

A Path Simulator Focusing on Time Consumption

Based on the Transport Network and the Data of Public Traffic
Vehicles in Shanghai

Name:张博伦

Partner:王鼎居

(our part 1,2,4,5,7 are the same)

【Abstract】 With the rapid growth of every aspect of our society, people's schedule is getting tighter and tighter. Thus, there is a will for the masses to spend less time on road and get the precise time consumption, in order to form a neat schedule with little time being wasted. What we want to construct is a program, which will return the minimized time consumption and the best route. This project centers on the mission of simulate the transport condition of any time and find a route which costs minimize time to reach the destination. The time we are simulating is not only present, but also for the future. The project can be mainly divided into four parts, including data acquiring, data processing, map building, the usage and comparison of graph algorithms.

【Key Words】 Trajectory big data; Routing

1. Introduction

With the rapid growth of every aspect of our society, people's schedule is getting tighter and tighter. Thus, there is a will for the masses to spend less time on road and get the precise time consumption, in order to form a neat schedule with little time being wasted. What we want to construct is a program, which will return the minimized time consumption and the best route. When searching for related papers, we found that most routing program is focused on the traffic between cities, rather than in cities and is hard to guide our daily commute. To solve this problem, our project centers on the mission of simulate the transport condition of any time and find a route which costs minimize time to reach the destination. The time we are simulating is not only present, but also for the future (for example, to schedule a conference on a workday's morning, the final time consumption will surely surpass the number revealed on the navigation app because of the rushing hour). The project is finished by Python and can be mainly divided into four parts, including data acquiring, data processing, map building, the usage and comparison of graph algorithms.

2. Related knowledge

2.1 All modules, packages and Library being used

2.1.1 OS

'os' is the attributive of the word 'operating_system'. It supplies random interfaces between various Python programs and operating systems. The os module can easily interact with the operating system, and greatly enhance the portability of code.

2.1.2 NumPy

"Numeric Python" is a library consisting of multidimensional array objects and a collection of routines for processing array.

2.1.3 Pandas

Python Data Analysis Library is a tool, origin from NumPy. It is used to solve data analysis tasks. Pandas incorporates a large number of libraries and some standard data models, providing the tools needed to operate large data sets efficiently. Pandas provides a large number of functions and methods that enable us to process data quickly and conveniently. You will soon find that it is one of the important factors that make Python a powerful and efficient data analysis environment.

2.1.4 Osm2gmns

Developed by ASU trans+ai lab team of Arizona State University, Osm2gmns can process connected network and simplify the network or build module automatically. For data inputs, it supports POI, which can be download from websites. With random network included, unified format, it is widely used to build all kinds of models.

2.1.5 OSMnx

OSMnx is a package built by geopandas, network and matplotlib. The capacity of OSMnx and Osm2gmns is alike, but the latter is good at routing.

2.1.6 Datetime

Datetime is a time related module, it can process year, month, day, hour, minute and second, return to the gap between two spots and return to the information of the time zone

2.1.7 Shapely

Shapely uses Python's ctypes module to perform set theory analysis and operation on planar features. The functions used are from the GEOS library GEOS is the migration of Java topology Suite (JTS) and the geometry engine of PostGIS spatial extension of PostgreSQL RDBMS. Point, LineString, LinearRing, Polygon... is included.

2.1.8 GC

Generational garbage collection can significantly reduce the memory usage and accelerate the program.

2.2 Trajectory big data

Thanks to the rapid development of spatial positioning technology and sensor networks, spatial positioning sensors have been widely used in aircraft, ships, cars, and handheld devices, generating, and accumulating a large amount of space-time trajectory data of moving targets. Trajectory data is the trajectory generated by moving objects in geographical space, which is usually represented by some columns of spatial points with time order. Generally, the elements of PI include: positioning point ID, trajectory ID, longitude, latitude, height, speed, time, etc.

