# Week 3 Tasks Sub-list 3

Student Name: Zhangli Wang
Student ID: 10868661
Student Central Name: n67315zw

In [1]:

```python
# packages
import pandas as pd
import numpy as np
import datetime
import copy
from itertools import combinations
```

In [2]:

```python
# Read dataset
df = pd.read_csv('./Datasets/rawpvr_2018-02-01_28d_1083 TueFri.csv')
```

*Explanation*:
Read the original csv file 'rawpvr_2018-02-01_28d_1083 TueFri.csv' into a dataframe.

In [3]:

```python
# slice to acquire the 'Date' column
date = df.loc[:, 'Date']
# type conversion: string > timestamp
date = pd.to_datetime(date)
# calculate days of the week
days_of_week = [item.dayofweek for item in date]
# type conversion: list > series
days_of_week = pd.Series(np.array(days_of_week).T)
# update the column 'Flags' in the original dataframe
df.loc[:, 'Flags'] = days_of_week.map(lambda x: x+1, na_action='ignore')
# update the column 'Flag Text' in the original dataframe
df.loc[:, 'Flag Text'] = days_of_week.map(lambda x: 'Tuesday' if x == 1 else ('Friday' if x
 == 4 else ''), na_action='ignore')
# add a new column 'Time' to store the time of the date
df = df.assign(Time = pd.to_datetime(df['Date']).dt.time)
# add a new column 'DateWOTime' to store the date without the time
df = df.assign(DateWOTime = pd.to_datetime(df['Date']).dt.date)
# transform the type of data in the 'Date' column of the dataset: string > timestamp
df = df.assign(Date = pd.to_datetime(df['Date']))
# add a new column 'DateHour' to store the date in the format of '%year%-%month%-%day%
%hour%:00:00'
df['DateHour'] = df['Date'].apply(lambda x: datetime.datetime(x.year, x.month, x.day,
x.hour, 0, 0))
# add a new column 'Hour' to store the hour in the format of '%hour%:00:00'
df['Hour'] = df['Date'].apply(lambda x: datetime.time(x.hour, 0, 0))
```

*Explanation*:
This code block is the data preparation operations for the input dataset.

- Update the column `Flag` and `Flag Text` with the day of the week id and the day of the week, respectively.
    - Create a series of the days of week.

        ```
        # slice to acquire the 'Date' column
        date = df.loc[:, 'Date']
        # type conversion: string > timestamp
        date = pd.to_datetime(date)
        # calculate days of the week
        days_of_week = [item.dayofweek for item in date]
        # type conversion: list > series
        days_of_week = pd.Series(np.array(days_of_week).T)
        ```

        Firstly, the column `Date` is sliced from the original dataframe and converted to the type timestamp. Then the days of week is calculated and its type is converted from list to series for adding columns in the dataframe.
    - Update the columns `Flag` and `Flag Text`.

        ```
        # update the column 'Flags' in the original dataframe
        df.loc[:, 'Flags'] = days_of_week.map(lambda x: x+1, na_action='ignore')
        # update the column 'Flag Text' in the original dataframe
        df.loc[:, 'Flag Text'] = days_of_week.map(lambda x: 'Tuesday' if x == 1 else ('Friday' if x == 4 else ''), na_action='ignore')
        ```

        Use map functions to map the series to the new columns in the dataframe with specified requirements.
- Add a new column `Time` in the dataframe to store the time in dates

    ```
    # add a new column `time` in the dataframe to store the time of the date
    df = df.assign(Time = pd.to_datetime(df['Date']).dt.time)
    ```

- Add a new 'DateWOTime' to store the date without the time

    ```
    # add a new column `DateWOTime` in the dataframe to store the date without the time
    df = df.assign(DateWOTime = pd.to_datetime(df['Date']).dt.date)
    ```

- Change the datatype of `Date` from `string` to `timestamp`

    ```
    # transform the type of data in the 'Date' column of the dataset: string > timestamp
    df = df.assign(Date = pd.to_datetime(df['Date']))
    ```

    This operation is used for filtering the dataframe according to dates in the future and is a necessary step for continuing the very next data preparasion operasion.
- Add a new column `DateHour` to store the date in the format of '%year%-%month%-%day% %hour%:00:00'

    ```
    # add a new column 'DateHour' to store the date in the format of '%year%-%month%-%day% %hour%:00:00'
    df['DateHour'] = df['Date'].apply(lambda x: datetime.datetime(x.year, x.month, x.day, x.hour, 0, 0))
    ```

    This operation is useful for filtering the dataframe according to one-hour time intervals.

The data preparasion operations mainly focus on altering and adding variaty of the time attribute for the convenience of filtering the dataframe accordign to time. In the future, for each sub-task, more data preparasion operations are executed and they mainly focus on filtering to a new dataframe which is to meet the constraints for achieving the final goals.

# Task5

## Task5.1

---

**Solution**:

In [4]:

```python
# filter the dataset with regard to 'Flag Text'
df1 = df.loc[df['Flag Text'] == 'Tuesday']
# filter the dataset with regard to 'Time'
start_time = datetime.time(7, 0, 0)
end_time = datetime.time(19, 0, 0)
df1 = df1.loc[(df1['Time'] > start_time) & (df1['Time'] < end_time)]
```

*Explanation*:

This code block is the data preparation operations for Task5.1 that filters the dataframe to a new dataframe `df1` according to constraints of the dataset in Task5.1.

