

19장 실전 데이터 분석 프로젝트

- 서울시에서 공개하는 업무추진비 데이터를 활용해 데이터 분석을 수행한다.

19.1 데이터 분석 프로세스

- 1. 주제 선정
 - 데이터 분석의 목표를 명확히 하고 이로부터 주제를 선정한다. 즉, 데이터 분석을 통해 얻고 싶은 결과가 무엇인지를 설정하는 것이다.
- 2. 데이터 수집
 - 주제가 선정되면 주제에 맞는 데이터를 수집해야 한다.
 - 국내외 무료 공개 데이터 저장소
 - 공공 데이터 포털(<https://www.data.go.kr>)
 - 서울시 정보소통광장(<https://opengov.seoul.go.kr/>)
 - 국가 통계 포털(<http://kosis.kr/>)
 - 세계 은행 오픈 데이터(<https://data.worldbank.org/>)
 - FRD(<https://fred.stlouisfed.org/>)
 - 캐글 사이트(<https://www.kaggle.com>)
 - U.S. Government's open data(<https://www.data.gov/>)
- 3. 데이터 처리
 - 데이터가 수집되면 다음은 데이터 분석이 편리하도록 데이터를 처리하는 과정을 거친다.
 - 데이터 정제: 가공되지 않은 초기 데이터에서 부정확한 데이터를 찾아내고 수정하거나 제거하는 절차이다.
 - 데이터 타이딩: 데이터 분석을 위해 데이터를 알아보기 쉽고 처리하기 편하도록 구조화하는 절차이다. 보통 데이터 정제 이후에 이뤄지는데, 이 과정을 거친 데이터를 깔끔한 데이터(Tidy data)라고 한다.
- 4. 데이터 분석
 - 앞 단계에서 처리된 데이터를 이용해 다양한 기법으로 데이터를 분석한다. 대표적인 데이터 분석 방법은 통계적 분석 방법이며 최근에는 머신러닝을 활용한 데이터 분석도 많이 수행한다.
 - 데이터 분석 이후에는 분석 결과를 효과적으로 전달하기 위해 다양한 시각화 방법을 이용한다.
- 5. 정보 도출
 - 데이터 분석 과정으로 도출된 결과에서 의미를 발견하는 단계이다. 이 단계에서 앞의 모든 과정이 올바르게 이뤄졌는지 검증하는 과정이 필요하다.

19.2 데이터 획득, 처리, 시각화 심화

19.2.1 깃허브에서 파일 내려받기

- requests를 이용해 깃허브에서 파일을 내려받고 저장하는 코드를 작성한다.

```
[ch19_project/ex01_github.py]
```

```

01 import requests
02
03 # 깃허브의 파일 URL
04 url = 'https://github.com/wikibook/python-for-data-analysis-rev/raw/master/readme.txt'
05
06 # URL에 해당하는 파일을 내려받음
07 r = requests.get(url)
08
09 # 파일을 저장할 폴더와 파일명을 지정
10 file_name = './ch19_project/readme.txt'
11
12 # 내려받은 파일을 지정한 폴더에 저장
13 with open(file_name, 'wb') as f:
14     f.write(r.content)

```

19.2.2 데이터에서 결측치 확인 및 처리

- 수집된 데이터에는 다양한 이유로 데이터 값이 누락될 수 있다. 이처럼 누락된 데이터를 결측치(Missing data)라고 하며, 데이터에 결측치가 있다면 데이터를 분석할 때 결과가 왜곡되거나 데이터 분석이 불가능할 수도 있다. 따라서 데이터에서 어떠한 값이 결측치인지 아는 것이 중요하다.
- Pandas의 DataFrame 형식으로 가져온 데이터의 경우 `isna()` 혹은 `isnull()`을 통해 결측치 여부를 판단할 수 있다.
 - 18: `isnull()`을 이용해 결측치 여부를 확인한다.
 - 29: `isnull()`의 결과에 `sum()`을 수행하면 결측치 개수를 셀 수 있다.

[ch19_project/ex02_missing_data.py]

```

01 # 결측치 확인
02 import pandas as pd
03
04 data_file = "./ch19_project/missing_data_test.csv"
05
06 df = pd.read_csv(data_file, encoding="cp949", index_col="연도")
07
08 print(df)
09 '''
10          제품1  제품2  제품3  제품4
11 연도
12 2015   250.0  150    NaN    NaN
13 2016   200.0  160   170.0    NaN
14 2017   150.0  200   100.0  150.0
15 2018   120.0  230   130.0  170.0
16 2019    NaN  250   140.0    NaN
17 '''
18
19 print(df.isnull())
20 '''
21          제품1  제품2  제품3  제품4
22 연도
23 2015   False  False   True   True
24 2016   False  False  False   True
25 2017   False  False  False  False
26 2018   False  False  False  False
27 2019    True  False  False   True
28 '''
29
30 print(df.isnull().sum())

```

```

31 '''
32 제품1    1
33 제품2    0
34 제품3    1
35 제품4    3
36 dtype: int64
37 '''
38
39
40 # 결측치 처리
41 print(df.drop(index=[2019]))
42 '''
43          제품1  제품2   제품3   제품4
44 연도
45 2015  250.0  150    NaN    NaN
46 2016  200.0  160   170.0    NaN
47 2017  150.0  200   100.0  150.0
48 2018  120.0  230   130.0  170.0
49 '''
50
51 print(df.drop(columns=['제품3', '제품4']))
52 '''
53          제품1  제품2
54 연도
55 2015  250.0  150
56 2016  200.0  160
57 2017  150.0  200
58 2018  120.0  230
59 2019    NaN  250
60 '''
61
62 print(df.drop(index=[2018, 2019], columns=['제품3', '제품4']))
63 '''
64          제품1  제품2
65 연도
66 2015  250.0  150
67 2016  200.0  160
68 2017  150.0  200
69 '''
70
71
72 print(df.dropna()) # df.dropna(axis=0)도 결과는 같습니다.
73 '''
74          제품1  제품2   제품3   제품4
75 연도
76 2017  150.0  200   100.0  150.0
77 2018  120.0  230   130.0  170.0
78 '''
79
80 print(df.dropna(axis=0, subset=['제품1']))
81 '''
82          제품1  제품2   제품3   제품4
83 연도
84 2015  250.0  150    NaN    NaN
85 2016  200.0  160   170.0    NaN
86 2017  150.0  200   100.0  150.0
87 2018  120.0  230   130.0  170.0
88 '''
89
90 print(df.dropna(axis=1))
91 '''
92          제품2
93 연도
94 2015  150
95 2016  160

```

```

96 2017 200
97 2018 230
98 2019 250
99 ...
100
101 print(df.fillna(0))
102 ...
103      제품1  제품2  제품3  제품4
104 연도
105 2015 250.0 150 0.0 0.0
106 2016 200.0 160 170.0 0.0
107 2017 150.0 200 100.0 150.0
108 2018 120.0 230 130.0 170.0
109 2019 0.0 250 140.0 0.0
110 ...
111
112 print(df.fillna(method='bfill'))
113 ...
114      제품1  제품2  제품3  제품4
115 연도
116 2015 250.0 150 170.0 150.0
117 2016 200.0 160 170.0 150.0
118 2017 150.0 200 100.0 150.0
119 2018 120.0 230 130.0 170.0
120 2019 NaN 250 140.0 NaN
121 ...
122
123 print(df.fillna(method='ffill'))
124 ...
125      제품1  제품2  제품3  제품4
126 연도
127 2015 250.0 150 NaN NaN
128 2016 200.0 160 170.0 NaN
129 2017 150.0 200 100.0 150.0
130 2018 120.0 230 130.0 170.0
131 2019 120.0 250 140.0 170.0
132 ...
133
134 values = {'제품1': 100, '제품4': 400}
135 print(df.fillna(value=values))
136 ...
137      제품1  제품2  제품3  제품4
138 연도
139 2015 250.0 150 NaN 400.0
140 2016 200.0 160 170.0 400.0
141 2017 150.0 200 100.0 150.0
142 2018 120.0 230 130.0 170.0
143 2019 100.0 250 140.0 400.0
144 ...

