

```

1 import os, shutil
2 import glob
3 import pandas as pd
4 import time
5 import multiprocessing as mp
6 import gc
7
8
9 #tbl_temp = pd.DataFrame() #정제된 CSV 컨테이너 iteration 발생 때 마다 저장 하기 위함.
10 tbl_temp_final = pd.DataFrame() # 마지막 일별 데이터를 concatenate함 (일별 컨테이너)
11 tbl_final = pd.DataFrame() # 마지막 일별 데이터를 concatenate함 (일별 컨테이너)
12 path = os.getcwd()
13 print("폴더 경로 입력(예: c:\자료): ")
14 folderloc = input()
15
16 path_Done= folderloc + "\\\" + "Done"
17 #path_temp= folderloc + "\\\" + "temp"
18 isExist = os.path.exists(path_Done)
19 #printing if the path exists or not
20 if isExist==True:
21     print("이미 Done 폴더가 있습니다.")
22 isExist = os.path.exists(path_Done)
23 if not isExist:
24     os.makedirs(path_Done)
25     print("Done폴더 생성")
26
27 #path_final_temp= "\\nt1.koti.re.kr\\13 모빌리티빅데이터분석팀\\02. 개인별 폴더\\12. 이준성\\분석과정및결과\" + "\\\" + "final_temp"
28 path_final_temp= folderloc + "\\\" + "final_temp"
29 path_final_temp_sido_time_dow_age_tf_base= path_final_temp + "\\\" + "tbl_sido_time_dow_age_tf_base"
30 path_final_temp_tbl_ch_sido_tf_result= path_final_temp + "\\\" + "tbl_ch_sido_tf_result"
31
32 isExistfin = os.path.exists(path_final_temp)
33 if isExistfin==True:
34     print("이미 final_temp 폴더가 있습니다.")
35 isExistfin = os.path.exists(path_final_temp)
36 if not isExistfin:
37     os.makedirs(path_final_temp)
38     os.makedirs(path_final_temp_sido_time_dow_age_tf_base)
39     os.makedirs(path_final_temp_tbl_ch_sido_tf_result)
40
41     print("final_temp폴더(하위포함) 생성")
42
43 #path_final= "\\nt1.koti.re.kr\\13 모빌리티빅데이터분석팀\\02. 개인별 폴더\\12. 이준성\\분석과정및결과\" + "\\\" + "final"
44 path_final= folderloc + "\\\" + "final"
45
46 #path_final= "\\nt1.koti.re.kr\\13 모빌리티빅데이터분석팀\\02. 개인별 폴더\\12. 이준성\\분석과정및결과\" + "\\\" + "final"
47 isExistfin = os.path.exists(path_final)
48 if isExistfin==True:
49     print("이미 final 폴더가 있습니다.")
50 isExistfin = os.path.exists(path_final)
51 if not isExistfin:
52     os.makedirs(path_final)

```

This program is a script translated from R to Python.

**[R script has only Preprocessing]**

It's important to consider that there are significant differences between the workstation previously used with R and the PC (assigned to me) now using Python.

[Workstation -> PC]

CPU (i9-10900 -> i7-6700)

RAM (128GB -> 32GB)

GPU (RTX 3090 -> 960GTX)

Considerations made during the translation from R to Python include:

Analysis should not halt if Python crashes.

If a crash occurs, resume from the save point.

Use minimal RAM, preferably through CLI-based operations (avoid using libraries like PyQt5).

Graphs are outputted in Excel following the existing reporting format (Matplotlib is not used).

Continuous analysis from day 1 to day 365 without interruption.

