it on the validation group, and seeing my model predicts correctl
y or not.

An assumption for the temporal proximity is that if there are more acc
idents within a timeframe, there will be less roads to drive on, and t
hus more traffic overall. If we neglected this feature, we miss out on
overall traffic conditions during the time of the taxi ride.

An assumption for the spatial proximity is that if there are less than
20 accidents that occured in the vicinity of the pickup, the most like

ly the pickup is in Manhattan, and vice versa. If we neglected this fe
ature, we will not be able to approximate the location of the taxi rid
e and thus its average duration (unless we are clustering spatially th
e locations of each trip with the lat/long).

## Part 3 Exports

We are not requiring you to export anything from this notebook, but you may find it useful to do so. There is a space below for you to export anything you wish.

```
In [33]: Path("data/part3").mkdir(parents=True, exist_ok=True)
         data_file = Path("data/part3", "data_part3.hdf") # Path of hdf file
         ...
```

Out[33]: Ellipsis

```
In [34]: merged.to_hdf(data_file, "collisions")
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:1996: PerformanceWarning:
your performance may suffer as PyTables will pickle object types that
it cannot
map directly to c-types [inferred_type->date,key->block3_values] [item
s->['date', 'DATE']]

  return pytables.to_hdf(path_or_buf, key, self, **kwargs)
```

## Part 3 Conclusions

We merged the NYC Accidents dataset with our NYC Taxi dataset, conditioning on temporal and spatial locality. We explored potential features by visualizing the relationship between number of accidents and the average duration of a ride.

**Please proceed to part 4 where we will be engineering more features and building our models using a processing pipeline.**

## Submission

You're almost done!

Before submitting this assignment, ensure that you have:

1. Restarted the Kernel (in the menubar, select Kernel→Restart & Run All)
2. Validated the notebook by clicking the "Validate" button.

Then,

1. **Submit** the assignment via the Assignments tab in **Datahub**

2. **Upload and tag** the manually reviewed portions of the assignment on **Gradescope**

In [ ]:

Before you turn in the homework, make sure everything runs as expected. To do so, select **Kernel**→**Restart & Run All** in the toolbar above. Remember to submit both on **DataHub** and **Gradescope**.

Please fill in your name and include a list of your collaborators below.

```
In [7]: NAME = "William Sheu"
        COLLABORATORS = ""
```

# Project 2: NYC Taxi Rides

## Part 4: Feature Engineering and Model Fitting

In this final part of the project, you will finally build a regression model that attempts to predict the duration of a taxi ride from all other available information.

You will build this model using a processing pipeline and submit your results to Kaggle. We will first walk you through a generic example using the data we saved from Part 1. Please carefully follow these steps as you will need to repeat this for your final model. After, we give you free reign and let you decide how you want to define your final model.

```
In [8]: import os
        import pandas as pd
        import numpy as np
        import sklearn.linear_model as lm
        import matplotlib.pyplot as plt
        import seaborn as sns
        from pathlib import Path
        from sqlalchemy import create_engine
        from sklearn.model_selection import cross_val_score, train_test_split,

        sns.set(style="whitegrid", palette="muted")

        plt.rcParams['figure.figsize'] = (12, 9)
        plt.rcParams['font.size'] = 12

        %matplotlib inline
```

### Training and Validation

The following code loads the training and validation data from part 1 into a Pandas DataFrame.

In [9]:
```python
# Run this cell to load the data.
data_file = Path("./", "cleaned_data.hdf")
train_df = pd.read_hdf(data_file, "train")
val_df = pd.read_hdf(data_file, "val")
```

## Testing

Here we load our testing data on which we will evaluate your model.

In [10]:
```python
test_df = pd.read_csv("./proj2_test_data.csv")
test_df['tpep_pickup_datetime'] = pd.to_datetime(test_df['tpep_pickup_d
test_df.head()
```

Out[10]:

| | record_id | VendorID | tpep_pickup_datetime | passenger_count | trip_distance | pickup_longitude |
|---|---|---|---|---|---|---|
| **0** | 10000 | 1 | 2016-01-02 01:45:37 | 1 | 1.20 | -73.982224 |
| **1** | 19000 | 2 | 2016-01-02 03:05:16 | 1 | 10.90 | -73.999977 |
| **2** | 21000 | 1 | 2016-01-02 03:24:36 | 1 | 1.80 | -73.986618 |
| **3** | 23000 | 2 | 2016-01-02 03:47:38 | 1 | 5.95 | -74.002922 |
| **4** | 27000 | 1 | 2016-01-02 04:36:44 | 1 | 1.60 | -73.986366 |