```

19.2.3 데이터의 요약 및 재구성

- CSV 데이터 파일(total_sales_data.csv)을 DataFrame 형식으로 읽어오는 코드는 다음과 같다.

```
[ch19_project/ex03_info.py]
```

```

01 import pandas as pd
02
03 data_file = "./ch19_project/total_sales_data.csv"
04

```

```

05 df_sales = pd.read_csv(data_file)
06 print(df_sales)
07 '''
08     매장명  제품종류  모델명  판매  재고
09 0  A  스마트폰  S1    1    2
10 1  A  스마트폰  S2    2    5
11 2  A   TV  V1    3    5
12 3  B  스마트폰  S2    4    6
13 4  B  스마트폰  S1    5    8
14 5  B   TV  V1    6    9
15 6  C  스마트폰  S2    2    4
16 7  C   TV  V1    3    6
17 8  C   TV  V2    7    9
18 '''
19
20 print(df_sales.info())
21 '''
22 <class 'pandas.core.frame.DataFrame'>
23 RangeIndex: 9 entries, 0 to 8
24 Data columns (total 5 columns):
25  #   Column  Non-Null Count  Dtype
26  ---  ---
27 0  매장명      9 non-null      object
28 1  제품종류    9 non-null      object
29 2  모델명      9 non-null      object
30 3  판매        9 non-null      int64
31 4  재고        9 non-null      int64
32 dtypes: int64(2), object(3)
33 memory usage: 488.0+ bytes
34 None
35 '''
36
37 print(df_sales['매장명'].value_counts())
38 '''
39 C    3
40 B    3
41 A    3
42 Name: 매장명, dtype: int64
43 '''
44
45 print(df_sales['제품종류'].value_counts())
46 '''
47 스마트폰    5
48 TV          4
49 Name: 제품종류, dtype: int64
50 '''
51
52
53 print(df_sales.pivot_table(index=["매장명", "제품종류", "모델명"],
54                             values=["판매", "재고"], aggfunc='sum'))
55 '''
56             재고  판매
57 매장명 제품종류 모델명
58 A   TV   V1      5    3
59   스마트폰 S1      2    1
60   S2      5    2
61 B   TV   V1      9    6
62   스마트폰 S1      8    5
63   S2      6    4
64 C   TV   V1      6    3
65   V2      9    7
66   스마트폰 S2      4    2
67 '''
68
69 print(df_sales.pivot_table(index=["매장명"], columns = ["제품종류"],

```

```

70         values=["판매", "재고"], aggfunc='sum'))
71     '''
72     재고      판매
73     제품종류 TV 스마트폰 TV 스마트폰
74     매장명
75     A      5      7      3      3
76     B      9     14      6      9
77     C     15      4     10      2
78     '''
79
80     print(df_sales.pivot_table(index=["매장명"], columns=["제품종류"],
81         values=["판매", "재고"], aggfunc='count'))
82     '''
83     재고      판매
84     제품종류 TV 스마트폰 TV 스마트폰
85     매장명
86     A      1      2      1      2
87     B      1      2      1      2
88     C      2      1      2      1
89     '''

```

19.2.4 워드 클라우드를 이용한 데이터 시각화

- 워드 클라우드는 텍스트 데이터에서 출현 빈도가 높은 단어는 크게 표시하고 출현 빈도가 낮은 단어는 작게 표시하는 방법으로 데이터를 시각화한다.
- 워드 클라우드 설치

```

D:\dev\workspace\python>pip install wordcloud
Collecting wordcloud
  Downloading wordcloud-1.8.1-cp38-cp38-win_amd64.whl (155 kB)
    |#####| 155 kB 6.8 MB/s
...(생략)...
Installing collected packages: wordcloud
Successfully installed wordcloud-1.8.1

```

- 다음은 텍스트 파일로 워드 클라우드 이미지를 생성해서 화면에 출력하는 코드이다.

```

[ch19_project/ex04_wordcloud.py]

01 from wordcloud import WordCloud
02 import matplotlib.pyplot as plt
03
04 file_name = './ch19_project/littleprince_djvu.txt'
05
06 with open(file_name) as f: # 파일을 읽기 모드로 열기
07     text = f.read() # 파일의 내용 읽어오기
08
09 # 워드 클라우드의 이미지를 생성합니다.
10 wordcloud_image = WordCloud().generate(text)
11
12 # 생성한 워드 클라우드 이미지를 화면에 표시합니다.
13 plt.imshow(wordcloud_image, interpolation='bilinear')
14 plt.axis("off")
15 plt.show()
16
17

```


부서별 사용내역 기간별 이용내역

부서: 서울시본청 서울특별시시장

범위: 2021년 10월 간행 2021년10월 다운로드 2016년 10월 이전자료는 **기간별 이용내역**에서 확인하세요

서울특별시시장 기관운영, 시책추진 업무추진비 사용내역 엑셀 다운로드

연번	구분	부서명	사용일시	사용장소	사용목적	사용금액	사용자 및 인원	결제방법	비고
1		행정국 총무과(서울특별시시장)	2021-10-01 13:15	서울시청 간담회장(서울 중구 세종대로 110)	시장현안 관련 의견수렴 간담회	150,000	시장 외 5명	카드	시책
2	현업-우수부서 격려 등	행정국 총무과(서울특별시시장)	2021-10-01 16:02	피자헛(서울특별시 중구 수표로 84)	현안 관련 격려부서 격려	362,800	자치행정과	계로책이	기관
3		행정국 총무과(서울특별시시장)	2021-10-01 19:55	한양동, 코프(Cof)(서울특별시 중구 경동길 21-15)	시장 추진 관련 자문 간담회	173,000	시장 외 5명	카드	시책

■ 다음과 같은 세부 주제를 선정해서 데이터 분석을 진행한다.

- 연도별 추이 분석
- 월별 집행금액 분석
- 부서별 집행 내역 분석
- 요일별 및 시간대별 집행 내역 분석

19.3.2 데이터 수집

■ 서울시 행정정보 공개 깃허브(<https://github.com/seoul-opengov/opengov>)에 있는 업무추진비 집행내역을 가져온다.