The data spans three years, with daily data ranging from 6-8GB. For a year's worth, it's 2.5TB-3TB (many files exceed 7GB).

```

60 gc.enable()
61
62 #loop_chk = 0
63 for f in csv_files:
64     tbl_group_sido_tf_all_base = pd.DataFrame() #정제된 CSV 컨테이너 iteration 발생 때 마다 저장 하기 위함.
65     tbl_group_ch_sido_tf_result = pd.DataFrame()
66     tbl_group_sido_time_dow_age_tf_base = pd.DataFrame()
67     tbl_tf_all_base = pd.DataFrame()
68
69     #loop_chk = loop_chk + 1
70     # print the location and filename
71     print('파일 경로:', f)
72     print('파일 이름:', f.split("\\")[1])
73     tbl_chunk = pd.read_csv(f, chunksize=600000)
74     start_time = time.time() # 알고리즘 시작 시간
75     for chunk in tbl_chunk:
76         chunk['o_sido'] = chunk['o_sido'].astype("Int64") # o_sido float 문제 해결
77         chunk.loc[~(chunk['traffic_time_all_mean']>=5), 'traffic_all']=None
78         chunk.loc[~(chunk['traffic_time_all_mean']<180), 'traffic_all']=None
79         chunk.loc[~(chunk['traffic_time_30_mean']>=5), 'traffic_30']=None
80         chunk.loc[~(chunk['traffic_time_30_mean']<180), 'traffic_30']=None
81         chunk.loc[~(chunk['traffic_time_60_mean']>=5), 'traffic_60']=None
82         chunk.loc[~(chunk['traffic_time_60_mean']<180), 'traffic_60']=None
83
84         chunk.drop(chunk[(chunk['stay_traffic_volume'] == None) & (chunk['traffic_all'] == None) & (chunk['traffic_30'] == None) & (chunk['traffic_60'] == None)].index, inplace = True) # 같은 행에 여러 변수가 NULL일 경우 알림 -> 최종 결과값
85         chunk.drop(chunk[(chunk['stay_traffic_volume'] == 0) & (chunk['traffic_all'] == 0) & (chunk['traffic_30'] == 0) & (chunk['traffic_60'] == 0)].index, inplace = True) # 같은 행에 여러 변수가 NULL일 경우 알림 -> 최종 결과값: 같음
86
87         #tbl.loc[tbl['stay_traffic_volume'] == "", 'stay_traffic_volume'] = 3
88         chunk['stay_traffic_volume'] = chunk['stay_traffic_volume'].replace('', 3)
89         chunk['stay_traffic_volume'] = chunk.stay_traffic_volume.astype(float)
90         chunk['stay_traffic_volume'] = (chunk['stay_traffic_volume'] * chunk['stay_traffic_volume_cor_index'])
91         #tbl.loc[tbl['stay_traffic_volume'], 'stay_traffic_volume'] = (tbl['stay_traffic_volume'] * tbl['stay_traffic_volume_cor_index'])
92
93         chunk['TT_V_all'] = (chunk['stay_traffic_volume'] * chunk['traffic_all'])
94         chunk['TT_T_all'] = (chunk['traffic_time_all_mean'] * chunk['TT_V_all'])
95         chunk['TT_V_30'] = (chunk['stay_traffic_volume'] * chunk['traffic_30'])
96         chunk['TT_T_30'] = (chunk['traffic_time_30_mean'] * chunk['TT_V_30'])
97         chunk['TT_V_60'] = (chunk['stay_traffic_volume'] * chunk['traffic_60'])
98         chunk['TT_T_60'] = (chunk['traffic_time_60_mean'] * chunk['TT_V_60'])
99         tbl_tf_all_base = chunk[['o_sido', 'd_sido', 'o_type', 'd_type', 'o_time', 'd_time', 'o_date', 'o_dow', 'age', 'stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']] #시역별 동행량 비교 표 3-9를 제외한
100         tbl_group_sido_tf_all_base = pd.concat([tbl_group_sido_tf_all_base, tbl_tf_all_base])
101     del tbl_chunk
102     del chunk

```

When using Pandas on Windows to load a file of around 7GB into 32GB of RAM, the system, including background processes, ends up utilizing more than 32GB of RAM. Consequently, tools like Pyarrow were deemed unnecessary luxuries. Since there was no urgent need for speed (I could just set it running and leave it to finish), I opted for a method to address the issue of exceeding the 32GB RAM limit by using chunksize. I designed the logic to read and analyze the CSV file in chunks. Then, I analyzed each chunk, dropped the leftover scraps of tables to free up memory, and repeated this process until all analyzed files were concatenated together. As an extra measure, I considered the possibility of Python crashing due to this approach, so I attempted to save the analysis results in chunks to files. However, this method proved too time-consuming and was therefore excluded.

```

104 tbl_group_sido_time_dow_age_tf_base = tbl_group_sido_tf_all_base.groupby(['o_sido','d_sido','o_type','d_type','o_time','d_time','o_dow','age'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']]
105 tbl_group_sido_time_dow_age_tf_base.to_csv(path_final_temp + "\\ " + "tbl_sido_time_dow_age_tf_base\\tbl_sido_time_dow_age_tf_base_" + (f.split("\\")[-1]).replace("zip", "csv"))
106 del tbl_group_sido_time_dow_age_tf_base
107 del tbl_group_sido_tf_all_base
108 del tbl_tf_all_base
109
110 print("파일 로딩 및 정제에 걸린 시간: " + str(round((time.time() - start_time),2)) + "분")
111 shutil.move(folderloc+"\\ "+(f.split("\\")[-1]), folderloc + "\\ " + "Done" + "\\ " + (f.split("\\")[-1]))
112 print("작업한 파일[" + f.split("\\")[-1] + "]을 작업완료 폴더(Done) 안으로 옮김")
113
114 #####
115
116 print("<마지막 연산 시작>")
117
118 tbl_final = pd.DataFrame() # [일별 컨테이너 포맷] 메모리 관리
119 csv_temp_final_files = glob.glob(os.path.join(path_final_temp_sido_time_dow_age_tf_base, "*.csv"))
120 for fin in csv_temp_final_files:
121     tbl_temp_final = pd.DataFrame()
122     tbl_temp_final = pd.read_csv(fin, engine="pyarrow")
123     tbl_temp_final = tbl_temp_final.groupby(['o_sido','d_sido','o_type','d_type','o_time','d_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
124     tbl_final = pd.concat([tbl_final, tbl_temp_final], ignore_index=True)
125     del tbl_temp_final
126     print("시도별 OD " + fin + " 로딩 완료")
127
128 ##### Epoch