3. My works (张博伦)

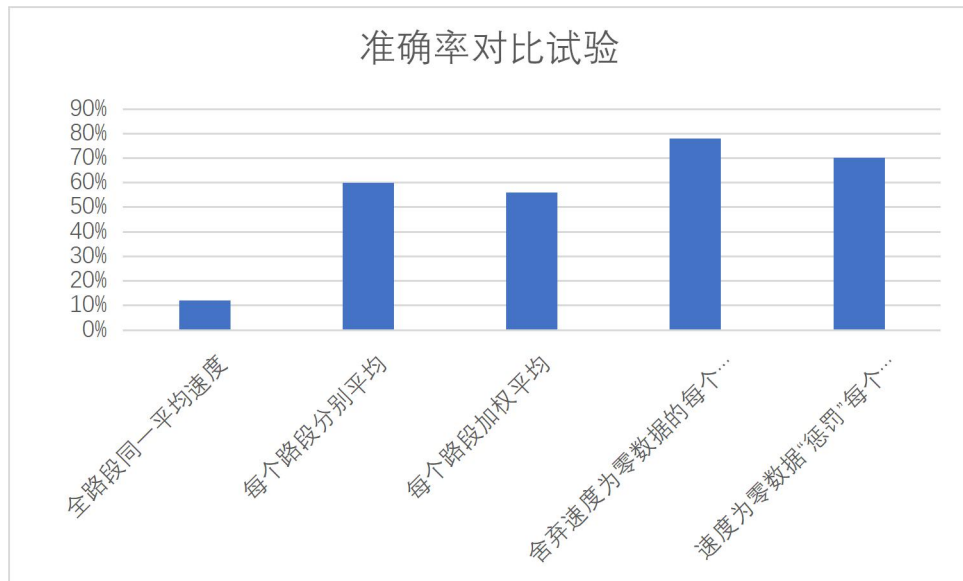
3.1 Road Data process

Using osmnx to get the road data of shanghai ,I successfully transformed the road network data into a graph that NetworkX can handle. According to the computing power of the computer and the need to reduce the test cost, I determined the test navigation range to be 10 kilometers around Minhang Campus of Shanghai Jiao Tong University (in fact, it can be extended to Shanghai).

3.2 Weight adjustment

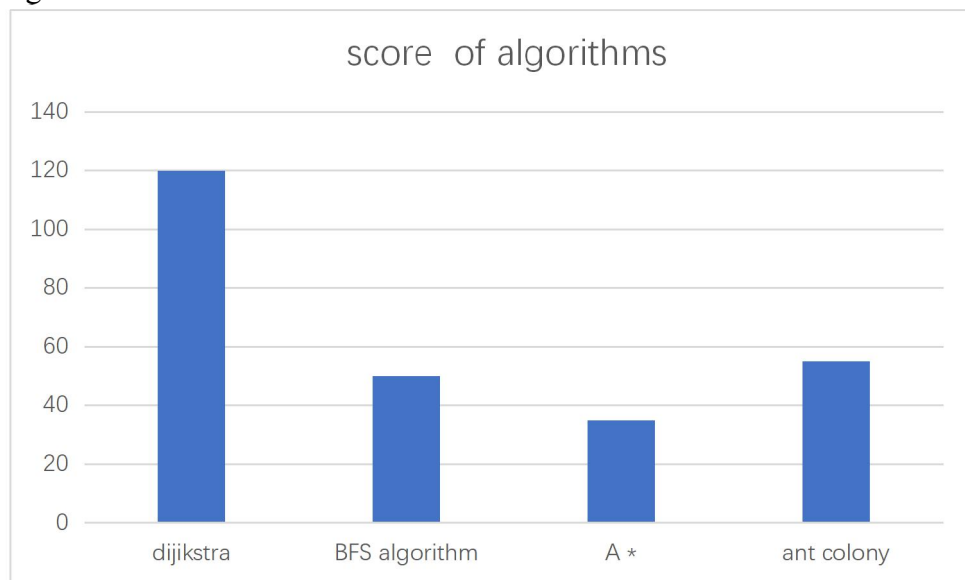
After several adjustment observation and comparison tests, it is determined that 11.11m/s is the average road speed without involving data, which is the most consistent with the actual situation