[ch19_project/ex05_데이터수집.py]

```
01 # 데이터 수집
02 import glob
03 import requests
04 import os
05 import pathlib
06
07 #인자: 확장자, 연도, 내려받을 폴더
08
09
10 def get_seoul_expense_list(extension, year, data_folder):
11
12     # 깃허브의 데이터 위치 지정
13     # ex) 'https://github.com/seoul-opengov/opengov/raw/master/expense_list2017/'
14     expense_list_year_url = 'https://github.com/seoul-opengov/opengov/raw/master/expense_list' + \
15         str(year) + '/'
16
17     # 데이터를 내려받을 폴더 지정
18     # ex) 'C:/myPyCode/data/seoul_expense/2017/'
19     expense_list_year_dir = data_folder + str(year) + '/'
20
21     # 내려받을 폴더가 없다면 폴더 생성
22     if(os.path.isdir(expense_list_year_dir)):
23         print("데이터 폴더({0})가 이미 있습니다. {0}년 데이터의 다운로드를 시작합니다.".format(year))
24     else:
25         print("데이터 폴더({0})가 없어서 생성했습니다. {0}년 데이터의 다운로드를 시작합니
26 다.".format(year))
27         # 폴더 생성
28         pathlib.Path(expense_list_year_dir).mkdir(parents=True, exist_ok=True)
29
30     # 지정한 폴더로 1월 ~ 12월 업무추진비 파일을 다운로드
31     for k in range(12):
32         file_name = '{0}{1:02d}_expense_list.{2}'.format(year, k+1, extension)
```



```

33     url = expense_list_year_url + file_name
34     print(url)
35     r = requests.get(url)
36     with open(expense_list_year_dir + file_name, 'wb') as f:
37         f.write(r.content)
38
39
40 # 내려받을 업무추진비 데이터의 파일 형식을 지정
41 extension = "csv"
42
43 # 내려받을 업무추진비 데이터의 연도를 지정
44 year = 2017
45
46 # 내려받을 업무추진비 데이터의 폴더를 지정
47 data_folder = './ch19_project/seoul_expense/'
48
49 # 함수를 실행
50 get_seoul_expense_list(extension, year, data_folder)
51
52
53 path_name = './ch19_project/seoul_expense/2017/' # 폴더 이름
54
55 # 지정 폴더에서 파일명에 list.csv가 포함된 파일만 지정
56 file_name_for_glob = path_name + "*list.csv"
57
58 csv_files = []
59 for csv_file in glob.glob(file_name_for_glob):
60     # 반환값에서 폴더는 제거하고 파일 이름만 추출
61     csv_files.append(csv_file.split("\\")[-1])
62
63 print("[폴더 이름]", path_name) # 폴더명 출력
64 print("* CSV 파일:", csv_files)
65 '''
66 [폴더 이름] ./ch19_project/seoul_expense/2017/
67 * CSV 파일: ['201701_expense_list.csv', '201702_expense_list.csv', '201703_expense_list.csv',
68 '201704_expense_list.csv', '201705_expense_list.csv', '201706_expense_list.csv',
69 '201707_expense_list.csv', '201708_expense_list.csv', '201709_expense_list.csv',
70 '201710_expense_list.csv', '201711_expense_list.csv', '201712_expense_list.csv']
71 '''
72
73
74 data_folder = './ch19_project/seoul_expense/'
75
76 years = [2017, 2018, 2019] # 다운로드받을 연도를 지정
77
78 extension = "csv"
79 # extension = "xlsx"
80 # extension = "xml"
81
82 for year in years:
83     get_seoul_expense_list(extension, year, data_folder)
84
85 print("모든 데이터를 다운로드 받았습니다.")
86
87
88 data_folder = './ch19_project/seoul_expense/'
89
90 years = [2017, 2018, 2019] # 다운로드받을 연도를 지정
91
92 for year in years:
93     path_name = data_folder + str(year) + "/" # 연도별 폴더명을 지정
94
95     # 지정 폴더에서 파일명에 list.csv가 포함된 파일만 지정
96     file_name_for_glob = path_name + "*list.csv"
97

```

```

98     csv_files = []
99     for csv_file in glob.glob(file_name_for_glob):
100         # 반환값에서 폴더는 제거하고 파일명만 추출
101         csv_files.append(csv_file.split("\\")[-1])
102
103     print("[폴더 이름]", path_name) # 폴더명 출력
104     print("* CSV 파일:", csv_files)
105     ...
106 [폴더 이름] ./ch19_project/seoul_expense/2017/
107 * CSV 파일: ['201701_expense_list.csv', '201702_expense_list.csv', '201703_expense_list.csv',
108 '201704_expense_list.csv', '201705_expense_list.csv', '201706_expense_list.csv',
109 '201707_expense_list.csv', '201708_expense_list.csv', '201709_expense_list.csv',
110 '201710_expense_list.csv', '201711_expense_list.csv', '201712_expense_list.csv']
111 [폴더 이름] ./ch19_project/seoul_expense/2018/
112 * CSV 파일: ['201801_expense_list.csv', '201802_expense_list.csv', '201803_expense_list.csv',
113 '201804_expense_list.csv', '201805_expense_list.csv', '201806_expense_list.csv',
114 '201807_expense_list.csv', '201808_expense_list.csv', '201809_expense_list.csv',
115 '201810_expense_list.csv', '201811_expense_list.csv', '201812_expense_list.csv']
116 [폴더 이름] ./ch19_project/seoul_expense/2019/
117 * CSV 파일: ['201901_expense_list.csv', '201902_expense_list.csv', '201903_expense_list.csv',
118 '201904_expense_list.csv', '201905_expense_list.csv', '201906_expense_list.csv',
119 '201907_expense_list.csv', '201908_expense_list.csv', '201909_expense_list.csv',
120 '201910_expense_list.csv', '201911_expense_list.csv', '201912_expense_list.csv']
121 '''

```

19.3.3 데이터 처리

- 자신이 직접 설계해서 만든 데이터가 아니고 외부에서 데이터를 가져온 경우 데이터가 원하는 구조가 아닐 수 있다. 이 경우 데이터 정제와 정리 과정이 필요하다.

(1) 수집된 데이터 파일의 구조 분석

- 데이터 수집 이후에는 수집된 데이터에 문제가 없는지 검토하는 과정이 필요하다.
 - 08~15: 텍스트 파일을 읽는 코드를 이용해 처음 세 줄을 읽어서 출력한다.
 - 16~31: 데이터의 첫째 줄에는 열 이름이 있고, 둘째 줄부터는 데이터 값이 있다.
 - 34~39: 첫째 줄에 있는 열 이름의 개수와 둘째 줄과 셋째 줄에 있는 데이터 값의 개수가 각각 몇 개인지 살펴본다.
 - 46~56: "시정 충남, 전북 현장 방문 관련 업무협약"과 같이 따옴표 안의 콤마를 제거한다.

[ch19_project/ex06_데이터처리.py의 일부]

```

01 # 수집된 데이터 파일의 구조 분석
02 from datetime import datetime
03 import os
04 import pandas as pd
05 import glob
06 data_file = './ch19_project/seoul_expense/2017/201701_expense_list.csv'
07
08 with open(data_file, encoding='utf-8') as f:
09     line1 = f.readline()
10     line2 = f.readline()
11     line3 = f.readline()
12
13     print(line1)
14     print(line2)
15     print(line3)

```

```

16  '''
17  nid,title,url,dept_nm_lvl_1,dept_nm_lvl_2,dept_nm_lvl_3,dept_nm_lvl_4,dept_nm_lvl_5,exec_yr,exec_month,
18  expense_budget,expense_execution,category,dept_nm_full,exec_dt,exec_loc,exec_purpose,target_nm,payment_
19  method,exec_amount
20
21  11430252,"2017년 1월 장애인복지정책과 업무추진비 집행내역",http://opengov.seoul.go.kr/public/11430252,
22  서
23  울시본청,복지본부,장애인복지정책과,,,2017,1,,,,"복지본부 장애인복지정책과","2017-01-26 13:10","동해일식
24  (중구 무교동)","기본소득과 장애인복지 논의간담회","장애인복지정책팀장 외 2명",카드,76000
25
26  11430252,"2017년 1월 장애인복지정책과 업무추진비 집행내역",http://opengov.seoul.go.kr/public/11430252,
27  서
28  울시본청,복지본부,장애인복지정책과,,,2017,1,,,,"복지본부 장애인복지정책과","2017-01-25 22:41","김앤장 (
29  중구 무교로)","장애인단체 활동지원 논의간담회","장애인복지정책과장 외 3명",카드,102000
30
31  '''
32
33
34  line1_len = len(line1.split(','))
35  line2_len = len(line2.split(','))
36  line3_len = len(line3.split(','))
37
38  print("[각 줄의 데이터값의 개수]")
39  print("첫째 줄:{}, 둘째 줄:{}, 셋째 줄:{}".format(line1_len, line2_len, line3_len))
40  '''
41  [각 줄의 데이터값의 개수]
42  첫째 줄:20, 둘째 줄:20, 셋째 줄:20
43  '''
44
45
46  def get_value_count(line):
47
48      line_rep_list = []
49      for k, x in enumerate(line.split('')):
50          if(k % 2 != 0):
51              x = x.replace(',', ' ')
52              line_rep_list.append(x)
53
54      line_rep_str = ''.join(line_rep_list)
55
56      return len(line_rep_str.split(','))
57
58
59  line1_len = get_value_count(line1)
60  line2_len = get_value_count(line2)
61  line3_len = get_value_count(line3)
62
63  print("[각 줄의 데이터값의 개수]")
64  print("첫째 줄:{}, 둘째 줄:{}, 셋째 줄:{}".format(line1_len, line2_len, line3_len))
65  '''
66  [각 줄의 데이터값의 개수]
67  첫째 줄:20, 둘째 줄:20, 셋째 줄:20
68  '''

