





- ✓ 데이터분석을 위한 문제 정의
- ☑ 연령별 인구구조가 유사한 지역을 찾기 위한 알고리즘
- ☑ 입력 받은 지역과 연령별 인구구조가 가장 유사한 지역을 찾기



학습목표

- ☑ 데이터분석을 위한 문제를 정의할 수 있다.
- ☑ 연령별 인구구조가 유사한 지역을 찾기 위한 알고리즘을 이해할 수 있다.
- ☑ 패키지 numpy의 배열 array를 사용해 문자열을 정수로 변환할 수 있다.
- ☑ 입력 받은 지역과 연령별 인구구조가 가장 유사한 지역을 찾기 위한 코드를 구현할 수 있다.

LESSON 01

명확한 문제 정의



☑ 명확한 문제 정의



→ 데이터 살펴보며 질문하기 (1/2)

❤ 데이터를 살펴보는 방법

- ☑ 엑셀 같은 스프레드시트 프로그램 등을 활용해 데이터 자세히 살펴보기
- ☑ 데이터가 담고 있는 내용(또는 담고 있지 않은 내용)
- 데이터가 기록된 기간은 언제부터 언제까지인지
- ☑ 어떤 형태로 시각화해보면 어떤 정보들을 알 수 있을지 생각해보기
- ☑ 데이터를 보며 궁금한 내용들 자유롭게 질문하기



→ 데이터 살펴보며 질문하기 (2/2)



- 전국에서 영유아들이 가장 많이 사는 지역은 어디일까?
- ☑ 보통 학군이 좋다고 알려진 지역에는 청소년들이 많이 살까?
- ☑ 광역시 데이터를 10년 단위로 살펴보면 청년 비율이 줄고 있다는 사실을 알 수 있을까?
- 서울에서 지난 5년간 인구가 가장 많이 증가한 구는 어디일까?
- ☑ 우리 동네의 인구 구조와 가장 비슷한 동네는 어디일까?

영유아들?

청소년들?

청년비율?

인구 구조와 가장 비슷한 동네?

●문제를 명확히 정의할 필요가 있어 보입니다.



⊸ 질문을 명확한 문제로 정의하기

- **❷** 예를 들어 아래와 같이 문제를 좀 더 명확하게 정의할 수 있습니다.
 - ☑ 전국에서 영유아들이 가장 많이 사는 지역은 어디일까?
 - ★ 전국에 있는 읍면동 중 만 0세 이상 6세 이하의 인구 비율이 높은 상위 10곳은?
 - 우리 동네의 인구 구조와 가장 비슷한 동네는 어디일까?
 - ◆ 전국에서 우리 동네의 연령별 인구 구조와 가장 형태가 비슷한 지역은 어디일까?

위 문제를 해결하기 위해서, 알고리즘을 어떻게 설계해야 할까요?



⊸ 알고리즘 설계하기

② [문제] 전국에서 우리 동네의 연령별 인구 구조와 가장 형태가 비슷한 지역은 어□일까?

- ☑ Step 1) 데이터를 읽어온다.
- ☑ Step 2) 궁금한 지역의 이름을 입력 받는다.
- ☑ Step 3) 궁금한 지역의 인구 구조를 저장한다.
- ☑ Step 4) 궁금한 지역의 인구 구조와 가장 비슷한 인구 구조를 가진 지역을 찾는다.
- ☑ Step 5) 가장 비슷한 곳의 인구 구조와 궁금한 지역의 인구 구조를 시각화한다.

numpy를 활용하여 궁금한 지역의 인구 데이터를 출력하는 코드를 작성해 보겠습니다!