```

The analyzed chunk data is systematically saved into tables and then stored back into files. By repeating this process 365 to 366 times (for leap years), an entire year's worth of data undergoes preprocessing! As a result, a daily 7GB file is reduced to 300MB to 500MB after preprocessing. Therefore, preprocessed data for 365 days amounts to approximately 15-20GB, significantly reducing the original size from 3TB to 15-20GB.

Once preprocessing is completed, the script moves on to the post-processing stage, referencing the saved file locations.

It's noteworthy that although Pyarrow didn't support chunking at the time of script creation, rendering it unsuitable for 7GB files, it could still be used for files sized between 300-500MB. Therefore, for the sake of time efficiency, Pyarrow was utilized for preprocessing the files of this size.

```

129 # 시도별 od
130 final= tbl_final.groupby(['o_sido','d_sido','o_type', 'd_type', 'o_time', 'd_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
131 tbl_final['TT_T_all'] = (tbl_final['TT_T_all'] / tbl_final['TT_V_all'])
132 tbl_final['TT_T_30'] = (tbl_final['TT_T_30'] / tbl_final['TT_V_30'])
133 tbl_final['TT_T_60'] = (tbl_final['TT_T_60'] / tbl_final['TT_V_60'])
134 tbl_final_seoul_o = tbl_final.loc[(tbl_final['o_sido'] == 11), ['o_sido', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].apply(
135 tbl_final_seoul_od = pd.concat([tbl_final_seoul_o, tbl_final.loc[(tbl_final['d_sido'] == 11), ['d_sido', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
136 del tbl_final_seoul_o
137 tbl_final_Busan_o = tbl_final.loc[(tbl_final['o_sido'] == 21), ['o_sido', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].apply(
138 tbl_final_Busan_od = pd.concat([tbl_final_Busan_o, tbl_final.loc[(tbl_final['d_sido'] == 21), ['d_sido', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
139 del tbl_final_Busan_o
140 tbl_final_Daegu_o = tbl_final.loc[(tbl_final['o_sido'] == 22), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].apply(
141 tbl_final_Daegu_od = pd.concat([tbl_final_Daegu_o, tbl_final.loc[(tbl_final['d_sido'] == 22), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
142 del tbl_final_Daegu_o
143 tbl_final_Incheon_o = tbl_final.loc[(tbl_final['o_sido'] == 23), ['o_sido', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
144 tbl_final_Incheon_od = pd.concat([tbl_final_Incheon_o, tbl_final.loc[(tbl_final['d_sido'] == 23), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
145 del tbl_final_Incheon_o
146 tbl_final_Gwangju_o = tbl_final.loc[(tbl_final['o_sido'] == 24), ['o_sido', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
147 tbl_final_Gwangju_od = pd.concat([tbl_final_Gwangju_o, tbl_final.loc[(tbl_final['d_sido'] == 24), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
148 del tbl_final_Gwangju_o
149 tbl_final_Daejeon_o = tbl_final.loc[(tbl_final['o_sido'] == 25), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
150 tbl_final_Daejeon_od = pd.concat([tbl_final_Daejeon_o, tbl_final.loc[(tbl_final['d_sido'] == 25), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
151 del tbl_final_Daejeon_o
152 tbl_final_Ulsan_o = tbl_final.loc[(tbl_final['o_sido'] == 26), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].apply(
153 tbl_final_Ulsan_od = pd.concat([tbl_final_Ulsan_o, tbl_final.loc[(tbl_final['d_sido'] == 26), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
154 del tbl_final_Ulsan_o
155 tbl_final_Sejong_o = tbl_final.loc[(tbl_final['o_sido'] == 29), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].apply(
156 tbl_final_Sejong_od = pd.concat([tbl_final_Sejong_o, tbl_final.loc[(tbl_final['d_sido'] == 29), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
157 del tbl_final_Sejong_o
158 tbl_final_Gyeonggi_o = tbl_final.loc[(tbl_final['o_sido'] == 31), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
159 tbl_final_Gyeonggi_od = pd.concat([tbl_final_Gyeonggi_o, tbl_final.loc[(tbl_final['d_sido'] == 31), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
160 del tbl_final_Gyeonggi_o
161 tbl_final_Gangwon_o = tbl_final.loc[(tbl_final['o_sido'] == 32), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
162 tbl_final_Gangwon_od = pd.concat([tbl_final_Gangwon_o, tbl_final.loc[(tbl_final['d_sido'] == 32), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
163 del tbl_final_Gangwon_o
164 tbl_final_Chungbuk_o = tbl_final.loc[(tbl_final['o_sido'] == 33), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
165 tbl_final_Chungbuk_od = pd.concat([tbl_final_Chungbuk_o, tbl_final.loc[(tbl_final['d_sido'] == 33), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
166 del tbl_final_Chungbuk_o
167 tbl_final_Chungnam_o = tbl_final.loc[(tbl_final['o_sido'] == 34), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
168 tbl_final_Chungnam_od = pd.concat([tbl_final_Chungnam_o, tbl_final.loc[(tbl_final['d_sido'] == 34), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
169 del tbl_final_Chungnam_o
170 tbl_final_Jeonbuk_o = tbl_final.loc[(tbl_final['o_sido'] == 35), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
171 tbl_final_Jeonbuk_od = pd.concat([tbl_final_Jeonbuk_o, tbl_final.loc[(tbl_final['d_sido'] == 35), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
172 del tbl_final_Jeonbuk_o
173 tbl_final_Jeonnam_o = tbl_final.loc[(tbl_final['o_sido'] == 36), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
174 tbl_final_Jeonnam_od = pd.concat([tbl_final_Jeonnam_o, tbl_final.loc[(tbl_final['d_sido'] == 36), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
175 del tbl_final_Jeonnam_o
176 tbl_final_Gyeongbuk_o = tbl_final.loc[(tbl_final['o_sido'] == 37), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
177 tbl_final_Gyeongbuk_od = pd.concat([tbl_final_Gyeongbuk_o, tbl_final.loc[(tbl_final['d_sido'] == 37), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
178 del tbl_final_Gyeongbuk_o
179 tbl_final_Gyeongnam_o = tbl_final.loc[(tbl_final['o_sido'] == 38), ['o_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].groupby(['o_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60','TT_T_60']].appl
180 tbl_final_Gyeongnam_od = pd.concat([tbl_final_Gyeongnam_o, tbl_final.loc[(tbl_final['d_sido'] == 38), ['d_sido', 'stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido'])[['stay_traffic volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply(
181 del tbl_final_Gyeongnam_o

