

서울특별시 자치구별 학구열 분석 & 서울런 교육 사업

헬로월드조
정가연 김지은 배지환 박수현 허지원

CONTENTS

STEP 1

주제 배경

STEP 2

사용한 데이터 소개

STEP 3

전처리 과정

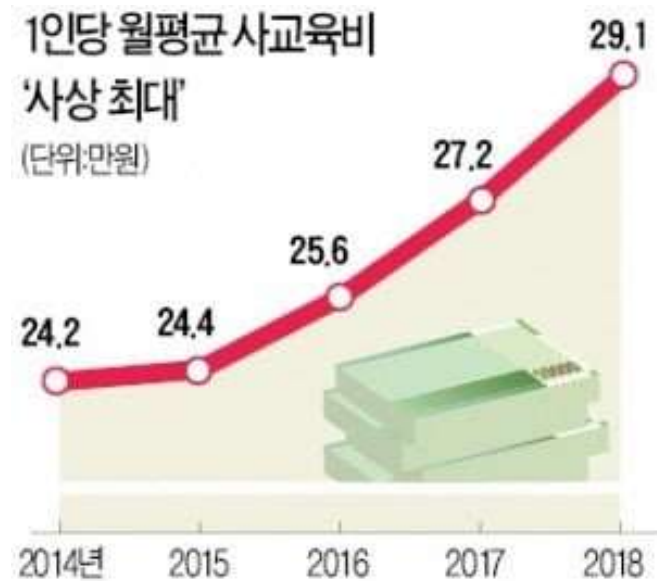
STEP 4

분석 과정

STEP 5

결론 및 활용방안

주제 배경



출처 : 한경닷컴

우리나라가 교육에 대해 매우 민감하고 예민하다는 것을 여러 매체들을 통해서 알 수 있다

사용 데이터 소개

메인 데이터	서브 데이터	외부 데이터	제작 데이터
서울특별시 학원 교습소정보.csv	서울시 가구형태별 가구 및 가구원 (동별) 통계.csv	2021 서울시 학급당 학생 수 통계. csv	helloworld_data_set.xlsx
서울특별시 학교 기본정보.csv	서울시 주민등록인구 (동별) 통계.csv		
서울특별시 유치원 일반현황.csv	서울특별시 유치원 통학차량 현황.csv		
서울시 주민등록연앙인구(연령별동별) 통계.csv	자치구별 1인당 지역내총생산 및 수준지수.csv		

사용 데이터 소개

메인 데이터	서브 데이터	외부 데이터	제작 데이터
서울특별시 학원 교습소정보.csv	서울시 가구형태별 가구 및 가구원 (동별) 통계.csv	2021 서울시 학급당 학생 수 통계. csv	helloworld_data_set.xlsx
서울특별시 학교 기본정보.csv	서울시 주민등록인구 (동별) 통계.csv		
서울특별시 유치원 일반현황.csv	서울특별시 유치원 통학차량 현황.csv		
서울시 주민등록연앙인구(연령별동별) 통계.csv	자치구별 1인당 지역내총생산 및 수준지수.csv		

사용 데이터 소개

메인 데이터	서브 데이터	외부 데이터	제작 데이터
서울특별시 학원 교습소정보.csv	서울시 가구형태별 가구 및 가구원 (동별) 통계.csv	2021 서울시 학급당 학생 수 통계. csv	helloworld_data_set.xlsx
서울특별시 학교 기본정보.csv	서울시 주민등록인구 (동별) 통계.csv		
서울특별시 유치원 일반현황.csv	서울특별시 유치원 통학차량 현황.csv		
서울시 주민등록연앙인구(연령별동별) 통계.csv	자치구별 1인당 지역내총생산 및 수준지수.csv		

사용 데이터 소개

메인 데이터	서브 데이터	외부 데이터	제작 데이터
서울특별시 학원 교습소정보.csv	서울시 가구형태별 가구 및 가구원 (동별) 통계.csv	2021 서울시 학급당 학생 수 통계. csv	helloworld_data_set.xlsx
서울특별시 학교 기본정보.csv	서울시 주민등록인구 (동별) 통계.csv		
서울특별시 유치원 일반현황.csv	서울특별시 유치원 통학차량 현황.csv		
서울시 주민등록연앙인구(연령별동별) 통계.csv	자치구별 1인당 지역내총생산 및 수준지수.csv		

전처리 과정

Main Data 1

01

```
1 # 도로명상세주소 + 도로명주소 + 교습계열명 + 교습과정명 + 일시수용능력인원합계 결측치 포함 행 제거
2 academy = academy.dropna(subset=['도로명상세주소', '교습계열명', '일시수용능력인원합계'])

1 academy['행정구역명'] = academy['행정구역명'].fillna(academy['도로명주소'].str.split(' ')[1])
2 academy['인양수강료내용'] = academy['인양수강료내용'].fillna('일수없음')
3 academy['기숙사학원여부'] = academy['기숙사학원여부'].fillna('일수없음')
4 academy['휴원시작일자'] = academy['휴원시작일자'].fillna(0)
5 academy['휴원종료일자'] = academy['휴원종료일자'].fillna(0)

1 # '서울특별시마포구'같이 외여쓰기가 안 되어 있는 데이터가 있어 '서울특별시'를 제거하고 앞의 공백을 제거
2 academy['도로명주소'] = academy['도로명주소'].str.replace('서울특별시', '')
3 academy['도로명주소'] = academy['도로명주소'].str.strip().tolist()
```

결측치 제거

02

```
1 academy.duplicated().sum() # 중복값 1개 확인
1
1 academy.shape # 중복값 제거 전 행, 열 개수
(24984, 16)

1 academy.drop_duplicates(inplace=True) # 중복값 제거

1 academy.shape # 중복값 제거 후 행, 열 개수
(24983, 16)
```

중복값 제거

03

```
1 academy['휴원종료일자'].unique() # 이상치 999999999

1 # 이?미용 -> 이 미용, 보습?논술 -> 보습 논술 / 데이터 출처 찾아서 ?를 ' '로 수정
2 academy['교습과정명'] = academy['교습과정명'].str.replace('?', ' ')
3 academy['교습계열명'] = academy['교습계열명'].str.replace('?', ' ')

1 # 일시수용능력인원합계 99999 이상치를 보고 10000 이하 / 99999 이상 학원/교습소들은 존재안, 정확으로 수업/병행 인원
2인 = academy['일시수용능력인원합계'] < 9999
3 academy['일시수용능력인원합계'] = academy['일시수용능력인원합계'].where(2, 0)
4
5 # 인원합계 50000 이상을 이상치로 보고 해당 행 삭제(50000명 이상인 학원/교습소 = 정략, 온라인 학원)
6 is_n = academy['인원합계'] < 50000
7 academy['인원합계'] = academy['인원합계'].where(6, np.nan)
8 academy = academy.dropna()

1 # 등록일자, 개설일자 Dtype object에서 datetime으로 변경
2 academy['등록일자'] = pd.to_datetime(academy['등록일자'], format='%Y%m%d')
3 academy['개설일자'] = pd.to_datetime(academy['개설일자'], format='%Y%m%d')

1 # 전처리와 이상치 처리 후 인덱스 번호 재설정
2 academy.reset_index(drop=True, inplace=True)
```

이상치 처리

전처리 과정

Main Data 2

01

```
1 # 해당 열들은 대부분의 값이 비어있음 -> 유용하지 않은 열이므로 삭제
2 # 고등학교구분명, 특수목적고등학교계열명, 계열명, 학과명
3 del school['고등학교구분명']
4 del school['특수목적고등학교계열명']
5 del school['계열명']
6 del school['학과명']
7 del school['도로명상세주소']

1 school['주소구'] = school['도로명주소'].str.split().str[1] # 자치구만 뽑아내기
2 school['주소구']

1 # 결측치 행 제거
2 school = school.dropna(subset=['도로명주소', '주소구'])
```

결측치 제거

02

```
1 school.duplicated().sum()
```

```
1 school.shape
```

```
(3867, 15)
```

```
1 school.drop_duplicates(inplace=True)
```

```
1 school.shape
```

```
(1446, 15)
```

중복값 제거

03

#이상치 없음

이상치 처리

전처리 과정

Main Data 3

01

[illegible]

결측치 제거

02

```
In [14]: kd.duplicated().sum()
```

```
Out[14]: 0
```

중복값 제거

03

```
1 #이상치 처리
2 kd['혼합학급수']=kd['혼합학급수'].clip(0,10)
3 kd['혼합유아수']=kd['혼합유아수'].clip(0, 44)
```