- Filter the dataset with regard to `Flag Text` that the day of the week is Tuesday.

  ```python
  # filter the dataset with regard to 'Flag Text'
  df1 = df.loc[df['Flag Text'] == 'Tuesday']
  ```

- Filter the dataframe with regard to `Time` that the time is between 7:00 and 19:00.

  ```python
  # filter the dataset with regard to 'Time'
  start_time = datetime.time(7, 0, 0)
  end_time = datetime.time(19, 0, 0)
  df1 = df.loc[(df['Time'] > start_time) & (df['Time'] < end_time)]
  ```

In [5]:

```python
num_cells = len(df1)     # number of cells

# number of non-empty cells of the column 'Gap (s)'
num_non_e_cells_gap = len(df1['Gap (s)'].notna()[df1['Gap (s)'].notna() == True])
# column completeness of the column 'Gap (s)'
column_completeness_gap = 100 * num_non_e_cells_gap / num_cells

# number of non-empty cells of the column 'Headway (s)'
num_non_e_cells_headway = len(df1['Headway (s)'].notna()[df1['Headway (s)'].notna() ==
True])
# column comleteness of the column 'Headway (s)'
column_completeness_headway = 100 * num_non_e_cells_headway / num_cells

print('column completeness gap: ' + str(column_completeness_gap))
print('column completeness headway: ' + str(column_completeness_headway))
```

```
column completeness gap: 98.03654443753885
column completeness headway: 98.96532007458049
```

*Explanation*:

This code block is calculation of column completeness for the columns `Gap (s)` and `Headway (s)`. Firstly, the length of the columns in `df1` are calculated by `num_cells = len(df1)`. Then, the number of non-empty cells of each column are calculated like

```
# number of non-empty cells of the column 'Gap (s)'
num_non_e_cells_gap = len(df1['Gap (s)'].notna()[df1['Gap (s)'].notna() == True])
```

. Finally, the the values of column completeness are calculated like

```
# column completeness of the column 'Gap (s)'
column_completeness_gap = 100 * num_non_e_cells_gap / num_cells
```

.

**Result**:

- column completeness gap: 98.03654443753885
- column completeness headway: 98.96532007458049

**Interpretation of the result**:

The completeness values of the columns 'Gap (s)' and 'Headway (s)' are both over 98.0, which means that over 98 percent of the values are not empty. So for both the two columns, most of the records are filled with data. But there still exists a small proportion of missing values.

## Task5.2

---

**Solution**:

In [6]:

```
1  # initialize dataframe for storing interpolated dataframe
2  df2 = df1.copy(deep=True)
3  # filter with regard to 'Lane Name'
4  df2 = df2.loc[df2['Lane Name'] == 'NB_MID']
```

*Explanation*:

This code block is the data preparation operations for Task5.2. Because the dataset to be used in this task also requires the same constraints as in Task5.1, a new dataframe `df2` is created from the dataframe used in Task5.1 as initialization.

```
df2 = df1.copy(deep=True)
```

- Filter the dataframe with regard to `Lane Name` that the lane name is 'NB_MID'

```
# filter with regard to 'Lane Name'
df2 = df2.loc[df2['Lane Name'] == 'NB_MID']
```

In conclusion, now the dataframe has additional restrictions that the `Flag Text` is 'Tuesday', the `Time` is ranged from 7:00 to 19:00, and the `Lane Name` is 'NB_MID'. In the next code block, data manipulation on `df2` can be done with the preparation.

```
 1  medians = dict()    # key: column name, value: medians of the column
 2  for column_name in ['Gap (s)', 'Headway (s)']:
 3        # record unique hour
 4        unique_hours = df2['Hour'].unique()
 5
 6        # calculate median values for all the time intervals
 7        medians_col = dict()    # key: hour, value: median value
 8        for unique_hour in unique_hours:
 9              # filter the dataframe with regard to date hour
10              filtered_df2 = df2.loc[df2['Hour'] == unique_hour]
11
12              # get the values of the column
13              values = filtered_df2[column_name]
14
15              # record the median
16              medians_col.update({unique_hour: values.median()})
17
18        # record the medians of the column
19        medians.update({column_name: medians_col})
20
21        # find the series of 'Date' of dataframe if the column_name is null
22        ss_date_missing = df2.loc[df2[column_name].isna()]['Date']
23
24        # update the cell of the dataframe with the median value
25        for item in ss_date_missing:
26              # get the row of the dataframe where the value is missing
27              tempdf = df2.loc[df2['Date'] == item]
28
29              # get the corresponding median value
30              median_value = medians[column_name].get(tempdf['Hour'].iloc[0])
31
32              # interpolate value to the dataframe
33              df2.loc[df2['Date'] == item, column_name] = median_value
34
35  # print results
36  # print(medians)
37
38  # clear memory
39  # del df1
40  # del df2
```

*Explanation*:

This code block is the data manipulation on `df2`. Overall, there are two processes explained below.

- Find median values of the columns `Gap (s)` and `Headway (s)` for all one-hour intervals.
  - Firstly we initialize the object to store all the median values. It is a dictionary storing names of the columns `Gap (s)` and `Headway (s)` as two keys. And their values are dictionaries that stores median values of corresponding columns with regard to one-hour time intervals in the day.

    ```
    medians = dict()  # key: column name, value: medians of the column
    ```

  - Then we iterate over the two columns `Gap (s)` and `Headway (s)` to execute the same operations to them. Following mentioned operations are executed in this iteration.

