

```
In [1]: import numpy as np
import pandas as pd

import requests
import time
from bs4 import BeautifulSoup
from selenium import webdriver
from selenium.webdriver.common.by import By
from selenium.webdriver.common.keys import Keys
import json
import googlemaps
import pprint

import matplotlib.pyplot as plt
import seaborn as sns

sns.set()

#----- 차트 관련 속성 (한글처리, 그리드) -----
plt.rcParams['font.family'] = 'Malgun Gothic'
plt.rcParams['axes.unicode_minus'] = False

#----- 주피언, 출력결과 넓이 늘리기 -----
# from IPython.core.display import display, HTML
from IPython.display import display, HTML

display(HTML("<style>.container{width:100% !important;}</style>"))
pd.set_option('display.max_rows', 100)
pd.set_option('display.max_columns', 100)
pd.set_option('max_colwidth', None)

import warnings
warnings.filterwarnings(action='ignore')
```

```
In [2]: # #----- 크롬 옵션 객체 생성
# # options = webdriver.ChromeOptions()
# # options.add_argument("window-size=1000x800") # 화면크기(윈도우화면)
# # user_agent = "Mozilla/5.0 (Windows NT 6.0; WOW64; AppleWebKit/537.36 (KHTML, like Gecko) Chrome/37.0.2484.0 Safari/537.36 "
# # options.add_argument('user-agent' + user_agent)
# # options.add_argument('headless') # headless 모드 실행
# # options.add_argument("disable-opts")
# # options.add_argument("disable-infobars")
# # options.add_argument("--disable-extensions")
# # options.add_argument("--mute-audio") mouse
# # options.add_argument('--blink-settings=imagesEnabled=false') #브라우저에서 이미지 로딩을 하지 않습니다.
# # options.add_argument('incognito') #익스플로러 모드의 브라우저가 실행됩니다.
# # options.add_argument("--start-maximized")
# # driver = webdriver.Chrome('.\chromedriver_102.0.5895.27.exe', options=options)

# #----- 크롬 드라이버 로드 110.0.5481.177
# # https://chromedriver.chromium.org/downloads
# # https://chromedriver.storage.googleapis.com/index.html?path=110.0.5481.77/
# #-----
# driver = webdriver.Chrome('.\DATA_COLLECTION\chromedriver_110.exe')
# driver.get("https://api.visitjeju.net/api/contents/list?_siteId=jeju&locale=kr&device=pc&cacheId=cache000000002&tag=&sorting=reviewcnt+desc,&title_kr=as&regionId=4&region2cd=4&pageSize=12&page=1")

# #----- 스크롤 다운
# # driver.find_element(By.TAG_NAME, "body").send_keys(Keys.END)

# htmlstr = driver.page_source
# htmlstr = htmlstr.replace("\n", "").replace("<br>", "")
# print(htmlstr)
```