```

However, there was a lot of computation to be done during the post-processing stage. Unfortunately, the elements to be analyzed were scattered all over the place, making it difficult to use for loops or regular expressions. So, I just resorted to using Ctrl + C and Ctrl + V.

```

185 del tbl_final_Jeju_o
186 tbl_final_sido_concat_od = pd.concat([tbl_final_Seoul_od,tbl_final_Busan_od,tbl_final_Daegu_od,tbl_final_Incheon_od,tbl_final_Gwangju_od,tbl_final_Daejeon_od,tbl_final_Ulsan_od,tbl_final_Sejong_od,tbl_final_Gyeonggi_od,tbl_final_Gangwon_od,
187 with pd.ExcelWriter(path_final + "\\\"+ 'final_tbl_sido_tf_base.xlsx') as excel: # 시도별 OD
188     final.to_excel(excel, sheet_name="All_sido")
189     tbl_final_sido_concat_od.to_excel(excel, sheet_name="All_sido_concat")
190     tbl_final_Seoul_od.to_excel(excel, sheet_name="서울")
191     tbl_final_Busan_od.to_excel(excel, sheet_name="부산")
192     tbl_final_Daegu_od.to_excel(excel, sheet_name="대구")
193     tbl_final_Incheon_od.to_excel(excel, sheet_name="인천")
194     tbl_final_Gwangju_od.to_excel(excel, sheet_name="광주")
195     tbl_final_Daejeon_od.to_excel(excel, sheet_name="대전")
196     tbl_final_Ulsan_od.to_excel(excel, sheet_name="울산")
197     tbl_final_Sejong_od.to_excel(excel, sheet_name="세종")
198     tbl_final_Gyeonggi_od.to_excel(excel, sheet_name="경기")
199     tbl_final_Gangwon_od.to_excel(excel, sheet_name="강원")
200     tbl_final_Chungbuk_od.to_excel(excel, sheet_name="충북")
201     tbl_final_Chungnam_od.to_excel(excel, sheet_name="충남")
202     tbl_final_Jeonbuk_od.to_excel(excel, sheet_name="전북")
203     tbl_final_Jeonnang_od.to_excel(excel, sheet_name="전남")
204     tbl_final_Gyeongbuk_od.to_excel(excel, sheet_name="경북")
205     tbl_final_Gyeongnam_od.to_excel(excel, sheet_name="경남")
206     tbl_final_Jeju_od.to_excel(excel, sheet_name="제주")
207
208 print("final_tbl_sido_tf_base 연산 완료")
209 del tbl_temp_final
210 del tbl_final
211 del csv_temp_final_files
212
213 tbl_final = pd.DataFrame() # (일별 컨테이너 포맷)
214 csv_temp_final_files = glob.glob(os.path.join(path_final_temp_sido_time_dow_age_tf_base, "*.csv"))
215 for fin in csv_temp_final_files:
216     tbl_temp_final = pd.read_csv(fin, engine="pyarrow")
217     tbl_final = pd.concat([tbl_final, tbl_temp_final], ignore_index=True) #요일별 출퇴근