LESSON 02

알고리즘 코드 표현



₩ 알고리즘 코드 표현



→ 데이터 읽어 오기

```
import csv

f = open('age.csv', encoding='cp949')
data = csv.reader(f)

header = next(data) # 헤터 제거
print(header)

for row in data:
    print(row)
    break

f.close()
```



→ Step 1) 데이터 읽어 오기(실행결과)

'2023년09월 계 3세'. '2023년09월 계 4세'. '2023년09월 계 8세'. '2023년09월 계 9세'. 세', '2023년09월 계 12세', '2023년09월 계 13세', '2023년09월 계 14세', '2023년09월 | 16세|, | '2023년09월_계_17세|, | '2023년09월_계_18세|, | '2023년09월_계_19세|, | '2023년09월 '2023년09월 계 29세'. 109월 계 31세'. '2023년09월_계_32세', '2023년09월_계_33세', '2023년09월_계_34세', '2023년09월_계_38세', 23년09월 계 41세', '2023년09월 계 42세', '2023년09월 계 43세', 계 60세', '2023년09월 계 61세', '2023년09월 계 62세', '2023년09월 계 63세', '2023년09월 계 월_계_65세', '2023년09월_계_66세', '2023년09월_계_67세', '2023년09월_계_68세', '2023년09월_계_72세', '2023년09월_계_73세' '2023년09월 계 77세'. '2023년09월_계_78세' 23년09월 계 80세'. '2023년09월 계 81세'. '2023년09월 계 82세'. '2023년09월 계 83세'. '2023년09월_계_85세', '2023년09월_계_86세', '2023년09월_계_87세', '2023년09월_계_88세', '2023년09월_계_89 세', '2023년09월_계_90세', '2023년09월_계_91세', '2023년09월_계_92세', '2023년09월_계_93세', '2023년09월_계 94세', '2023년09월_계_95세', '2023년09월_계_96세', '2023년09월_계_97세', '2023년09월_계_98세', '2023년09월 계 99세', '2023년09월 계 100세 이상' (1100000000)', '9,407,540', '9,407,540', '38,101', '41,599', '43,518', '44,893', 1,107', '56,215', '64,176', '66,545', '66,588', '68,791', '73,452', '73,667', '70,811', '70,182', '75,957' '76,605', '71,528', '72,968', '85,145', '88,965', '95,128', '114,431', '130,725', '134,291', '147,512', '15 7,700', '164,260', '165,915', '169,505', '169,721', '167,914', '154,238', '145,524', '139,088', '132,125', '129,417', '129,449', '128,551', '132,545', '141,799', '149,998', '151,352', '151,696', '141,877', |5', '134,730', '132,627', '140,149', '152,989', '154,654', '157,451', '166,643', '157,897', '159,261', '151 094', '139,108', '134,184', '138,870', '130,717', '141,481', '135,134', '151,553', '151,136', '137,487', '12 9.042'. '128.405'. '117.703'. '123.829'. '101.401'. '92.397'. '96.553'. '70.066'. '78.008'. '78.177'. '76.17 1', '75,999', '56,031', '55,814', '53,247', '54,490', '56,569', '44,137', '37,671', '34,905', '29,591', '25, 447', '21,166', '18,331', '14,079', '11,132', '8,984', '6,319', '5,441', '4,430', '3,158', '2,286', '1,306' '916', '703', '1,505'l

문자열 자료형이고 숫자 사이에 쉼표 (,)가 있습니다.

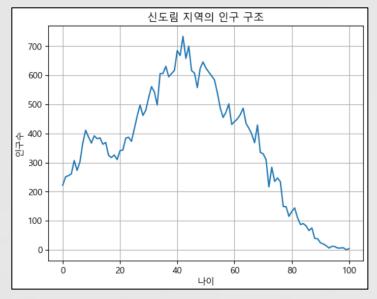


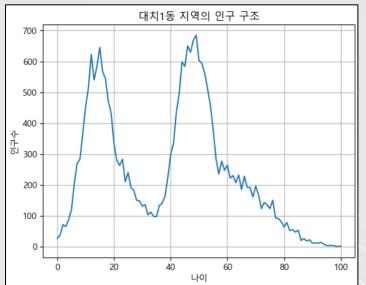
⊸ 궁금한 지역의 인구 구조 시각화하기

```
신도림 지역의 인구 구조
import csv
import numpy as np
import matplotlib.pyplot as plt
f = open('age.csv', encoding='cp949')
data = csv.reader(f)
next (data)
                                                                                        300 -
|name = input('인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : ')
for row in data :
   row[2:] = map(lambda s: s.replace(',', ''), row[2:])
   if name in row[0] :
       home = np.array(row[3:], dtype = int)
       break
                                                                                                    대치1동 지역의 인구 구조
f.close()
plt.rc('font', family ='Malgun Gothic')
plt.plot(home)
plt.title(f'{name} 지역의 인구 구조')
plt.grid(True)
                                                                                      음
300·
plt.xlabel('나이')
plt.ylabel('인구수')
plt.show()
                                                                                        100
인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : 신도림
```