이상치 처리

전처리 과정

Main Data 4

01

```
1 #결측치 중앙값으로 대체
2 pop_age['0~4세'] = pop_age['0~4세'].fillna(pop_age['0~4세'].median())
3 pop_age['0~4세'] = pop_age['0~4세'].astype('int')
4 pop_age['5~9세'] = pop_age['5~9세'].fillna(pop_age['5~9세'].median())
5 pop_age['5~9세'] = pop_age['5~9세'].astype('int')
6 pop_age['10~14세'] = pop_age['10~14세'].fillna(pop_age['10~14세'].median())
7 pop_age['10~14세'] = pop_age['10~14세'].astype('int')
8 pop_age['15~19세'] = pop_age['15~19세'].fillna(pop_age['15~19세'].median())
9 pop_age['15~19세'] = pop_age['15~19세'].astype('int')
10 pop_age['20~24세'] = pop_age['20~24세'].fillna(pop_age['20~24세'].median())
11 pop_age['20~24세'] = pop_age['20~24세'].astype('int')
12 pop_age['30~34세'] = pop_age['30~34세'].fillna(pop_age['30~34세'].median())
13 pop_age['30~34세'] = pop_age['30~34세'].astype('int')
14 pop_age['45~49세'] = pop_age['45~49세'].fillna(pop_age['45~49세'].median())
15 pop_age['45~49세'] = pop_age['45~49세'].astype('int')
16 pop_age['50~54세'] = pop_age['50~54세'].fillna(pop_age['50~54세'].median())
17 pop_age['50~54세'] = pop_age['50~54세'].astype('int')
32 pop_age['100세 이상*'] = pop_age['100세 이상*'].fillna(0)
33 pop_age['100세 이상*'] = pop_age['100세 이상*'].astype('int')
```

결측치 제거

02

```
1 #중복값 삭제
2 #자치구열에 함께 행 삭제
3 all = pop_age[pop_age['자치구'].str.contains('합계')].index
4 pop_age.drop(all, inplace=True)
```

중복값 제거

03

#이상치 없음

이상치 처리

전처리 과정

Sub Data 1

01

```
4 #첫번째 행 삭제
5 family = family.drop([0,1])
6
7 #송가구를 수치형으로 변경
8 family['송가구'] = family['송가구'].astype('int')
9
10 #1인가구를 수치형으로 변경
11 family['1인가구'] = family['1인가구'].astype('int')
12
13 #2인가구를 수치형으로 변경
14 family['2인가구'] = family['2인가구'].st
15 try:
16     family['2인가구'] = family['2인가구']
17 except:
18     family['2인가구'] = family['2인가구'].astype('float')
19
20 #3인가구를 수치형으로 변경
21 family['3인가구'] = family['3인가구'].astype('int')
22
23 #4인가구를 수치형으로 변경
24 family['4인가구'] = family['4인가구'].astype('float')
25
26 #5인가구를 수치형으로 변경
27 family['5인가구'] = family['5인가구'].astype('float')
28
29 #6인가구를 수치형으로 변경
30 family['6인가구'] = family['6인가구'].astype('float')
31
32 family['4인가구'] = family['4인가구'].astype('float')
33 except:
34     family['7인미상가구'] = family['7인미상가구'].astype('float')
```

결측치 제거

1 #결측치 확인 -> 없음
2 family.isnull().sum()

02

중복값 처리
family = family.drop_duplicates()

중복값 제거

03

```
#03 결측치 확인
family['자치구'].value_counts()

자치구
송파구 28
강남구 23
관악구 22
성북구 21
강서구 21
노원구 20
강동구 19
서초구 19
영등포구 19
양천구 19
중도구 18
성동구 18
은평구 17
마포구 17
구로구 17
송파구 17
흥안구 17
송구 16
홍안구 16
광안구 16
도봉구 15
동대문구 15
서대문구 15
강북구 14
금천구 11
Name: 자치구, dtype: int64
```

이상치 처리

전처리 과정

Sub Data 2

01

```
1 #첫번째,두번째,세번째 행 삭제
2 pop_dong = pop_dong.drop([0,1,2])
3 pop_dong.reset_index(drop=True)
```

```
1 pop_dong['세대'] = pop_dong['세대'].str.replace('.', '')
2 pop_dong['세대'].unique()
3 try:
4     pop_dong['세대'] = pop_dong['세대'].astype('int')
5 except:
6     pop_dong['세대'] = pop_dong['세대'].astype('float')
7
8 pop_dong['종인구'] = pop_dong['종인구'].str.replace('.', '')
9 pop_dong['종인구'].unique()
10 try:
11     pop_dong['종'] = pop_dong['종'].astype('int')
12 except:
13     pop_dong['종'] = pop_dong['종'].astype('float')
14
15 pop_dong['남자인구'] = pop_dong['남자인구'].str.replace('.', '')
16 pop_dong['남자인구'].unique()
17 try:
18     pop_dong['남자인구'] = pop_dong['남자인구'].astype('int')
19 except:
20     pop_dong['남자인구'] = pop_dong['남자인구'].astype('float')
21
```

1 #결측치 확인 -> 없음
2 pop_dong.isnull().sum()

결측치 제거

02

```
1 # 중복값 처리
2 pop_dong = pop_dong.drop_duplicates()
```

중복값 제거

03

```
1 #이상치 없음
2 pop_dong['자치구'].value_counts()
```

```
1 #이상치 없음
2 pop_dong['동'].value_counts()
```

이상치 처리

전처리 과정

Sub Data 3

01

```
1 #결측치 확인 -> 없음  
2 kinder_bus.isnull().sum()
```

결측치 제거

02

```
1 # 중복값 처리  
2 kinder_bus = kinder_bus.drop_duplicates()
```

중복값 제거

03

```
1 #이상치 없음  
2 kinder_bus['교육지원청명'].value_counts()  
3 kinder_bus['유치원명'].value_counts()  
4 kinder_bus['설립유형'].value_counts()  
5 kinder_bus['주소'].value_counts()  
6 kinder_bus.describe().T
```

이상치 처리

전처리 과정

Sub Data 4

01

```
1 #결측치 확인 -> 없음  
2 grdp.isnull().sum()
```

결측치 제거

02

```
1 # 중복값 처리  
2 grdp = grdp.drop_duplicates()
```

중복값 제거

03

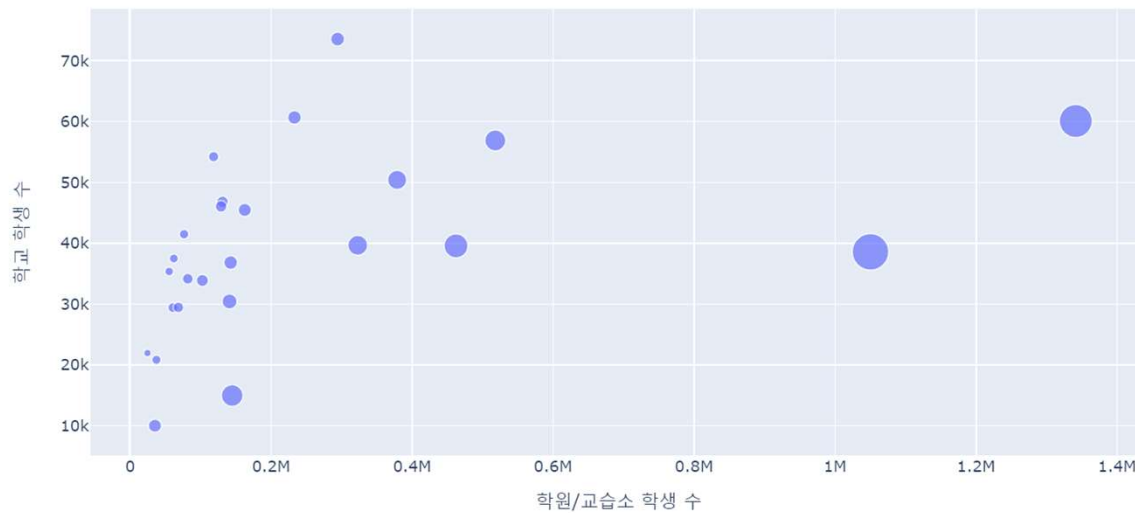
```
1 #이상치 확인 -> 없음  
2 grdp.value_counts()
```