Similarly, in several comparison tests, I compared the influence of the same average speed of the whole section, the average of each section separately, the weighted average of each section, the average of each section with zero speed data abandoned, and the average of each section with zero speed data "punished" the average of each section separately on the accuracy of results. Finally, it was determined that the method of "discarding the data with zero speed and averaging each section separately" was optimal, and the weight of the data after processing was assigned to the edge of the graph



3.3 Shortest path algorithm selection

In figure on the processing of the shortest path, my team members of ant colony algorithm is proposed, and I suggest using dijkstra algorithm, through the contrast experiments, I also participated in comparison with A * algorithm and BFS algorithm, according to the processing speed from champion to the fourth are the 40,20,10,5 points respectively, processing accuracy first to fourth place, respectively, 100, 50, 25,10 points, and finally dijkstra algorithm wins with the highest score, so dijkstra algorithm is selected



3.4 Accuracy evaluation

Call the Amap interface and compare the output route and time with the result we get .*The route is consistent and the time difference within three minutes is considered consistent .*

3.5 visualization

Call OpenStreetMap, output road network diagram, shortest path diagram, parametric graph theory diagram, and output the shortest time

3.6 Program running and testing

4. Process of building project

4.1 Data acquiring

Sources from Shanghai public data update platform: <https://data.sh.gov.cn/>. The taxi data and road-net data is downloaded at this website.

4.2 Data processing(brief)

In this period, we are going to introduce the method we use briefly, more information is attached with the code.

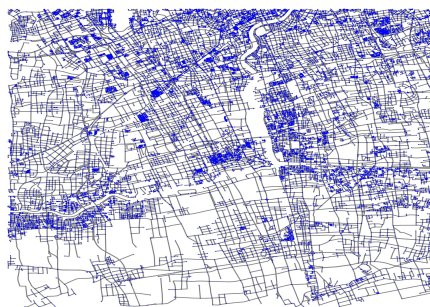
4.2.1 Road net

While abstracting route planning problem as a TSP problem. Running the program, with all road included, time consumption is disastrous (over 2h).

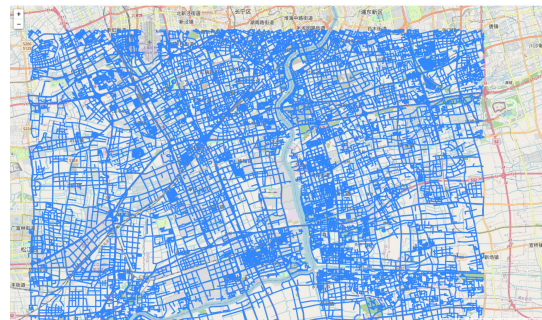
. For the massive trajectory data set, there are a large number of redundant points. When dealing TSP problems

4.2.1.1 Osm2gmns

Using osm2gmns to simplify the road net, return 'node.csv' and 'edge.csv'. The running time of the program is proportional to the length of 'edge.csv', thus the running time of the program is significantly.



Result



For contrast

4.2.1.2 K-means

Considering the feature of crossings: most of them have familiar distance, when searching for the same cluster, it is easy to get stuck at locally optimal value. More annoyingly, the constant k is hard to select, so we abandon this method.

4.2.2 Taxi data

Due to the running time of the program and the limitation of the performance of laptop, we decide to upgrade the data every 15min, so we place those data in a new folder. When shifting the data, we preserve the time, speed, and position of normal taxis and save the data as 'shanghai_taxi_spd_calculation_dict_byTime.npy'.