In [11]:
```python
test_df.describe()
```

Out[11]:

| | record_id | VendorID | passenger_count | trip_distance | pickup_longitude | pickup_lat |
|---|---|---|---|---|---|---|
| **count** | 1.377400e+04 | 13774.000000 | 13774.000000 | 13774.000000 | 13774.000000 | 13774.00 |
| **mean** | 3.465950e+07 | 1.536082 | 1.663642 | 2.954688 | -72.953619 | 40.18 |
| **std** | 2.015133e+07 | 0.498714 | 1.311739 | 3.704427 | 8.628431 | 4.75 |
| **min** | 1.000000e+04 | 1.000000 | 0.000000 | 0.000000 | -77.039436 | 0.00 |
| **25%** | 1.719975e+07 | 1.000000 | 1.000000 | 1.000000 | -73.992058 | 40.73 |
| **50%** | 3.457400e+07 | 2.000000 | 1.000000 | 1.700000 | -73.981846 | 40.75 |
| **75%** | 5.216875e+07 | 2.000000 | 2.000000 | 3.157500 | -73.967119 | 40.76 |
| **max** | 6.940400e+07 | 2.000000 | 6.000000 | 104.800000 | 0.000000 | 40.86 |

# Modeling

We've finally gotten to a point where we can specify a simple model. Remember that we will be fitting our model on the training set we created in part 1. We will use our validation set to evaluate how well our model might perform on future data.

### Reusable Pipeline

Throughout this assignment, you should notice that your data flows through a single processing pipeline several times. From a software engineering perspective, this should be sufficient motivation to abstract parts of our code into reusable functions/methods. We will now encapsulate our entire pipeline into a single function `process_data_gm` . gm is shorthand for "guided model".

In [12]:
```python
# Copied from part 2
def haversine(lat1, lng1, lat2, lng2):
    """
    Compute haversine distance
    """
    lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
    average_earth_radius = 6371
    lat = lat2 - lat1
    lng = lng2 - lng1
    d = np.sin(lat * 0.5) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(l
    h = 2 * average_earth_radius * np.arcsin(np.sqrt(d))
    return h

# Copied from part 2
def manhattan_distance(lat1, lng1, lat2, lng2):
    """
    Compute Manhattan distance
    """
    a = haversine(lat1, lng1, lat1, lng2)
    b = haversine(lat1, lng1, lat2, lng1)
    return a + b

# Copied from part 2
def bearing(lat1, lng1, lat2, lng2):
    """
    Compute the bearing, or angle, from (lat1, lng1) to (lat2, lng2).
    A bearing of 0 refers to a NORTH orientation.
    """
    lng_delta_rad = np.radians(lng2 - lng1)
    lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
    y = np.sin(lng_delta_rad) * np.cos(lat2)
    x = np.cos(lat1) * np.sin(lat2) - np.sin(lat1) * np.cos(lat2) * np.
    return np.degrees(np.arctan2(y, x))

# Copied from part 2
def add_time_columns(df):
    """
    Add temporal features to df
    """
    df.is_copy = False # propogate write to original dataframe
    df.loc[:, 'month'] = df['tpep_pickup_datetime'].dt.month
    df.loc[:, 'week_of_year'] = df['tpep_pickup_datetime'].dt.weekofyea
    df.loc[:, 'day_of_month'] = df['tpep_pickup_datetime'].dt.day
    df.loc[:, 'day_of_week'] = df['tpep_pickup_datetime'].dt.dayofweek
    df.loc[:, 'hour'] = df['tpep_pickup_datetime'].dt.hour
    df.loc[:, 'week_hour'] = df['tpep_pickup_datetime'].dt.weekday * 24
    return df

# Copied from part 2
def add_distance_columns(df):
    """
    Add distance features to df
    """
    df.is_copy = False # propogate write to original dataframe
    df.loc[:, 'manhattan'] = manhattan_distance(lat1=df['pickup_latitud
                                        lng1=df['pickup_longitu
```

```
                                        lat2=df['dropoff_latitu(
                                        lng2=df['dropoff_longit(

        df.loc[:, 'bearing'] = bearing(lat1=df['pickup_latitude'],
                                       lng1=df['pickup_longitude'],
                                       lat2=df['dropoff_latitude'],
                                       lng2=df['dropoff_longitude'])
        df.loc[:, 'haversine'] = haversine(lat1=df['pickup_latitude'],
                                       lng1=df['pickup_longitude'],
                                       lat2=df['dropoff_latitude'],
                                       lng2=df['dropoff_longitude'])
        return df

    def select_columns(data, *columns):
        return data.loc[:, columns]
```

```
In [13]:  def process_data_gm1(data, test=False):
              X = (
                  data

                  # Transform data
                  .pipe(add_time_columns)
                  .pipe(add_distance_columns)

                  .pipe(select_columns,
                        'pickup_longitude',
                        'pickup_latitude',
                        'dropoff_longitude',
                        'dropoff_latitude',
                        'manhattan',
                      )
              )
              if test:
                  y = None
              else:
                  y = data['duration']

              return X, y
```

We will use our pipeline defined above to pre-process our training and test data in exactly the same way. Our functions make this relatively easy to do!