이상치 처리

분석 과정

주제 1 - 자치구별 학교 다니는 학생수 비율, 자치구별 학원 다니는 학생수 비율

학원/교습소에 다니는 학생 수와 학교 학생 수와의 상관관계



1) 자치구별 grdp 총 수와 자치구별 학교수의 상관관계

1. 자치구별 유치원, 초등학교, 중학교, 고등학교에 다니는 학생 수 산출
2. 자치구별 연령별 총 인구 수 산출
3. 연령별로 학교에 다니는 학생 수 / 전체 학생 수(교육열)를 구함
4. 교육열을 정렬한 뒤 막대그래프로 시각화
5. 연령별 총 학생 수와 학교에 다니는 학생 수 간의 상관관계를 나타냄

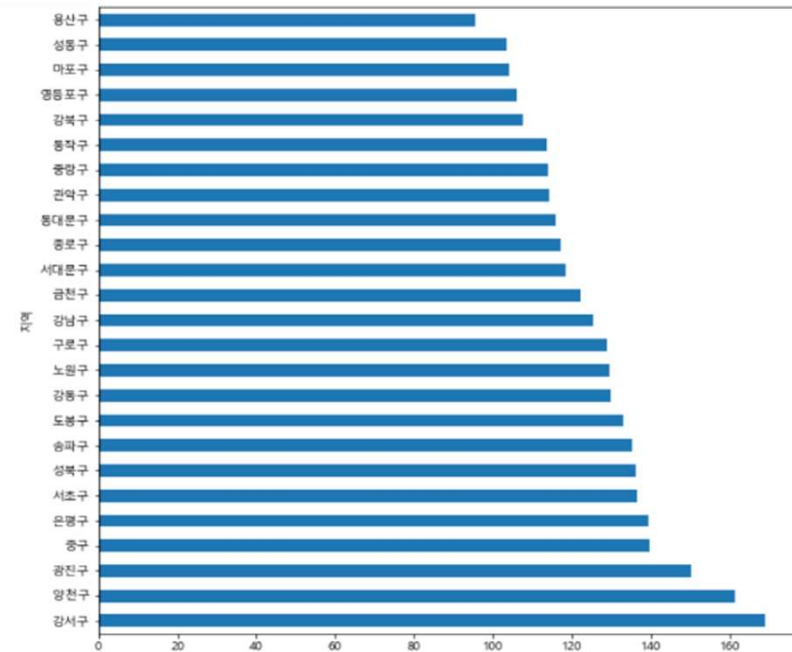
분석 과정

주제 1 - 자치구별 학교 다니는 학생수 비율, 자치구별 학원 다니는 학생수 비율

2) 유·초·중·고별 교육열 시각화



유치원교육열

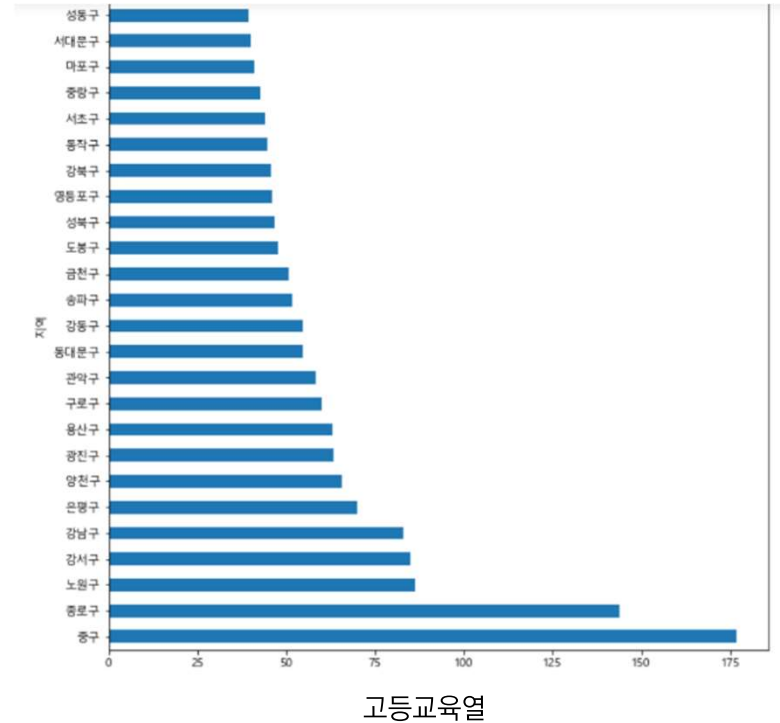
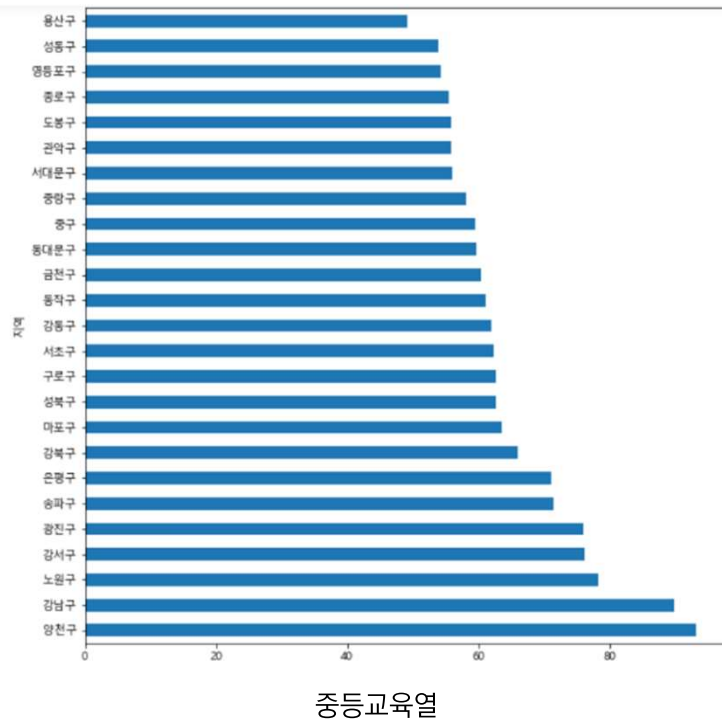


초등교육열

분석 과정

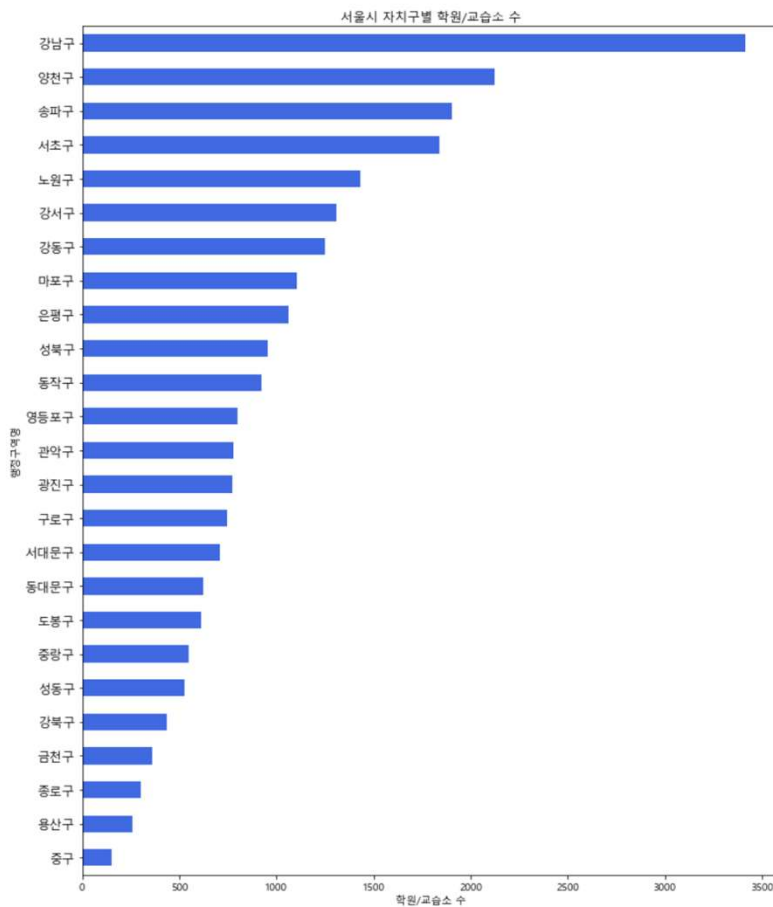
주제 1 - 자치구별 학교 다니는 학생수 비율, 자치구별 학원 다니는 학생수 비율

