



ESC 20-1 Final Project

Week1

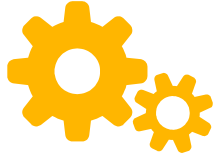
Team 2

김윤환 조인식 신혜연 백채빈 이상완 손지우

“ Our Goal

서울시 가구원의 행복도
예측 모델 만들기

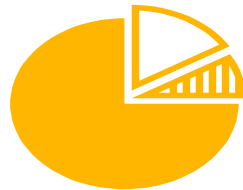
Contents



**Pre-processing
Visualization**

**Feature
Extraction**

Seoul Happiness 2014



S
Survey
2014
U
L

2014 서울서베이 (도시정책지표 조사표 - 가구원용)

조사원이 작성하는 한입니다.	일련번호 □□□□□□□□	행정구역코드 □□□□□□□□	가구원 코드 □□□□□□□□	출생년도 □□□□□□□□	성 명 □□□□□□□□	
조사주관 서울특별시 행정국 정보공개정책과	조사기관 서울특별시정보관리데이터서비스팀 (02)3488-2715 (서울특별시청 02)8188-0011	제 20111 호				

■ 서울 서베이 교환 연서도

문1. 귀하는 서울에 거주하면서 서울이 2014년 느끼는지요?

- ① 서울에서 태어나서 서울이 코랄같이 느껴진다
- ② 서울에서 태어났으나 서울이 코랄 같지 않다
- ③ 서울에서 태어나서 서울이 코랄하지만 잘 모르니 서울이 코랄 같지
않게 된다
- ④ 서울에서 태어나서 코랄하고 서울이 코랄하지도 코랄하지도 않음

■ 서울축제에 대한 연서도 및 참여자 만족도

문2. 귀하는 서울시 및 구청에서 개최하고 있는 2014년 축제에
얼마나 관심이 있습니까?

- (예 : H1 Seoul 포스트타워, 한강시민공원, 불꽃축제 등)
- ① 매우 관심있다 ② 약간 관심있다 ③ 보통이다
 - ④ 약간 관심없다 ⑤ 전혀 관심없다

문2-1. 귀하는 서울에서 개최된 축제에 참여한 경험이 있으신지요?

- ① 있다 - (문2-2번)
- ② 없다 - (문3번)

문2-2. 귀하는 축제에 참여하신 후 어느 정도 만족하셨는지요?

- ① 매우 만족 ② 약간 만족 ③ 보통이다
- ④ 약간 불만족 ⑤ 매우 불만족

■ 사회사건

문3. 귀하는 지난 1년 동안 다음과 모양 또는 만화책 등에
참여한 적이 있습니까? 참여 경험에 있는 것을 모두
표기해 주십시오.

- ① 전통문화/민속제 ② 동물보호/동물복지
- ③ 지역발전/행정개혁 ④ 입주민/주민자치
- ⑤ 동호회 ⑥ 자원봉사단체
- ⑦ 시민단체 ⑧ 노조 및 직능 단체
- ⑧ 정당 ⑨ 종교단체
- ⑩ 기타(구체적으로 :)
- ⑪ 어느 모양이나 단체에도 참여된 적이 없다

■ 정책지수

문4. 귀하는 다음과 같은 소스로 행복이라고 생각하십니까? 가장 행복한 소스를
10점으로 가장 불행한 소스를 0점으로 하여 각 항목별 자신의
행복지수를 표시해 주십시오.

가장 불행한 소스	10점	9점	8점	7점	6점	5점	4점	3점	2점	1점	가장 행복한 소스
1) 자신의 건강상태.....	10	9	8	7	6	5	4	3	2	1	
2) 자신의 재정상태.....	10	9	8	7	6	5	4	3	2	1	
3) 우리 정치, 현 정부의 관계.....	10	9	8	7	6	5	4	3	2	1	
4) 가정생활.....	10	9	8	7	6	5	4	3	2	1	
5) 사회생활(학교, 종교, 사회, 레크리 등).....	10	9	8	7	6	5	4	3	2	1	

문4-1. 귀하는 지금 얼마나 행복하십니까? 가장 행복한 소스를 100
점으로, 가장 불행한 소스를 0점으로 하여, 귀하는 자신의 행복
지수를 어떻게 측정해 주십니까? 100점 중 1 점

■ 도시위험도

문5. 귀하는 서울에 거주하면서 다음과 같은 항목에 대해 어느 정도
위험하다고 생각하십니까? 각 항목에 대해 위험도를 평가해 주십시오.

매우 위험하다	약간 위험하다	보통 이다	별로 위험하지 않다	전혀 위험하지 않다
5	4	3	2	1
1) 화재나 홍수, 산사태 등의 재해로 인한 피해				
2) 밤늦게 일어나는 경우				
3) 강도, 소매치기, 성추행 등 다양한 범죄 피해				
4) 각종 불법행위(차량, 주차장 등)에 의한 피해				

■ 현대사회의 위험 요인에 따른 피해 정도

문6. 아래의 표에 나열되어 있는 위험들이 발생한다면, 이로 인해 예상
되는 피해가 얼마나 크다고 생각하십니까? 다음 중 귀하에서 생각
하시는 위험도에 대해 정도를 표시하는 것에 동의해 주십시오.

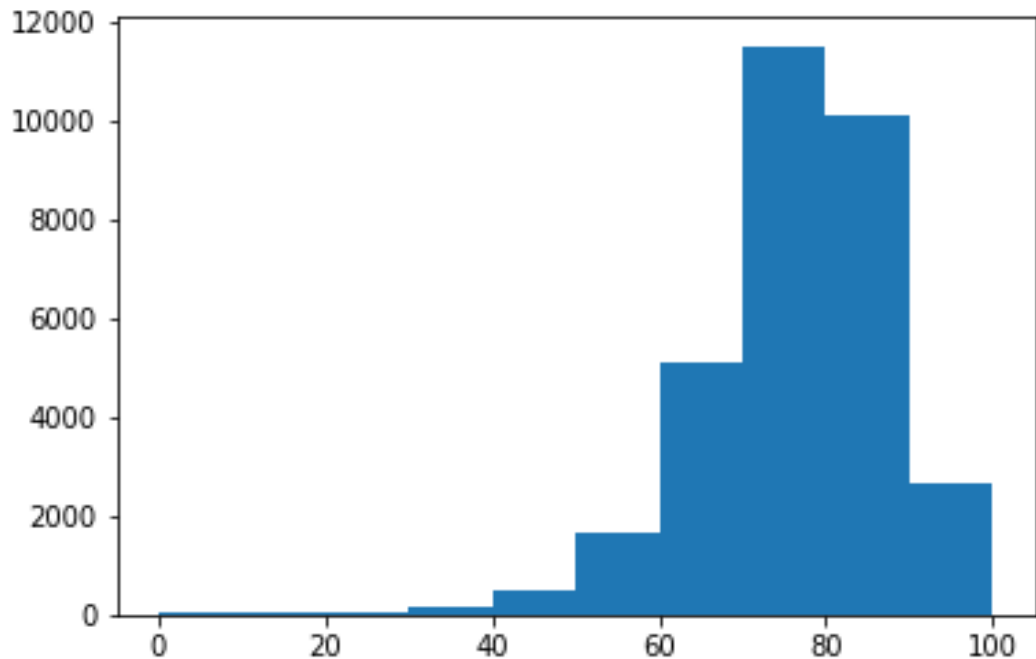
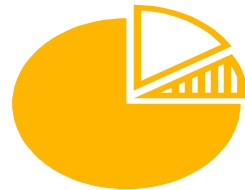
매우 크다	상당 크다	보통 이다	별로 크지 않다	전혀 크지 않다
5	4	3	2	1
1) 자연재해(태풍, 지진, 홍수 등)				
2) 교통사고				
3) 살인				
4) 화재(가정, 화재, 화재, 화재)				
5) 전염병(인사, 전염, 전염, 전염)				
6) 부정부패				
7) 폭력 범죄(강도, 폭도, 폭도, 폭도)				
8) 사회 갈등				
9) 경제위기(금융위기 등)				
10) 환경				
11) 연립정권 확대				
12) 컴퓨터 바이러스, 사이버범죄로 인한 혼란				
13) 핵전쟁(전쟁, 전쟁, 전쟁, 전쟁)				
14) 생활비(전, 생활비, 생활비, 생활비)				

문7. 귀하는 10년 전과 비교할 때 서울 시민이 오늘날 경험하는 위험
의 정도가 어떻게 변했다고 생각하십니까?