```

(2) 첫 번째 줄의 이름과 개수 변경

- 서울시의 업무추진비 CSV 파일에 첫 번째 줄에 누락된 열 이름이 있는 경우, 열 이름을 추가하고 알아보기 어려운 열 이름은 한글로 바꾼다.
 - 05: splitlines()는 문자열을 개행문자(\n)를 중심으로 나누는 함수이다.
 - 13: '\n'.join(lines)는 lines에 있는 리스트 요소를 개행문자로 연결해서 하나의 문자

열로 만들어준다.

- 89~106: 모든 데이터 파일에 대해 첫 번째 줄의 열 이름을 변경한 파일이 잘 생성됐는지 확인한다.
- 102: 파이썬 내장 모듈인 glob.glob() 함수는 인자로 받은 패턴과 이름이 일치하는 모든 파일과 디렉터리의 리스트를 반환한다. 패턴을 그냥 *라고 주면 모든 파일과 디렉터리를 볼 수 있다.

[ch19_project/ex06_데이터처리.py의 일부]

```
01 # 2. 첫 번째 줄의 열 이름과 개수 변경
02 def change_csv_file_first_line_value(old_file_name, new_file_name):
03     with open(old_file_name, encoding='utf-8') as f: # 파일을 읽기 모드로 열기
04         # 전체 데이터를 읽어서 한 줄씩 lines 리스트의 각 요소에 할당
05         lines = f.read().splitlines()
06
07     # 첫째 줄의 내용을 변경할 열 이름을 지정해서 변경
08     lines[0] = 'nid,제목,url,부서레벨1,부서레벨2,부서레벨3,부서레벨4,부서레벨5,\
09 집행연도,집행월,예산,집행,구분,부서명,집행일시,집행장소,집행목적,대상인원,결재방법,집행금액'
10
11     with open(new_file_name, 'w', encoding='utf-8') as f: # 파일을 쓰기 모드로 열기
12         # 리스트 내의 각 요소를 개행문자(\n)로 연결해서 파일로 저장
13         f.write('\n'.join(lines))
14
15
16 # 기존의 파일
17 old_file_name = './ch19_project/seoul_expense/2017/201701_expense_list.csv'
18
19 # 새로운 파일
20 new_file_name = './ch19_project/seoul_expense/2017/201701_expense_list_new.csv'
21
22 # 첫째 줄의 내용을 변경한 새로운 파일 생성
23 change_csv_file_first_line_value(old_file_name, new_file_name)
24
25
26 with open(new_file_name, encoding='utf-8') as f: # 파일을 읽기 모드로 열기
27     for k in range(3):
28         print(f.readline())
29     '''
30 nid,제목,url,부서레벨1,부서레벨2,부서레벨3,부서레벨4,부서레벨5,집행연도,집행월,예산,집행,구분,부서명,집
31 행일시
32 ,집행장소,집행목적,대상인원,결재방법,집행금액
33
34 11430252,"2017년 1월 장애인복지정책과 업무추진비 집행내역",http://opengov.seoul.go.kr/public/11430252,
35 서울시
36 본청,복지본부,장애인복지정책과,,2017,1,,,"복지본부 장애인복지정책과","2017-01-26 13:10","동해일식 (중
37 구 무
38 교동)","기본소득과 장애인복지 논의간담회","장애인복지정책팀장 외 2명",카드,76000
39
40 11430252,"2017년 1월 장애인복지정책과 업무추진비 집행내역",http://opengov.seoul.go.kr/public/11430252,
41 서울시
42 본청,복지본부,장애인복지정책과,,2017,1,,,"복지본부 장애인복지정책과","2017-01-25 22:41","김앤장 (중구
43 무교
44 로)","장애인단체 활동지원 논의간담회","장애인복지정책과장 외 3명",카드,102000
45
46 '''
47
48
49 # 인자: 연도, 데이터 파일이 있는 폴더
50 def change_year_csv_file_first_line_value(year, data_folder):
51
52     # 데이터 파일이 있는 폴더 지정
53     # ex) 'C:/myPyCode/data/seoul_expense/2017/'
54     expense_list_year_dir = data_folder + str(year) + '/'
```

```

55
56     extension = 'csv' # 확장자 이름
57
58     # 지정한 폴더에 있는 월별 업무추진비 파일에서 첫 번째 줄의 열 이름을 변경
59     for k in range(12):
60         # 기존의 파일 이름 지정
61         old_file_name = expense_list_year_dir + \
62             '{0}{1:02d}_expense_list.{2}'.format(year, k+1, extension)
63
64         # 새로운 파일 이름 지정
65         new_file_name = expense_list_year_dir + \
66             '{0}{1:02d}_expense_list_new.{2}'.format(year, k+1, extension)
67
68         # 첫째 줄의 내용을 변경한 새로운 파일 생성
69         change_csv_file_first_line_value(old_file_name, new_file_name)
70
71
72     data_folder = './ch19_project/seoul_expense/'
73
74     years = [2017, 2018, 2019] # 연도를 지정
75
76     for year in years:
77         print("{}년 데이터의 첫 번째 줄의 열 이름을 변경해서 새 파일에 저장합니다.".format(year))
78         change_year_csv_file_first_line_value(year, data_folder)
79
80     print("모든 데이터의 첫 번째 줄의 열 이름을 변경해서 새 파일로 저장했습니다.")
81     '''
82     2017년 데이터의 첫 번째 줄의 열 이름을 변경해서 새 파일에 저장합니다.
83     2018년 데이터의 첫 번째 줄의 열 이름을 변경해서 새 파일에 저장합니다.
84     2019년 데이터의 첫 번째 줄의 열 이름을 변경해서 새 파일에 저장합니다.
85     모든 데이터의 첫 번째 줄의 열 이름을 변경해서 새 파일로 저장했습니다.
86     '''
87
88
89     data_folder = './ch19_project/seoul_expense/'
90
91     years = [2017, 2018, 2019] # 연도를 지정
92
93     for year in years:
94         path_name = data_folder + str(year) # 폴더명을 지정
95         print("[폴더 이름]", path_name) # 폴더명 출력
96
97         new_csv_files = []
98
99         # 지정 폴더에서 파일명에 _new.csv가 포함된 파일만 지정
100        file_name_for_glob = path_name + "/*_new.csv"
101
102        for new_csv_file in glob.glob(file_name_for_glob):
103            # 반환값에서 폴더는 제거하고 파일 이름만 추출
104            new_csv_files.append(new_csv_file.split("\\")[-1])
105
106        print("* 새롭게 생성된 CSV 파일:", new_csv_files)
107        '''
108        [폴더 이름] ./ch19_project/seoul_expense/2017
109        * 새롭게 생성된 CSV 파일: ['201701_expense_list_new.csv', '201702_expense_list_new.csv',
110        '201703_expense_list_new.csv', '201704_expense_list_new.csv', '201705_expense_list_new.csv',
111        '201706_expense_list_new.csv', '201707_expense_list_new.csv', '201708_expense_list_new.csv',
112        '201709_expense_list_new.csv', '201710_expense_list_new.csv', '201711_expense_list_new.csv',
113        '201712_expense_list_new.csv']
114        [폴더 이름] ./ch19_project/seoul_expense/2018
115        * 새롭게 생성된 CSV 파일: ['201801_expense_list_new.csv', '201802_expense_list_new.csv',
116        '201803_expense_list_new.csv', '201804_expense_list_new.csv', '201805_expense_list_new.csv',
117        '201806_expense_list_new.csv', '201807_expense_list_new.csv', '201808_expense_list_new.csv',
118        '201809_expense_list_new.csv', '201810_expense_list_new.csv', '201811_expense_list_new.csv',
119        '201812_expense_list_new.csv']