2) 유·초·중·고별 교육열 시각화



분석 과정

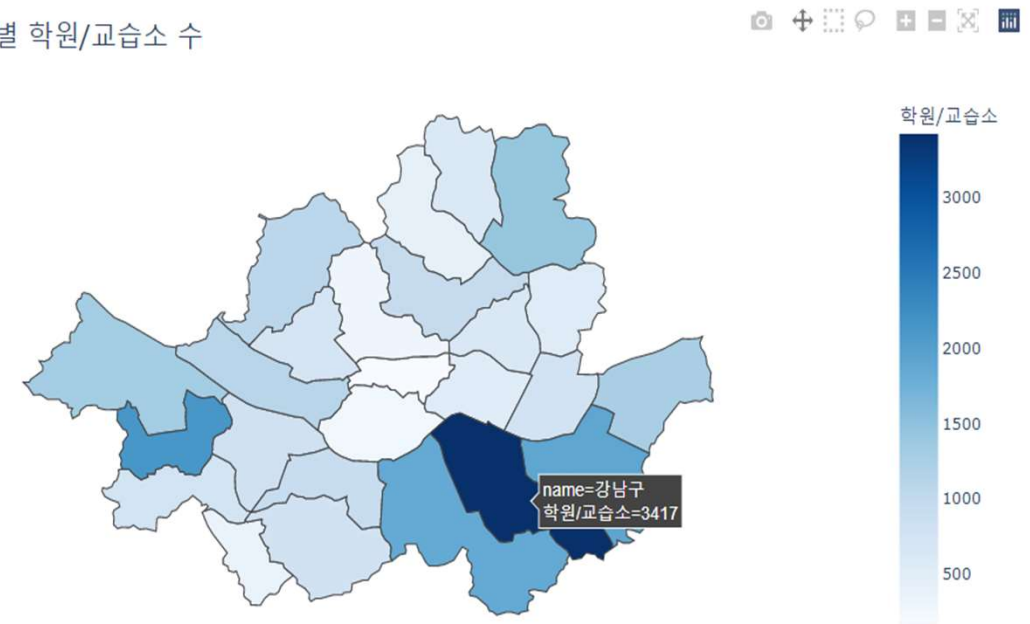
주제 2 – 강남 8학군은 실제로 학구열이 높을까?



1) 자치구별 학원/교습소 수 비교

1. 자치구별 학원/교습소 수 합계 산출

서울시 자치구별 학원/교습소 수



분석 과정

주제 2 – 강남 8학군은 실제로 학구열이 높을까?

2) 학군별 학원/교습소 수 & 학교 수 비교

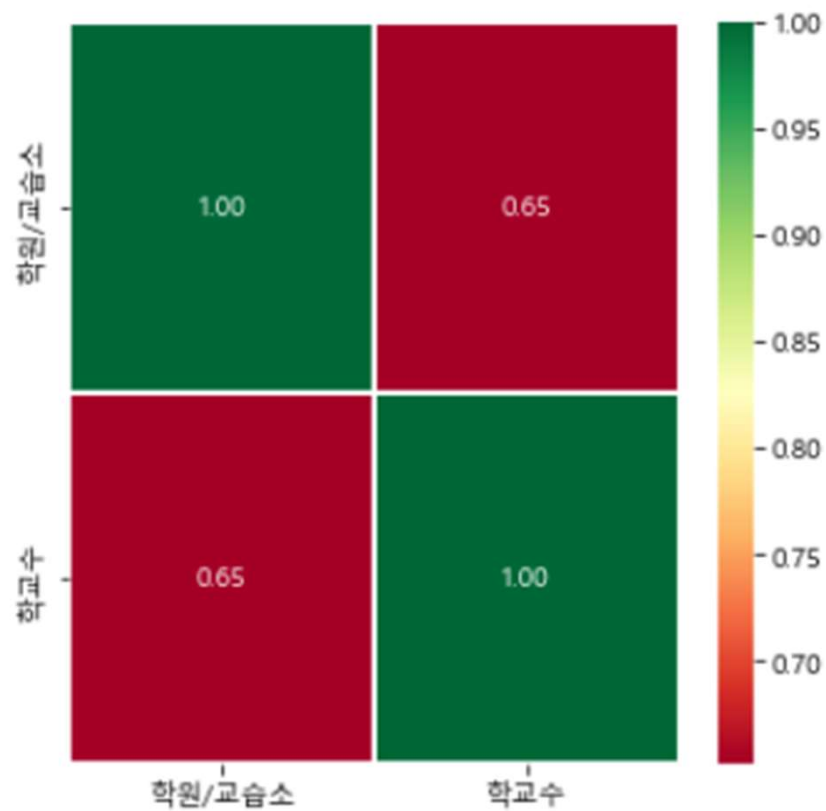
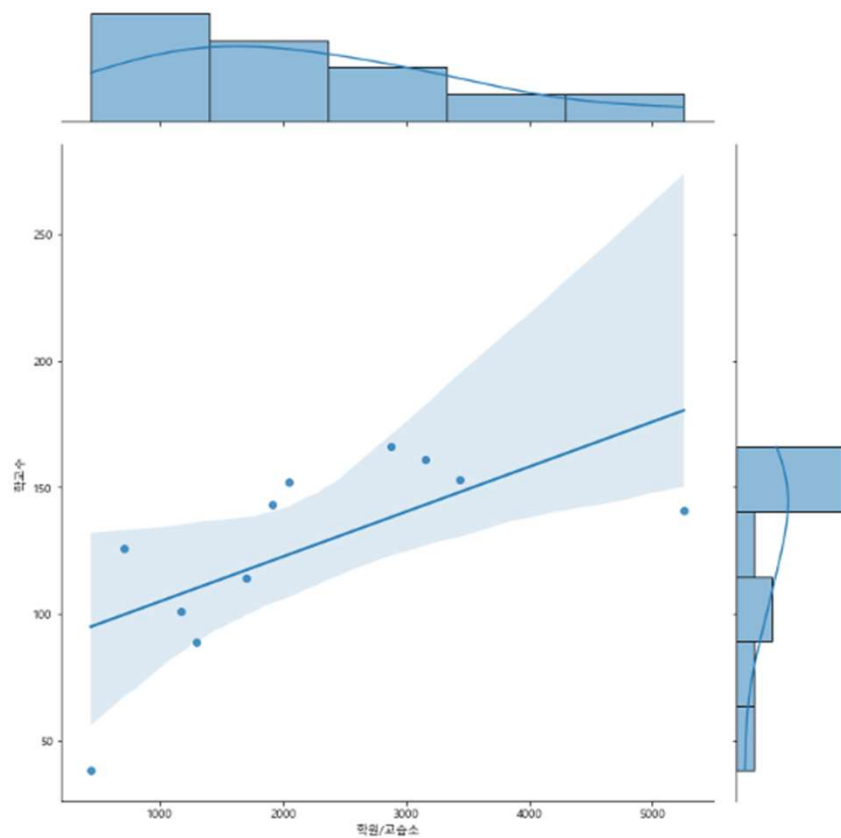
1. 자치구별 해당하는 학군 열 추가 (강남구, 서초구 > 8학군)
2. 학군별 학원 수 합계 산출

서울시 학군별 학원/교습소 & 학교 수



분석 과정

주제 2 – 강남 8학군은 실제로 학구열이 높을까?



분석 과정

주제 2 – 강남 8학군은 실제로 학구열이 높을까?