- ① 위험이 매우 커졌다 ② 위험이 상당히 커졌다
- ③ 10년 전과 비슷하다 ④ 위험이 약간 줄었다
- ⑤ 위험이 많이 줄었다

5

Y



- 0-100 주관적 행복점수

- Mean 72.68

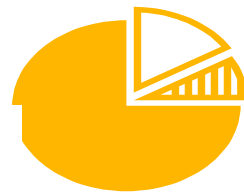
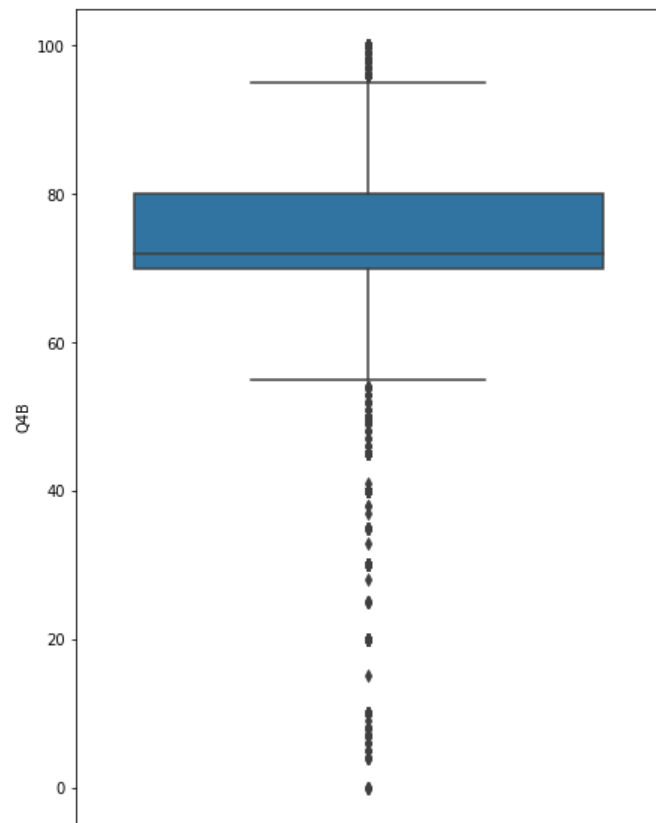
Y



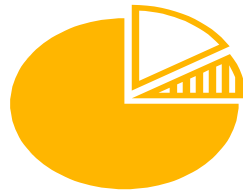
```
y.describe()
```



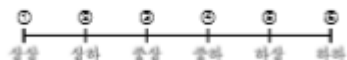
```
count    31848.000000  
mean      72.680482  
std       11.822053  
min        0.000000  
25%       70.000000  
50%       72.000000  
75%       80.000000  
max      100.000000  
Name: Q4B, dtype: float64
```



X variables



문27. 귀하의 정치경제사회적인 위치는 어느 계층에 속한다고 생각하십니까?



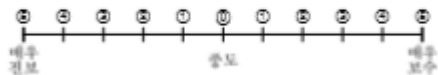
문11. 귀하가 거주하는 지역의 녹지(공원, 숲 등)에 대해 얼마나 만족하십니까?

- ⑤ 매우 만족 ④ 약간 만족 ③ 보통이다
② 약간 불만족 ① 매우 불만족



■ 개인성향

문41. 귀하는 어느 정도 보수적 또는 진보적이라고 생각하십니까?



■ 교통 이용 만족도

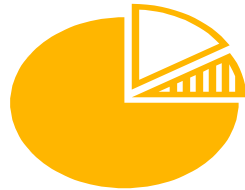
문45. 귀하는 서울의 교통수단에 대해 어느 정도 만족하십니까?
교통수단의 쾌적성, 정시성, 안전성 등을 전반적으로 생각하여
다음 각 항목에 대한 귀하의 생각을 말씀해 주십시오.

매우 만족	약간 만족	보통 이다	약간 불만족	매우 불만족	이용하지 않음
5-----4-----3-----2-----1					9

1) 버스(시내버스, 마을버스, 화석버스 등)	
2) 지하철	
3) 택시	



서울시민 고향 인식도
서울축제에 대한 인지도
사회자본
도시위험도
문화활동
환경
여가활동형태
운동
평생교육
자원봉사활동
계층인식
사회적 신뢰
교통이용만족도
은퇴시기
...



변수삭제

1)정보 없음

Ex)FC 가구원코드

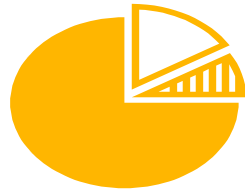
2)겹치는 정보

Ex) 출생년도-나이

지역분류 변수 다수 존재

3)행복지수와 관련 없는 변수

Ex) 여가활동의 종류/ 공공시설 불편이유/ 마을공동체관련 질문



변수변환

1)카테고리가 너무 많다

자가/전세/반전세/월세/상속 -> 자가/전세/기타

단체활동경험 -> 단체활동 개수

박물관,전시회 방문횟수/비용 -> 문화체험 연간 방문횟수

1순위 /2순위 -> 1순위만 선택

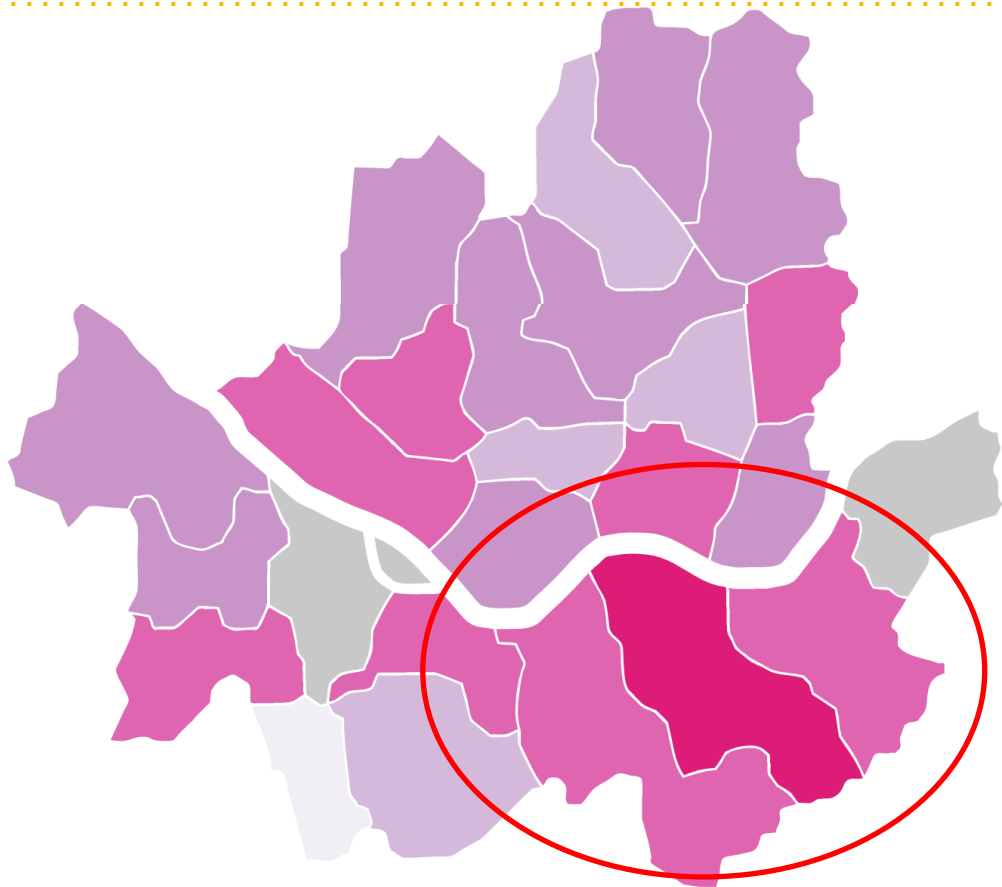
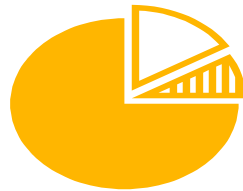
응답내 의미가 겹친다 -> 묶자

