

Big Data Analytic HW1

0556555

Used Tools:

Since the Taxi data is quite huge (10 million data per month). Weka may not be a good solution for this since the I/O overhead is more time-costly than other naive language.

I choose **python** as my processing tool to run this analytic. Also, several useful libraries were imported to accelerate my calculation, such as **pandas, pyplot**.

Data Set: Yellow trip data 2016_10

Q1.What regions have most pickups and drop-offs?

At first, I wanted to use DBSCAN to cluster the pickup and dropoff, but the 10million dataset killed my memories many times. I choose to use the data after July, since the location changed into LocationID. Then just need to calculate the occurrence for every IDs, we could get the answer for this question.

Most Pickups: ID 237 Manhattan Upper East Side South

Most Dropoffs: ID 236 Manhattan Upper East Side North

Q2. When are the peak hours and off-peak hours in taking taxi?

The same way as q1. Calculate the occurrence for every single hour. I use pandas date time format extraction to extract the day-hour for specific analytic.

Peak hour: 19 p.m.

Off-Peak hour: 5 a.m.

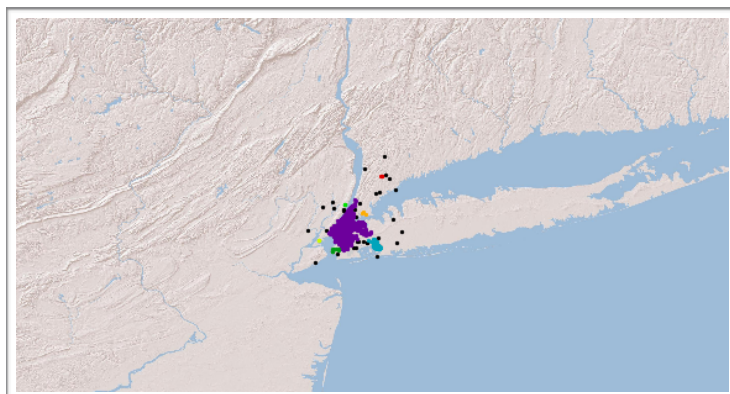
Q3. What differences exist between short and long distance trips of taking taxi?

I took the average of trip_distance as a milestone to distinguish between short and long distance trip.

I noticed that the **average tip_amount** performs better in long trip distance. And people in long trip distance tends more to **pay in credit card.** (71% vs 65%)

Note

At first, I tended to use mpl_toolkit with basemap to draw the cluster on the Geography map. But the memory issue to use DBSCAN in sklearn stopped me, I can only show the 60000 points result in the end.



Figures.