Visit Jeju

- API : <https://m.visitjeju.net/kr/visitJejuApi>
- 여행 : cate1cd=cate0000000002 (/datasets_jeju02.여행.json) 360 rows
- 쇼핑 : cate1cd=cate0000000003 (/datasets_jeju03.쇼핑.json) 146 rows
- 숙박 : cate1cd=cate0000000004 (/datasets_jeju04.숙박.json) 360 rows
- 음식 : cate1cd=cate0000000005 (/datasets_jeju05.음식.json) 360 rows
- 여행 : https://api.visitjeju.net/api/contents/list?_siteId=jeju&locale=kr&device=pc&cate1cd=cate0000000002&tag=&sorting=reviewcnt+desc,&title_kr=as&C2%AEion1cd=%C2%AEion2cd=&pageSize=360&page=1
- 쇼핑 : https://api.visitjeju.net/api/contents/list?_siteId=jeju&locale=kr&device=pc&cate1cd=cate0000000003&tag=&sorting=reviewcnt+desc,&title_kr=as&C2%AEion1cd=%C2%AEion2cd=&pageSize=360&page=1
- 숙박 : https://api.visitjeju.net/api/contents/list?_siteId=jeju&locale=kr&device=pc&cate1cd=cate0000000004&tag=&sorting=reviewcnt+desc,&title_kr=as&C2%AEion1cd=%C2%AEion2cd=&pageSize=360&page=1
- 음식 : https://api.visitjeju.net/api/contents/list?_siteId=jeju&locale=kr&device=pc&cate1cd=cate0000000005&tag=&sorting=reviewcnt+desc,&title_kr=as&C2%AEion1cd=%C2%AEion2cd=&pageSize=360&page=1

```
In [3]: col_list = ["contentsid", "alltag", "label", "title", "address", "tag", "introduction",
               "readcnt", "likecnt", "reviewcnt", "markcnt", "snssharecnt", "schedulecnt", "visitcnt", "eveipcnt",
               "latitude", "longitude", "phonoeno", "sbst", "img", "thumb"]

file_list = ["02.여행", "03.쇼핑", "04.숙박", "05.음식"]

df_list = []

for fname in file_list :
    with open(r'./datasets_jeju/{fname}.json', "r", encoding='utf-8') as f:
        data_txt = f.read()
        #print(data_txt)
        dic = json.loads(data_txt)
        print ( len(dic["items"]) )

#-----
data_list = []
for item in dic["items"]:
    data_list.append([
        item["contentsid"], #pk
        item["alltag"], #text-----
        item["contentscd"]['label'],
        item["title"], #text-----
        item["address"], #text-----
        item["tag"], #text-----
        item["introduction"], #text-----
        item["readcnt"], #687790,
        item["likecnt"], #222,
        item["reviewcnt"], #389,
        item["markcnt"], #3979,
        item["snssharecnt"], #881,
        item["schedulecnt"], #9,
        item["visitcnt"], #21,
        item["eveipcnt"], #9,
        item["latitude"], #33.462147,
        item["longitude"], #126.936424,
        item["phonoeno"], #7064-793-8959,
        item["sbst"], #text-----
        item["repPhoto"]['photoId']["imgpath"],
        item["repPhoto"]['photoId']["thumbnailipath"],
    ])
df = pd.DataFrame(data_list, columns=col_list)
df.to_csv('df_list.csv')
df_list.append(df)

for key in dic["items"][0].keys():
    print(f'{key}',',', end=" ")

print("")
df = pd.concat(df_list)
print(df.shape)
df.to_csv(r'./datasets_jeju/data.csv', index=False)
```

```
In [4]: null_idx = df[df['latitude'].isna()].index.values
null_idx
```

```
Out[4]: array([ 41, 262], dtype=int64)
```

```
In [5]: print(df.shape)
df = df.drop(null_idx, axis=0)
print(df.shape)
(1226, 21)
```

```
In [6]: df['eveipcnt'] = df['eveipcnt'].fillna(0)
df['eveipcnt'] = df['eveipcnt'].astype('int')
df[['alltag', 'phonoeno', 'sbst']] = df[['alltag', 'phonoeno', 'sbst']].fillna('')
```

```
In [7]: df['tag_orig'] = df['tag']
df['tag'] = df['tag_orig'].apply(lambda x: x.split(",")[0] if len(x) > 0 else cate)
```

```
In [8]: # 중복제거
df.drop_duplicates(['contentsid'], keep='first', inplace=True, ignore_index=True)
```

```
In [9]: df.to_csv(r'./datasets_jeju/data.csv', index=False)

• https://github.com/sqalchemy/sqalchemy/issues/4265
• https://cx-oracle.readthedocs.io/en/latest/user\_guide/sql\_execution.html
```

```
In [10]: df.select_dtypes('object').columns
```

```
Out[10]: Index(['contentsid', 'alltag', 'label', 'title', 'address', 'tag',