2) 이어지는 문항

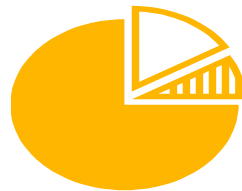
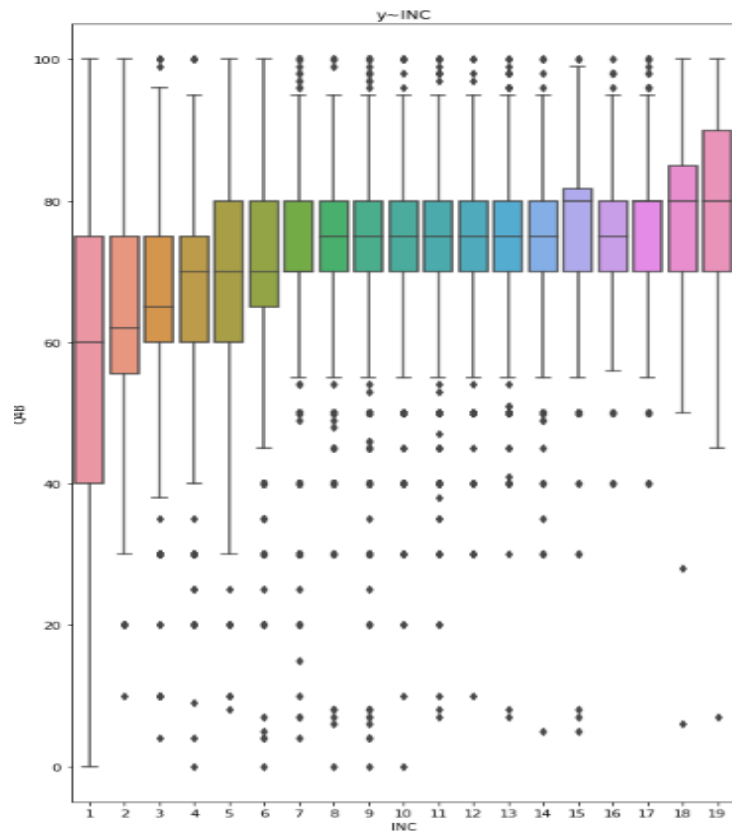
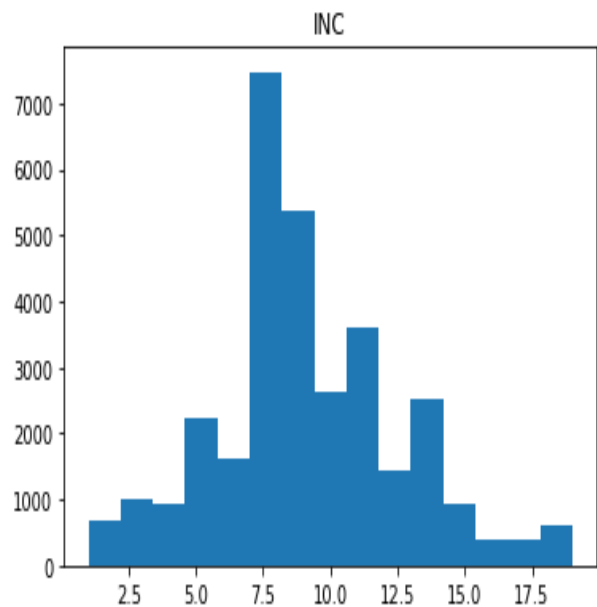
축제경험있다-> 3번문항으로 가시오

-> 관련변수 합쳐서 카테고리화

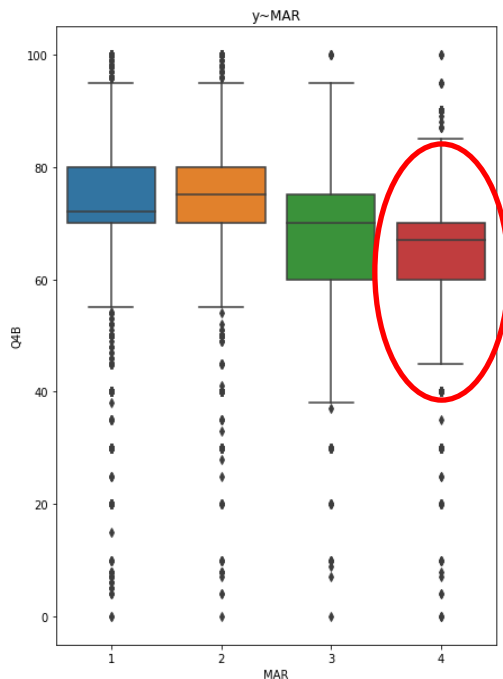
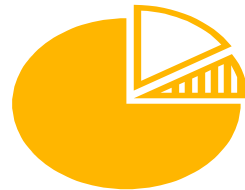
행정구역



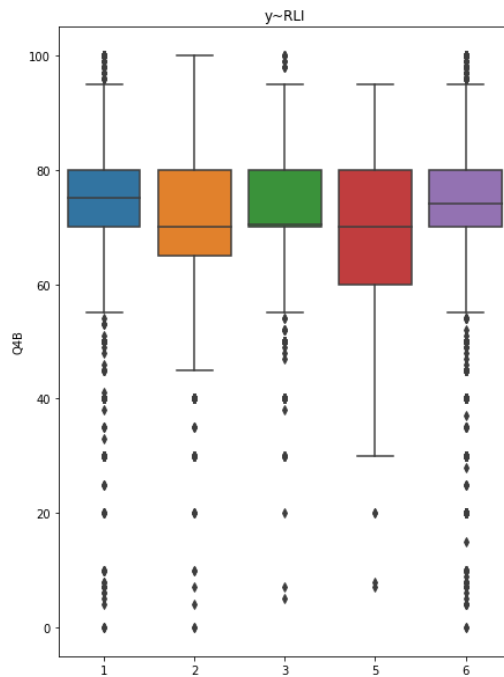
소득



혼인상태/ 종교

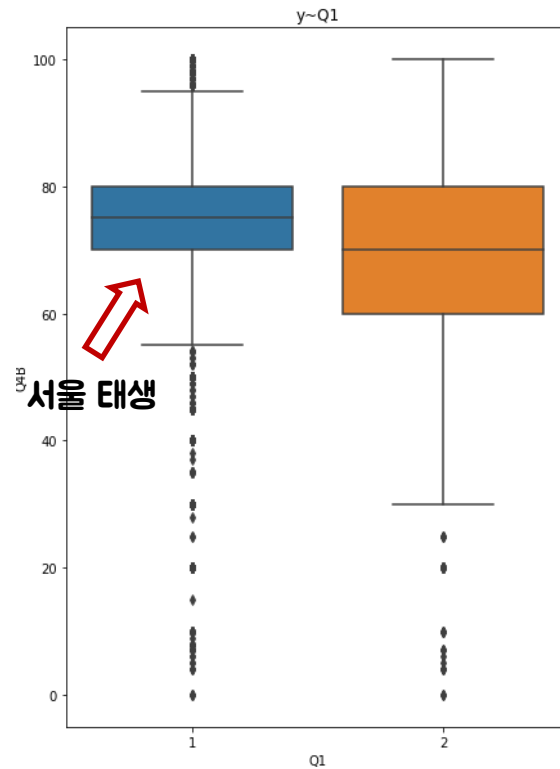
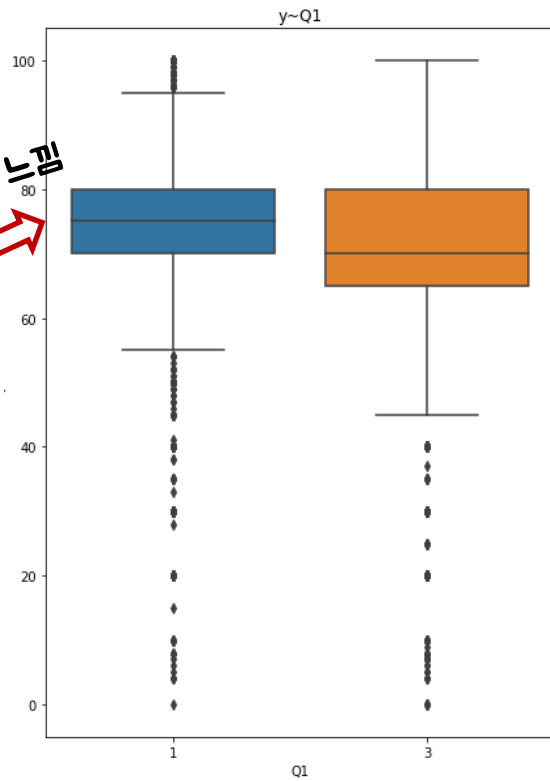


기혼 미혼 이혼 사별

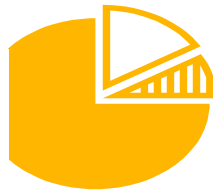


개신교 불교 천주교 기타 무교

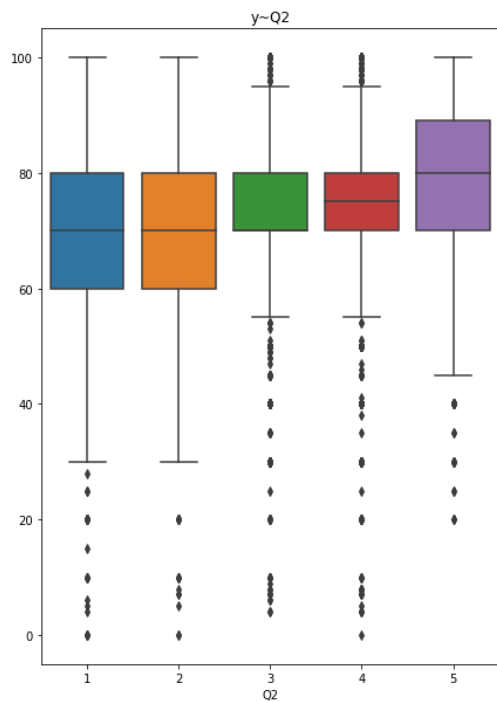
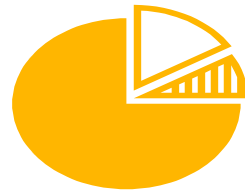
서울을 고향이라 느낌