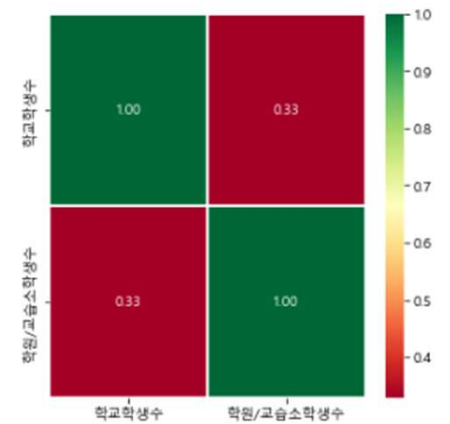
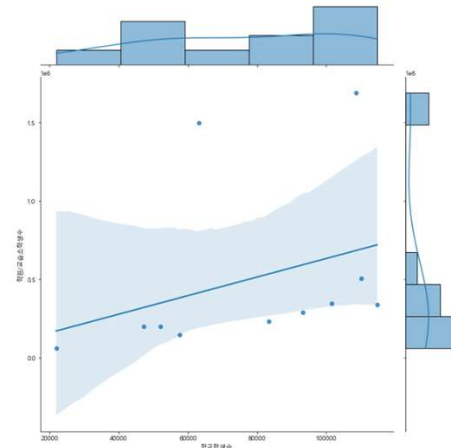
학군별 학원/교습소 학생 수 / 학교 학생 수 비율



3) 학군별 학원/교습소 학생 수 & 학교 학생 수 비교

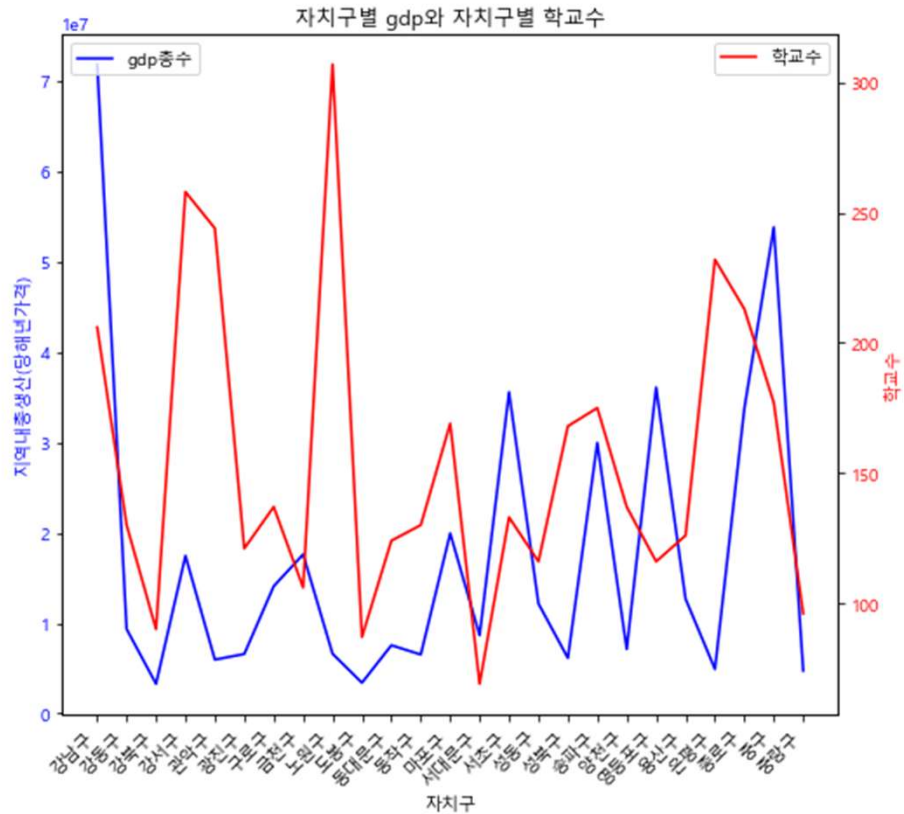
1. 온라인, 원격 키워드 포함한 학원명 행 제거
(실제 등원하는 학생 수 데이터가 필요)
2. 학군별 학원/교습소 학생 수 합계 산출

서울시 학군별 학원/교습소 & 학교 학생 수



분석 과정

주제 3 - grdp가 높은 자치구 학구열 높을까?

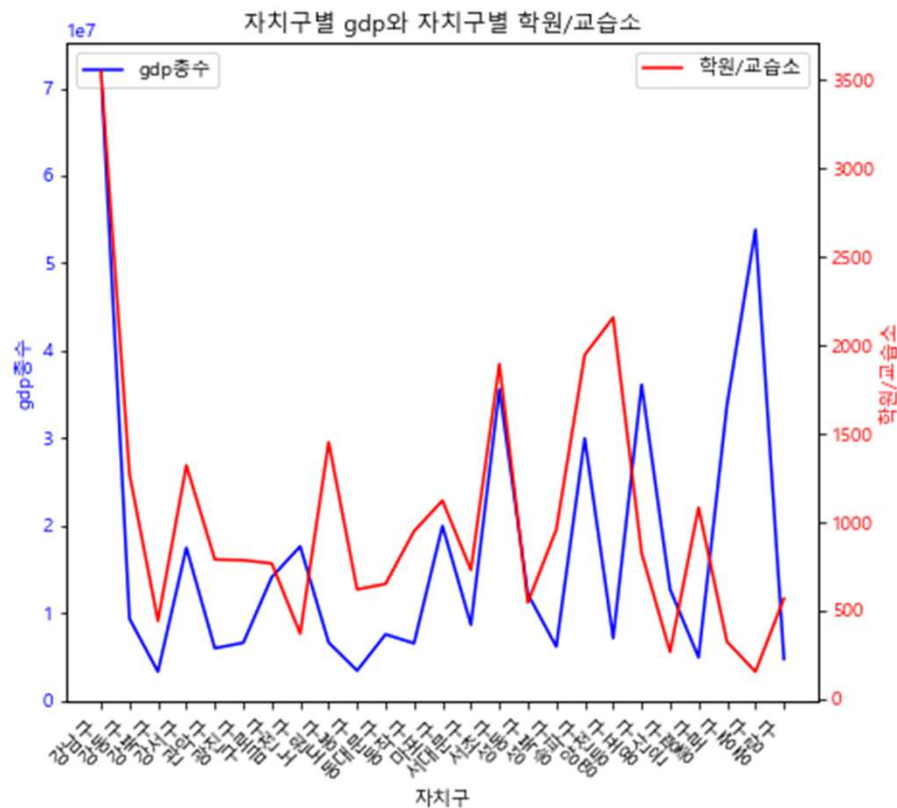


1) 자치구별 grdp 총 수와 자치구별 학교수의 상관관계

1. grdp 데이터에서 자치구별로 총가구수의 개수(sum) 구해서 새로운 데이터프레임 생성
2. school 데이터에서 자치구별 학교 개수 구해서 새로운 데이터 프레임 생성
3. 둘의 데이터를 자치구 기준으로 concat 해준 뒤, 새로운 데이터 프레임 생성

분석 과정

주제 3 - grdp가 높은 자치구 학구열 높을까?