    ```
    for column_name in ['Gap (s)', 'Headway (s)']:
    ```

- To acheive the goal, we may want to iteratively filter `df2` according to each one-hour time interval between 7:00 and 19:00 and calculate median values. So to begin with we find all unique one-hour time intervals in `df2`.

```
# record unique hour
unique_hours = df2['Hour'].unique()
```

- Then we execute operations mentioned in the previous bullet point and the median values for the column are calculated and updated to `medians`.

```
# calculate median values for all the time intervals
medians_col = dict()  # key: date hour, value: median value
for unique_hour in unique_hours:
    # filter the dataframe with regard to hour
    filtered_df2 = df2.loc[df2['Hour'] == unique_hour]

    # get the values of the column
    values = filtered_df2[column_name]

    # record the median
    medians_col.update({unique_hour: values.median()})

# record the medians of the column
medians.update({column_name: medians_col})
```

- Update the cells with missing values with the corresponding median values.
  - To acheive this goal, we need to find and record all the rows that have missing values. We can only take the column `Date` for future comparison as the dates are unique.

```
# find the series of 'Date' of dataframe if the column_name is null
ss_date_missing = df2.loc[df2[column_name].isna()]['Date']
```

  - Now we can update the values by iteratively find a missing row in `df2` and update it with the corresponding value.

```
# update the cell of the dataframe with the median value
for item in ss_date_missing:
    # get the row of the dataframe where the value is missing
    tempdf = df2.loc[df2['Date'] == item]

    # get the corresponding median value
    median_value = medians[column_name].get(tempdf['Hour'].iloc[0])

    # interpolate value to the dataframe
    df2.loc[df2['Date'] == item, column_name] = median_value
```

In the iteration, the first step is to get the row that the date (unique key) in `df2` is matched with the one in `ss_date_missing`. The second step is to find the corresponding median value in `medians` according to the time of the row. The final step is updating the row.

**Result of median values**:

{'Gap (s)':
  {datetime.time(7, 0): 1.834,
   datetime.time(8, 0): 2.1020000000000003,
   datetime.time(9, 0): 2.08,
   datetime.time(10, 0): 2.5835,
   datetime.time(11, 0): 2.7785,
   datetime.time(12, 0): 2.795,
   datetime.time(13, 0): 2.7560000000000002,
   datetime.time(14, 0): 2.8049999999999997,
   datetime.time(15, 0): 2.577,
   datetime.time(16, 0): 2.3225,
   datetime.time(17, 0): 2.187,
   datetime.time(18, 0): 2.228},
 'Headway (s)':
  {datetime.time(7, 0): 2.722,
   datetime.time(8, 0): 3.1910000000000003,
   datetime.time(9, 0): 2.72,
   datetime.time(10, 0): 3.0,
   datetime.time(11, 0): 3.21,
   datetime.time(12, 0): 3.16,
   datetime.time(13, 0): 3.102,
   datetime.time(14, 0): 3.133,
   datetime.time(15, 0): 2.92,
   datetime.time(16, 0): 2.853,
   datetime.time(17, 0): 2.9065,
   datetime.time(18, 0): 2.79}}

**Result of screenshot of the updated dataframe**:

| | Date | Lane | Lane Name | Direction | Direction Name | Speed (mph) | Headway (s) | Gap (s) | Flags | Flag Text | Time | DateWOTime | DateHour | Hour |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 68517 | 2018-02-06 07:00:01.160 | 2 | NB_MID | 1 | North | 27.962 | 1.720 | 1.834 | 2 | Tuesday | 07:00:01.160000 | 2018-02-06 | 2018-02-06 07:00:00 | 07:00:00 |
| 68518 | 2018-02-06 07:00:03.000 | 2 | NB_MID | 1 | North | 27.962 | 2.480 | 1.992 | 2 | Tuesday | 07:00:03 | 2018-02-06 | 2018-02-06 07:00:00 | 07:00:00 |
| 68521 | 2018-02-06 07:00:04.020 | 2 | NB_MID | 1 | North | 27.962 | 1.520 | 0.856 | 2 | Tuesday | 07:00:04.020000 | 2018-02-06 | 2018-02-06 07:00:00 | 07:00:00 |
| 68524 | 2018-02-06 07:00:05.020 | 2 | NB_MID | 1 | North | 29.205 | 1.226 | 0.680 | 2 | Tuesday | 07:00:05.020000 | 2018-02-06 | 2018-02-06 07:00:00 | 07:00:00 |
| 68529 | 2018-02-06 07:00:07.020 | 2 | NB_MID | 1 | North | 29.825 | 2.137 | 1.694 | 2 | Tuesday | 07:00:07.020000 | 2018-02-06 | 2018-02-06 07:00:00 | 07:00:00 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 496010 | 2018-02-27 18:58:51.060 | 2 | NB_MID | 1 | North | 35.417 | 1.200 | 0.725 | 2 | Tuesday | 18:58:51.060000 | 2018-02-27 | 2018-02-27 18:00:00 | 18:00:00 |
| 496029 | 2018-02-27 18:59:26.080 | 2 | NB_MID | 1 | North | 41.010 | 35.200 | 34.941 | 2 | Tuesday | 18:59:26.080000 | 2018-02-27 | 2018-02-27 18:00:00 | 18:00:00 |
| 496050 | 2018-02-27 18:59:48.040 | 2 | NB_MID | 1 | North | 34.798 | 22.600 | 22.382 | 2 | Tuesday | 18:59:48.040000 | 2018-02-27 | 2018-02-27 18:00:00 | 18:00:00 |
| 496051 | 2018-02-27 18:59:50.000 | 2 | NB_MID | 1 | North | 33.554 | 1.800 | 1.343 | 2 | Tuesday | 18:59:50 | 2018-02-27 | 2018-02-27 18:00:00 | 18:00:00 |
| 496061 | 2018-02-27 18:59:58.060 | 2 | NB_MID | 1 | North | 43.495 | 7.600 | 7.307 | 2 | Tuesday | 18:59:58.060000 | 2018-02-27 | 2018-02-27 18:00:00 | 18:00:00 |