서울 태생

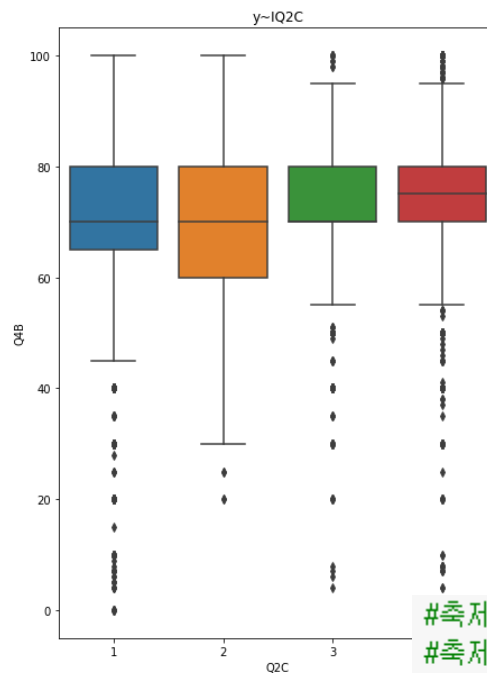


축제관련



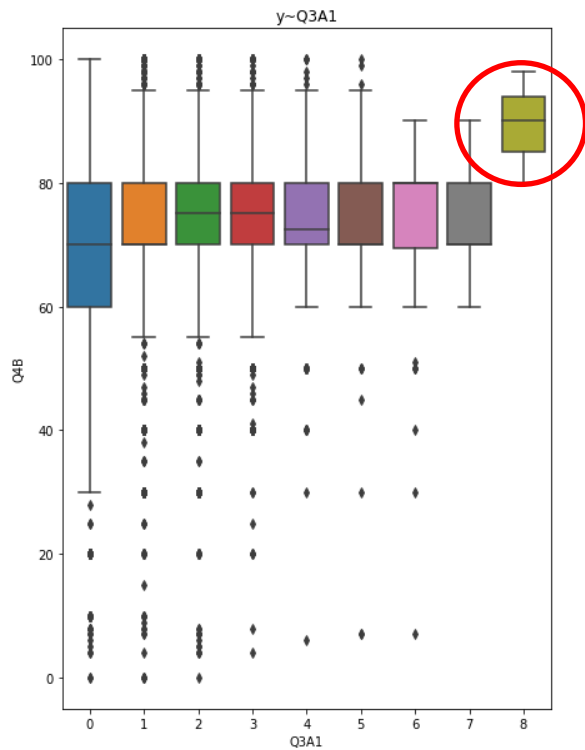
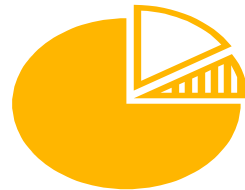
←————→

불만족 만족

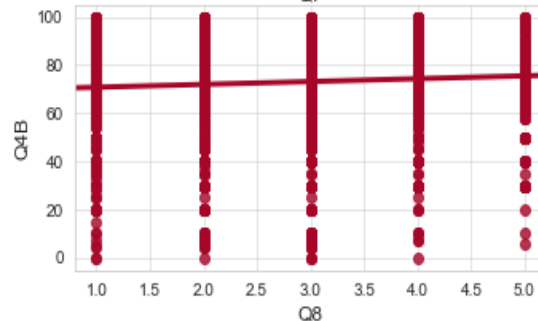
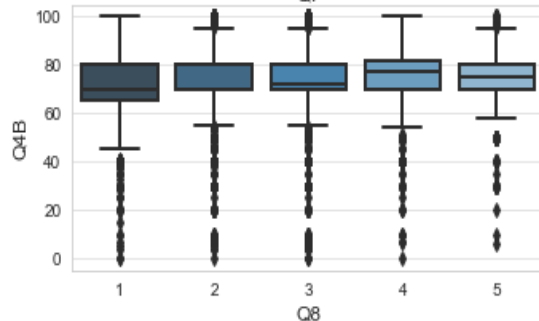
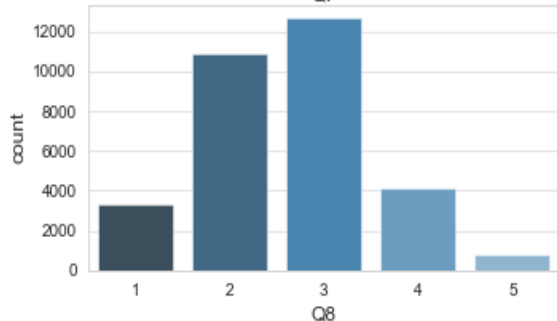
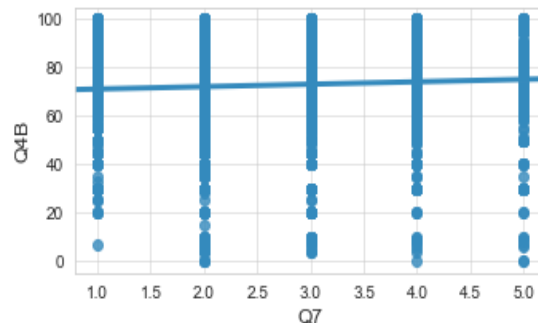
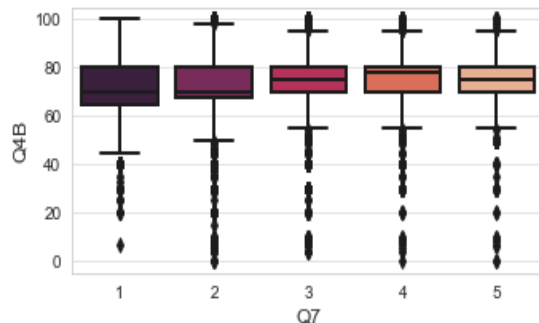
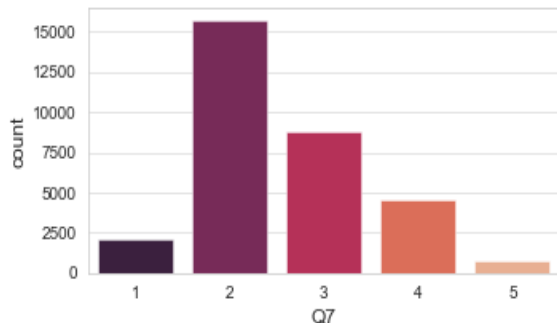
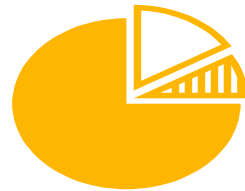


#축제경험있다(1)-만족(4)->4
 #축제경험있다(1)-보통(3)->3
 #축제경험있다(1)-불만족(2)->2
 #축제경험없다(2)->1

단체활동 개수

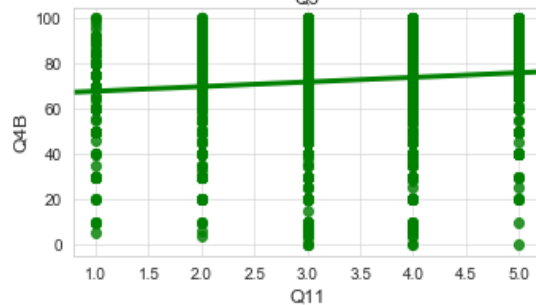
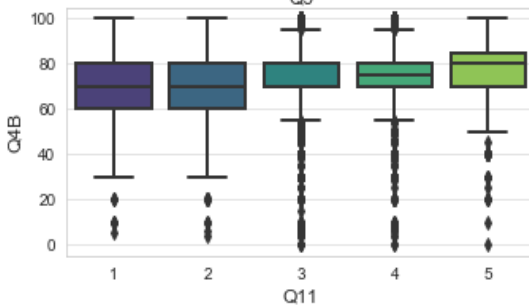
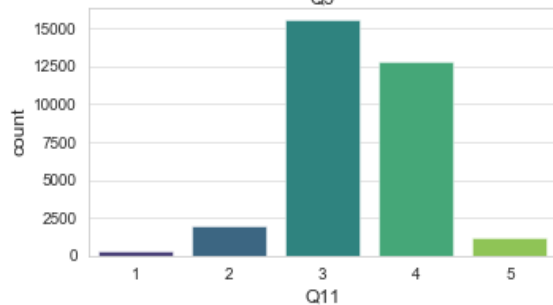
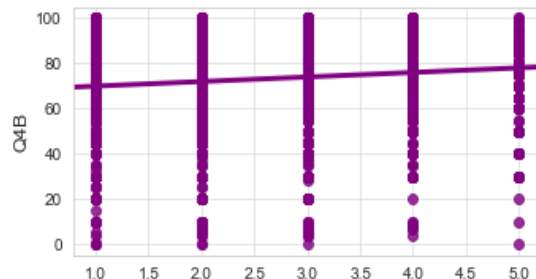
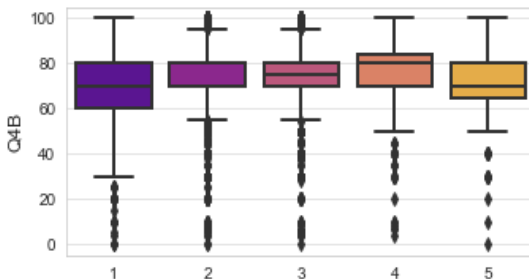
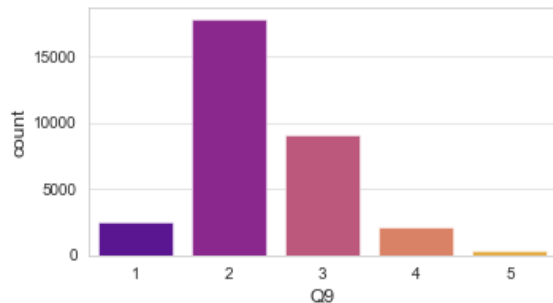
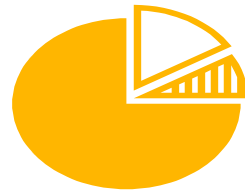