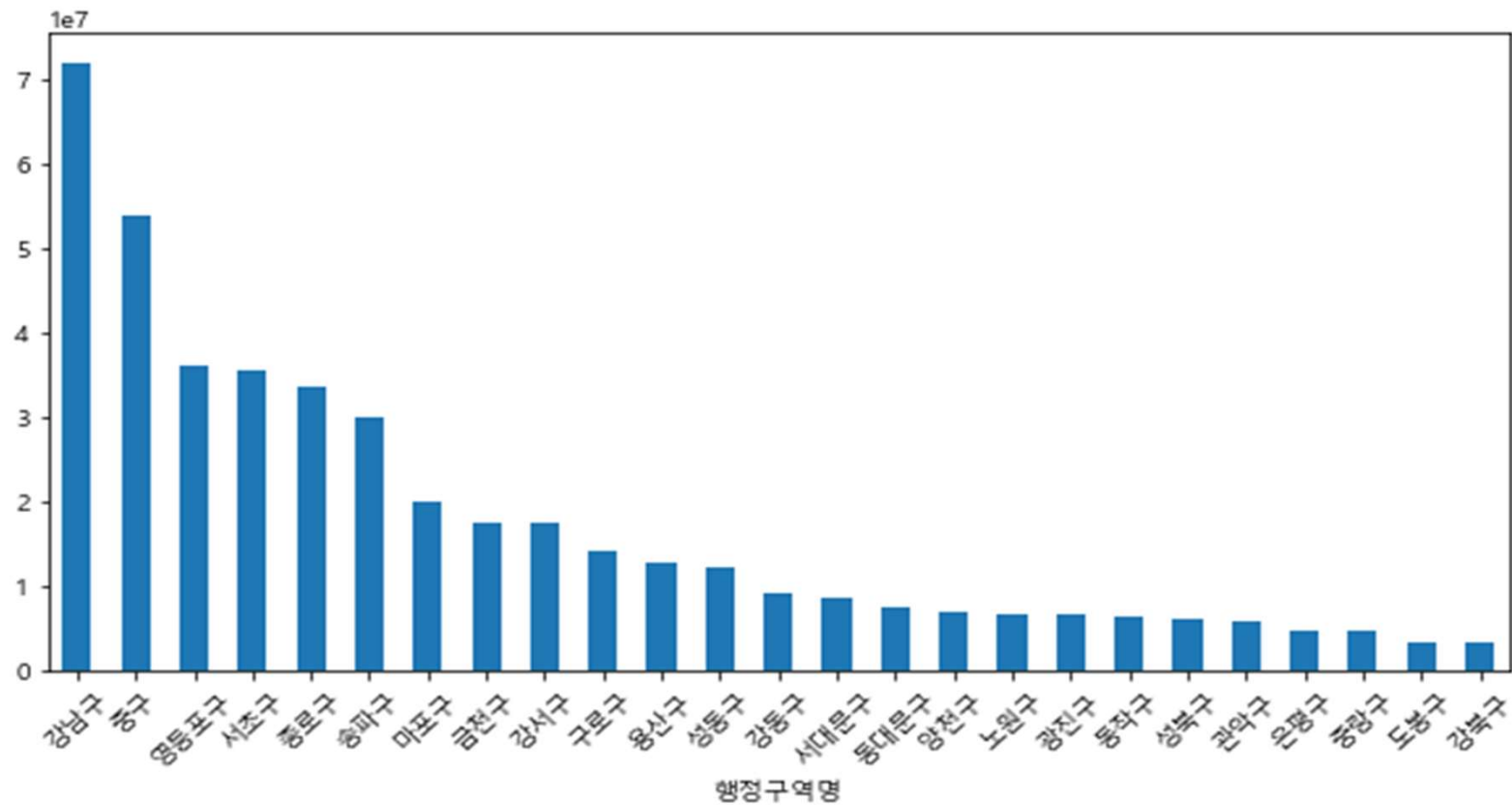
2) 자치구별 grdp 총수와 자치구별 학원 및 교습소의 상관관계

1. grdp 데이터에서 자치구별로 총가구수의 개수(sum) 구해서 새로운 데이터프레임 생성
2. academy 데이터에서 자치구별 학원/교습소 개수 구해서 새로운 데이터프레임 생성
3. 둘의 데이터를 자치구 기준으로 concat 해준 뒤, 새로운 데이터프레임 생성

분석 과정

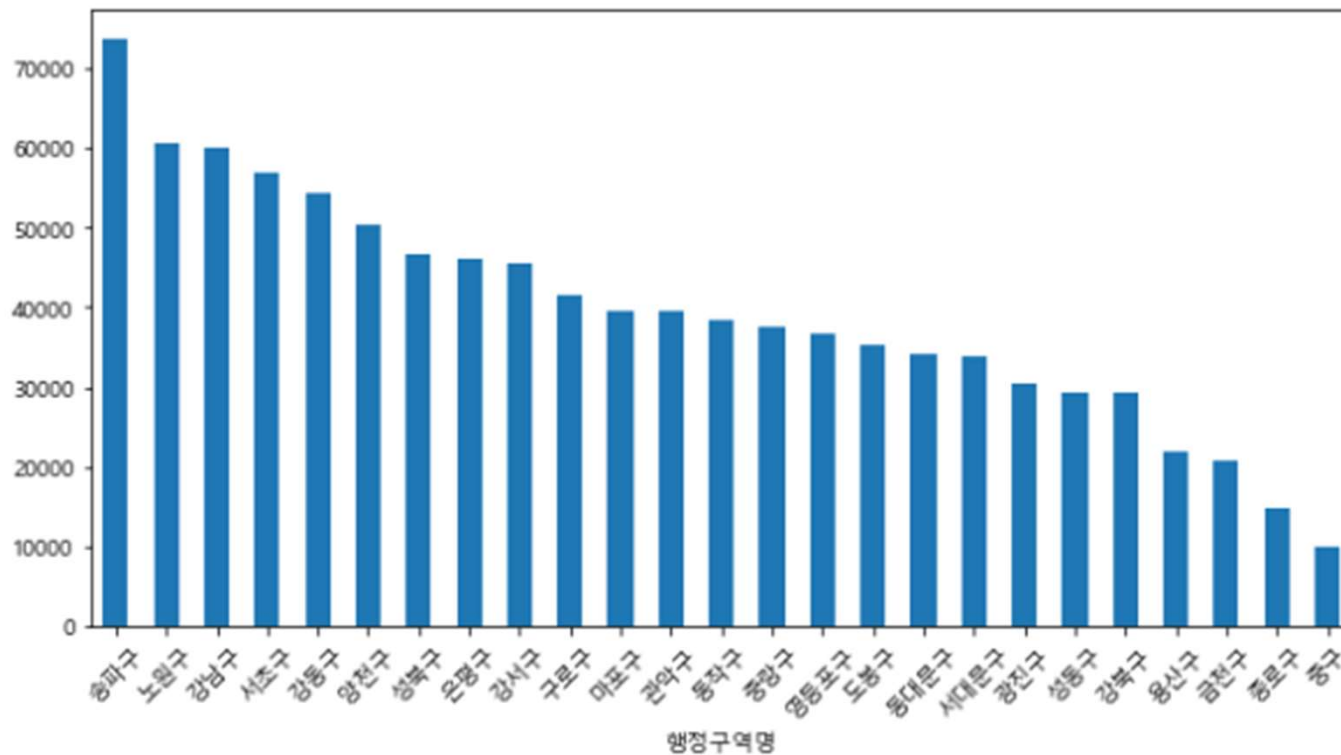
주제 4- 서울 외곽 지역에 학원 수가 많은 이유는?

1) 자치구별 grdp 지수 분포



분석 과정

주제 4- 서울 외곽 지역에 학원 수가 많은 이유는?