However, when I try to attach the speed value to every edge, no matter what border value is chosen, the average speed is always equal or missing (for a more precise view, please open the 'spd' folder, the number after k is the maximum value of the eccentricity of the oval, which stand for the range of the road). After the observation towards 'node.csv', 'edge.csv' and the real map, I found out that edge is not well-related with the real map. Also, because the imprecisely recorded GPS data (GPS has a precision of 15m, when cars travel through high buildings, the average error will increase significantly), taxis can't be divided into related clusters by covering them with oval.

Finally, to separate the taxi data, we decide to use OSMnx to attach speed data to node.

4.2.3 Calculate time consumption for each edge

First, we attach all speed data to its nearest edges (with the format of 'list') and calculate average speed by function '.avg()'. Then divide the length of the road with average speed, forming a map with time as its weight.

4.3 Map building(visualization)

Using `osmnx.graph_from_address` ('place', `dist = length`, `network_type=` in['all', 'cars', etc.]) to get street data and `osmnx.plot_graph`(graph, `show=True`) . The result is as follows (an example of finding the shortest route from the campus in Minhang to the campus in Xuhui)

4.4 The usage and comparison of graph algorithms

4.4.1 Ant Colony algorithm

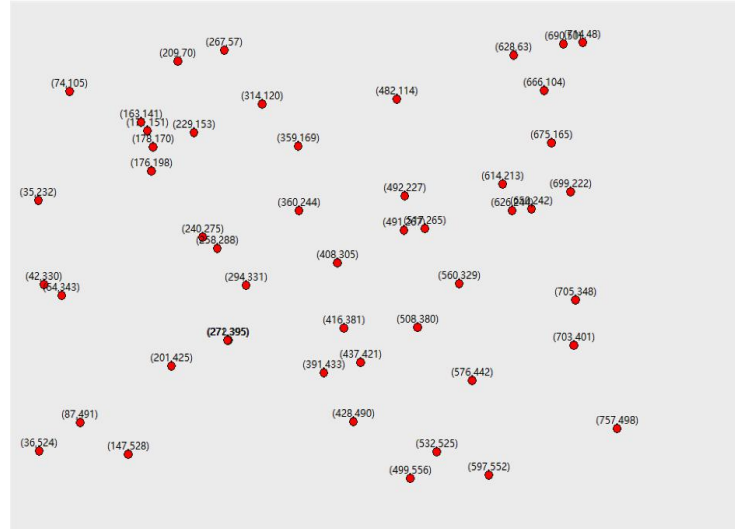
4.4.1.1 Definition

Ant colony algorithm is inspired by the research on the foraging behavior of real ant colonies. Biological research shows that a group of cooperative ants can find the shortest path between food and nest, while a single ant cannot. Biologists have found through a lot of careful observation and research that the behavior of ant individuals is interactive. In the process of movement, an ant can leave a substance called pheromone on the path it passes through, and this substance is precisely the carrier of information transmission and communication between ant individuals. Ants can perceive this substance when moving, and are used to tracking this substance to crawl. Of course, pheromones will be released during crawling. The thicker the pheromone trace on a road, the higher the probability that other ants will follow and crawl this path, so the pheromone trace on this path will be strengthened. Therefore, the collective behavior of ant colonies composed of a large number of ants will show a positive feedback phenomenon. The more ants walk along a certain path, the more likely the latecomers are to choose the path. It is through this indirect communication mechanism that ant individuals achieve the goal of collaborative search for the shortest path.