위험사회에 관한 인식



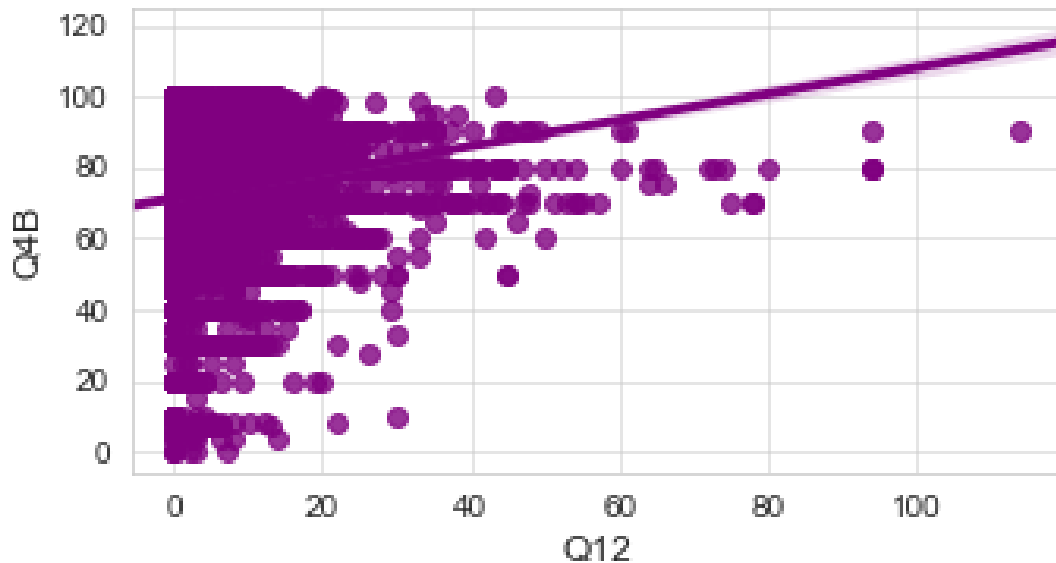
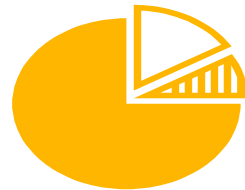
위험 증가

스트레스와 녹지만족도



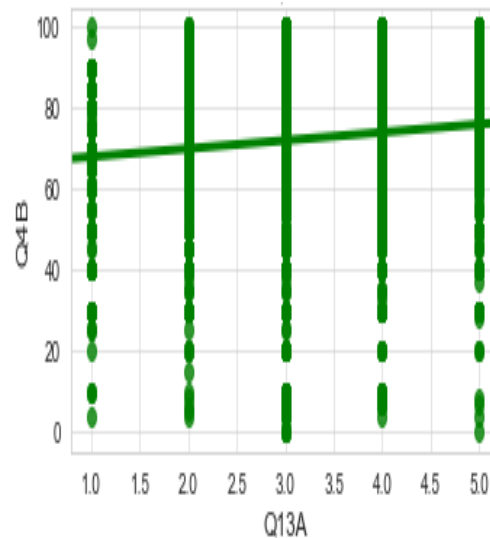
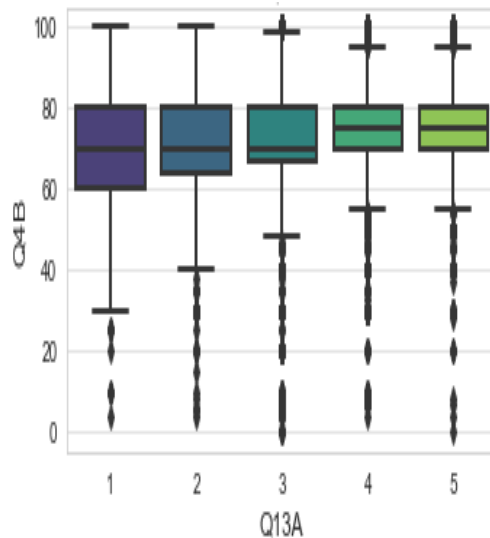
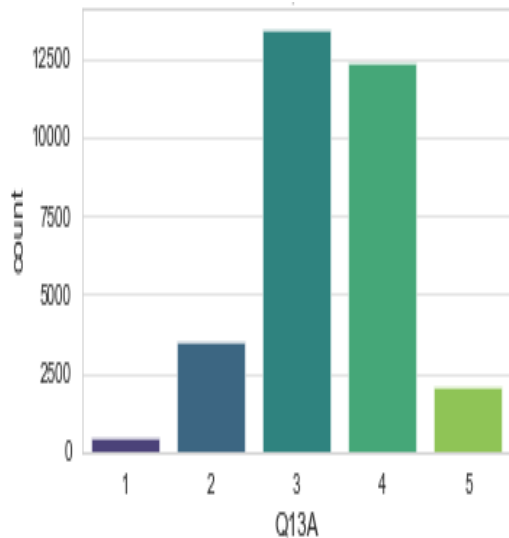
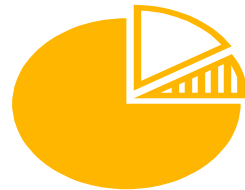
스트레스가 없고
녹지만족함

문화경험수



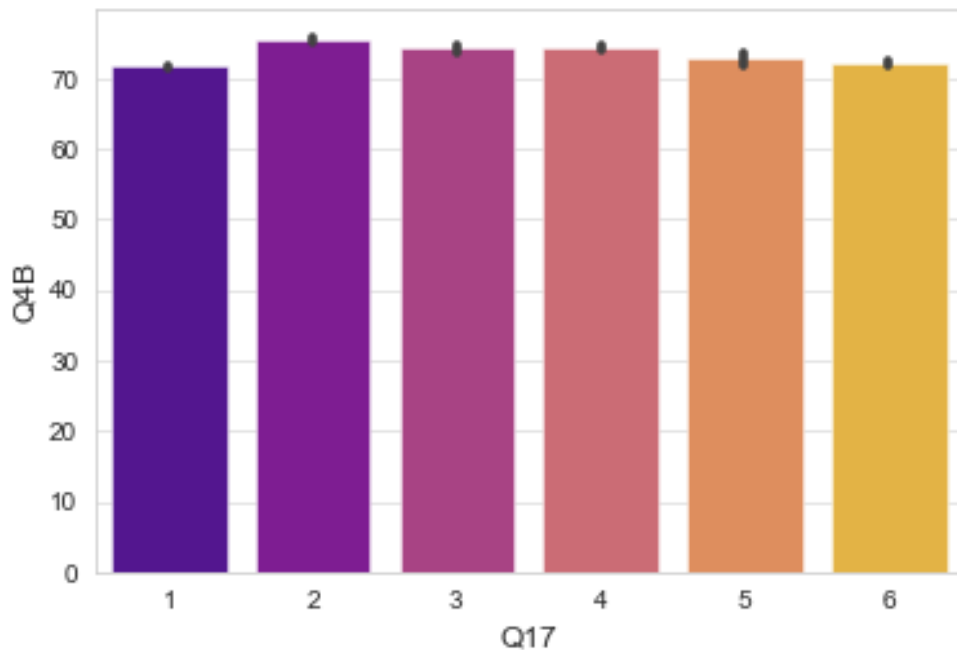
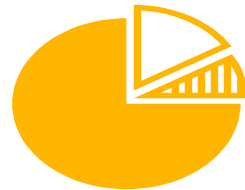
1년간 문화경험 총 횟수