              'introduction', 'phonoeno', 'sbst', 'img', 'thumb', 'tag_orig'],
              dtype='object')
```

```
In [11]: from sqlalchemy import create_engine
import sqlalchemy as sa
engine = create_engine("oracle+cx_oracle://ai:0000@localhost:1521/XE")
df.to_sql("JEJU_TRAVEL", engine,
         if_exists="replace", #, fail , append
         # index=True,
         # index_label = 'contentsid',
         dtype={'latitude': sa.FLOAT(), 'longitude': sa.FLOAT()},
         # 'contentsid':sa.String(4000), 'alltag':sa.String(4000), 'label':sa.String(4000), 'title':sa.String(4000), 'address':sa.String(4000), 'tag':sa.String(4000),
         # 'phonoeno':sa.String(4000), 'sbst':sa.String(4000), 'img':sa.String(4000), 'thumb':sa.String(4000), 'tag_orig':sa.String(4000), #'introduction':sa.String(4000),
         )

Out[11]: 1213
```

```
In [12]: from sqlalchemy import create_engine, text
engine = create_engine("oracle+cx_oracle://ai:0000@localhost:1521/XE")
df = pd.read_sql(text("SELECT * FROM JEJU_TRAVEL"), con = engine.connect())
print(df.shape)
df = df.set_index('index')
print(df.info())
# df.head(2)

(1213, 23)
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1213 entries, 0 to 1212
Data columns (total 22 columns):
 # Column Non-Null Count Dtype
---
 0 contentsid 1213 non-null object
 1 alltag 1212 non-null object
 2 label 1213 non-null object
 3 title 1213 non-null object
 4 address 1213 non-null object
 5 tag 1213 non-null object
 6 introduction 1213 non-null object
 7 readcnt 1213 non-null int64
 8 likecnt 1213 non-null int64
 9 reviewcnt 1213 non-null int64
10 markcnt 1213 non-null int64
11 snssharecnt 1213 non-null int64
12 schedulecnt 1213 non-null int64
13 visitcnt 1213 non-null int64
14 eveipcnt 1213 non-null int64
15 latitude 1213 non-null float64
16 longitude 1213 non-null float64
17 phonoeno 1176 non-null object
18 sbst 1143 non-null object
19 img 1213 non-null object
20 thumb 1213 non-null object
21 tag_orig 1213 non-null object
dtypes: float64(2), int64(8), object(12)
memory usage: 218.0+ KB
None
```

```
In [13]: df = pd.read_csv("./datasets_jeju/data.csv")
df['label'].value_counts().index.values
```

```
Out[13]: array(['음식집', '관광지', '숙박', '쇼핑'], dtype=object)
```

```
In [14]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1213 entries, 0 to 1212
Data columns (total 22 columns):
 # Column Non-Null Count Dtype
---
 0 contentsid 1213 non-null object
 1 alltag 1212 non-null object
 2 label 1213 non-null object
 3 title 1213 non-null object
 4 address 1213 non-null object
 5 tag 1213 non-null object
 6 introduction 1213 non-null object
 7 readcnt 1213 non-null int64
 8 likecnt 1213 non-null int64
 9 reviewcnt 1213 non-null int64
10 markcnt 1213 non-null int64
11 snssharecnt 1213 non-null int64
12 schedulecnt 1213 non-null int64
13 visitcnt 1213 non-null int64
14 eveipcnt 1213 non-null int64
15 latitude 1213 non-null float64
16 longitude 1213 non-null float64
17 phonoeno 1176 non-null object
18 sbst 1143 non-null object
19 img 1213 non-null object
20 thumb 1213 non-null object
21 tag_orig 1213 non-null object
dtypes: float64(2), int64(8), object(12)
memory usage: 208.6+ KB
```

```
In [15]: df[df['contentsid'].isin(['CNTS_0000000000022562', 'CNTS_000000000019005', 'CNTS_0000000000081471', 'CNT_0000000000561100' ])] title.values

Out[15]: array(['마노큰별장', '9기진옥조', '양가형제', '포도호텔'], dtype=object)
```

```
In [16]: df.isna().sum()

Out[16]: contentsid 0
alltag 1
label 0
title 0
address 0
tag 0
introduction 0
readcnt 0
likecnt 0
reviewcnt 0
markcnt 0
snssharecnt 0
schedulecnt 0
visitcnt 0
eveipcnt 0
latitude 0
longitude 0
phonoeno 37
sbst 70
img 0
thumb 0
tag_orig 0
dtype: int64
```

```
In [ ]:
```

```
In [17]: df[df['label']=="음식집"][['title', 'address', 'likecnt', 'reviewcnt', 'eveipcnt', 'thumb']]

Out[17]:
```