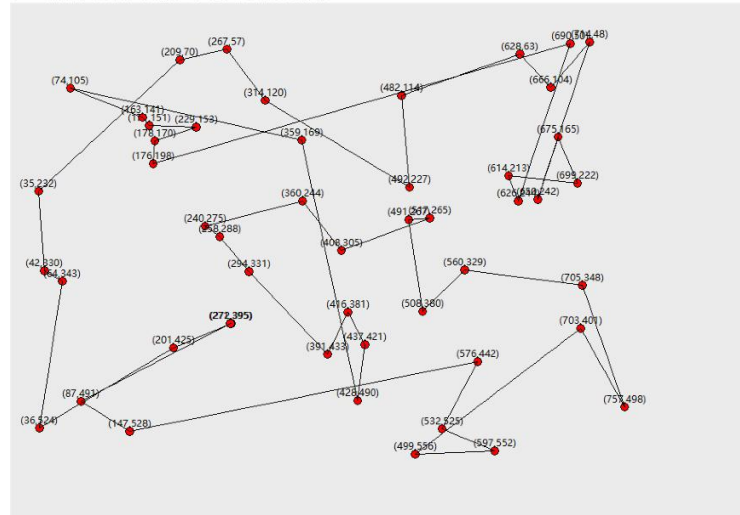
4.4.1.2 Problem revealed while trying to use ACO algorithm

After the simplification of road net, there are still 1348 nodes left. The ant needs to ergodic every edge for several times to form the best result. Let us take a simple *Traveling Salesman Problem (TSP)* as example: the overall length reach 3687 in the 155th running.

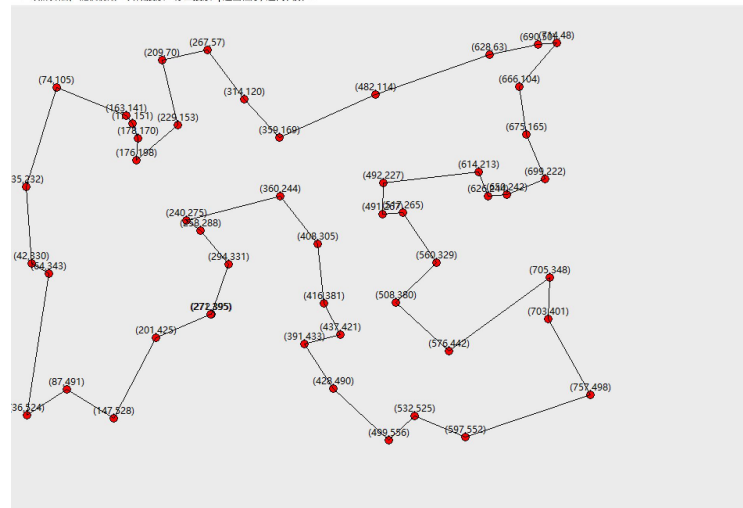
TSP蚁群算法(n:初始化 e:开始搜索 s:停止搜索 q:退出程序)