# Task6

**Solution**:

In [8]:

```python
# read the dataset of the site 1415
dfb = pd.read_csv('./Datasets/rawpvr_2018-02-01_28d_1415 TueFri.csv')
```

In [9]:

```python
# apply the same data preparation executions to dataframe b

# slice to acquire the 'Date' column
date = dfb.loc[:, 'Date']
# type conversion: string > timestamp
date = pd.to_datetime(date)
# calculate days of the week
days_of_week = [item.dayofweek for item in date]
# type conversion: list > series
days_of_week = pd.Series(np.array(days_of_week).T)
# update the column 'Flags' in the original dataframe
dfb.loc[:, 'Flags'] = days_of_week.map(lambda x: x+1, na_action='ignore')
# update the column 'Flag Text' in the original dataframe
dfb.loc[:, 'Flag Text'] = days_of_week.map(lambda x: 'Tuesday' if x == 1 else ('Friday' if
    x == 4 else ''), na_action='ignore')
# add a new column 'Time' to store the time of the date
dfb = dfb.assign(Time = pd.to_datetime(dfb['Date']).dt.time)
# add a new column 'DateWOTime' to store the date without the time
dfb = dfb.assign(DateWOTime = pd.to_datetime(dfb['Date']).dt.date)
# transform the type of data in the 'Date' column of the dataset: string > timestamp
dfb = dfb.assign(Date = pd.to_datetime(dfb['Date']))
# add a new column 'DateHour' to store the date in the format of '%year%-%month%-%day%
    %hour%:00:00'
dfb['DateHour'] = dfb['Date'].apply(lambda x: datetime.datetime(x.year, x.month, x.day,
    x.hour, 0, 0))
```

*Explanation*:
In this code block, we simply execute the same data preparation operations for `dfb` as the ones we did for `df` for consistency and future usage.

In [10]:

```python
# create a alias for the dataframe of the site 1083 for this task
dfa = df
# use a variable to reference the two dataframes
dfs = {'1083': dfa,
       '1415': dfb}
```

*Explanation*:
Firstly an alias of the dataframe for the site 1083 is created ( `dfa` ), this is for standardize the names for the dataframes in avoidance of generating potential inconsistency in the future. Then the dataframes are referenced together in one variable ( `dfs` ) for the convenience of coding.

```
1  # filter the two dataframes
2
3  # create a copy of the old dataframes
4  dfs1 = copy.deepcopy(dfs)
5  for key, df in dfs1.items():
6      # filter the dataframe with regard to 'Time'
7      start_time = datetime.time(17, 0, 0)
8      end_time = datetime.time(18, 0, 0)
9      df = df.loc[(df['Time'] > start_time) & (df['Time'] < end_time)]
10
11     # filter the dataframe with regard to 'Direction'
12     df = df.loc[df['Direction'] == 1]
13
14     # filter the dataframe with regard to 'Flag Text'
15     df = df.loc[df['Flag Text'] == 'Friday']
16
17     # record
18     dfs1[key] = df
```

*Explanation*:

This code block is the data preparation operations for Task6. Because the operations are specificly focused on Task6, we created a new version `dfs1` of the dataframe `dfs` similar to the things we did before `dfs1 = copy.deepcopy(dfs)` .\

Then we iterate on `dfs1` for each dataframe and perform the following operations.

- Filter the dataframe with regard to 'Time' that the time is between 17:00 and 18:00.

    ```
    # filter the dataframe with regard to 'Time'
    start_time = datetime.time(17, 0, 0)
    end_time = datetime.time(18, 0, 0)
    df = df.loc[(df['Time'] > start_time) & (df['Time'] < end_time)]
    ```

- Filter the dataframe with regard to 'Direction' that the direction is heading north.

    ```
    # filter the dataframe with regard to 'Direction'
    df = df.loc[df['Direction'] == 1]
    ```

- Filter the dataframe with regard to 'Flag Text' that the day of the week is Friday.

    ```
    # filter the dataframe with regard to 'Flag Text'
    df = df.loc[df['Flag Text'] == 'Friday']
    ```

```
1  # record unique lane names
2  unique_lane_names = dict()
3  for key, df in dfs1.items():
4      # find unique lane names for the dataframe
5      unique_lane_name = df['Lane Name'].unique()
6
7      # record
8      unique_lane_names.update({key: unique_lane_name})
```

*Explanation*: This code block is an intermediate process to record unique lane names `unique_lane_names` of the two dataframes for the usage in the next code block.

**Result of** `unique_lane_names` :

```
{'1083': array(['NB_OS', 'NB_MID', 'NB_NS'], dtype=object),
 '1415': array(['NE_NS', 'NE_OS'], dtype=object)}
```