```

```

120 [폴더 이름] ./ch19_project/seoul_expense/2019
121 * 새롭게 생성된 CSV 파일: ['201901_expense_list_new.csv', '201902_expense_list_new.csv',
122 '201903_expense_list_new.csv', '201904_expense_list_new.csv', '201905_expense_list_new.csv',
123 '201906_expense_list_new.csv', '201907_expense_list_new.csv', '201908_expense_list_new.csv',
124 '201909_expense_list_new.csv', '201910_expense_list_new.csv', '201911_expense_list_new.csv',
125 '201912_expense_list_new.csv']
126 '''

```

(3) 데이터의 구조 및 결측치 살펴보기

- CSV 파일을 pandas의 DataFrame 형식으로 가져와서 데이터의 구조를 살펴보고 빠진 데이터 값(결측치)이 있는지도 살펴본다.
 - 68~79: isna()를 통해 결측치가 몇 개 있는지 확인한다.
 - 82~84: 결측치가 많은 열은 제거하고 나머지 열의 데이터만 이용한다.
 - 113~134: 원하는 열의 데이터만 선택해서 새로운 파일로 저장한다.
 - 137~144: 2017년, 2018년, 2019년에 해당하는 모든 데이터를 읽어서 필요 없는 열은 제거하고 새로운 파일 이름(연도_expense_list_tidy.csv)으로 저장한다.

[ch19_project/ex06_데이터처리.py의 일부]

```

01 # 3. 데이터의 구조 및 결측치 살펴보기
02 expense_list2016_dir = './ch19_project/seoul_expense/2017/'
03 file_name = "201701_expense_list_new.csv"
04
05 df = pd.read_csv(expense_list2016_dir + file_name)
06 print(df.head(2))
07 '''
08          nid          제목 ... 결제방법   집행금액
09 0  11430252  2017년 1월 장애인복지정책과 업무추진비 집행내역 ...   카드   76000
10 1  11430252  2017년 1월 장애인복지정책과 업무추진비 집행내역 ...   카드  102000
11
12 [2 rows x 20 columns]
13 '''
14
15
16 year = 2017
17 expense_list_year_dir = './ch19_project/seoul_expense/' + str(year) + '/'
18
19 df_year = pd.DataFrame()
20 for k in range(12):
21
22     # 파일 이름 지정
23     file_name = "{0}{1:02d}_expense_list_new.csv".format(year, k+1)
24
25     # pandas DataFrame 형식으로 csv 데이터 불러오기
26     df_month = pd.read_csv(expense_list_year_dir + file_name)
27
28     # df_year에 df_month를 세로 방향으로 추가해서 다시 df_year에 할당
29     # 통합된 dataframe의 순서대로 index를 할당하기 위해서 `ignore_index = True` 옵션 지정
30     df_year = df_year.append(df_month, ignore_index=True)
31
32 print(df_year.head(2))
33 '''
34          nid          제목 ... 결제방법   집행금액
35 0  11430252  2017년 1월 장애인복지정책과 업무추진비 집행내역 ...   카드   76000
36 1  11430252  2017년 1월 장애인복지정책과 업무추진비 집행내역 ...   카드  102000
37
38 [2 rows x 20 columns]
39 '''

```

```

40
41 print(df_year.tail(2))
42 '''
43          nid                      제목 ... 결제방법   집행금액
44  70130  14292506  2017년 12월 사업소_은평병원_원무과 업무추진비 내역 ...   카드   820000
45  70131  14292506  2017년 12월 사업소_은평병원_원무과 업무추진비 내역 ...   카드  440000
46
47 [2 rows x 20 columns]
48 '''
49
50 print(df_year.info())
51 '''
52 <class 'pandas.core.frame.DataFrame'>
53 RangeIndex: 70132 entries, 0 to 70131
54 Data columns (total 20 columns):
55  #   Column  Non-Null Count  Dtype
56  ---  ---
57  0    nid      70132 non-null  int64
58  1    제목      70132 non-null  object
59  2    url       70132 non-null  object
60  ... (생략) ...
61  18   결제방법   69929 non-null  object
62  19   집행금액   70132 non-null  int64
63 dtypes: float64(2), int64(4), object(14)
64 memory usage: 10.7+ MB
65 None
66 '''
67
68 print(df_year.isna().sum())
69 '''
70 nid          0
71 제목          0
72 url           0
73 부서레벨1      0
74 부서레벨2     58
75 ... (생략) ...
76 결제방법      203
77 집행금액       0
78 dtype: int64
79 '''
80
81
82 df_year_drop = df_year.drop(columns=['nid', 'url', '부서레벨3', '부서레벨4', '부서레벨5',
83                                     '예산', '집행', '구분'])
84 print(df_year_drop.head(2))
85 '''
86          제목  부서레벨1  부서레벨2  집행연도   ...   집행목적
87 대상
88 인원 결제방법   집행금액
89 0  2017년 1월 장애인복지정책과 업무추진비 집행내역  서울시본청  복지본부  2017   ...   기본소득과 장애인
90 복지 논
91 의간담회 장애인복지정책팀장 외 2명   카드   76000
92 1  2017년 1월 장애인복지정책과 업무추진비 집행내역  서울시본청  복지본부  2017   ...   장애인단체 활동지
93 원 논
94 의간담회 장애인복지정책과장 외 3명   카드  102000
95
96 [2 rows x 12 columns]
97 '''
98
99
100 year = 2017
101 expense_list_year_dir = './ch19_project/seoul_expense/' + str(year) + '/'
102
103 expense_list_tidy_file = "{}_expense_list_tidy.csv".format(year)
104 df_year_drop.to_csv(expense_list_year_dir +

```

```

105         expense_list_tidy_file, index=False)
106
107
108 file_name = expense_list_year_dir + expense_list_tidy_file
109 print(file_name)
110 os.path.isfile(file_name)
111
112
113 def select_columns_save_file(year, data_folder, drop_columns_list):
114
115     expense_list_year_dir = data_folder + str(year) + '/'
116     expense_list_tidy_file = "{}_expense_list_tidy.csv".format(year)
117     df_year = pd.DataFrame()
118
119     for k in range(12):
120         # 파일 이름 지정
121         file_name = "{}{0:02d}_expense_list_new.csv".format(year, k+1)
122
123         # aDtaFrame 형식으로 csv 데이터 불러오기
124         df_month = pd.read_csv(expense_list_year_dir + file_name)
125
126         # fd_year에 df_month를 새로 추가해서 다시 df_year에 할당
127         # 통합된 adtaFrame의 순서대로 index를 할당하기 위해서 `ignore_index = True` 옵션 지정
128         df_year = df_year.append(df_month, ignore_index=True)
129
130     df_year_drop = df_year.drop(columns=drop_columns_list)
131     new_file_name = expense_list_year_dir + expense_list_tidy_file
132     df_year_drop.to_csv(new_file_name, index=False)
133
134     print("==> {} 파일을 생성했습니다.".format(expense_list_tidy_file))
135
136
137 data_folder = './ch19_project/seoul_expense/'
138 years = [2017, 2018, 2019]
139 drop_columns_list = ['nid', 'url', '부서레벨3', '부서레벨4', '부서레벨5', '예산', '집행', '구분']
140
141 for year in years:
142     print("{}년 데이터를 정리해서 저장하고 있습니다.".format(year))
143     select_columns_save_file(year, data_folder, drop_columns_list)
144 print("모든 연도의 데이터를 정리해서 파일로 저장했습니다.")
145 '''
146 ./ch19_project/seoul_expense/2017/2017_expense_list_tidy.csv
147 2017년 데이터를 정리해서 저장하고 있습니다.
148 ==> 2017_expense_list_tidy.csv 파일을 생성했습니다.
149 2018년 데이터를 정리해서 저장하고 있습니다.
150 ==> 2018_expense_list_tidy.csv 파일을 생성했습니다.
151 2019년 데이터를 정리해서 저장하고 있습니다.
152 ==> 2019_expense_list_tidy.csv 파일을 생성했습니다.
153 모든 연도의 데이터를 정리해서 파일로 저장했습니다.
154 '''
155
156 years = [2017, 2018, 2019]
157
158 for year in years:
159
160     expense_list_year_dir = data_folder + str(year) + '/'
161     expense_list_tidy_file = "{}_expense_list_tidy.csv".format(year)
162
163     file_name = expense_list_year_dir + expense_list_tidy_file
164     print(file_name, "==> ", end="")
165     print(os.path.isfile(file_name))
166
167
168 def get_file_info(year, data_folder):
169     expense_list_year_dir = data_folder + str(year) + '/'