TSP蚁群算法(n:随机初始 e:开始搜索 s:停止搜索 q:退出程序) 迭代次数: 3



TSP蚁群算法(n:随机初始 e:开始搜索 s:停止搜索 q:退出程序) 迭代次数: 331

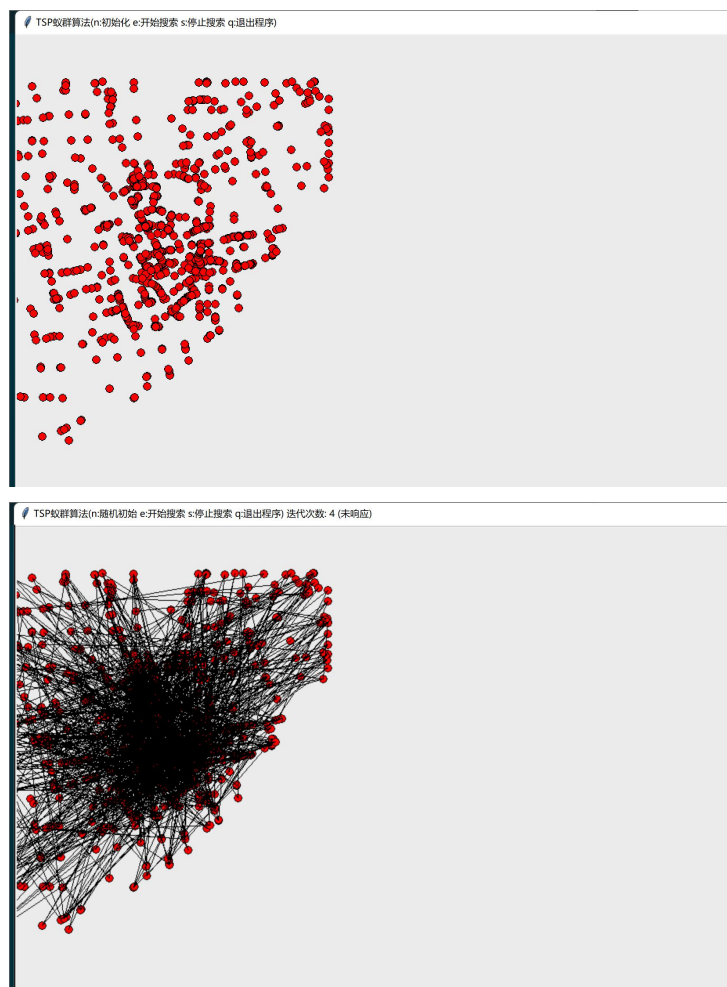


When the number of nodes is n , the number of all edges possible will reach $\frac{(n-1)!}{2}$, with the number of 1348, the number of possible edges

is 8.438299E3631. Unfortunately, for the nodes that are not linked with an edge, we can only set a high distance (eg. 100km), but delete it, for the ant does not know those two nodes are not linked.

For the net of Shanghai: the total length is dropping, but in a relatively low speed and the time consumption is unbearable.

迭代次数: 1 最佳路径总距离: 134670
迭代次数: 2 最佳路径总距离: 134644
迭代次数: 34 最佳路径总距离: 115858
迭代次数: 72 最佳路径总距离: 11473
迭代次数: 73 最佳路径总距离: 11473
迭代次数: 74 最佳路径总距离: 11473



The low productiveness while dealing TSP can somehow reflect the situation on dealing routing problem with an original a. Due to the limitation of computer performance, we stopped trying to solve routing problem with cute little ants and only give some ways for betterment:

- (1) Inspired by the function 'has_path' in networkx Set the track for the ant, so that we can reduce the number of edges significantly. We can define a function to test whether two node is reachable.

- (2) Reduce some unimportant nodes after the designation of the start and end, which may reduce the number of edges by reducing the nodes.

4.4.2 Dijkstra algorithm

4.4.2.1 Definition

The aim of this algorithm is to find the shortest path from one vertex to other vertices in weighted graph. The main feature of Dijkstra algorithm is that it starts from the starting point and adopts the strategy of greedy algorithm. Each time, it traverses the adjacent nodes of the vertices that are closest to the starting point and have not been visited until it extends to the end point.