In [14]:
```python
# Train
X_train, y_train = process_data_gm1(train_df)
X_val, y_val = process_data_gm1(val_df)
guided_model_1 = lm.LinearRegression(fit_intercept=True)
guided_model_1.fit(X_train, y_train)

# Predict
y_train_pred = guided_model_1.predict(X_train)
y_val_pred = guided_model_1.predict(X_val)
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  return object.__setattr__(self, name, value)
```

Here, `y_val` are the correct durations for each ride, and `y_val_pred` are the predicted durations based on the 7 features above ( `vendorID` , `passenger_count` , `pickup_longitude` , `pickup_latitude` , `dropoff_longitude` , `dropoff_latitude` , `manhattan` ).

In [15]:
```python
assert 600 <= np.median(y_train_pred) <= 700
assert 600 <= np.median(y_val_pred) <= 700
```

The resulting model really is a linear model just like we saw in class, i.e. the predictions are simply generated by the product $\Phi\theta$. For example, the line of code below generates a prediction for $x_1$ by computing $\phi_1^T\theta$. Here `guided_model_1.coef_` is $\theta$ and `X_train.iloc[0, :]` is $\phi_1$.

Note that unlike in class, here the dummy intercept term is not included in $\Phi$.

In [16]:
```python
X_train.iloc[0, :].dot(guided_model_1.coef_) + guided_model_1.intercept
```

Out[16]: 558.751330511368

We see that this prediction is exactly the same (except for possible floating point error) as generated by the `predict` function, which simply computes the product $\Phi\theta$, yielding predictions for every input.

In [17]:
```python
y_train_pred[0]
```

Out[17]: 558.75133051135344

In this assignment, we will use Mean Absolute Error (MAE), a.k.a. mean L1 loss, to measure the quality of our models. As a reminder, this quantity is defined as:

$$MAE = \frac{1}{n}\sum_i |y_i - \hat{y}_i|$$

Why may we want to use the MAE as a metric, as opposed to Mean Squared Error (MSE)? Using our domain knowledge that most rides are short in duration (median is roughly 600 seconds), we know that MSE is susceptible to outliers. Given that some of the outliers in our dataset are quite extreme, it is probably better to optimize for the majority of rides rather than for the outliers. You may want to remove some of these outliers later on.

In [18]:
```python
def mae(actual, predicted):
    """
    Calculates MAE from actual and predicted values
    Input:
      actual (1D array-like): vector of actual values
      predicted (1D array-like): vector of predicted/fitted values
    Output:
      a float, the MAE
    """

    mae = np.mean(np.abs(actual - predicted))
    return mae
```

In [19]:
```python
assert 200 <= mae(y_val_pred, y_val) <= 300
print("Validation Error: ", mae(y_val_pred, y_val))
```

Validation Error:  266.136130855

Side note: scikit-learn also has tools to compute mean absolute error ( `sklearn.metrics.mean_absolute_error` ). In fact, most metrics that we have discussed in this class can be found as part of the `sklearn.metrics` module (https://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics). Some of these may come in handy as part of your feature engineering!

# Visualizing Error

You should be getting between 200 and 300 MAE, which means your model was off by roughly 3-5 minutes on trips of average length 12 minutes. This is fairly decent performance given that our basic model uses only using the pickup/dropoff latitude and manhattan distance of the trip. 3-5 minutes may seem like a lot for a trip of 12 minutes, but keep in mind that this is the *average* error. This metric is susceptible to extreme outliers, which exist in our dataset.

Now we will visualize the residual for the validation set. We will plot the following:

1. Distribution of residuals
2. Average residual grouping by ride duration

```
In [20]:  # Distribution of residuals
          plt.figure(figsize=(8,4))
          sns.distplot(np.abs(y_val - y_val_pred))
          plt.xlabel('residual')
          plt.title('distribution of residuals');
```
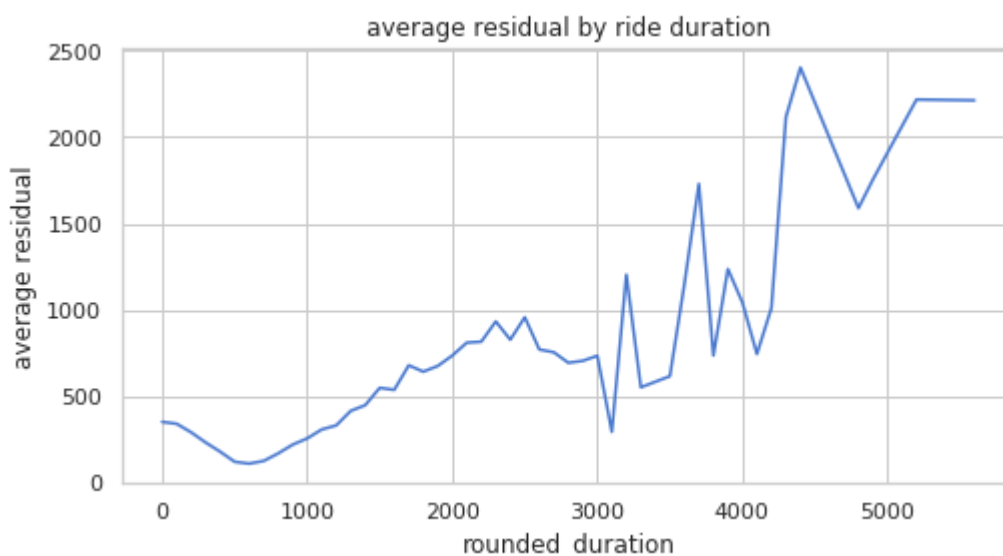


distribution of residuals

```
In [21]:  # Average residual grouping by ride duration
          val_residual = X_val.copy()
          val_residual['duration'] = y_val
          val_residual['rounded_duration'] = np.around(y_val, -2)
          val_residual['residual'] = np.abs(y_val - y_val_pred)
          tmp = val_residual.groupby('rounded_duration').mean()
          plt.figure(figsize=(8,4))
          tmp['residual'].plot()
          plt.ylabel('average residual')
          plt.title('average residual by ride duration');
```



average residual by ride duration

In the first visualization, we see that most of the residuals are centered around 250 seconds ~ 4 minutes. There is a minor right tail, suggesting that we are still unable to accurately fit some outliers in our data. The second visualization also suggests this, as we see the average residual

increasing as a somewhat linear function of duration. But given that our average ride duration is roughly 600-700 seconds, it seems that we are indeed optimizing for the average ride because the residuals are smallest around 600-700.

Keep this in mind when creating your final model! Visualizing the error is a powerful tool and may help diagnose shortcomings of your model. Let's go ahead and submit to kaggle, although your error on the test set may be higher than 300.

# Submission to Kaggle

The following code will write your predictions on the test dataset to a CSV, which you can submit to Kaggle. You may need to modify it to suit your needs, but we recommend you make a copy and preserve the original function.

Remember that if you've performed transformations or featurization on the training data, you must also perform the same transformations on the test data in order to make predictions. For example, if you've created features for the columns `pickup_datetime` or `pickup_latitude` on the training data, you must also extract the same features in order to use scikit-learn's `.predict(...)` method.

In [22]:
```python
from datetime import datetime
def generate_submission(test, predictions, force=False):
    if force:
        if not os.path.isdir("submissions"):
            os.mkdir("submissions")
        submission_df = pd.DataFrame({
            "id": test_df.index.values,
            "duration": predictions,
        },
            columns=['id', 'duration'])

        timestamp = datetime.isoformat(datetime.now()).split(".")[0]

        submission_df.to_csv(f'submissions/submission_{timestamp}.csv',

        print(f'Created a CSV file: submission_{timestamp}.csv')
        print('You may now upload this CSV file to Kaggle for scoring.'
```

In [23]:
```python
X_test, _ = process_data_gm1(test_df, True)
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  return object.__setattr__(self, name, value)
```

```
In [24]: assert list(X_train.columns) == list(X_test.columns), "Different columns
         submission_predictions = (guided_model_1
                                     .fit(X_train, y_train)
                                     .predict(X_test))
         submission_predictions = submission_predictions.astype(int)
         submission_predictions[submission_predictions < 0] = 0
         generate_submission(test_df, submission_predictions, True)
```

```
Created a CSV file: submission_2018-12-05T10:15:06.csv
You may now upload this CSV file to Kaggle for scoring.
```

```
In [25]: # Check your submission
         assert isinstance(submission_predictions, np.ndarray), "Submission not a
         assert all(submission_predictions >= 0), "Duration must be non-negative
         assert issubclass(submission_predictions.dtype.type, np.integer), "Secor
```

## Your Turn!

Now it's your turn! Draw upon everything you have learned this semester to find the best features to help your model accurately predict the duration of a taxi ride.

You may use whatever method you prefer in order to create features. You may use features that we created and features that you discovered yourself from any of the 2 datasets. However, we want to make it fair to students who are seeing these techniques for the first time. As such, you are only allowed regression models and their regularized forms. This means no random forest, k-nearest-neighbors, neural nets, etc.