→ 궁금한 지역의 인구 구조 시각화하기







- → 궁금한 지역의 인구 구조와 가장 비슷한 인구 구조를 가진 지역 찾기
- ❤️ 내가 알고자 하는 궁금한 지역: A
- ❤️ 비교할 지역(= A지역을 제외한 나머지 지역): B

A와 B의 인구 구조가 비슷하다는 것을 어떻게 알 수 있을까요?



```
# 관심 지역과 '화촌면'의 연령별 인구를 그리기
import csv
import numpy as np
import matplotlib.pyplot as plt
                                                                                            • 관심지역과 '화촌면'의 인구 구조 시각화
f = open('age.csv', encoding='cp949')
data = csv.reader(f)
header = next(data)
name = input('인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : ')
for row in data :
   row[2:] = map(lambda s: s.replace(',', ''), row[2:])
   if name in row[0] :
      home = np.array(row[3:1, dtype = int)
   if '화촌면' in row[0] :
      homeB = np.array(row[3:], dtype = int)
      print(row[0])
                                                                                    신도림 지역의 인구 구조
f.close()
                                                                                                          화촌면
plt.rc('font', family ='Malgun Gothic')
plt.plot(home, label=name)
plt.plot(homeB, label='화촌면')
plt.title(name +' 지역의 인구 구조')
plt.grid(True)
                                                                    400
plt.xlabel('LHOI')
plt.ylabel('인구수')
                                                                    300
plt.legend()
                                                                    200
plt.show()
                                                                    100
인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : 신도림
강원특별자치도 홍천군 화촌면(5172031000)
                                                                                                               100
                                                                                           나이
```

인구수 차이가 아니라 인구 비율을 고려해 보겠습니다.



```
# 관심 지역과 '화촌면'의 연령별 인구비율 그리기
import csv
import numpy as np
import matplotlib.pyplot as plt
f = open('age.csv', encoding='cp949')
                                                                                              • 관심지역과 '화촌면'의 연령별비율 시각화
data = csv.reader(f)
header = next(data)
name = input('인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : ')
for row in data
   # row[2]를 나누기에 사용하므로 먼저 정수로 변환
   row[2:] = map(int, map(lambda s: s.replace(',', ''), row[2:]))
   if name in row[0] :
      home = np.array(row[3:], dtype = int) / row[2]
   if '화촌면' in row[0] :
      homeB = np.array(row[3:], dtype = int) / row[2]
                                                                               신도림 지역과 화촌면 지역의 인구 구조 비교
f.close()
                                                                      0.035
                                                                                                               신도림
# plt.style.use('ggplot')
                                                                                                               화촌면
                                                                      0.030
plt.rc('font', family ='Malgun Gothic')
                                                                      0.025
plt.plot(home, label=name)
plt.plot(homeB, label='화촌면')
                                                                   - 0.020
plt.title(f'{name} 지역과 화촌면 지역의 인구 구조 비교')
                                                                   녌 0.015
plt.grid(True)
plt.xlabel('나이')
plt.vlabel('연구비율')
                                                                      0.010
plt.legend()
                                                                      0.005
plt.show()
인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : 신도림
                                                                      0.000
```

