Midterm report

This semester, we have learnt how to use QGIS as the visualization tool and PostgreSQL as database. These tools are powerful in illustrating the geographic information from a visionary perspective. With the practice during the course, I managed to learn how to generate database from various methods; how to export SQL files; how to generate graph in QGIS; and most importantly, how to adjust the graph elements to optimize the reception experience of viewers.

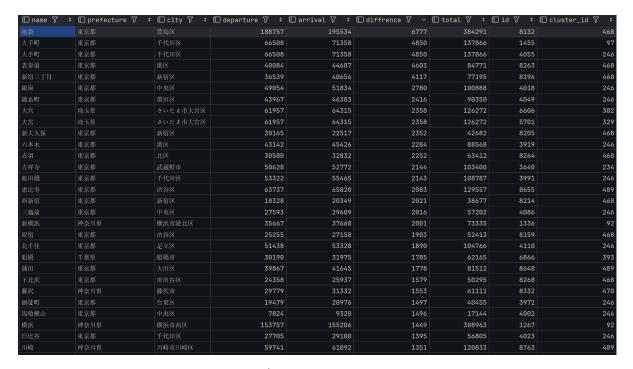
My report is a very simple one. Although I am not able to produce stunning graphs with my humble aesthetics ability, I managed to produce a decent graph with figures on the stations along the train lines.

Given limited data, my initial attempt was to calculate the ratio of empty seats pf each line. First of all, let me explain how the ratio of empty seats bothers me. Let us assume we have a line between a and b. And $a \to b$ is busy while $b \leftarrow a$ is empty. And if a train goes from a to b, since the seats are all occupied, we can say the energy is well spent because the train fulfills its job as a commuting tool as well as earns profit. Yet on its way back, this train is so empty that there is no economic profit in this run and thus the budget is wasted. I believe the occupation of seats is very important to measure if a line is energy efficient and therefore environmental friendly.

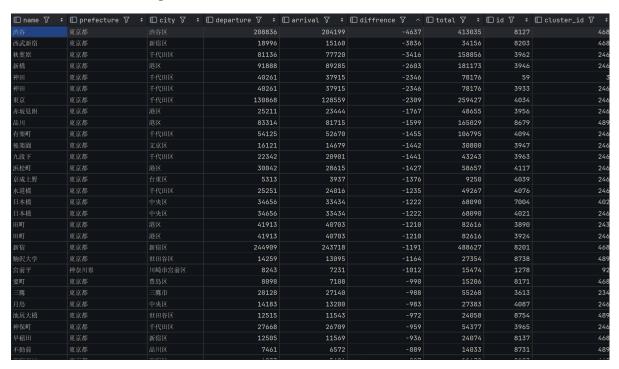
However, there is no way to measure the seat occupation directly from the data. First of all, we only get the information of departure and arrival in each station. Beside, some of the stations are shared by different lines. Therefore, the exact number of seat occupation is impossible to get regarding the missing data. So, I measured the difference between arrival and departure of each station instead.



The picture above is the part of my map1. In which, I found the station with highest influx(arrival-difference) is 池袋.



While the station with highest outflux is 渋谷.



However, I think it is not reasonable to infer any results from this data. First of all, these two stations has largest passenger flow among all the stations. And the larger the flow is, the more likely to have higher difference. Therefore, I managed to sort the stations by the ratio of influx/outflux to the total number of passengers. And map2 is the result. However, this data is not very convincing since I found some data is abnormal.

| □ name 🎖 : | ∷ □ prefecture 🎖 💠 | □city 7 ÷ | ☐ departure 🎖 💠 | □arrival 7 ÷ | \square diffrence $ abla$ ÷ | □ total 7 ÷ | □ratio 7 ~ |
|----------------|--------------------|--------------|-----------------|--------------|-------------------------------|-------------|-------------------------------------------|
| 白丸 | 東京都 | 奥多摩町 | 16 | | -16 | 16 | 1 |
| 竹沢 | 埼玉県 | 小川町 | 0 | 11 | 11 | 11 | 1 |
| 笹子 | 山梨県 | 大月市 | | | | | 1 |
| 大宝 | 茨城県 | 下妻市 | | 12 | 12 | 12 | 1 |
| 東浪見 | 千葉県 | 一宮町 | 17 | | -11 | 23 | 0.4782608695652173913 |
| 東成田 | 千葉県 | 成田市 | 131 | 367 | 236 | 498 | 0.47389558232931726908 |
| 初狩 | | 大月市 | 19 | 38 | 19 | 57 | 0.333333333333333333333 |
| 合戦場 | 栃木県 | 栃木市 | 22 | 11 | -11 | 33 | 0.333333333333333333333 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 武蔵横手 | 埼玉県 | | 38 | 63 | 25 | 101 | 0.24752475247524752475 |
| 川崎新町 | 神奈川県 | 川崎市川崎区 | 363 | 230 | -133 | 593 | 0.22428330522765598651 |
| 相模金子 | 神奈川県 | 大井町 | 62 | 40 | -22 | 102 | 0.21568627450980392157 |
| 塔ノ沢 | 神奈川県 | 箱根町 | 15 | 23 | | 38 | 0.21052631578947368421 |
| 柴又 | 東京都 | 葛飾区 | 1502 | 985 | -517 | 2487 | 0.20788098110172899075 |
| 府中競馬正門前 | 東京都 | 府中市 | 47 | 31 | -16 | 78 | 0.20512820512820512821 |
| 鳥沢 | 山梨県 | 大月市 | 121 | 81 | -40 | 202 | 0.1980198019801980198 |
| 柳生 | 埼玉県 | | 83 | 56 | -27 | 139 | 0.19424460431654676259 |
| 松田 | 神奈川県 | 松田町 | 51 | 75 | 24 | 126 | 0.19047619047619047619 |
| 水郷 | 千葉県 | 香取市 | 19 | 13 | -6 | 32 | 0.1875 |
| 葛生 | 栃木県 | 佐野市 | 62 | 43 | -19 | 105 | 0.18095238095238095238 |
| 柳小路 | 神奈川県 | 藤沢市 | 807 | 564 | -243 | 1371 | 0.17724288840262582057 |
| 多々良 | 群馬県 | 館林市 | 28 | 20 | -8 | 48 | 0.1666666666666666666 |
| 田島 | 栃木県 | 佐野市 | 21 | 15 | -6 | 36 | 0.166666666666666666 |
| 大平下 | 栃木県 | 栃木市 | 63 | 46 | -17 | 109 | 0.15596330275229357798 |
| 大田郷 | 茨城県 | | 56 | 41 | -15 | 97 | 0.15463917525773195876 |
| Michael en alc | Advisor to CER | Acto Los mor | *** | 100 | | 222 | ^ 4 - 0 - 0 - 0 - 0 - 0 - 0 - 0 - 0 - 0 - |