**Here are some ideas to improve your model:**

- **Data selection**: January 2016 was an odd month for taxi rides due to the blizzard. Would it help to select training data differently?
- **Data cleaning**: Try cleaning your data in different ways. In particular, consider how to handle outliers.
- **Better features**: Explore the 2 datasets and find what features are most helpful. Utilize external datasets to improve your accuracy.
- **Regularization**: Try different forms of regularization to avoid fitting to the training set. Recall that `Ridge` and `Lasso` are the names of the classes in `sklearn.linear_model` that combine `LinearRegression` with regularization techniques.
- **Model selection**: You can adjust parameters of your model (e.g., the regularization parameter) to achieve higher accuracy. GridSearchCV (http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html) may be helpful.
- **Validation**: Recall that you should use cross-validation to do feature and model selection properly! Otherwise, you will likely overfit to your training data.

There's many things you could try that could help your model. We have only suggested a few. Be creative and innovative! Please use `proj2_extras.ipynb` for all of your extraneous work. Note that you will be submitting `proj2_extras.ipynb` and we will be grading it. Please properly comment and format this notebook!

Once you are satisfied with your results, answer the questions in the Deliverables section. You may want to read this section in advance so you have an idea of what we're looking for.

# Deliverables

# Feature/Model Selection Process

Let's first look at selection of better features. In this following cell, describe the process of choosing good features to improve your model. You should use at least 3-4 sentences each to address the follow questions. Backup your responses with graphs supporting your claim (you can save figures and load them, no need to add the plotting code here). Use these questions to concisely summarize all of your extra work!

## Question 1a

How did you find better features for your model?

```
In [26]:  q1a_answer = r"""

I started by looking at the train_df and seeing which columns would be u

"""
          print(q1a_answer)
          # YOUR CODE HERE
          #raise NotImplementedError()
```

```
I started by looking at the train_df and seeing which columns would be
useful in determining duration. I also tried combining multiple column
s together, thinking that the resulting column would be useful in dete
rmining duration.
```

## Question 1b

What did you try that worked / didn't work?

In [27]:
```python
q1b_answer = r"""

I decided to try out using fare_amount, since it would make sense for d

"""
print(q1b_answer)
# YOUR CODE HERE
#raise NotImplementedError()
```

I decided to try out using fare_amount, since it would make sense for
duration to scale positively with the fare_amount. This worked out ver
y well, and reduced my error significantly. Additionally, I added the
hour as a feature, as at some hours, there is consistently more traffi
c than other hours. I also added month as a feature, since I thought t
he month would greatly affect duration of a ride. Both of these featur
es slightly reduced my error.

## Question 1c

What was surprising in your search for good features?

In [28]:
```python
q1c_answer = r"""

I was expecting the fare_amount to be a great feature, but it exceeded

"""
print(q1c_answer)
# YOUR CODE HERE
#raise NotImplementedError()
```

I was expecting the fare_amount to be a great feature, but it exceeded
my expectations, as it more than halved my error. I also added month a
s a feature, which slightly lowered my error, and I had to import 5 ad
ditional months worth of data.

## Question 2

Just as in the guided model above, you should encapsulate as much of your workflow into
functions as possible. Define `process_data_fm` and `final model` in the cell below. In
order to calculate your final model's MAE, we will run the code in the cell after that.

**Note:** You *MUST* name the model you wish to be evaluated on `final_model`. This is what we will be using to generate your predictions. We will take the state of `final_model` right after executing the cell below and run the following code:

```
# Load in test_df, solutions
X_test, _ = process_data_fm(test_df, True)
submission_predictions = final_model.predict(X_test)
# Generate score for autograding
```

We encourage you to conduct all of your exploratory work in `proj2_extras.ipynb`, which will be graded for 10 points.

```python
In [29]:  train_extra_df = pd.read_hdf(Path("data", "data_extra.hdf"), "train_ext

          AVERAGE_LONG_DROP=np.mean(train_extra_df['dropoff_longitude'])
          AVERAGE_LAT_DROP=np.mean(train_extra_df['dropoff_latitude'])
          AVERAGE_LONG_PICK=np.mean(train_extra_df['pickup_longitude'])
          AVERAGE_LAT_PICK=np.mean(train_extra_df['pickup_latitude'])

          def process_data_fm(data, test=False):
              # Put your final pipeline here
              if test:
                  data1 = data.copy()
              else:
                  data1 = data.copy()[(data['duration'] < 10000)]
              if test:
                  data1.loc[(data1['dropoff_longitude'] > -70), 'dropoff_longitude
                  data1.loc[(data1['dropoff_latitude'] < 35), 'dropoff_latitude']
                  data1.loc[(data1['pickup_longitude'] > -70), 'pickup_longitude'
                  data1.loc[(data1['pickup_latitude'] < 35), 'pickup_latitude'] =
              X = (data1.pipe(add_time_columns).pipe(add_distance_columns)
                   .pipe(select_columns,
                         'pickup_longitude',
                         'pickup_latitude',
                         'dropoff_longitude',
                         'dropoff_latitude',
                         'manhattan',
                         'fare_amount',
                         'hour',
                         'month'
                        )
                  )
              if test:
                  y = None
              else:
                  y = data1['duration']

              return X, y

          X_train2, y_train2 = process_data_fm(train_extra_df)
          final_model = lm.LinearRegression(fit_intercept=True)
          final_model.fit(X_train2, y_train2)

          # YOUR CODE HERE
          #raise NotImplementedError()
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  return object.__setattr__(self, name, value)
```

```
Out[29]:  LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=
          False)
```