전체 인구수가 다르더라도 인구 구조를 비교할 수 있을 것 같습니다.



```
import csv
import numpy as np
import matplotlib.pyplot as plt
f = open('age.csv', encoding='cp949')
data = csv.reader(f)
                                                아무것도 출력되지 않습니다!
header = next(data)
print(type(data))
name = input('인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : ')
for row in data :
   # row[2]를 나누기에 사용하므로 먼저 정수로 변환
   row[2:] = map(int, map(lambda s: s.replace(',', ''), row[2:]))
   if name in row[0] :
      home = np.array(row[3:], dtype = int) / row[2]
# print(home)
# 반복자는 한번 반복한 이후 더 이상 반복을 못함
for row in data:
   print(row)
f.close()
<class ' csv.reader'>
인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : 신도림
```



```
import csv
import numpy as np
import matplotlib.pyplot as plt
f = open('age.csv', encoding='cp949')
data = csv.reader(f)
header = next(data)
                         다시 반복을 하기 위해 리스트로 변환해서 작업
print(type(data))
data = list(data) # 奉카
print(type(data))
|name = input('인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : ')
for row in data:
   # row[2]를 나누기에 사용하므로 먼저 정수로 변환
   row[2:] = map(int, map(lambda s: s.replace(',', ''), row[2:]))
   if name in row[0] :
      home = np.array(row[3:], dtype = int) / row[2]
for row in data :
   print(row)
f.close()
```



- ⊸ 궁금한 지역의 인구 구조와 가장 비슷한 인구 구조를 가진 지역 찾기
- **❷** 리스트의 항목을 수정해 저장
 - ☑ 원래의 리스트에 반영

```
d = [['1', '2'], ['10', '20', '30'], ['11', '22', '33']]
print(d) # 수정 전 출력

# 리스트의 내부 항목을 수정해서 대입하면 바로 리스트가 수정됨
for item in d:
    for i in range(len(item)):
        item[i] = int(item[i])

print(d) # 수정 후 출력

[['1', '2'], ['10', '20', '30'], ['11', '22', '33']]
[[1, 2], [10, 20, 30], [11, 22, 33]]
```



⊸ 가장 비슷한 곳의 인구 구조와 궁금한 지역의 인구 구조 시각화하기

알 알고리즘 버전 ("차이의 제곱의 합"이 가장 작은 지역)

- ☑ ⓐ 전국의 모든 지역 중 한 곳 (B)을 선택한다.
- ☑ ⓑ 궁금한 지역 A의 0세 인구 비율에서 B의 0세 인구 비율을 뺀다.
- ☑ ⑥ ⑥를 "100세 이상 인구수"에 해당하는 값까지 반복한 후 각각의 차이의 제곱의 합을 구한다.
- ☑ ⓓ 전국의 모든 지역에 대해 반복하며 그 차이가 가장 작은 지역을 찾는다.

| | 0세 | 1세 | 2세 | 100세이상 |
|-----------|-------------|-------------|-------------|-----------------|
| A지역(away) | 인구수/지역전체인구수 | 인구수/지역전체인구수 | 인구수/지역전체인구수 | 인구수/지역전체인구수 |
| | | | | |
| B지역(home) | 인구수/지역전체인구수 | 인구수/지역전체인구수 | 인구수/지역전체인구수 | 인구수/지역전체인구수 |

np.sum((home - away) ** 2)



⊸ 가장 비슷한 곳의 인구 구조와 궁금한 지역의 인구 구조 시각화하기

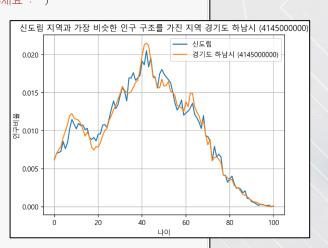
- ❤️ 알고리즘 버전 ("차이의 제곱의 합"이 가장 작은 지역)
 - ☑ ② 전국의 모든 지역 중 한 곳 (B)을 선택한다.
 - ☑ ⓑ 궁금한 지역 A의 0세 인구 비율에서 B의 0세 인구 비율을 뺀다.
 - ☑ ⑥ ⑥를 "100세 이상 인구수"에 해당하는 값까지 반복한 후 각각의 차이의 제곱의 합을 구한다.
 - ☑ ⓓ 전국의 모든 지역에 대해 반복하며 그 차이가 가장 작은 지역을 찾는다.

```
Iname = input('인구 구조가 알고 싶은 지역의 이름(옵면동 단위)을 입력해주세요 : ')
for row in data :
   # row[2]를 나누기에 사용하므로 먼저 정수로 변환
   row[2:] = map(int, map(lambda s: s.replace(',', ''), row[2:]))
    if name in row[0]:
       home = np.array(row[3:], dtype = int) / row[2]
for row in data :
    if name not in row[0]:
       if row[2] == 0:
           continue
       away = np.array(row[3:], dtype = int) / row[2]
       s = np.sum((home - away) ** 2)
       if s < min val:
           min_val = s
           result_name = row[0]
           result = away
```