```

The first-stage post-processing resulted in a file named 'Commute by City' that needed to be split and saved. I divided it into sheets within one Excel file.

However, after the first-stage post-processing was completed... there was more computation to be done. There are various other elements that need to be extracted from the preprocessed file, such as commuting by day of the week, commuting by hour, commuting time by age group, and others. Different data computation formulas are involved in various scenarios, so I began the computation using the preprocessed file.



```

processing the preprocessed files for each name as many times as there are names in the 'final_tbl_' name base.
214 csv_temp_final_files = glob.glob(os.path.join(path_final_temp_sido_time_dow_age_tf_base, "*.csv"))
215 for fin in csv_temp_final_files:
216     tbl_temp_final = pd.read_csv(fin, engine="pyarrow")
217     tbl_final = pd.concat([tbl_final, tbl_temp_final], ignore_index=True)  #요일별 출퇴근
218 ##### Epoch
219 #tbl_final = tbl_final.loc[(tbl_final['o_type'] == "H") & (tbl_final['d_type'] == "C") | (tbl_final['o_type'] == "C") & (tbl_final['d_type'] == "H"), ['o_dow', 'o_time', 'o_type', 'd_type', 'stay_traffic_volume', 'TT_V_all', 'TT_T_all',
220 tbl_final=tbl_final.groupby(['o_sido', 'd_sido', 'o_dow', 'o_time', 'd_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
221 tbl_final['TT_T_all'] = (tbl_final['TT_T_all'] / tbl_final['TT_V_all'])
222 tbl_final['TT_T_30'] = (tbl_final['TT_T_30'] / tbl_final['TT_V_30'])
223 tbl_final['TT_T_60'] = (tbl_final['TT_T_60'] / tbl_final['TT_V_60'])
224 tbl_final = tbl_final.reset_index()
225 tbl_final_sun_o = tbl_final.loc[(tbl_final['o_dow'] == 1), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
226 tbl_final_sun_od = pd.concat([tbl_final_sun_o, tbl_final.loc[(tbl_final['o_dow'] == 1), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
227 tbl_final_mon_o = tbl_final.loc[(tbl_final['o_dow'] == 2), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
228 tbl_final_mon_od = pd.concat([tbl_final_mon_o, tbl_final.loc[(tbl_final['o_dow'] == 2), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
229 tbl_final_tue_o = tbl_final.loc[(tbl_final['o_dow'] == 3), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
230 tbl_final_tue_od = pd.concat([tbl_final_tue_o, tbl_final.loc[(tbl_final['o_dow'] == 3), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
231 tbl_final_wed_o = tbl_final.loc[(tbl_final['o_dow'] == 4), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
232 tbl_final_wed_od = pd.concat([tbl_final_wed_o, tbl_final.loc[(tbl_final['o_dow'] == 4), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
233 tbl_final_thu_o = tbl_final.loc[(tbl_final['o_dow'] == 5), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
234 tbl_final_thu_od = pd.concat([tbl_final_thu_o, tbl_final.loc[(tbl_final['o_dow'] == 5), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
235 tbl_final_fri_o = tbl_final.loc[(tbl_final['o_dow'] == 6), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
236 tbl_final_fri_od = pd.concat([tbl_final_fri_o, tbl_final.loc[(tbl_final['o_dow'] == 6), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
237 tbl_final_sat_o = tbl_final.loc[(tbl_final['o_dow'] == 7), ['o_dow', 'o_sido', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['o_sido', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].apply
238 tbl_final_sat_od = pd.concat([tbl_final_sat_o, tbl_final.loc[(tbl_final['o_dow'] == 7), ['o_dow', 'd_sido', 'd_time', 'stay_traffic_volume', 'TT_V_all', 'TT_V_30', 'TT_V_60']].groupby(['d_sido', 'd_time'])[['stay_traffic_volume', 'TT_V_all',
239
240 tbl_final_dow_concat = pd.concat([tbl_final_sun_od, tbl_final_mon_od, tbl_final_tue_od, tbl_final_wed_od, tbl_final_thu_od, tbl_final_fri_od, tbl_final_sat_od])
241
242 with pd.ExcelWriter(path_final + "\\\"+ 'final_tbl_dow_base.xlsx') as excel_dow: # 시도별 OD
243     final.to_excel(excel_dow, sheet_name="All_sido")
244     tbl_final_dow_concat.to_excel(excel_dow, sheet_name="요일별 통합")
245     tbl_final_sun_od.to_excel(excel_dow, sheet_name="일")
246     tbl_final_mon_od.to_excel(excel_dow, sheet_name="월")
247     tbl_final_tue_od.to_excel(excel_dow, sheet_name="화")
248     tbl_final_wed_od.to_excel(excel_dow, sheet_name="수")
249     tbl_final_thu_od.to_excel(excel_dow, sheet_name="목")
250     tbl_final_fri_od.to_excel(excel_dow, sheet_name="금")
251     tbl_final_sat_od.to_excel(excel_dow, sheet_name="토")
252
253 print("final_tbl_dow_base 연산 완료")
254 del tbl_temp_final
255 del tbl_final
256 del csv_temp_final_files
257