```



```

170     expense_list_tidy_file = "{}_expense_list_tidy.csv".format(year)
171
172     path_file_name = expense_list_year_dir + expense_list_tidy_file
173     print(path_file_name)
174     result = os.path.isfile(path_file_name)
175
176     # 파일 수정 시간
177     modified_time = datetime.fromtimestamp(os.path.getmtime(path_file_name))
178
179     # 파일 생성 시간
180     created_time = datetime.fromtimestamp(os.path.getctime(path_file_name))
181
182     # 파일 크기
183     file_size = os.path.getsize(path_file_name)
184
185     if(result == True):
186         print("[생성한 CSV 데이터 파일의 정보]")
187         print('* 폴더 위치 :', expense_list_year_dir)
188         print('* 파일 이름 :', expense_list_tidy_file)
189         print('* 수정 시간 :', modified_time.strftime('%Y-%m-%d %H:%M:%S'))
190         print('* 생성 시간 :', created_time.strftime('%Y-%m-%d %H:%M:%S'))
191         print('* 파일 크기 : {0:,} 바이트'.format(file_size))
192
193
194     data_folder = './ch19_project/seoul_expense/'
195     years = [2017, 2018, 2019]
196
197     for year in years:
198
199         get_file_info(year, data_folder)
200         print("")
201     '''
202     [생성한 CSV 데이터 파일의 정보]
203     * 폴더 위치 : ./ch19_project/seoul_expense/2017/
204     * 파일 이름 : 2017_expense_list_tidy.csv
205     * 수정 시간 : 2021-11-21 11:22:37
206     * 생성 시간 : 2021-11-20 23:27:26
207     * 파일 크기 : 21,505,565 바이트
208
209     ./ch19_project/seoul_expense/2018/2018_expense_list_tidy.csv
210     [생성한 CSV 데이터 파일의 정보]
211     * 폴더 위치 : ./ch19_project/seoul_expense/2018/
212     * 파일 이름 : 2018_expense_list_tidy.csv
213     * 수정 시간 : 2021-11-21 11:22:38
214     * 생성 시간 : 2021-11-20 23:32:30
215     * 파일 크기 : 23,012,590 바이트
216
217     ./ch19_project/seoul_expense/2019/2019_expense_list_tidy.csv
218     [생성한 CSV 데이터 파일의 정보]
219     * 폴더 위치 : ./ch19_project/seoul_expense/2019/
220     * 파일 이름 : 2019_expense_list_tidy.csv
221     * 수정 시간 : 2021-11-21 11:22:39
222     * 생성 시간 : 2021-11-20 23:32:31
223     * 파일 크기 : 24,459,405 바이트
224     '''

```

19.3.4 데이터 분석

- 2017년, 2018년, 2019년 서울시의 업무추진비 데이터를 가져와 정제한 후 필요한 데이터만 남겨서 CSV 파일로 저장했다. 이제 이 파일을 이용해 서울시의 업무추진비가 어떻게 집행됐는지 분석해 본다.

- 07~19: 정제 후 깔끔하게 정리된 2017년, 2018년, 2019년치 데이터 파일을 DataFrame 형식으로 가져와 하나로 통합한다.
- 22: df_expense_all의 전체 데이터 구조를 확인하기 위해 info()를 사용한다.
- 23~44: 3년간의 업무추진비 전체 데이터 파일에서 데이터를 읽어오다 보니 데이터 개수가 아주 많다.
- 46~59: df_expense_all 변수의 앞쪽 데이터 일부를 확인해 본다.
- 61~74: df_expense_all 변수의 뒤쪽 데이터 일부도 확인해 본다.

[ch19_project/ex07_데이터분석.py의 일부]

```

01 # 데이터 분석
02 from wordcloud import WordCloud
03 import matplotlib
04 import matplotlib.pyplot as plt
05 import pandas as pd
06
07 data_folder = './ch19_project/seoul_expense/'
08 years = [2017, 2018, 2019]
09
10 df_expense_all = pd.DataFrame()
11
12 for year in years:
13     expense_list_year_dir = data_folder + str(year) + '/'
14     expense_list_tidy_file = "{}_expense_list_tidy.csv".format(year)
15
16     path_file_name = expense_list_year_dir + expense_list_tidy_file
17
18     df_expense = pd.read_csv(path_file_name)
19     df_expense_all = df_expense_all.append(df_expense, ignore_index=True)
20
21
22 print(df_expense_all.info())
23 '''
24 <class 'pandas.core.frame.DataFrame'>
25 RangeIndex: 216557 entries, 0 to 216556
26 Data columns (total 12 columns):
27 #   Column   Non-Null Count  Dtype
28 ---  ---
29 0   제목      216557 non-null object
30 1   부서레벨1 216557 non-null object
31 2   부서레벨2 216273 non-null object
32 3   집행연도  216557 non-null int64
33 4   집행월    216557 non-null int64
34 5   부서명    216478 non-null object
35 6   집행일시  216557 non-null object
36 7   집행장소  214401 non-null object
37 8   집행목적  216535 non-null object
38 9   대상인원  215535 non-null object
39 10  결제방법  216354 non-null object
40 11  집행금액  216557 non-null int64
41 dtypes: int64(3), object(9)
42 memory usage: 19.8+ MB
43 None
44 '''
45
46 print(df_expense_all.head(2))
47 '''
48                                     제목   부서레벨1  부서레벨2   집행연도   ...   집행목적
49 대상
50 인원 결제방법   집행금액
51 0   2017년 1월 장애인복지정책과 업무추진비 집행내역   서울시본청   복지본부   2017   ...   기본소득과 장애인
52 복지 논

```

```

53  의간담회 장애인복지정책팀장 외 2명 카드 76000
54  1 2017년 1월 장애인복지정책과 업무추진비 집행내역 서울시본청 복지본부 2017 ... 장애인단체 활동지
55  원 논
56  의간담회 장애인복지정책과장 외 3명 카드 102000
57
58  [2 rows x 12 columns]
59  '''
60
61  print(df_expense_all.tail(2))
62  '''
63
64  제목 부서레벨1 부서레벨2 ... 대상인원 결제방
65  법
66  집행금액
67  216555 2019년 12월 사업소 서울시립대학교 중앙도서관 사서과 업무추진비 - 전체 사업소 서울시립대학
68  교 ...
69  사서과장 외 8명 카드 23940
70  216556 2019년 12월 사업소 서울시립대학교 중앙도서관 사서과 업무추진비 - 전체 사업소 서울시립대학
71  교 ...
72  이용자 80명 카드 53420
73  [2 rows x 12 columns]
74  '''

```

(1) 연도별 추이 분석

- 2017년부터 2019년까지 연도별 전체 집행 금액을 비교해서 추이 분석을 해 본다.
 - 02: 연도별 업무추진비 집행 횟수를 알고 싶다면 '집행연도' 열에 대해 value_counts()를 수행한다.
 - 12~20: 앞의 데이터를 막대 그래프로 시각화한다.