문화환경 만족도



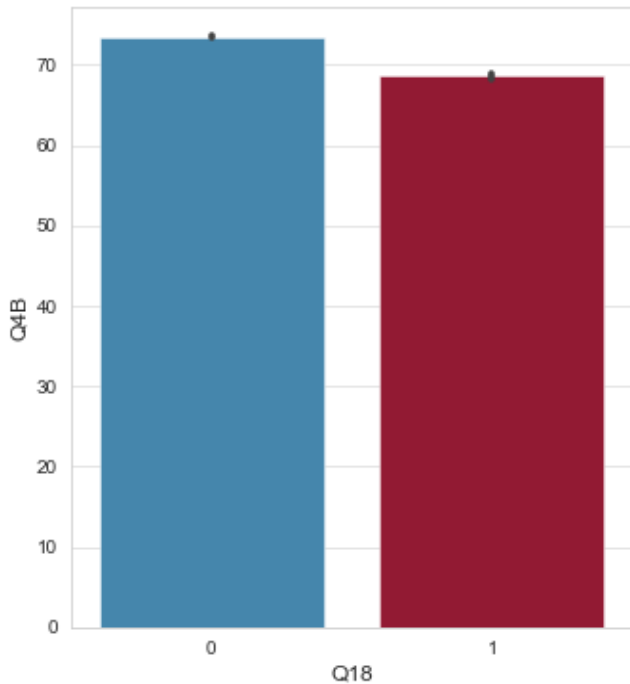
문화 시설 만족

취미

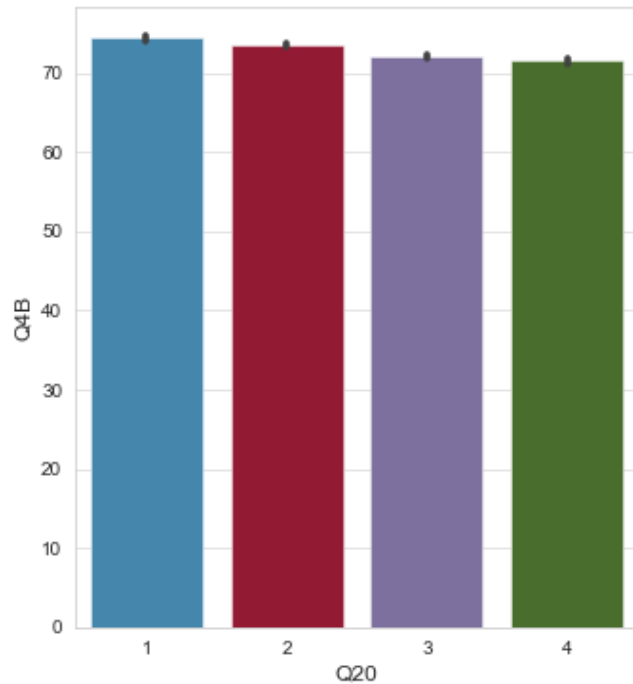


- 1: TV, 컴퓨터
- 2: 문화, 창작 활동
- 3: 운동
- 4: 여행
- 5: 사회봉사
- 6: 기타

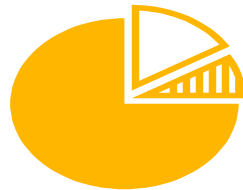
기타



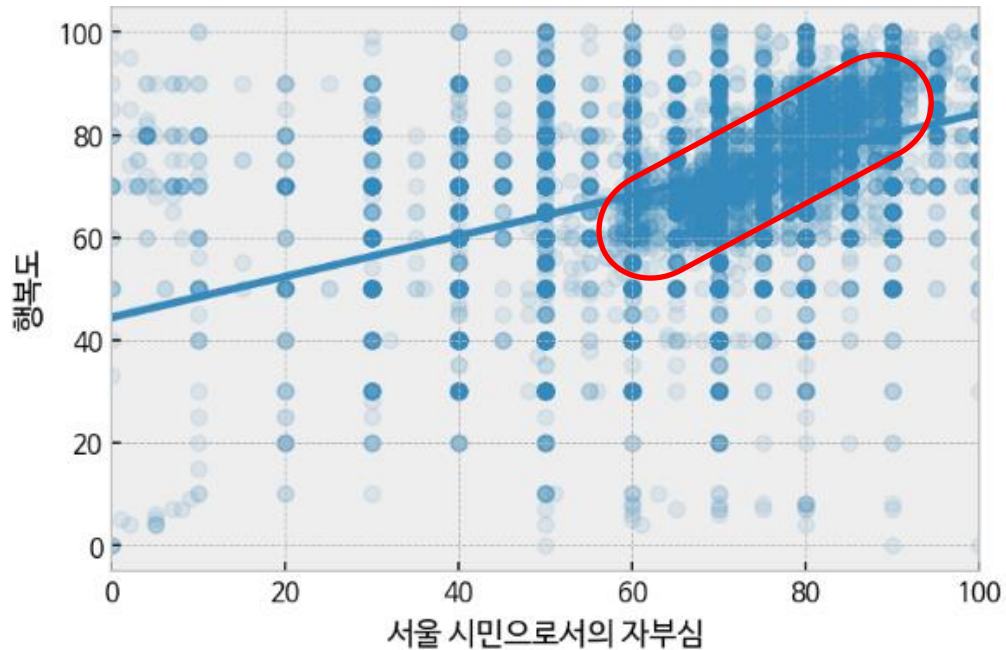
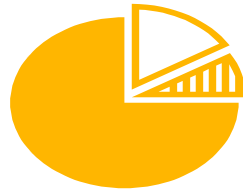
1 : 혼자 취미활동
0 : 혼자가 아님



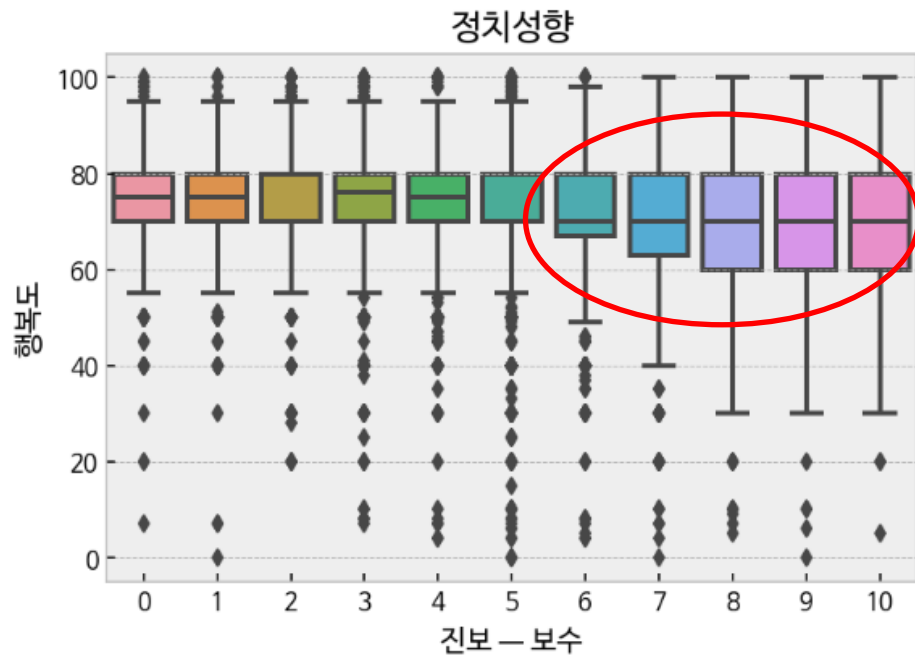
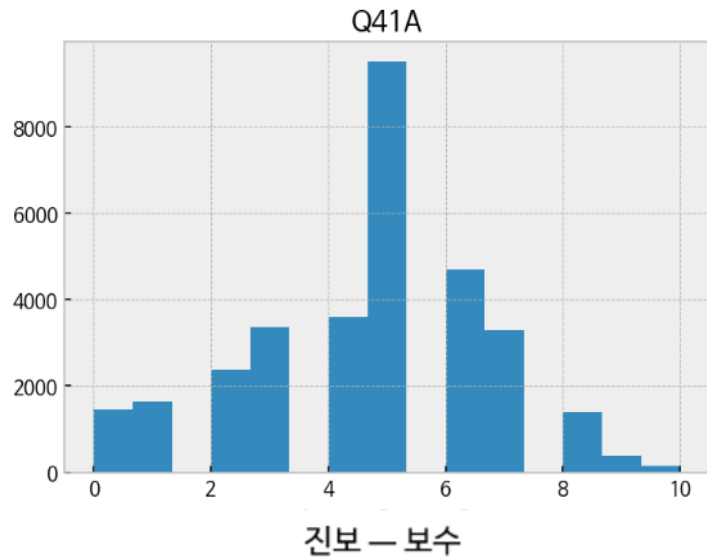
→ 4로 갈수록 운동량 적음