In [30]:
```python
# Feel free to change this cell
X_test, _ = process_data_fm(test_df, True)
final_predictions = final_model.predict(X_test)
final_predictions = final_predictions.astype(int)
generate_submission(test_df, final_predictions, False) # Change to true
print(final_predictions)
```

```
[ 409 1861  544 ...,  964  632  437]

/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  return object.__setattr__(self, name, value)
```

## Question 3

The following hidden cells will test your model on the test set. Please do not delete any of them if you want credit!

In [31]:
```python
# NO TOUCH
```

In [32]:
```python
# NOH
```

In [33]:
```python
# STAHP
```

In [34]:
```python
# NO MOLESTE
```

In [35]:
```python
# VA-T'EN
```

In [36]:
```python
# NEIN
```

In [37]:
```python
# PLSNO
```

In [38]:
```python
# THIS SPACE IS NOT YOURS
```

In [39]:
```python
# TAWDEETAW
```

In [40]:
```python
# MAU LEN
```

In [41]:
```python
# ALMOST
```

In [42]:
```python
# TO
```

In [43]:
```python
# THE
```

```
In [44]: # END
```

```
In [45]: # Hmph
```

```
In [46]: # Good riddance
```

```
In [47]: generate_submission(test_df, submission_predictions, True)
```

```
Created a CSV file: submission_2018-12-05T10:15:07.csv
You may now upload this CSV file to Kaggle for scoring.
```

This should be the format of your CSV file.
Unix-users can verify it running `!head submission_{datetime}.csv` in a jupyter notebook cell.

```
id,duration
id3004672,965.3950873305439
id3505355,1375.0665915134596
id1217141,963.2285454171943
id2150126,1134.7680929570924
id1598245,878.5495792656438
id0668992,831.6700312449248
id1765014,993.1692116960185
id0898117,1091.1171629594755
id3905224,887.9037911118357
```

Kaggle link: https://www.kaggle.com/t/f8b3c6acc3a045cab152060a5bc79670 (https://www.kaggle.com/t/f8b3c6acc3a045cab152060a5bc79670)

# Submission

You're almost done!

Before submitting this assignment, ensure that you have:

1. Restarted the Kernel (in the menubar, select Kernel→Restart & Run All)
2. Validated the notebook by clicking the "Validate" button.

Then,

1. **Submit** the assignment via the Assignments tab in **Datahub**
2. **Upload and tag** the manually reviewed portions of the assignment on **Gradescope**

```
In [ ]:
```

Before you turn in the homework, make sure everything runs as expected. To do so, select **Kernel→Restart & Run All** in the toolbar above. Remember to submit both on **DataHub** and **Gradescope**.

Please fill in your name and include a list of your collaborators below.

In [1]:
```python
NAME = "William Sheu"
COLLABORATORS = ""
```

# Project 2: NYC Taxi Rides

## Extras

Put all of your extra work in here. Feel free to save figures to use when completing Part 4.

In [2]:
```python
import os
import pandas as pd
import numpy as np
import sklearn.linear_model as lm
import matplotlib.pyplot as plt
import seaborn as sns
import zipfile
from utils import fetch_and_cache
from pathlib import Path
from sqlalchemy import create_engine
from datetime import datetime
from sklearn.model_selection import cross_val_score, train_test_split,
```

```
In [3]: test_df = pd.read_csv("./proj2_test_data.csv")
        DB_URI = "sqlite:////srv/db/taxi_2016_student_small.sqlite"
        TABLE_NAME = "taxi"
        query = """    SELECT *
            FROM (
            SELECT *
            FROM (
        SELECT *, julianday(tpep_dropoff_datetime) - julianday(tpep_pickup_date
        FROM (
                    SELECT *
                    FROM taxi
                    WHERE tpep_pickup_datetime
                        BETWEEN '2016-01-01' AND '2016-07-01'
                        AND record_id % 100 == 0
                    ORDER BY tpep_pickup_datetime
                    )
        WHERE duration < 0.1157407
                    )
            WHERE (
                    pickup_longitude <= -73.75 AND
                    pickup_longitude >= -74.03 AND
                    dropoff_longitude <= -73.75 AND
                    dropoff_longitude >= -74.03 AND
                    pickup_latitude <= 40.85 AND
                    pickup_latitude >= 40.63 AND
                    dropoff_latitude <= 40.85 AND
                    dropoff_latitude >= 40.63
                    )
                    )
            WHERE (passenger_count > 0)"""
        sql_engine = create_engine(DB_URI)
        processed_df = pd.read_sql_query(query, sql_engine)
        processed_df['tpep_pickup_datetime'] = pd.to_datetime(processed_df['tpe
        processed_df['tpep_dropoff_datetime'] = pd.to_datetime(processed_df['tp
        processed_df['duration'] = processed_df['duration']*86400
```