I got some stations with result ratio of 1 or 0. I was a little bit confused about the result. And I assume the result is caused by a small amount of data from abnormal stations. So, I normalized the data by lift the bar by 20. And I got the third data.

| □ name 🎖 💠 | \square prefecture $ abla$ ÷ | □city 7 ÷ | □ departure 🎖 💠 | □ arrival 7 ÷ | \square diffrence $ abla$ ÷ | □ total 7 ^ | □ratio \7 ÷ |
|----------------|--------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------|---------------|-------------------------------|-------------|-----------------------------------------|
| 東成田 | 千葉県 | 成田市 | 131 | 367 | 236 | 498 | 0.47389558232931726908 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 栄町 | 千葉県 | 千葉市中央区 | 85 | 142 | 57 | 227 | 0.2511013215859030837 |
| 武蔵横手 | 埼玉県 | 日高市 | 38 | 63 | 25 | 101 | 0.24752475247524752475 |
| | 神奈川県 | 松田町 | 51 | 75 | 24 | 126 | 0.19047619047619047619 |
| 新高島 | 神奈川県 | 横浜市西区 | 2046 | 2741 | 695 | 4787 | 0.14518487570503446835 |
| 葭川公園 | 千葉県 | 千葉市中央区 | 319 | 413 | 94 | 732 | 0.12841530054644808743 |
| 神泉 | 東京都 | 渋谷区 | 3218 | 4108 | 890 | 7326 | 0.12148512148512148512 |
| 扇町 | 神奈川県 | 川崎市川崎区 | 76 | 96 | 20 | 172 | 0.11627906976744186047 |
| 扇町 | 神奈川県 | 川崎市川崎区 | 76 | 96 | 20 | 172 | 0.11627906976744186047 |
| 石下 | 茨城県 | 常総市 | 89 | 111 | 22 | 200 | 0.11 |
| 新整備場 | 東京都 | 大田区 | 1239 | 1544 | 305 | 2783 | 0.10959396334890406037 |
| 東武竹沢 | 埼玉県 | 小川町 | 97 | 120 | 23 | 217 | 0.10599078341013824885 |
| 武蔵白石 | 神奈川県 | 川崎市川崎区 | 339 | 405 | 66 | 744 | 0.08870967741935483871 |
| 新芝浦 | 神奈川県 | 横浜市鶴見区 | 175 | 209 | 34 | 384 | 0.08854166666666666667 |
| 馬喰横山 | 東京都 | 中央区 | 7824 | 9320 | 1496 | 17144 | 0.08726084927671488567 |
| 昭和 | 神奈川県 | 川崎市川崎区 | 137 | 162 | 25 | 299 | 0.08361204013377926421 |
| 新千葉 | 千葉県 | 千葉市中央区 | 421 | 497 | 76 | 918 | 0.08278867102396514161 |
| | 千葉県 | | 174 | 205 | | 379 | 0.08179419525065963061 |
| 小絹 | 茨城県 | つくばみらい市 | 202 | 237 | 35 | 439 | 0.07972665148063781321 |
| 南新宿 | 東京都 | 渋谷区 | 1140 | 1335 | 195 | 2475 | 0.07878787878787878788 |
| 二重橋前 | 東京都 | 千代田区 | 6500 | 7595 | 1095 | 14095 | 0.07768712309329549486 |
| 京成高砂 | 東京都 | 葛飾区 | 6366 | 7410 | 1044 | 13776 | 0.07578397212543554007 |
| 石岡 | 茨城県 | 石岡市 | 510 | 592 | 82 | 1102 | 0.07441016333938294011 |
| 松が谷 | 東京都 | 八王子市 | 530 | 612 | 82 | 1142 | 0.07180385288966725044 |
| 大崎広小路 | 東京都 | 品川区 | 2263 | 2609 | 346 | 4872 | 0.0710180623973727422 |
| dett plan dela | ±12 → 407 | And the last of th | 2012 | 40740 | 4750 | 40001 | 0.0000000000000000000000000000000000000 |

Now the data start to make sense. And it is also shown on map3.

In conclusion, although the desired ratio of empty seats is not acquired. But, I think the difference of influx/outflux can be used as a side-factor indicating the tend of utilization of space of each train. So far, I think in order to infer the empty seats ratio, we might need more data on the extent of crowdedness or more specific data on passenger counts of each line at various points of the time. Also, there should require more advanced statistic tools to make a reasonable prediction.