[ch19_project/ex07_데이터분석.py의 일부]

```

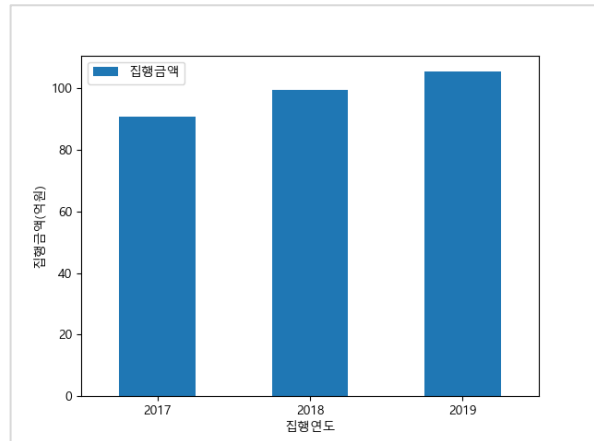
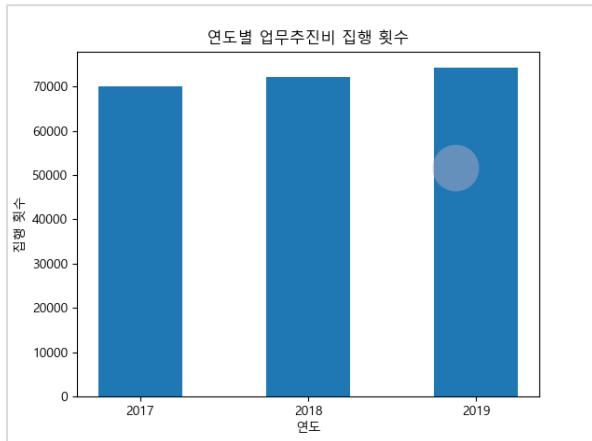
01  # 1. 연도별 추이 분석
02  year_expense = df_expense_all['집행연도'].value_counts()
03  print(year_expense)
04  '''
05  2019    74207
06  2018    72218
07  2017    70132
08  Name: 집행연도, dtype: int64
09  '''
10
11  matplotlib.rcParams['font.family'] = 'Malgun Gothic'
12  matplotlib.rcParams['axes.unicode_minus'] = False
13
14  plt.bar(year_expense.index, year_expense.values,
15          tick_label=year_expense.index, width=0.5)
16  plt.title("연도별 업무추진비 집행 횟수")
17  plt.xlabel("연도")
18  plt.ylabel("집행 횟수")
19  plt.show()
20
21
22  year_total = pd.pivot_table(df_expense_all, index=[
23                          '집행연도'], values=['집행금액'], aggfunc=sum)
24  print(year_total)
25  '''
26
27  집행금액
28  집행연도
29  2017    9076941387

```

```

29 2018 9937556542
30 2019 10532330632
31 '''
32
33 eok_won = 1000000000 # 억원
34 (year_total/eok_won).plot.bar(rot=0) # 'rot = 각도'로 xtick 회전 각도를 지정
35 plt.ylabel('집행금액(억원)')
36 plt.show()

```



(2) 월별 집행금액 분석

- 이번에는 월별로 합계를 구해서 1월에서 12월까지 월별로 업무추진비 집행금액의 변화를 살펴본다.
 - 02: pivot_table()을 이용해 집행월별 집행금액의 합을 구한다.
 - 22~23: 3년간의 데이터에서 월별 집행금액만 추출한다.

[ch19_project/ex07_데이터분석.py의 일부]

```

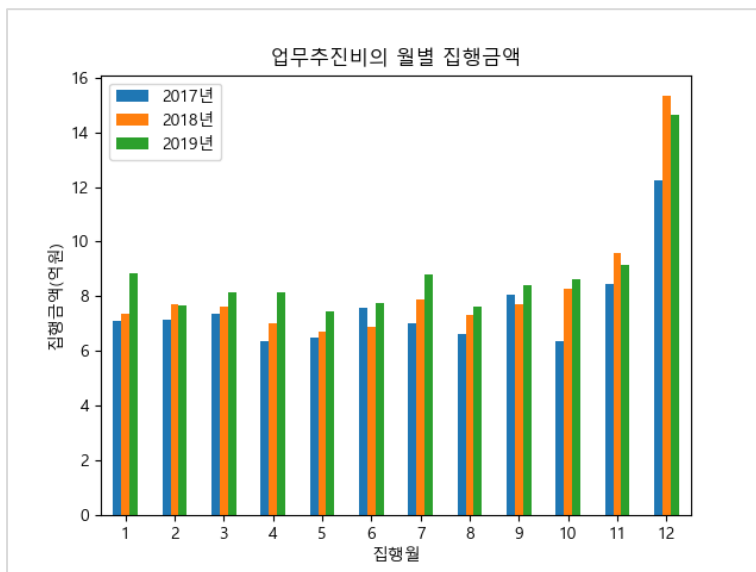
01 # 2. 월별 집행금액 분석
02 month_total = pd.pivot_table(df_expense_all, index=['집행월'], values=['집행금액'],
03                               aggfunc=sum)
04 print(month_total)
05 '''
06             집행금액
07 집행월
08 1      2328469179
09 2      2250971737
10 3      2314589911
11 4      2153704599
12 5      2063883588
13 6      2224855495
14 7      2372256669
15 8      2153716469
16 9      2417217365
17 10     2326108698
18 11     2719921484
19 12     4221133367
20 '''
21
22 year_month_total = pd.pivot_table(df_expense_all, index=['집행월'], columns=['집행연도'],
23                                   values=['집행금액'], aggfunc=sum)
24 print(year_month_total)
25 '''

```

```

26             집행금액
27 집행연도      2017      2018      2019
28 집행월
29 1      710368860  735587570  882512749
30 2      712679864  769360005  768931868
31 3      737250454  761059010  816280447
32 4      635265805  703781418  814657376
33 5      647582378  669044701  747256509
34 6      758257342  690652154  775945999
35 7      701604626  788926477  881725566
36 8      661174850  730290532  762251087
37 9      806170700  769404957  841641708
38 10     637219943  827022975  861865780
39 11     843619171  960310221  915992092
40 12     1225747394 1532116522 1463269451
41 '''
42
43 eok_won = 1000000000 # 억원
44
45 (year_month_total/eok_won).plot.bar(rot=0)
46 plt.ylabel('집행금액(억원)')
47 plt.title("업무추진비의 월별 집행금액")
48 plt.legend(['2017년', '2018년', '2019년'])
49 plt.show()

```



(3) 부서별 집행 내역 분석

■ 이번에는 부서별로 업무 추진비 집행 내역을 분석해 본다.