	title	address	likecnt	reviewcnt	eveipcnt	thumb
103	마노큰별장	서귀포시 한림면 덕우리 2952	13	30	5	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/88853437-8884-495d-b43d-9897917014ce.jpg
110	카친옥조	제주시 한림읍 한라리 958-1	2	28	5	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/b319ed56-2466-48d9-aae1-7a55bb6b970c.jpg
142	비디리	제주특별자치도 서귀포시 대포동 2384	9	20	4	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/9b20bc0e-5c29-4236-b3b3-6a67ea7f6996.jpg
232	양가형제 제주특별자치도 제주시 한림면 형우리 746-8		3	9	4	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/5c3d90f-ed82-456e-a84e-e4a8bc87d211.jpg
287	바다를책갈피로 제주특별자치도 제주시 한림면 만우리 2861-4		1	5	5	https://api.cdn.visitjeju.net/photomng/thumbnaipath/202111/15/01f0b09-5562-48ac-85ba-b6bb08487bce.jpg
...
1208	누리해방궁	제주특별자치도 제주시 물귀동 1819-3	1	1	5	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/f9b1579-1763-47f6-b022-90e714a65d.jpg
1209	달마야해물탕	제주특별자치도 제주시 조천읍 복촌리 1363-1	0	1	3	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/d4c4072f-6820-4c9e-a484-7acdc9b1892f.jpg
1210	달마해물탕	제주특별자치도 제주시 한림면 천리 1647-4	0	1	3	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/6569b984-3886-41dc-85cb-ad9cf711436f.jpg
1211	달곶관유식 제주특별자치도 제주시 애월읍 남원리 1249-4		0	1	3	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/89a124e6-b6f1-469b-9091-477e9160c869.jpg
1212	대관정식당	제주특별자치도 서귀포시 서귀동 297-1	1	1	4	https://api.cdn.visitjeju.net/photomng/thumbnaipath/201804/30/96c1f1c4-2aa3-452e-b0f1-8605885862df.jpg
361 rows × 6 columns						

상세보기 페이지

- https://api.visitjeju.net/api/node/tour/contents/read.json?d=CONT_0000000005003498,_siteId=jeju&locale=kr&device=pc&cacheTime=60

주변 관광지/맛집/숙소

- https://api.visitjeju.net/api/bigdata/list?_siteId=jeju&locale=kr&device=pc&distance=3&lat=33.462147&lng=126.936424&date=20230302&gender=&years=

댓글 크롤링

- https://www.visitjeju.net/ko/detail/view?contentsId=CONT_0000000005003498
- https://api.visitjeju.net/api/node/tour/contents/read.json?d=CONT_0000000005003498,_siteId=jeju&locale=kr&device=pc&cacheTime=60

```
In [ ]:
```

```
In [ ]:
```

```
In [20]: mm = ["CONT_000000000500349", "CONT_000000000500477", "CNTS_0000000001105"]
print(" ".join(mm))

CONT_000000000500349,CONT_000000000500477,CNTS_0000000001105
```