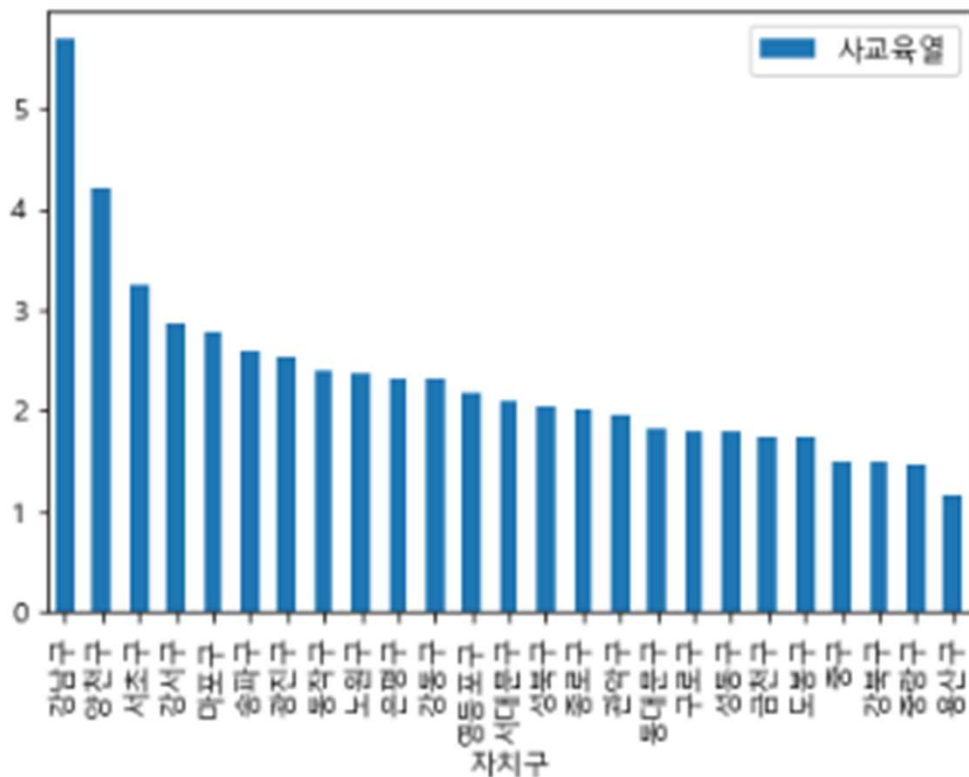


2) 자치구별 학령인구(5세~19세) 분포도

1. 주민등록연앙인구 데이터에서 학령인구의 열만 추출
2. 자치구별 학령인구 총 계 산출

분석 과정

주제 5 - 서울특별시의 교육사각지대와 서울런 사업



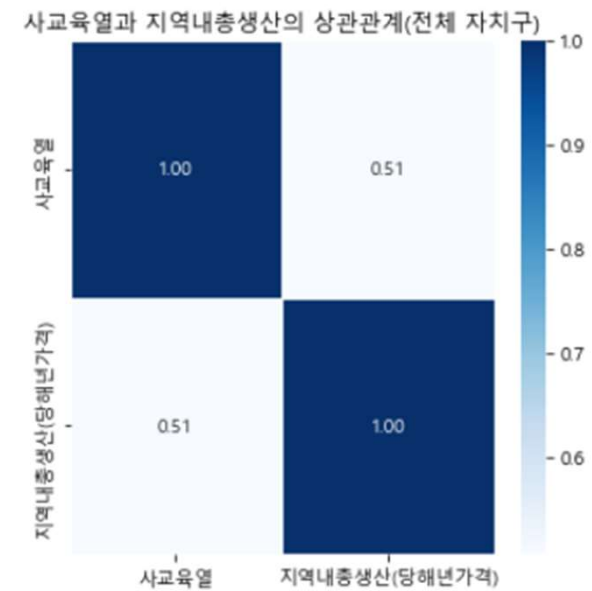
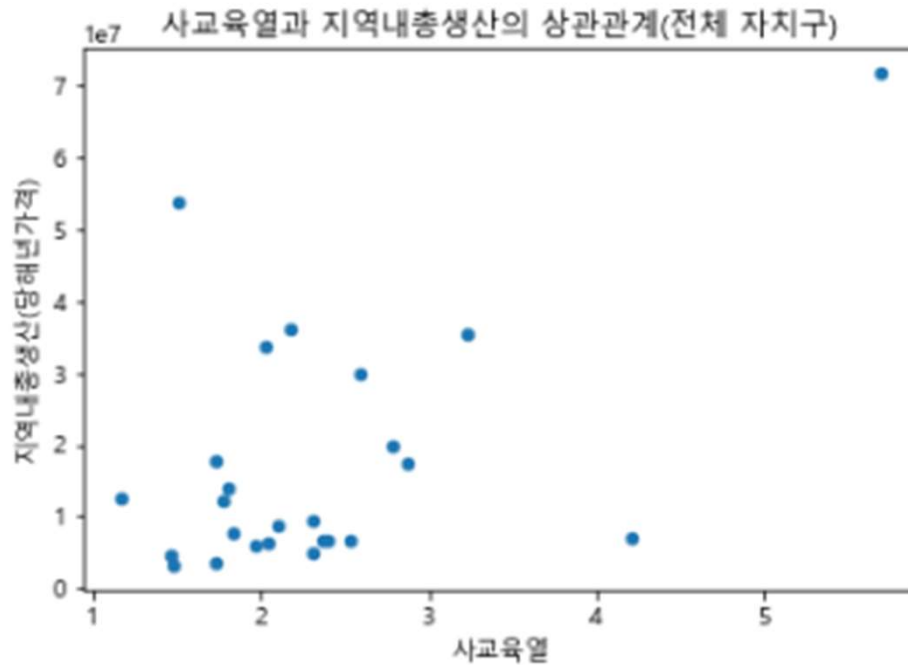
1) 사교육열 시각화

1. 서울특별시 학원 교습소 정보data에서 자치구별로 학원의 개수를 구한다
2. 주민등록연령인구(연령별동별)통계data에서 5~9세, 10~14세, 15~19세 column을 합쳐서 학령인구 column 생성
3. 새롭게 만든 학령인구 column을 자치구별로 나눠서 구한다.
4. 자치구별로 학원의 개수/학령인구*100을 계산하고 이를 사교육열로 정의한다.

분석 과정

주제 5 - 서울특별시의 교육사각지대와 서울런 사업

2) 사교육열과 지역내총생산의 상관관계(전체 자치구)

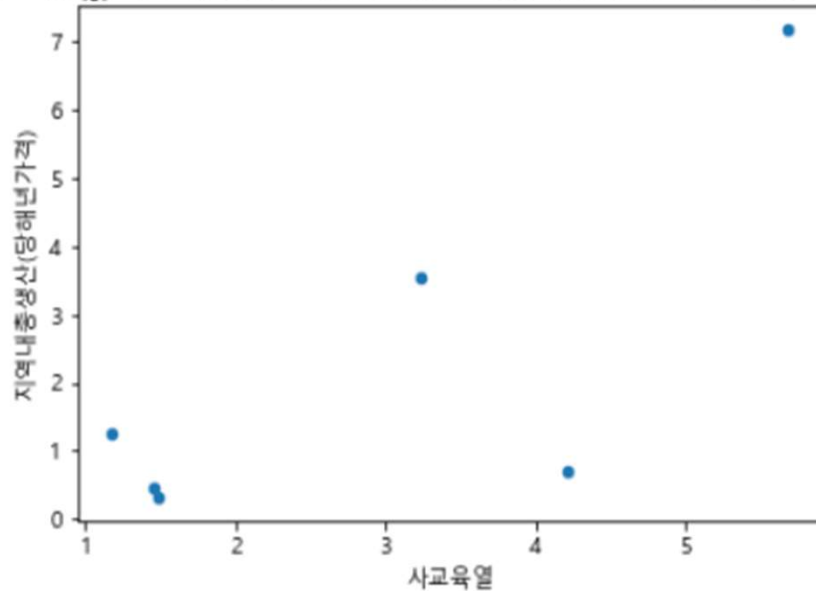


분석 과정

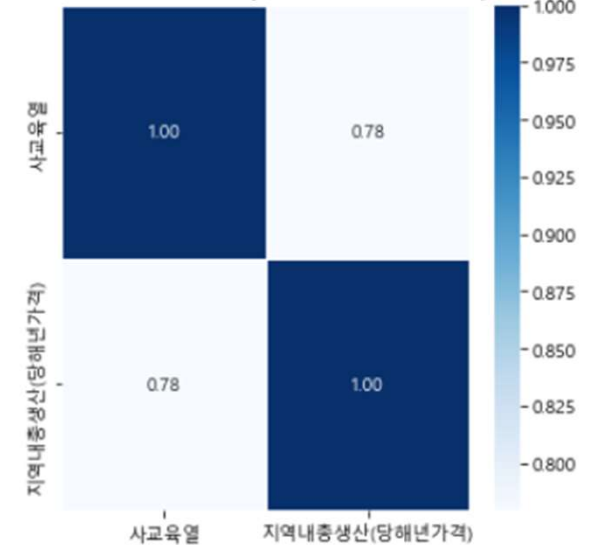
주제 5 - 서울특별시의 교육사각지대와 서울런 사업

3) 사교육열과 지역내총생산의 상관관계 (grdp 상위, 하위 10% 자치구)

사교육열과 지역내총생산의 상관관계(지역내총생산 상위10%, 하위10% 자치구)



#사교육열과 지역내총생산의 상관관계(지역내총생산 상위10%, 하위10% 자치구)



결론

주제 1 - 자치구별 학교 다니는 학생수 비율, 자치구별 학원 다니는 학생수 비율

1) 유치원교육열

-높음: 노원구, 종로구, 성북구 / -낮음: 서초구, 강북구, 성동구

2) 초등교육열

-높음: 강서구, 양천구, 광진구 / -낮음: 용산구, 성동구, 마포구

3) 중등교육열

-높음: 양천구, 강남구, 노원구 / -낮음: 용산구, 성동구, 영등포구

4) 고등교육열

-높음: 중구, 종로구, 노원구 / -낮음: 성동구, 서대문구, 마포구

→ 교육열이 낮은 자치구는 전 연령에서 동일하게 낮은 교육열을 보이지만
교육열이 높은 자치구는 연령에 따라서 바뀌는 것을 알 수 있음.