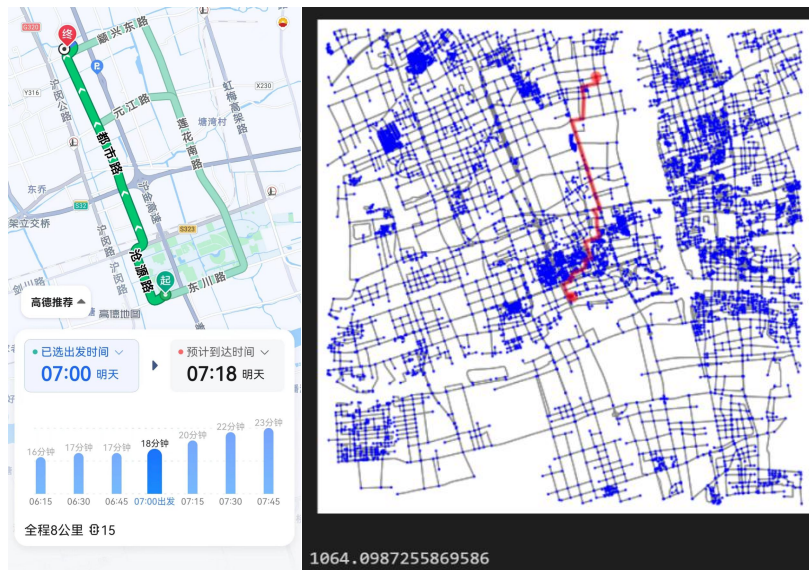
4.4.2.2 Advantage of Dijkstra algorithm

Dijkstra algorithm, or greedy algorithm is easy to understand and implement. When dealing with numerous node, greedy algorithm only needs to select the nearest node but consider the influence from other aspects. Thus, although it might not return the best answer, it is reliable and costs acceptable time.

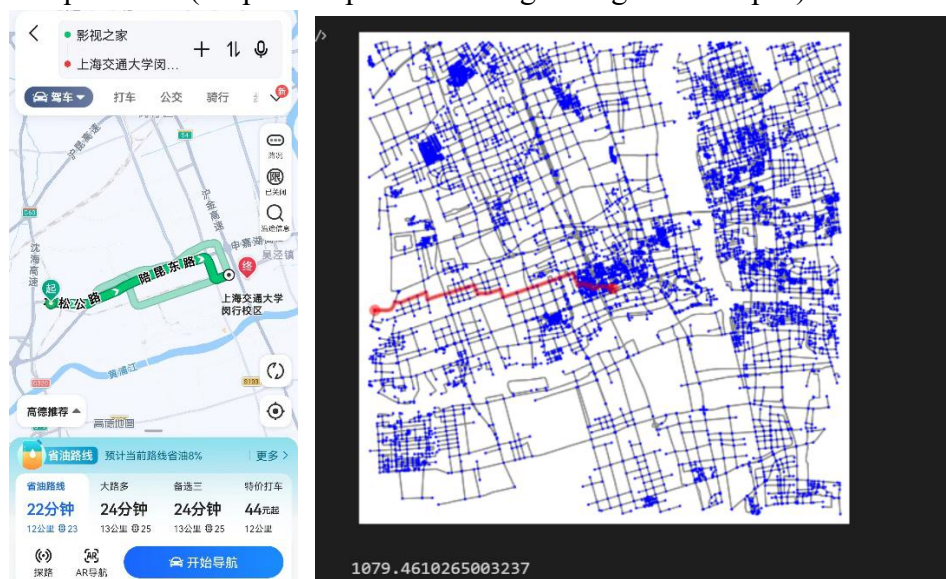
5. Result analysis

(Among the many matches, we selected two simple examples)

We set the starting time as 7 in the morning. The app ‘Gaode Navigator’ shows that it may cost 16 minutes, equals 1080 seconds, slightly longer than what we have predicted(1064second) two routes are quite alike(despite the part of striking through the campus) . The differences of time is less than 2 minute, which is a common error in the real world



The app 'Gaode Navigator' shows that it may cost 22 minutes, equals the range of 1290seconds to 1350seconds, shorter than what we have predicted(1079second) two routes are quite alike(despite the part of striking through the campus).



6. Individual thinking of boundedness

6.1 Data

Taxi data cannot fully reflect road conditions and commuting time. In order to obtain the most real results, it is necessary to understand the distribution of red street lights, traffic accidents, road construction and other factors

6.2 Standard of judging

Can only compare the results with existing navigation software on the market. More realistic, more realistic criteria are needed

7. Problems and betterment

7.1 Relate to reality

During the process we found out while routing the path from campus in Minhang district to campus in Xuhui district, the final path has crossed the campus, which means the simplification program can't tell whether the road is usable for most drivers. For the same time, the data regression reshapes our 3-D world into a 2-D graph, the average speed, and the ability to turn to another road for highway and frontage must be different. So, the final route might not be suitable in the reality.

7.2 Time consumption

The program is in urge needs to be bettered for it costs too much time. We have already use gc. In data processing, but our laptop continued to break down during the process of running the code.