```

processing the preprocessed files for each name as many times as there are names in the 'final\_tbl\_' name base.

```

259 tbl_final = pd.DataFrame() # (일별 컨테이너 포맷)
260 csv_temp_final_files = glob.glob(os.path.join(path_final_temp_sido_time_dow_age_tf_base, "*.csv"))
261 for fin in csv_temp_final_files:
262     tbl_temp_final = pd.read_csv(fin, engine="pyarrow")
263     tbl_final = pd.concat([tbl_final, tbl_temp_final], ignore_index=True)
264 ##### Epoch
265 #tbl_final = tbl_final.loc[(tbl_final['o_type'] == "H") & (tbl_final['d_type'] == "C") | (tbl_final['o_type'] == "C") & (tbl_final['d_type'] == "H"), ['o_sido', 'o_type', 'd_type', 'stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30',
266 tbl_final= tbl_final.groupby(['o_sido', 'd_sido', 'o_type', 'd_type'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
267 tbl_final['TT_T_all'] = (tbl_final['TT_T_all'] / tbl_final['TT_V_all'])
268 #tbl_group.loc['TT_T_all'] = (tbl_df['TT_T_all'] / tbl_df['TT_V_all'])
269 tbl_final['TT_T_30'] = (tbl_final['TT_T_30'] / tbl_final['TT_V_30'])
270 #tbl_group.loc['TT_T_30'] = (tbl_df['TT_T_30'] / tbl_df['TT_V_30'])
271 tbl_final['TT_T_60'] = (tbl_final['TT_T_60'] / tbl_final['TT_V_60'])
272 #tbl_group.loc['TT_T_60'] = (tbl_df['TT_T_60'] / tbl_df['TT_V_60'])
273 tbl_final.to_excel(path_final + "\\\"+ 'final_tbl_sido_time_tf_base.xlsx')
274
275 print("final_tbl_sido_time_tf_base 연산 완료")
276 del tbl_temp_final
277 del tbl_final
278 del csv_temp_final_files
279
280
281 tbl_final = pd.DataFrame() # (일별 컨테이너 포맷)
282 csv_temp_final_files = glob.glob(os.path.join(path_final_temp_sido_time_dow_age_tf_base, "*.csv"))
283 for fin in csv_temp_final_files:
284     tbl_temp_final = pd.read_csv(fin, engine="pyarrow")
285     tbl_final = pd.concat([tbl_final, tbl_temp_final], ignore_index=True) #연령별 출퇴근 통행비율
286 ##### Epoch
287 #tbl_final = tbl_final.loc[(tbl_final['o_type'] == "H") & (tbl_final['d_type'] == "C") | (tbl_final['o_type'] == "C") & (tbl_final['d_type'] == "H"), ['age', 'o_time', 'o_type', 'd_type', 'stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30',
288 tbl_final= tbl_final.groupby(['o_type', 'd_type', 'age', 'o_time'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
289 tbl_final['TT_T_all'] = (tbl_final['TT_T_all'] / tbl_final['TT_V_all'])
290 #tbl_group.loc['TT_T_all'] = (tbl_df['TT_T_all'] / tbl_df['TT_V_all'])
291 tbl_final['TT_T_30'] = (tbl_final['TT_T_30'] / tbl_final['TT_V_30'])
292 #tbl_group.loc['TT_T_30'] = (tbl_df['TT_T_30'] / tbl_df['TT_V_30'])
293 tbl_final['TT_T_60'] = (tbl_final['TT_T_60'] / tbl_final['TT_V_60'])
294 #tbl_group.loc['TT_T_60'] = (tbl_df['TT_T_60'] / tbl_df['TT_V_60'])
295 tbl_final.to_excel(path_final + "\\\"+ 'final_tbl_age_base.xlsx')
296 print("final_tbl_age_base 연산 완료")
297 del tbl_temp_final
298 del tbl_final
299 del csv_temp_final_files
300
301
302     # tbl_group_ch_sido_tf_result = pd.concat([tbl_group_ch_sido_tf_result, tbl_ch_sido_tf_result])
303     # tbl_group_ch_sido_tf_result = tbl_group_ch_sido_tf_result.groupby(['o_sido', 'o_type', 'd_type'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
304     # tbl_group_ch_sido_tf_result.to_csv(path_final_temp + "\\\" + "tbl_ch_sido_tf_result\\tbl_ch_sido_tf_result_" + str(chunkcount) + "_" + (f.split("\\")[1]).replace("zip", "csv"))
305     # del tbl_group_ch_sido_tf_result
306