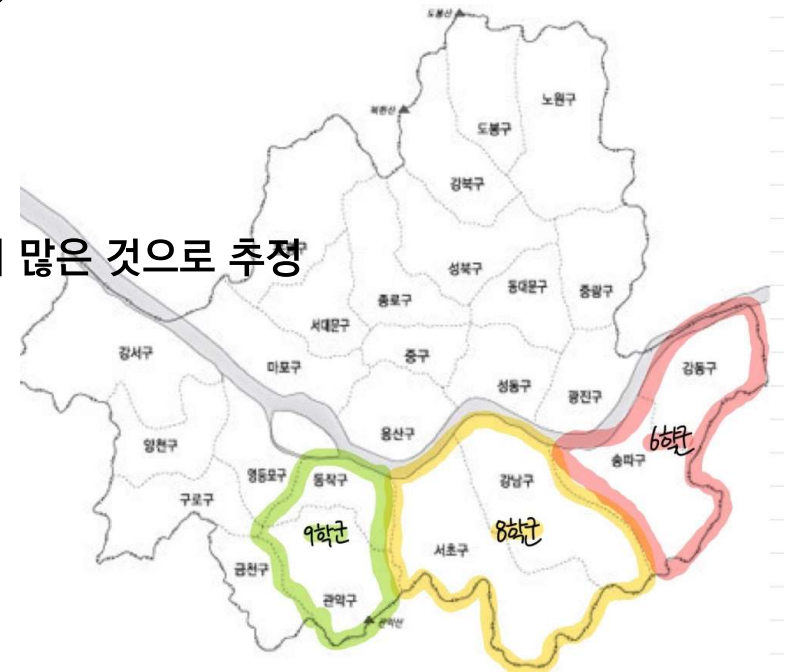
결론

주제 2 – 강남 8학군은 실제로 학구열이 높을까?

1) 자치구별 학원/교습소 수 비교

-학원/교습소 수 상위 5개 자치구 : 강남구, 양천구, 송파구, 서초구, 노원구, 강서구
→ 학원/교습소는 **6,8,9학군에 집중되어 있음**

-학원/교습소 수는 서울 외곽에 다수 위치
→ 인천과 붙어있는 강서구, 양천구에 인천 소재 학생 대상 학원이 많은 것으로 추정



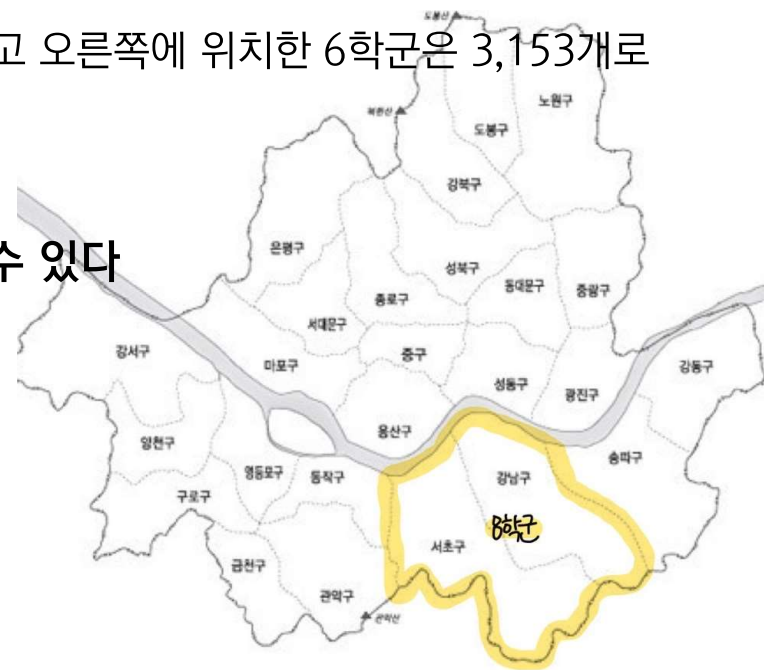
결론

주제 2 – 강남 8학군은 실제로 학구열이 높을까?

2) 학군별 학원/교습소 수 & 학교 수 비교

- 7~9학군을 제외한 나머지 학군의 학원/교습소 수와 학교 수 비율 비슷함
- 지리적으로 8학군 왼쪽에 위치한 9학군의 학원/교습소 수는 1,701개이고 오른쪽에 위치한 6학군은 3,153개로 8학군 5,258개와 큰 차이를 보이고 있다.

→ 다른 학군 자치구에서 8학군 소재 학원에 다니고 있음을 추측할 수 있다

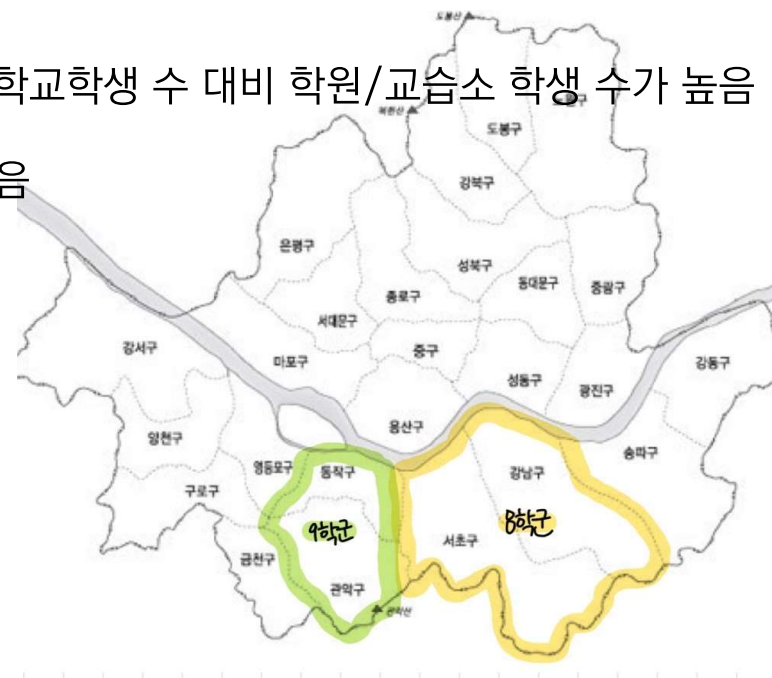


결론

주제 2 – 강남 8학군은 실제로 학구열이 높을까?

3) 학군별 학원/교습소 학생 수 & 학교 학생 수 비교

- 모든 학군에서 학교 학생 수보다 **학원/교습소 학생 수가 더 많음**
- 특히 **8, 9학군의 학원 학생 수가** 다른 학군에 비해 압도적으로 **많음**. 또한 학교학생 수 대비 학원/교습소 학생 수가 높음
- 타 학군 소재 학교 학생들이 8, 9 학군 소재 학원에 등원함을 유추할 수 있음



결론

주제 3 - grdp가 높은 자치구 학구열 높을까?

- 자치구별 grdp가 높은 곳일수록 학교, 학원 및 교습소 수도 많음
→ 강남구/양천구/서초구
- grdp가 높을수록 학구열이 높은 경향이 있음
- 학구열이 높은 7,8학군의 강남구, 서초구, 양천구는 입시, 검정 및 보습학원의 비율 ↑
중랑구는 예체능 학원의 비율 ↑

주제 4 - 서울 외곽 지역에 학원 수가 많은 이유는?

- 학령인구 수 상위 10개의 자치구 중 9군데가 외곽에 위치해 있으므로 외곽지역에 학령인구가 많음
- 외곽 지역의 grdp가 다른 지역의 grdp보다 상대적으로 높음

→ 즉, 소득수준이 높으면 학원에 많이 다니므로 소득이 높은 외곽지역에 학원 수가 많음

결론

주제 5 - 서울특별시의 교육사각지대와 서울런 사업

-서울런이란?

: 학습자원 접근이 어려운 청소년들을 대상으로 하는 교육 서비스 사업.

초중고 학생이 가장 선호하는 온라인 교육 콘텐츠를 종합반 수준으로 무제한 수강할 수 있도록 제공하고,

가입자 전원에게 멘토링 서비스를 제공하며 효율적인 학습 관리를 지원하는 서비스



결론

주제 5 - 서울특별시의 교육사각지대와 서울런 사업

- 자치구 전체를 살펴봤을 때 사교육열과 지역내총생산(grdp)의 상관관계가 약 50% 정도로 낮은 상관관계를 보임
 - 서울 지역내총생산 상위 10% 자치구와 하위 10% 자치구만 따로 뽑아서 상관관계를 구했더니 78%라는 값 도출
- 지역내총생산이 평균인 지역들은 지역내총생산과 사교육열 간의 상관관계가 크지 않았지만
지역내총생산이 매우 높은 지역일수록 사교육열이 더욱 높았으며,
지역내총생산이 매우 낮은 지역일수록 사교육열이 줄어드는 것을 확인함

결론

주제 5 - 서울특별시의 교육사각지대와 서울런 사업

지역내총생산이 가장 적은 강북구, 중랑구, 용산구는 교육사각지대로 볼 수 있으며 정부의 교육 지원이 필요하다고 판단
서울런의 가입자수는 계속 늘어날 전망이며 이를 통해 교육격차 해소를 기대해 볼 수 있다

[서울런 가입자 추이]