In [4]:
```python
# Copied from part 2
def haversine(lat1, lng1, lat2, lng2):
    """
    Compute haversine distance
    """
    lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
    average_earth_radius = 6371
    lat = lat2 - lat1
    lng = lng2 - lng1
    d = np.sin(lat * 0.5) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(l
    h = 2 * average_earth_radius * np.arcsin(np.sqrt(d))
    return h

# Copied from part 2
def manhattan_distance(lat1, lng1, lat2, lng2):
    """
    Compute Manhattan distance
    """
    a = haversine(lat1, lng1, lat1, lng2)
    b = haversine(lat1, lng1, lat2, lng1)
    return a + b

# Copied from part 2
def bearing(lat1, lng1, lat2, lng2):
    """
    Compute the bearing, or angle, from (lat1, lng1) to (lat2, lng2).
    A bearing of 0 refers to a NORTH orientation.
    """
    lng_delta_rad = np.radians(lng2 - lng1)
    lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
    y = np.sin(lng_delta_rad) * np.cos(lat2)
    x = np.cos(lat1) * np.sin(lat2) - np.sin(lat1) * np.cos(lat2) * np.
    return np.degrees(np.arctan2(y, x))

# Copied from part 2
def add_time_columns(df):
    """
    Add temporal features to df
    """
    df.is_copy = False # propogate write to original dataframe
    df.loc[:, 'month'] = df['tpep_pickup_datetime'].dt.month
    df.loc[:, 'week_of_year'] = df['tpep_pickup_datetime'].dt.weekofyea
    df.loc[:, 'day_of_month'] = df['tpep_pickup_datetime'].dt.day
    df.loc[:, 'day_of_week'] = df['tpep_pickup_datetime'].dt.dayofweek
    df.loc[:, 'hour'] = df['tpep_pickup_datetime'].dt.hour
    df.loc[:, 'week_hour'] = df['tpep_pickup_datetime'].dt.weekday * 24
    return df

# Copied from part 2
def add_distance_columns(df):
    """
    Add distance features to df
    """
    df.is_copy = False # propogate write to original dataframe
    df.loc[:, 'manhattan'] = manhattan_distance(lat1=df['pickup_latitud
                                    lng1=df['pickup_longitu
```

```
                                            lat2=df['dropoff_latitud
                                            lng2=df['dropoff_longitu

    df.loc[:, 'bearing'] = bearing(lat1=df['pickup_latitude'],
                                   lng1=df['pickup_longitude'],
                                   lat2=df['dropoff_latitude'],
                                   lng2=df['dropoff_longitude'])
    df.loc[:, 'haversine'] = haversine(lat1=df['pickup_latitude'],
                                       lng1=df['pickup_longitude'],
                                       lat2=df['dropoff_latitude'],
                                       lng2=df['dropoff_longitude'])
    return df

def select_columns(data, *columns):
    return data.loc[:, columns]

def mae(actual, predicted):
    """
    Calculates MAE from actual and predicted values
    Input:
      actual (1D array-like): vector of actual values
      predicted (1D array-like): vector of predicted/fitted values
    Output:
      a float, the MAE
    """
    mae = np.mean(np.abs(actual - predicted))
    return mae

def generate_submission(test, predictions, force=False):
    if force:
        if not os.path.isdir("submissions"):
            os.mkdir("submissions")
        submission_df = pd.DataFrame({
            "id": test_df.index.values,
            "duration": predictions,
        },
            columns=['id', 'duration'])

        timestamp = datetime.isoformat(datetime.now()).split(".")[0]

        submission_df.to_csv(f'submissions/submission_{timestamp}.csv',

        print(f'Created a CSV file: submission_{timestamp}.csv')
        print('You may now upload this CSV file to Kaggle for scoring.'
```

```
In [5]: train_df, val_df = train_test_split(processed_df, test_size=0.2, random

        AVERAGE_LONG_DROP=np.mean(train_df['dropoff_longitude'])
        AVERAGE_LAT_DROP=np.mean(train_df['dropoff_latitude'])
        AVERAGE_LONG_PICK=np.mean(train_df['pickup_longitude'])
        AVERAGE_LAT_PICK=np.mean(train_df['pickup_latitude'])

        def process_data_gm2(data, test=False):
            if test:
                data1 = data.copy()
            else:
                data1 = data.copy()[(data['duration'] < 10000)]
            if test:
                data1.loc[(data1['dropoff_longitude'] > -70), 'dropoff_longitud
                data1.loc[(data1['dropoff_latitude'] < 35), 'dropoff_latitude']
                data1.loc[(data1['pickup_longitude'] > -70), 'pickup_longitude'
                data1.loc[(data1['pickup_latitude'] < 35), 'pickup_latitude'] =
            X = (data1.pipe(add_time_columns).pipe(add_distance_columns)
                     .pipe(select_columns,
                             'pickup_longitude',
                             'pickup_latitude',
                             'dropoff_longitude',
                             'dropoff_latitude',
                             'manhattan',
                             'fare_amount',
                             'hour',
                             'month'
                          )
                )
            if test:
                y = None
            else:
                y = data1['duration']

            return X, y
```

```
In [6]: X_train, y_train = process_data_gm2(train_df)
        X_val, y_val = process_data_gm2(val_df)
        guided_model_2 = lm.LinearRegression(fit_intercept=True)
        guided_model_2.fit(X_train, y_train)
        y_val_pred = guided_model_2.predict(X_val)
        print(mae(y_val_pred, y_val))
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  return object.__setattr__(self, name, value)

175.869207619
```