In [13]:

```python
1  # find average travel time of the lanes
2  avg_times_lanes = dict()    # the average travel times of the two dataframes
3  for key, df in dfs1.items():
4      avg_times_lanes_df = list()    # the average speeds of the dataframe
5      for lane_name in unique_lane_names[key]:
6          # filter the dataframe with regard to 'Lane Name'
7          filtered_df = df.loc[df['Lane Name'] == lane_name]
8
9          # calculate the average speed
10         avg_speed = filtered_df['Speed (mph)'].mean()
11
12         # calculate the average travel time in minute by the equation 'time =
   distance/speed'
13         #  4.86km = 3.019864miles
14         avg_time_lane = 3.019864 / avg_speed * 60
15
16         # record
17         avg_times_lanes_df.append(avg_time_lane)
18     avg_times_lanes.update({key: avg_times_lanes_df})
19
20  # print result of average speeds in different lanes
21  print(avg_times_lanes)
22
23  # calculate the average value of the result
24  templ = list(map(lambda l: (len(l), sum(l)), avg_times_lanes.values()))
25  sum_templ = tuple(map(sum, zip(*templ)))    # * is for unpacking the list  cool
26  avg_time = sum_templ[1]/sum_templ[0]
27
28  # print result of average time
29  print(avg_time)
```

```
{'1083': [6.009038525021391, 6.295469340374707, 7.279743256362566], '1415': [7.51940
2241603255, 6.6714418953907915]}
6.7550190517505415
```

*Explanation*:

This code block is about finding the average travel times in different lanes and calculating their average value. The two processes are introduced below.

- Find the average travel times in different lanes.
  - To begin with, a dictionary `avg_times_lanes` is created `avg_times_lanes = dict()` to store the name of the sites as the keys and their corresponding average values with regard to different lanes.
  - Then the dataframes are iterated. In each iteration, we iteratively filter the dataframe with regard to the lane names.

```python
for lane_name in unique_lane_names[key]:
    # filter the dataframe with regard to 'Lane Name'
    filtered_df = df.loc[df['Lane Name'] == lane_name]
```

- Next, the average speed of the `filtered_df` is calculated `avg_speed = filtered_df['Speed (mph)'].mean()`.
- Next, the average time was calculated by the average speed and the given distance, we firstly changed the unit of the distance from 'km' to 'miles' `avg_time_lane = 3.019864 / avg_speed * 60`.
- Finally we record the results, following is the integral process.

```python
avg_times_lanes_df = list()  # the average speeds of the dataframe
for lane_name in unique_lane_names[key]:
  # filter the dataframe with regard to 'Lane Name'
  filtered_df = df.loc[df['Lane Name'] == lane_name]

  # calculate the average speed
  avg_speed = filtered_df['Speed (mph)'].mean()

  # calculate the average travel time in minute by the equation 'time = distance/speed'
  #  4.86km = 3.019864miles
  avg_time_lane = 3.019864 / avg_speed * 60

  # record
  avg_times_lanes_df.append(avg_time_lane)
avg_times_lanes.update({key: avg_times_lanes_df})
```

- Calculate the average value of the average travel times.
  - For calculating the average value of the result, firstly the number of elements and the sum of the values of the elements are calculated and stored in `templ`.

    ```python
    templ = list(map(lambda l: (len(l), sum(l)), avg_times_lanes.values()))

    [(3, 19.584251121758662), (2, 14.190844136994047)]
    ```

  - Then, a map function is used to calculate the sum of values in the tuples in the list and finally calculate the average travel time by dividing the sum of speeds to the number of elements.

    ```python
    sum_templ = tuple(map(sum, zip(*templ)))  # * is for unpacking the list  cool
    avg_time = sum_templ[1]/sum_templ[0]
    ```

**Results of average travel time in different lanes**:

```
{'1083': [6.009038525021391, 6.295469340374707, 7.279743256362566],
 '1415': [7.519402241603255, 6.6714418953907915]}
```

**Result of overall average travel time**:

- 6.7550190517505415 minutes

**Interpretation of the result**:
The result shows that according to this dataset, the average travel time of north lanes between the two sites in the time from 17:00 to 18:00 in Friday is around 6.755 minutes.

# Task7.1

**Solution (i)**:

I can describe the *Row Completeness* as a set of fractions with regard to combinations of columns that can possibly have missing values, where the numerators are the numbers of rows that are with missing values in certain columns and the denominators are the total number of rows.

$$RowCompleteness_1 = \{RowCompleteness_C : C \text{ is the set of all possible combinations of columns}\}$$
$$RowCompleteness_C = NumberOfMissingRows_C \times 100 \div NumberOfRows$$

The row completeness value indicate the proportion of the number of not-missing values to the total number of values.

Besides we can calculate the row completeness of the proportion of the rows that have any missing values to the total number of rows.

$$RowCompleteness_2 = NumberOfRowsWithMissingValues \times 100 \div NumberOfRows$$

**Solution (ii)**:

In [14]:

```
1  # filter the dataframe as we did before
2  df3 = df1.copy(deep=True)
```

*Explanation*:
This code block is the data preparation operations for Task7.1. Since the filtering requirements in this task are the same as in Task5.1, we can reference to the codes in the Task5.1 section.

In [15]:

```
1  # find the columns that contain missing values
2  ifna_cols = df3.isna().any()
3  missing_cols = list(ifna_cols[ifna_cols == True].keys())
4
5  # find all possible combinations in the missing columns
6  # reference: https://stackoverflow.com/questions/464864/how-to-get-all-possible-combinations-
   of-a-list-s-elements?rq=1
7  combs = sum([list(map(list, combinations(missing_cols, i))) for i in range(len(missing_cols)
   + 1)], [])
8  combs.remove([])
```

*Explanation*:
Before we start the manipulations, we should know which columns contain missing values and find all possible combinations of them. As we do not need the empty set as a combination, we simply remove it from our combinations.

```
1  # initialize row completeness object
2  RC_count = dict()    # key: combination of columns, value: count
3
4  for comb in combs:
5      # filter the dataframe with regard to combination of columns
6      filtered_df3 = df3[comb]
7
8      # calculate the number of rows that all cells are of missing values
9      count = len(filtered_df3.loc[filtered_df3.isnull().all(axis=1)].index)
10
11     # record
12     RC_count.update({str(comb): count})
13
14 # calculate the row completeness
15 total_count = len(df3.index)
16 RC1 = dict(zip(RC_count.keys(), map(lambda x: 100 * int(x) / total_count,
   RC_count.values())))
17
18 # calculate the the row completeness for the entire dataframe
19 count = len(df3.loc[df3.isnull().any(axis=1)].index)
20 RC2 = 100 * count / total_count
```

*Explanation*:

This code block is about calculating the row completeness.

