

COMS4060A/7056A: Assignment #2

Tim Bristow
tim@bristow.za.net

University of the Witwatersrand — May 26, 2022

Introduction

This assignment is based on content covering geospatial and time series data. You will be required to perform some basic data cleaning and exploration techniques on prescribed datasets.

The aim is to explore the dataset and make observations. There is no strict requirement on the format of your submission, but any answers should be reasoned, discussed, and relevant data or results should be provided to substantiate this. You might like to submit either a PDF or a Jupyter notebook, for example. You can use any programming language or tool you would like, however.

Useful Python packages

You might find the following package particularly useful:

- folium
- fiona
- shapely
- geopandas
- sklearn.cluster.DBSCAN

The following method might be useful as well (note this is for numpy, but you can create a standard Python version if you wish):

```
def haversine_np(lon1, lat1, lon2, lat2):  
    """  
    Calculate the great circle distance between two points  
    on the earth (specified in decimal degrees)  
    All args must be of equal length.  
    """  
    lon1, lat1, lon2, lat2 = map(np.radians, [lon1, lat1, lon2, lat2])  
  
    dlon = lon2 - lon1  
    dlat = lat2 - lat1  
  
    a = np.sin(dlat/2.0)**2 + np.cos(lat1) * np.cos(lat2) * np.sin(dlon/2.0)**2  
  
    c = 2 * np.arcsin(np.sqrt(a))  
    km = 6367 * c  
    return km
```

NYC Taxi Data

You will analyse a public dataset from Uber available on Kaggle, available here: <https://www.kaggle.com/c/nyc-taxi-trip-duration>. It is also available on Moodle for download.

Your primary dataset is one released by the NYC Taxi and Limousine Commission (TLC), which includes pickup time, geo-coordinates, number of passengers, and several other variables for 1.5 million trips between 2016-01-01 and 2016-06-30. Note that for this analysis, just use the *training* sample.

- `id` - a unique identifier for each trip
- `vendor_id` - a code indicating the provider associated with the trip record
- `pickup_datetime` - date and time when the meter was engaged
- `dropoff_datetime` - date and time when the meter was disengaged
- `passenger_count` - the number of passengers in the vehicle (driver entered value)
- `pickup_longitude` - the longitude where the meter was engaged
- `pickup_latitude` - the latitude where the meter was engaged
- `dropoff_longitude` - the longitude where the meter was disengaged
- `dropoff_latitude` - the latitude where the meter was disengaged
- `store_and_fwd_flag` - This flag indicates whether the trip record was held in vehicle memory before sending to the vendor because the vehicle did not have a connection to the server - Y=store and forward; N=not a store and forward trip
- `trip_duration` - duration of the trip in seconds

There is more up-to-date data available from TLC, but the datasets are large (10GB+ per year since 2018). They do include additional fields, however.

For this exercise, you can assume the centre of New York City is at the long/lat coordinates (40.716662, -74.009899).

1 Data Cleaning [3 marks]

There are several outliers in the data. Identify these and give justification for why you can remove them from the analysis. (Hint: look at trip duration, speed, distance, etc). [3 marks]

2 Feature generation [3 marks]

This can be done before or after the data cleaning step. Generate additional columns for at least these features (but you're welcome to add more!):

- Distance of trip
- Day of week
- Average speed of trip

[3 marks]

3 Time-based [14 marks]

Assume pickup time unless otherwise specified.

1. Which day of the week is the most popular? Show plots to motivate your answer. [2 marks]
2. What hour of the day is the most popular on each day? Plot a distributions of the hours and make observations and give possible suggestions for why the data looks like it does. [3 marks]
3. Investigate the differences between weekdays and weekends. What would account for this? [2 marks]
4. Look at how these patterns change on the major holidays (do they change?). Look at the following: St. Patrick's Day, Easter, Memorial Day, Valentine's Day, Martin Luther King Day. Make sure you use the correct dates for these for the relevant year. [5 marks]
5. How does the average speed of trips change throughout the day? What time of day are trips fastest? Show plots to motivate your answer. [2 marks]

4 Location clusters [12 marks]

4.1 Heatmaps

Produce a heatmap of all of the trip pickups over (do not do a scatter plot... there are 1.5 million data points and this will almost certainly crash your computer):

1. weekdays and weekends,
2. morning and evening (choose reasonable hours).

Comment on any findings you make. [4 marks]

4.2 Hotspots

If you were a taxi driver wanting to plan your evenings so that you could get the most trips, you would want to know where the popular areas are. Looking at the time periods 23:00 on a Friday evening to 02:00 on a Saturday morning, and between 17:00 and 20:00 on a Thursday, find hotspot locations (areas where there are a large number of trips happening). If you were to use k-means, you would define the number of clusters. However, here the number of clusters is not at all clear. DBSCAN (available in sklearn) determines this for you, and works well on spatial data. DBSCAN has two configurable parameters: ϵ - the maximum distance between any two points, and the minimum number of samples to determine a cluster. Your hotspot location might be defined as at least 15 pickups in that location in an hour, and locations might be required to be within 50 or 100 metres from each other (*motivate* your choice of parameters). Using DBSCAN, identify clusters and plot these on a map. How many clusters did you find? [8 marks]

The following code might be useful (but edit it to suit your own needs):

```
# calculate epsilon parameter using the user defined distance
kms_per_radian = 6371.0088
# The epsilon parameter is the max distance that points can be from each other to be
considered a cluster.
epsilon = max_distance / kms_per_radian
db = DBSCAN(eps=epsilon, min_samples=min_cars, algorithm='ball_tree',
            metric='haversine').fit(np.radians(coords))
```

5 Airports [9 marks]

Find out how long it takes, on average, to travel to JFK airport from the Empire State Building. Produce a plot showing the travel time by time of day. How does this compare with Newark Airport? Assume the following coordinates for the *centre point* of the locations (long, lat):

- JFK Airport: (40.647929, -73.777813)
- Empire State Building: (40.756724, -73.983806)
- Newark Airport: (40.689442, -74.173242)

Use a reasonable (and motivate!) radius around these locations when determining if a GPS coordinate is at that location. [9 marks]

6 Boroughs [11 marks]

You can find the shapefile containing NYC boroughs (basically neighbourhoods) here <https://data.cityofnewyork.us/City-Government/2010-Neighborhood-Tabulation-Areas-NTAs-/cpf4-rkhq>. It is also available for download on Moodle.

1. Using this shapefile find the neighbourhoods for the trip start and end locations (try geopandas, shapely, or fiona, for example). [3 marks]
2. Plot a choropleth of all pickups and all dropoffs in NYC. What do you notice about the difference in distribution? [2 marks]
3. Which boroughs have the most incoming trips and the most outgoing trips? [2 marks]
4. Which neighbourhood(s) is/are the quietest at night, between midnight and 5AM? (Not everyone wants to party). [2 marks]
5. Which neighbourhood(s) is/are the busiest at night, between midnight and 5AM? (Some people party, well, only Rod actually). [2 marks]

Submission

Total marks available: 52 (full marks is 50)

Work by yourself or in groups of up to four people. Submit your work to Moodle as a PDF or Jupyter notebook.

Deadline: 15 June 2022 (please contact me directly regarding extensions)