→ 가장 비슷한 곳의 인구 구조와 궁금한 지역의 인구 구조 시각화하기

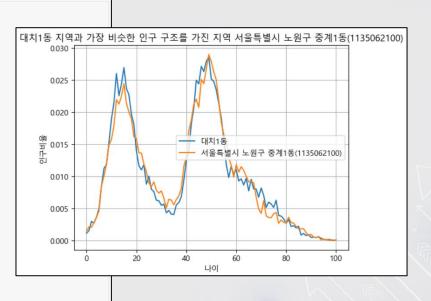
```
import csv
                                                plt.rc('font', family ='Malgun Gothic')
import numpy as np
                                                plt.plot(home, label=name)
import matplotlib.pyplot as plt
                                                plt.plot(result, label=result_name)
                                                plt.title(f'{name} 지역과 가장 비슷한 인구 구조를 가진 지역 {result_name}')
                                                plt.xlabel('나이'
f = open('age.csv', encoding='cp949')
                                                plt.ylabel('인구비율')
data = csv.reader(f)
                                                plt.grid(True)
header = next(data)
                                                plt.legend()
data = list(data) #추가
                                                plt.show()
                                                 인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : 신도림
min val = 1
result name = ''
result = 0
name = input('인구 구조가 알고 싶은 지역의 이름(옵면동 단위)을 입력해주세요 : ')
for row in data:
   # row[2]를 나누기에 사용하므로 먼저 정수로 변환
   row[2:] = map(int. map(lambda s: s.replace('.', ''), row[2:]))
    if name in row[0]:
       home = np.arrav(row[3:], dtype = int) / row[2]
for row in data:
    if name not in row[0] :
       if row[2] == 0:
           cont inue
       away = np.array(row[3:], dtype = int) / row[2]
       s = np.sum((home - away) ** 2)
       if s < min val:
           min val = s
           result name = row[0]
           result = awav
f.close()
```





→ 가장 비슷한 곳의 인구 구조와 궁금한 지역의 인구 구조 시각화하기

```
import csv
import numby as no
import matplotlib.pyplot as plt
f = open('age.csv', encoding='cp949')
data = csv.reader(f)
header = next(data)
data = list(data) #季가
min val = 1
result name = ''
result = 0
|name = input('인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : ')
for row in data :
   # row[2]를 나누기에 사용하므로 먼저 정수로 변환
   row[2:] = map(int, map(lambda s: s.replace(',', ''), row[2:]))
   if name in row[0]
       home = np.array(row[3:], dtype = int) / row[2]
for row in data :
   if name not in row[0]:
      if row[2] == 0:
          continue
       away = np.array(row[3:], dtype = int) / row[2]
       s = np.sum((home - away) ** 2)
       if s < min_val:
          min val = s
          result name = row[0]
          result = away
f.close()
plt.rc('font', family ='Malgun Gothic')
plt.plot(home, label=name)
plt.plot(result, label=result name)
plt.title(f'{name} 지역과 가장 비슷한 인구 구조를 가진 지역 {result_name}')
plt.xlabel('나이')
plt.ylabel('인구비율')
plt.grid(True)
plt.legend()
plt.show()
인구 구조가 알고 싶은 지역의 이름(읍면동 단위)을 입력해주세요 : 대치1동
```



SUMMARY

학습정긴





...

- ⊙ 데이터분석을 위한 문제를 보다 명확하게 정의
 - 코드로 구현이 가능하도록
- 🌣 csv.reader(f)의 결과
 - ≫ 클래스 _csv.reader로 반복자(iterator)
 - 한 번 반복하면 다시 반복이 불가능
- 🌣 차이 분석
 - '차이 값의 제곱의 합'이 작은 것이 차이가 작다
- 내장 함수 map을 사용한 '콤마 숫자 문자열' 숫자 변환
 - row[2:] = map(int, map(lambda s: s.replace(',', ''), row[2:]))
- 🧑 data의 자료 형 변환
 - np.array(data, dtype=int)