```

```

307 tbl_final = pd.DataFrame() # <일별 컨테이너 포맷>
308 csv_temp_final_files = glob.glob(os.path.join(path_final_temp_sido_time_dow_age_tf_base, "*.csv"))
309 for fin in csv_temp_final_files:
310     tbl_temp_final = pd.read_csv(fin, engine="pyarrow")
311     tbl_final = pd.concat([tbl_final, tbl_temp_final], ignore_index=True)
312 ##### Epoch
313 #tbl_final = tbl_final.loc[(tbl_final['o_type'] == "H") & (tbl_final['d_type'] == "C") | (tbl_final['o_type'] == "C") & (tbl_final['d_type'] == "H"), ['o_sido', 'o_type', 'd_type', 'stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30',
314 tbl_final = tbl_final.loc[(tbl_final['o_type'] == "H") & (tbl_final['d_type'] == "C") | (tbl_final['o_type'] == "C") & (tbl_final['d_type'] == "H"), ['o_sido', 'o_type', 'd_type', 'o_date', 'o_time', 'stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
315 tbl_final = tbl_final.groupby(['o_sido', 'o_type', 'd_type'])[['stay_traffic_volume', 'TT_V_all', 'TT_T_all', 'TT_V_30', 'TT_T_30', 'TT_V_60', 'TT_T_60']].apply(sum)
316 tbl_final['TT_T_all'] = (tbl_final['TT_T_all'] / tbl_final['TT_V_all'])
317 #tbl_group.loc['TT_T_all'] = (tbl_df['TT_T_all'] / tbl_df['TT_V_all'])
318 tbl_final['TT_T_30'] = (tbl_final['TT_T_30'] / tbl_final['TT_V_30'])
319 #tbl_group.loc['TT_T_30'] = (tbl_df['TT_T_30'] / tbl_df['TT_V_30'])
320 tbl_final['TT_T_60'] = (tbl_final['TT_T_60'] / tbl_final['TT_V_60'])
321 #tbl_group.loc['TT_T_60'] = (tbl_df['TT_T_60'] / tbl_df['TT_V_60'])
322 tbl_final.to_excel(path_final + "\\\"+ 'final_tbl_ch_sido_tf_result.xlsx')
323 print("final_tbl_ch_sido_tf_result <원산> <완료>")
324 del tbl_temp_final
325 del tbl_final
326 del csv_temp_final_files
327
328
329 #print(tbl_group) # <엑셀에서 열기 힘들 경우>
330 #tbl_final = tbl_final.reset_index(drop=False, inplace=True) #index <추가하고 싶을 시 사용>
331 print("<전체 작업에 걸린시간>: " + str(round((time.time() - start_time),2)) + "초")
332 os.startfile(path_final + "\\") # <폴더 위치 열기>
333 #os.startfile(path_final + "\\\"+ 'final.csv') # <변환된 파일 연결프로그램으로 열기>
334
335
336 #print(np.where((tbl_date_base['TT_T_30']) != (tbl_date_base['TT_T_60'])))
337
338
339 #print(tbl_df['traffic_time_all_mean'])
340
341 #tbl_grp = tbl.groupby(['o_sido', 'o_type', 'd_type', 'o_date'])
342 #print(tbl_grp)

```

I executed this program on the PC allocated by the company. Since there was only one PC assigned, I had to wait while performing this task without being able to do anything else. I couldn't even open Chrome; I just ran this program and if it crashed, I would troubleshoot, improve the code, and start again. Looking back at the code I wrote in November 2022, memories flood back even though it's been over a year and a half. Running computations on that computer took about three days. It had an i7-6700 processor, 32GB of RAM, and a 960GTX graphics card. I also wrote a CUDA version of the script, but it didn't work properly because the VRAM of the 960GTX was limited to 2GB. I couldn't even read the post-processed files properly. To optimize the computer, I completed the script using the concepts of Epoch and iteration that I learned during a year-long deep learning project in college. It took me 15 days to develop the script while experimenting with various CUDA implementations and other techniques. If ChatGPT had been available back then, I wonder if it would have taken only a day to accomplish the same tasks. ~~There's a lot to say about the original R code, but I'll omit that for now.~~