- First, we initialize a row completeness count dictionary that stores the combinations of columns as key and the count of the rows that cells in the combination of columns are missing.

```
RC_count = dict()    # key: combination of columns, value: count
```

- Then we iteratively filter `df3` with regard to the combination of columns `filtered_df3 = df3[comb]`, calculate the counts, and record them to `RC_count`.

```
for comb in combs:
    # filter the dataframe with regard to combination of columns
    filtered_df3 = df3[comb]

    # calculate the number of rows that all cells are of missing values
    count = len(filtered_df3.loc[filtered_df3.isnull().all(axis=1)].index)

    # record
    RC_count.update({str(comb): count})
RC_count
```

- Finally, we need to map `RC_count` according to the previously defined row completeness equation.

```
# calculate the row completeness
total_count = len(df3.index)
RC = dict(zip(RC_count.keys(), map(lambda x: 100 * int(x) / total_count, RC_count.val
ues())))
```

**Result of Row Completeness Calculation 1**:

{"['Speed (mph)']": 0.004474829086389061,
 "['Headway (s)']": 1.0346799254195151,
 "['Gap (s)']": 1.963455562461156,
 "['Speed (mph)', 'Headway (s)']": 0.0,
 "['Speed (mph)', 'Gap (s)']": 0.0,
 "['Headway (s)', 'Gap (s)']": 1.0346799254195151,
 "['Speed (mph)', 'Headway (s)', 'Gap (s)']": 0.0}

**Result of Row Completeness Calculation 2**:
1.967930391547545

# Discussion

### Interpretation of the Results

The result shows the values of row completeness considering all combinations of columns. Overall, the total rate of rows that have missing values is about 1.96793. For individual columns, the 'Speed (mph)' column has the minimum percentage of missing values, which is around 0.00447. And the column 'Gap (s)' has the highest proportion 1.96345. And we can infer from the result that for all the combinations that have the column 'Speed (mph)' and other columns, all the 'Speed (mph)' columns are filled with values. And in the combination ('Headway (s)', 'Gap (s)'), at least all the values in the column 'Headway (s)' are missing.

### Understanding of Completeness

The completeness is an important aspect in data quality. Data completeness can be defined that "the degree to which a data collection has values or all attributes of entities that are supposed to have values" [1]. In my interpretation, to see from the perspective of application, the completeness of data is that within the scale of the potential utilization of the data, values of objective attributes should be there for utilization. Hence completeness is inextricably linked with data quality. For implementation, in a dataset, the completeness can be demonstrated as the proportion of values that are stored to the number of possible records of the dataset. Or by Pipino, Lee, and Wang's definition [2], it is "a function of the missing values in a column of a table".

### Comparition of the two Completeness Formulas

I can compare the two formulas in terms of the difference between the representations of the completeness values. The row completeness demonstrates the completeness of values accross multiple columns, while the column completeness demonstrate the completeness of values in columns. Hence the row completeness formula can show the degree of completeness over the entire dataset, while the column completeness formula shows the degree of completeness in specific columns.For this dataset, because only three columns may have missing values, calculating the row completeness formula regarding combinations of columns are easy to implement, and the values calculated can provide a more thorough overview of the completeness than the ones calculated column completeness formula. We can obtain the overview of the rate of any values in these columns are missing as well as the inspect of a specific combination of columns that all the values inside them are missing. In the column completeness calculation, we can only inspect the value for an individual column.

[1] Li, C.: Computing complete answers to queries in the presence of limited access patterns. VLDB J. 12, 211–277 (2003)
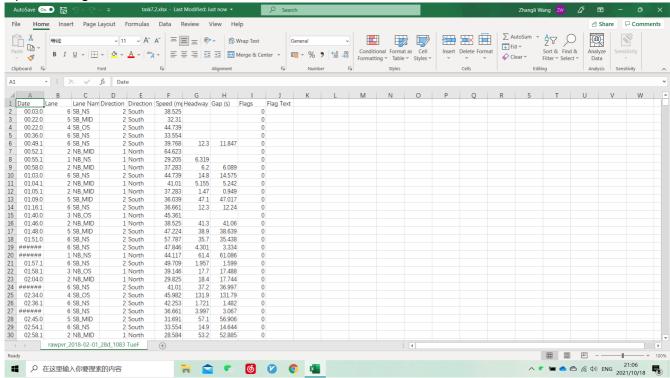[2] Pipino, L.L., Lee, Y.W., Wang, R.Y.: Data quality assessment. Commun. ACM 45, 211–218 (2002)