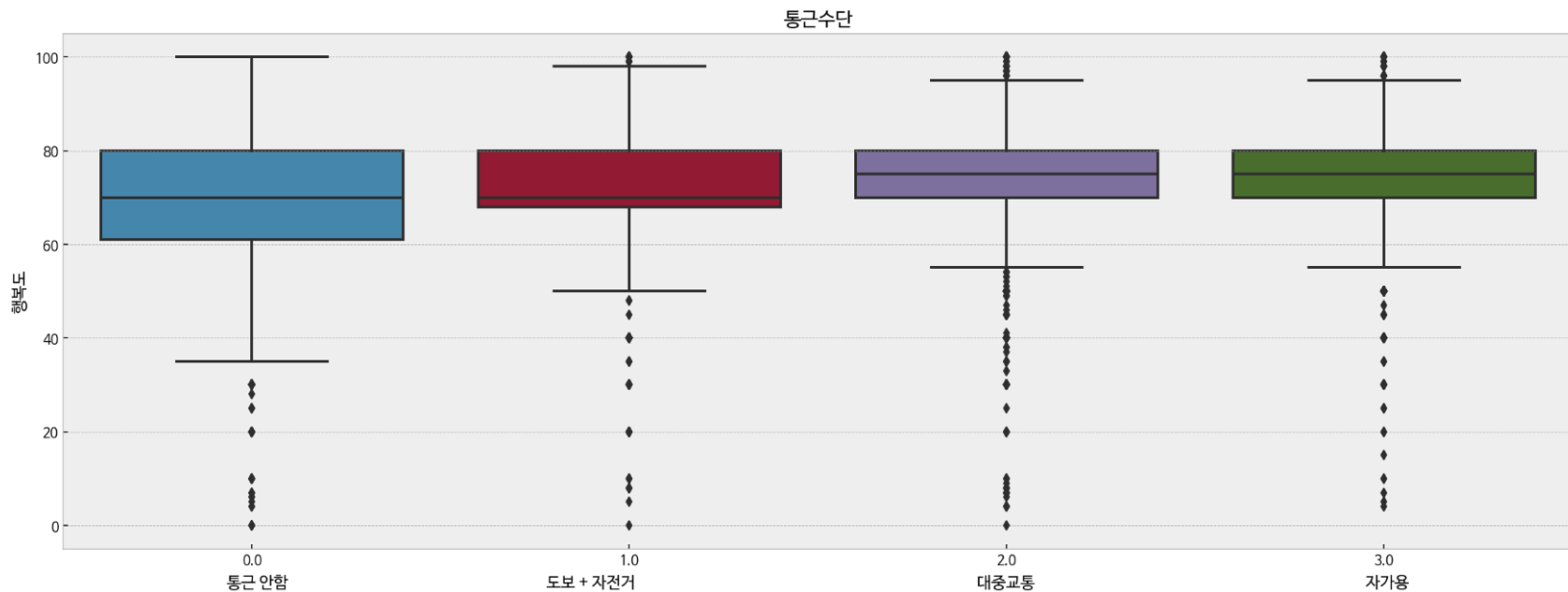
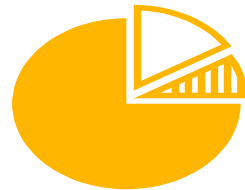
서울 시민 자부심



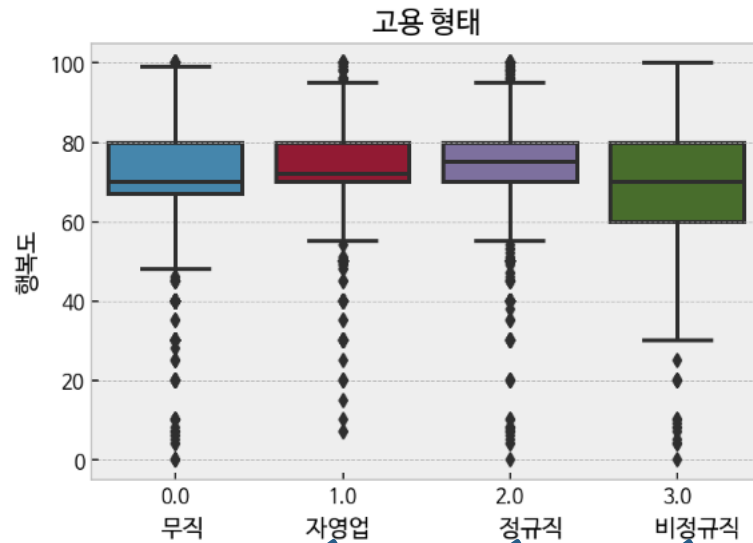
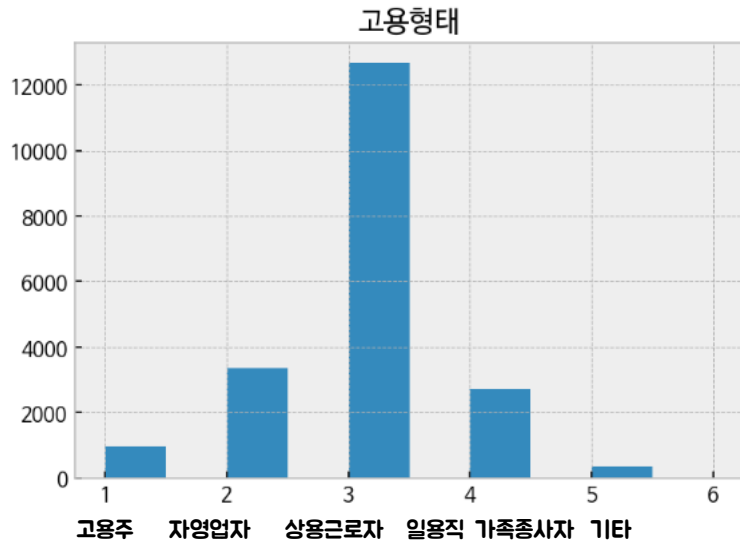
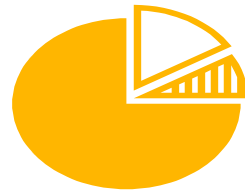
정치 성향



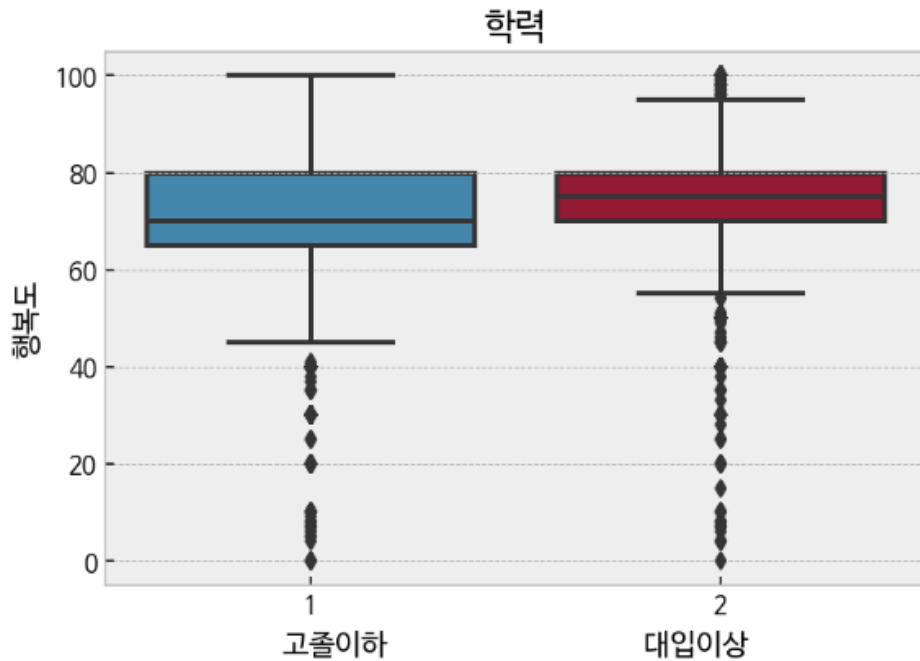
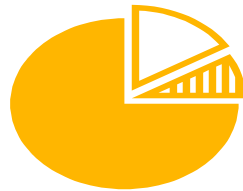
통근 수단



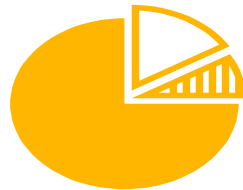
고용 형태



학력



자전거 이용



문44. 귀하는 현재 자전거를 이용 하고 계십니까?

- ① 이동(출근, 장보기, 업무 등) 수단으로 이용 → [문44-1로]
- ② 운동(레저 등) 수단으로 이용 → [문44-1로]
- ③ 이용하지 않음 → [문45로]

문44-1. 귀하는 자전거 이용 환경에 대해 전반적으로 얼마나 만족하십니까?