In [7]:
```python
print(guided_model_2.coef_)
train_df.head(10)
```

```
[-1466.7188094      923.73865708  -1740.60100377    107.43336519       2.502
18244
     61.37640329       2.88886216     10.32365152]
```

Out[7]:

| | record_id | VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip |
|---|---|---|---|---|---|---|
| 30390 | 15892900 | 1 | 2016-02-10 14:46:02 | 2016-02-10 15:13:52 | 1 | |
| 124167 | 59546000 | 2 | 2016-06-03 22:19:02 | 2016-06-03 22:25:05 | 1 | |
| 65114 | 29985300 | 2 | 2016-03-23 10:50:48 | 2016-03-23 10:59:22 | 1 | |
| 113090 | 54223000 | 2 | 2016-05-20 13:01:23 | 2016-05-20 13:08:20 | 1 | |
| 92178 | 44175100 | 1 | 2016-04-25 08:39:56 | 2016-04-25 08:42:36 | 1 | |
| 43260 | 19244300 | 1 | 2016-02-26 05:53:23 | 2016-02-26 05:55:53 | 1 | |
| 98668 | 47420300 | 1 | 2016-05-03 08:53:52 | 2016-05-03 09:05:52 | 1 | |
| 61205 | 27975300 | 1 | 2016-03-18 14:01:48 | 2016-03-18 14:09:26 | 1 | |
| 135630 | 64839900 | 1 | 2016-06-18 11:17:10 | 2016-06-18 11:34:00 | 1 | |
| 91719 | 43943300 | 1 | 2016-04-24 14:09:57 | 2016-04-24 14:28:49 | 3 | |

10 rows × 21 columns

In [8]:
```python
test_df['tpep_pickup_datetime'] = pd.to_datetime(test_df['tpep_pickup_da
test_df = test_df.pipe(add_distance_columns)
test_df[test_df['manhattan'] == 0].loc[:,['pickup_longitude', 'pickup_l
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/gener
ic.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will
be removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/gener
ic.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will
be removed in a future version.
  return object.__setattr__(self, name, value)
```

In [10]:
```python
X_test, _ = process_data_gm2(test_df, True)
final_predictions = guided_model_2.predict(X_test)
final_predictions = final_predictions.astype(int)
generate_submission(test_df, final_predictions, False) # Change to true
final_predictions
```

Created a CSV file: submission_2018-12-04T23:52:45.csv
You may now upload this CSV file to Kaggle for scoring.

/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generi
c.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will b
e removed in a future version.
  return object.__setattr__(self, name, value)

Out[10]: array([ 409, 1861,  544, ...,  964,  632,  437])

In [12]:
```python
data_file = Path("data", "data_extra.hdf") # Path of hdf file
train_df.to_hdf(data_file, "train_extra_df")
```

## Submission

You're almost done!

Before submitting this assignment, ensure that you have:

1. Restarted the Kernel (in the menubar, select Kernel→Restart & Run All)
2. Validated the notebook by clicking the "Validate" button.

Then,

1. **Submit** the assignment via the Assignments tab in **Datahub**
2. **Upload and tag** the manually reviewed portions of the assignment on **Gradescope**