# Task7.2

**Solution**:
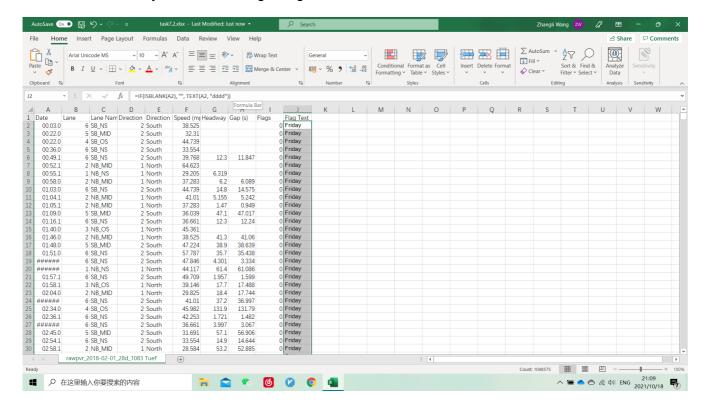
I use **Microsoft Excel** as the tool for this task.

1. Open the original dataset in Excel.



2. For the entire column of "Flag Text", apply the function

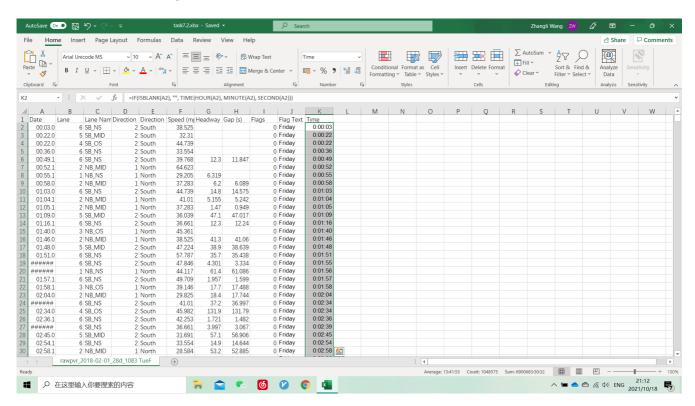    $$=IF(ISBLANK(A2), \text{""}, TEXT(A2, \text{"dddd"}))$$

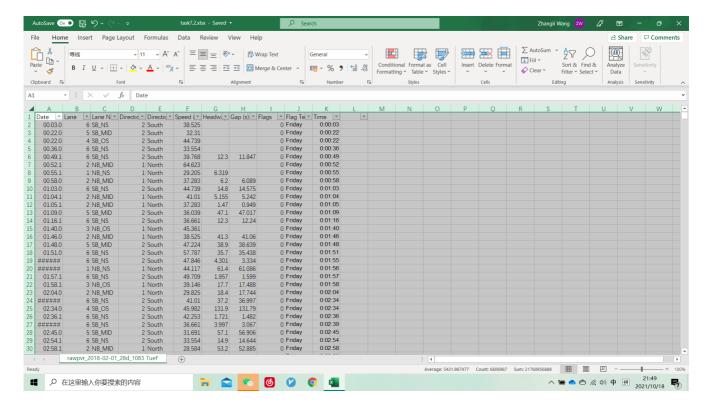, which calculates days of the week regarding the "Date" column.

3. Then we want to add a new column "Time" that stores the corresponding time in the "Date" column, apply this function

```
=IF(ISBLANK(A2), "", TIME(HOUR(A2), MINUTE(A2), SECOND(A2)))
```

to the column to extract hour, minute, and second from the corresponding date. Then, update the type of the column from default "General" to "Time".



4. Select the entire table, press `Ctrl` + `Shift` + `L` to initialize the filters. Then we can apply filtering to the table.

5. Filter the sheet that the time is between 17:00 and 18:00, the direction is 1, and the day of the week is Friday.

Customly filter the 'Time' to be between 17:00 and 18:00.



Filter the 'Direction Name' as 'North'.



Filter the 'Flag Text' as Friday.
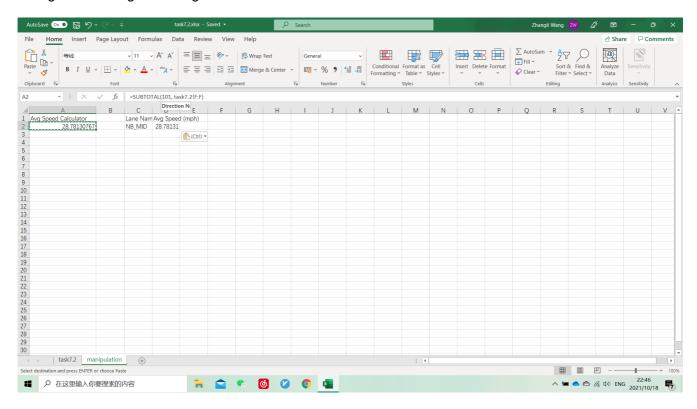
6. Then, filter the 'Lane Name' to 'NB_MID'.



7. Create a sheet named 'manipulation' for storing the results of calculating the average speeds and other results. Create a column named 'Avg Speed Calculator' as the label of calculating function
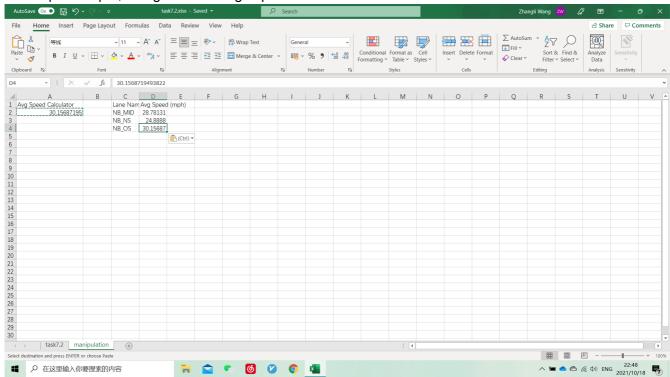
```
=SUBTOTAL(101, task7.2!F:F)
```

to calculate the average speed for 'NB_MID' according to the filtering. Note that the first arguement is the code of average function with ignoring hidden values, we use this one because filtering hide rows.