[ch19_project/ex07_데이터분석.py의 일부]

```

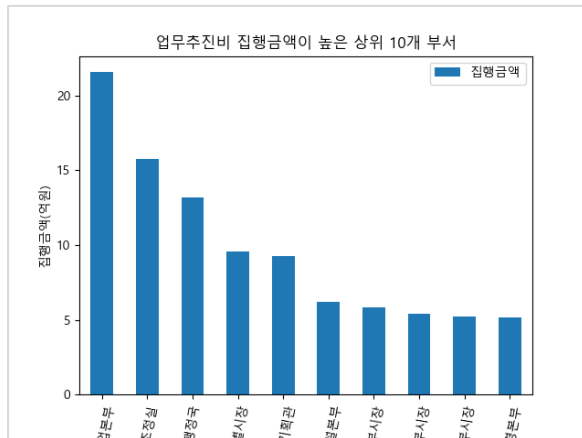
01 # 3. 부서별 집행 내역 분석
02 dept_level1_total = pd.pivot_table(df_expense_all, index=['부서레벨1'], values=['집행금액'],
03                                     aggfunc=sum)
04 print(dept_level1_total)
05 '''
06             집행금액
07 부서레벨1
08 사업소      6552128899

```

```

09 서울시본청      16606242519
10 소방재난본부(소방서)  5147645293
11 의회사무처      1240811850
12 ...
13
14
15 dept_level_2_total = pd.pivot_table(df_expense_all, index=['부서레벨2'], values=['집행금액'],
16                                     aggfunc=sum)
17 print(dept_level_2_total.head())
18 ...
19             집행금액
20 부서레벨2
21 119특수구조단  119225100
22 감사위원회    343281170
23 강남소방서    229660520
24 강동소방서    188773330
25 강북소방서    167700000
26 ...
27
28
29 dept_level_2_total_top10 = dept_level_2_total.sort_values(
30     by=['집행금액'], ascending=False)[0:10]
31 print(dept_level_2_total_top10)
32 ...
33             집행금액
34 부서레벨2
35 상수도사업본부  2156404778
36 기획조정실     1572753168
37 행정국         1320839804
38 서울특별시장   955448760
39 시민소통기획관  923338423
40 도시기반시설본부 620669144
41 경무부시장     581806882
42 행정1부시장    540457390
43 행정2부시장    522277598
44 기후환경본부   515222890
45 ...
46
47 eok_won = 1000000000 # 억원
48
49 (dept_level_2_total_top10/eok_won).plot.bar(rot=80)
50 plt.ylabel('집행금액(억원)')
51 plt.title("업무추진비 집행금액이 높은 상위 10개 부서")
52 plt.show()
53
54
55 korean_font_path = 'C:/Windows/Fonts/malgun.ttf' # 한글 폰트(맑은 고딕) 파일명
56
57 # 워드 클라우드 이미지 생성
58 wc = WordCloud(font_path=korean_font_path, background_color='white',
59               width=800, height=600)
60
61 frequencies = dept_level_2_total['집행금액'] # pandas의 Series 형식이 됨
62 wordcloud_image = wc.generate_from_frequencies(frequencies)
63
64 plt.figure(figsize=(12, 9))
65 plt.axis('off')
66 plt.imshow(wordcloud_image, interpolation='bilinear')
67 plt.show()

```



(4) 요일별 및 시간대별 집행 내역 분석

- 이번에는 '집행일시' 열에 있는 날짜와 시간을 이용해 요일별/시간대별로 업무추진비 집행 내역을 분석한다. 요일별/시간대별로 분석하면 무슨 요일, 무슨 시간대에 업무추진비를 많이 집행했는지 알 수 있다.

[ch19_project/ex07_데이터분석.py의 일부]

```
01 # 4. 요일별 및 시간대별 집행 내역 분석
02 print(df_expense_all['집행일시'].values)
03 '''
04 ['2017-01-26 13:10' '2017-01-25 22:41' '2017-01-24 12:35' ...
05 '2019-12-19 11:34' '2019-12-16 12:39' '2019-12-03 17:35']
06 '''
07
08 expense_date_time = pd.to_datetime(df_expense_all['집행일시'])
09 print(expense_date_time.values)
10 '''
11 ['2017-01-26T13:10:00.000000000' '2017-01-25T22:41:00.000000000'
12 '2017-01-24T12:35:00.000000000' ... '2019-12-19T11:34:00.000000000'
13 '2019-12-16T12:39:00.000000000' '2019-12-03T17:35:00.000000000']
14 '''
15
16 week_day_name = ["월", "화", "수", "목", "금", "토", "일"]
17
18 df_expense_all['집행일시_요일'] = [week_day_name[weekday]
19                                     for weekday in expense_date_time.dt.weekday]
20
21 df_expense_all['집행일시_시간'] = [hour for hour in expense_date_time.dt.hour]
22
23 print(df_expense_all.head(3))
24 '''
25 집행금액 집행일시_요일 집행일시_시간
26 0 2017년 1월 장애인복지정책과 업무추진비 집행내역 서울시본청 복지본부 2017 1 ... 장애인복지정
27 책팀장
28 외 2명 카드 76000 목 13
29 1 2017년 1월 장애인복지정책과 업무추진비 집행내역 서울시본청 복지본부 2017 1 ... 장애인복지정
30 책과장
31 외 3명 카드 102000 수 22
32 2 2017년 1월 장애인복지정책과 업무추진비 집행내역 서울시본청 복지본부 2017 1 ... 장애인복지
33 정책팀
34 장외7명 카드 80000 화 12
35
36 [3 rows x 14 columns]
37 '''
```

```

38
39
40 expense_weekday = df_expense_all['집행일시_요일'].value_counts()
41 print(expense_weekday)
42 '''
43 목      45683
44 화      43812
45 수      42343
46 금      41381
47 월      39498
48 토       2238
49 일       1602
50 Name: 집행일시_요일, dtype: int64
51 '''
52
53 expense_weekday = expense_weekday.reindex(index=week_day_name)
54 print(expense_weekday)
55 '''
56 월      39498
57 화      43812
58 수      42343
59 목      45683
60 금      41381
61 토       2238
62 일       1602
63 Name: 집행일시_요일, dtype: int64
64 '''
65
66 expense_weekday.plot.bar(rot=0)
67 plt.title("요일별 업무추진비 집행 횟수")
68 plt.xlabel("요일")
69 plt.ylabel("집행 횟수")
70 plt.show()
71
72
73 expense_hour_num = df_expense_all['집행일시_시간'].value_counts()
74 print(expense_hour_num)
75 '''
76 12      87518
77 20      23013
78 13      20990
79 19      16766
80 21      12210
81 11       8356
82 14       8311
83 15       7168
84 10       5824
85 18       5783
86 16       5169
87 0        4919
88 9        3486
89 17       2889
90 22       2563
91 8         875
92 7         412
93 23        128
94 1          44
95 6          42
96 3          27
97 4          26
98 5          19
99 2           19
100 Name: 집행일시_시간, dtype: int64
101 '''
102

```



```

103
104 work_hour = [(k+8) % 24 for k in range(24)]
105 expense_hour_num = expense_hour_num.reindex(index=work_hour)
106 print(expense_hour_num)
107 '''
108      8      875
109      9     3486
110     10     5824
111     11     8356
112     12    87518
113     13   20990
114     14     8311
115     15     7168
116     16     5169
117     17     2889
118     18     5783
119     19    16766
120     20   23013
121     21   12210
122     22   2563
123     23     128
124      0    4919
125      1      44
126      2      19
127      3      27
128      4      26
129      5      19
130      6      42
131      7     412
132 Name: 집행일시_시간, dtype: int64
133 '''
134
135
136 expense_hour_num.plot.bar(rot=0)
137 plt.title("시간별 업무추진비 집행 횟수")
138 plt.xlabel("집행 시간")
139 plt.ylabel("집행 횟수")
140 # plt.show()
141
142
143 expense_hour_total = pd.pivot_table(df_expense_all, index=['집행일시_시간'],
144                                     values=['집행금액'], aggfunc=sum)
145 print(expense_hour_total.head())
146 '''
147 집행일시_시간
148 0      842523116
149 1       7024161
150 2      2265190
151 3      7215762
152 4      5818431
153 '''
154
155
156 eok_won = 100000000 # 억원
157 expense_hour_total = expense_hour_total.reindex(index=work_hour)
158
159 (expense_hour_total/eok_won).plot.bar(rot=0)
160 plt.ylabel('집행금액(억원)')
161 plt.title("시간별대 업무추진비 집행금액")
162 plt.show()

```