# Execution

```
파일 정제에 걸린 시간: 6285.04초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221021_masking.csv
파일 이름: koti_final_all_20221021_masking.csv
파일 로딩에 걸린 시간: 6553.42초
koti_final_all_20221021_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 6658.57초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221022_masking.csv
파일 이름: koti_final_all_20221022_masking.csv
파일 로딩에 걸린 시간: 6909.51초
koti_final_all_20221022_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 6980.29초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221023_masking.csv
파일 이름: koti_final_all_20221023_masking.csv
파일 로딩에 걸린 시간: 7202.49초
koti_final_all_20221023_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 7252.06초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221024_masking.csv
파일 이름: koti_final_all_20221024_masking.csv
파일 로딩에 걸린 시간: 7518.24초
koti_final_all_20221024_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 7618.84초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221025_masking.csv
파일 이름: koti_final_all_20221025_masking.csv
파일 로딩에 걸린 시간: 7894.78초
koti_final_all_20221025_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 7997.4초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221026_masking.csv
파일 이름: koti_final_all_20221026_masking.csv
파일 로딩에 걸린 시간: 8282.7초
koti_final_all_20221026_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 8389.1초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221027_masking.csv
파일 이름: koti_final_all_20221027_masking.csv
파일 로딩에 걸린 시간: 8684.63초
koti_final_all_20221027_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 8786.88초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221028_masking.csv
파일 이름: koti_final_all_20221028_masking.csv
파일 로딩에 걸린 시간: 9098.19초
koti_final_all_20221028_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 9199.26초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221029_masking.csv
파일 이름: koti_final_all_20221029_masking.csv
파일 로딩에 걸린 시간: 9505.43초
koti_final_all_20221029_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 9568.56초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221030_masking.csv
파일 이름: koti_final_all_20221030_masking.csv
파일 로딩에 걸린 시간: 9820.06초
koti_final_all_20221030_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 9861.44초
파일 경로: C:\Users\Jun\Desktop\KOTI\자료분석\1달 자료\koti_final_all_20221031_masking.csv
파일 이름: koti_final_all_20221031_masking.csv
파일 로딩에 걸린 시간: 10175.54초
koti_final_all_20221031_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 10249.16초
<마지막 연산 시작>
전체 작업에 걸린 시간: 10249.45초
```

```
builder@ASUS-G14-2021: /mn  X + v

파일 이름: temp_koti_final_all_20221023_masking.csv
파일 로딩에 걸린 시간: 519.65초
temp_koti_final_all_20221023_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 519.75초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221024_masking.csv
파일 이름: temp_koti_final_all_20221024_masking.csv
파일 로딩에 걸린 시간: 549.42초
temp_koti_final_all_20221024_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 549.57초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221025_masking.csv
파일 이름: temp_koti_final_all_20221025_masking.csv
파일 로딩에 걸린 시간: 579.99초
temp_koti_final_all_20221025_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 580.15초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221026_masking.csv
파일 이름: temp_koti_final_all_20221026_masking.csv
파일 로딩에 걸린 시간: 610.75초
temp_koti_final_all_20221026_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 610.91초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221027_masking.csv
파일 이름: temp_koti_final_all_20221027_masking.csv
파일 로딩에 걸린 시간: 641.39초
temp_koti_final_all_20221027_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 641.56초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221028_masking.csv
파일 이름: temp_koti_final_all_20221028_masking.csv
파일 로딩에 걸린 시간: 669.91초
temp_koti_final_all_20221028_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 670.05초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221029_masking.csv
파일 이름: temp_koti_final_all_20221029_masking.csv
파일 로딩에 걸린 시간: 682.82초
temp_koti_final_all_20221029_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 682.92초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221030_masking.csv
파일 이름: temp_koti_final_all_20221030_masking.csv
파일 로딩에 걸린 시간: 692.17초
temp_koti_final_all_20221030_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 692.25초
파일 경로: /mnt/c/Users/Jun/Desktop/KOTI/자료분석/1달 자료/temp/temp_koti_final_all_20221031_masking.csv
파일 이름: temp_koti_final_all_20221031_masking.csv
파일 로딩에 걸린 시간: 704.0초
temp_koti_final_all_20221031_masking.csv에 대한 데이터 정제 작업 완료
파일 정제에 걸린 시간: 704.08초
<마지막 연산 시작>
전체 작업에 걸린 시간: 705.43초
```