And we also add two columns 'Avg Speed (mph)' and 'Lane Name' for storing the average speeds for lanes and the lane names, respectively. Then we execute the function, copy and paste the value by plain text to
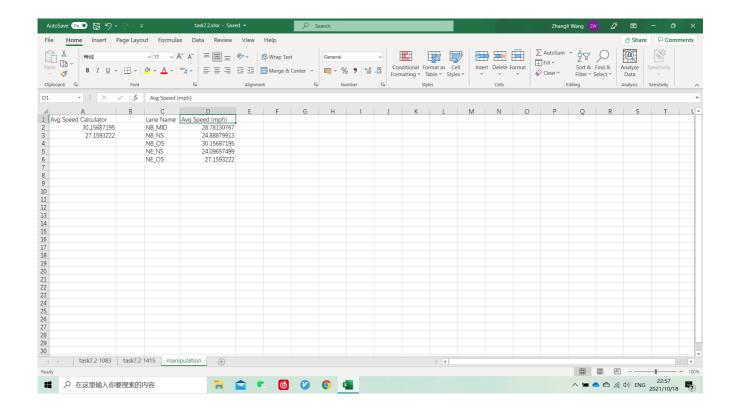
the corresponding place in the column 'Avg Speed (mph)'. We do this because the result of this function changes according to filtering.



8. Now repeat step 6, 7 to get the average speeds of all lanes.
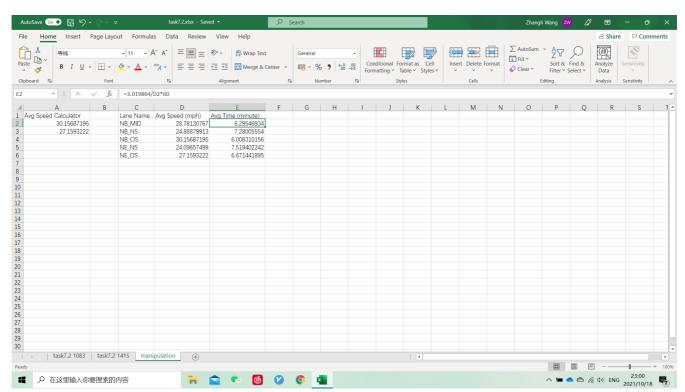


9. Load the dataset of site 1415 in a newly created sheet named `task7.2 1415` and repeat steps 1 to 8 for this sheet. After all steps are finished, we can acquire the average values of the five north lanes.
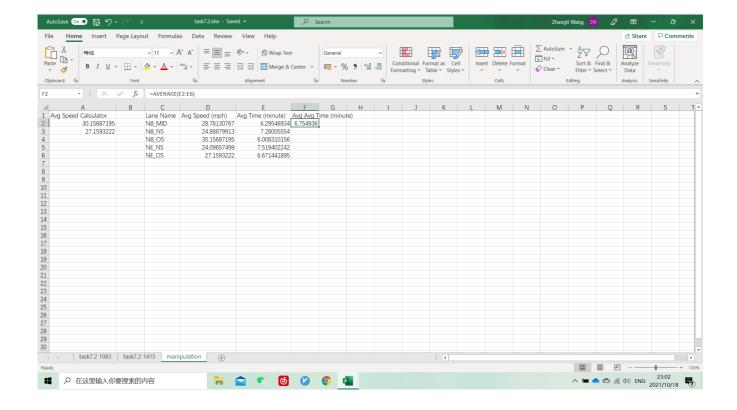
10. Now we can calculate the average time with the given distance and the speed with the function

   $=3.019864/D2*60$

.



Finally, we can compute the average over the average time by the function

   $=AVERAGE(E2:E6)$

.

# Comparition of the Two Technologies

Both of the two technologies have the ability to solve the tasks, they both have built-in functions for the same purposes in terms of data manipulation. And they can both manipulate the tabular dataset by filtering according to columns and aggregating according to rows.

To discuss about the differences, the manipulation logics are different. For python, the manipulation is realized by codes. For Excel, the manipulation is realized by the combination of tools in the interface, such as the tool for applying filters to the 'Lane Name' column, and the built-in functions. Python codes can be compiled statically while every operation in Excel is executed instantaneously.

For python, one advantage is that the process of manipulation is more interpretable and reviewable for a programmer. And using Jupyter Notebook to program this coursework can help me to divide operations into stages and add explanation everywhere I want. For example, for every task, I can split my manipulation operations in different code blocks by mini goals. Another advantage can be that the library pandas imported in python has efficient built-in functions such as `loc` that can automatically optimize the data manipulation operations. But the disadvantage is that bugs can possibly occur everywhere in my codes and the intermediate dataset cannot be entirely displayed during the manipulation.

For Excel, the advantage is that the interface is intuitive and can display the data in a proper tabular structure and the drawing of charts is also very convenient. And for manipulating an entire column, Excel can help me with it with its tools conveniently rather than coding for loops manually. But the disadvantages are also significant. Because most of my manipulations are functions, their dependencies result in everytime I change a part of the Excel file, all the functions are re-calculated, which is time consuming. Another disadvantage can be that using tools in the interface and built-in functions in combination can bring incompatibility in some cases. For example, for calculating the average speed for different lanes, I have to manually filter the table for several times and then record the values manually. When the dataset become larger and more complex, each operation executed in Excel can take much longer time due to the deep function dependencies and large quantity of calculations. And manually manipulate some variables can be cumbersome.