- ⑤ 매우 만족 ④ 약간 만족 ③ 보통이다
- ② 약간 불만족 ① 매우 불만족

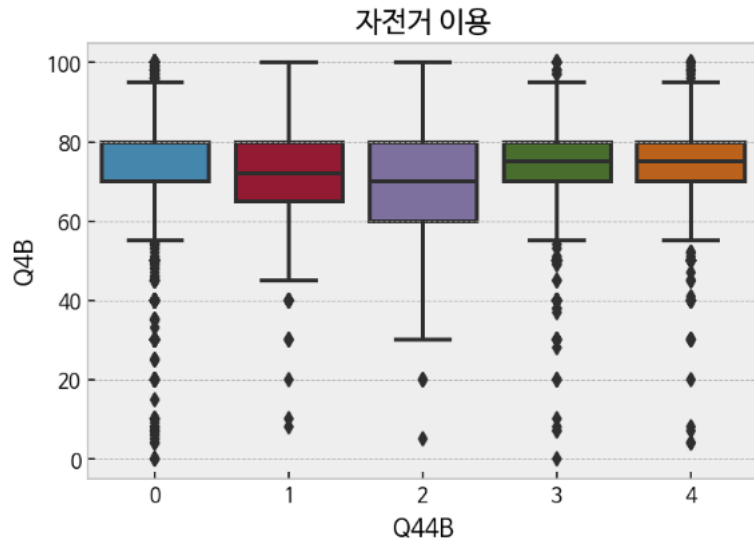
0. 자전거 이용 안함

1. 자전거 통근 + 환경만족

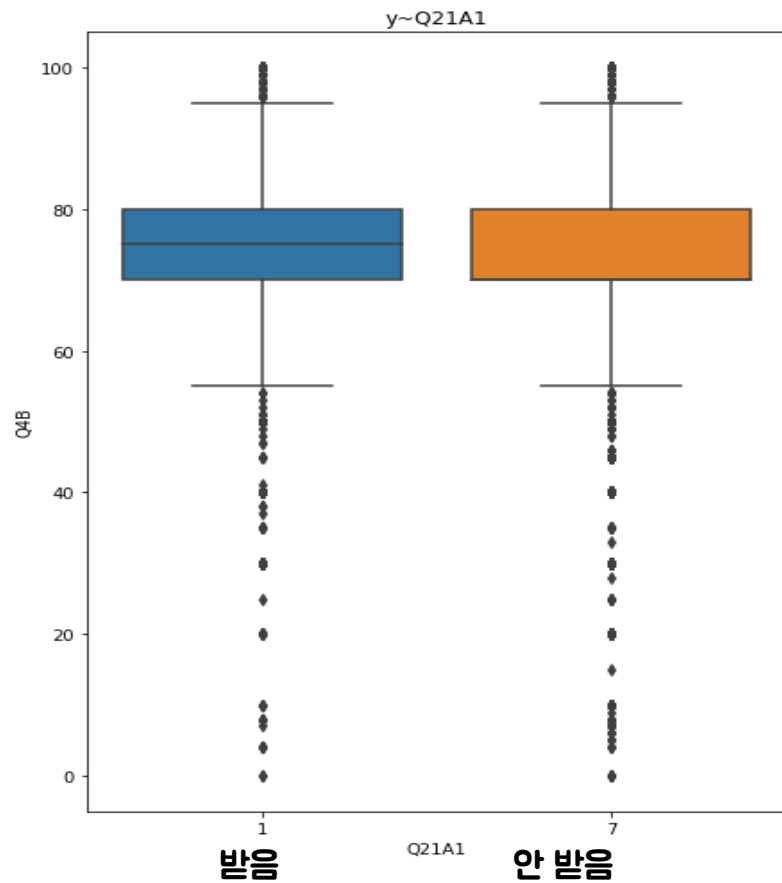
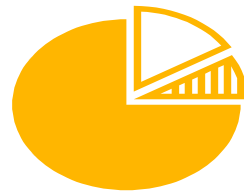
2. 자전거 통근 + 환경 불만족

3. 자전거 운동 + 환경 만족

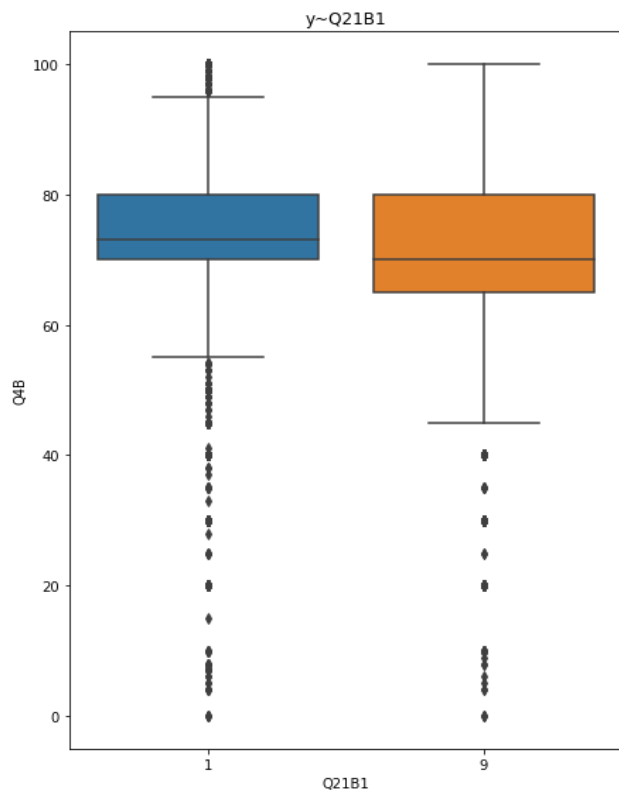
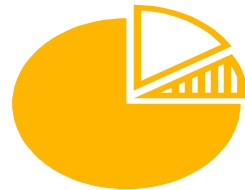
4. 자전거 운동 + 환경 불만족



이
호



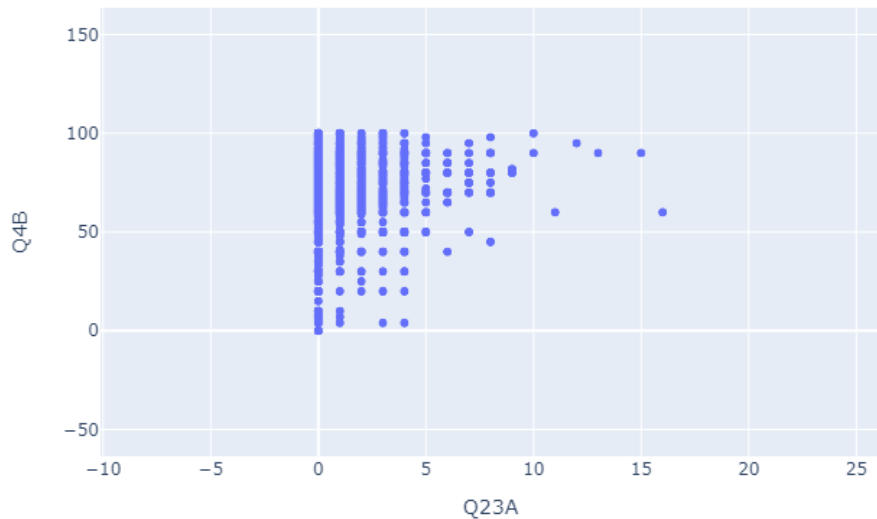
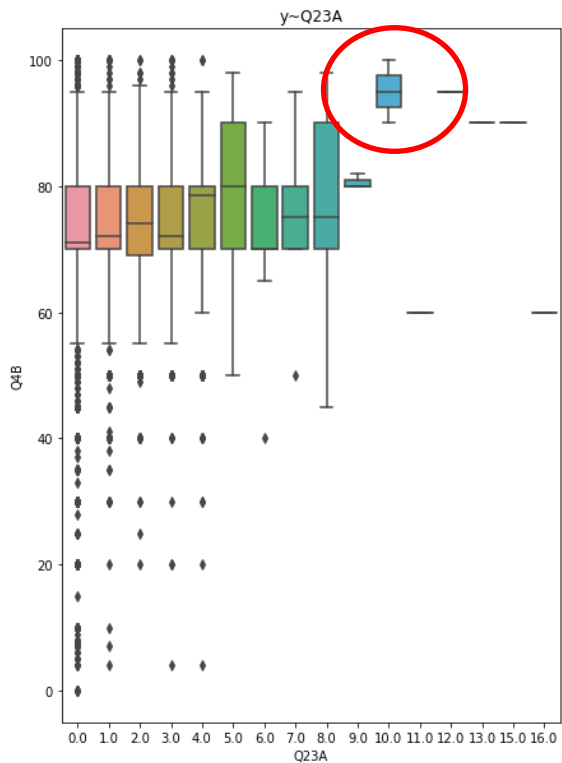
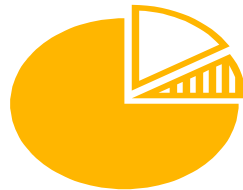
이
년



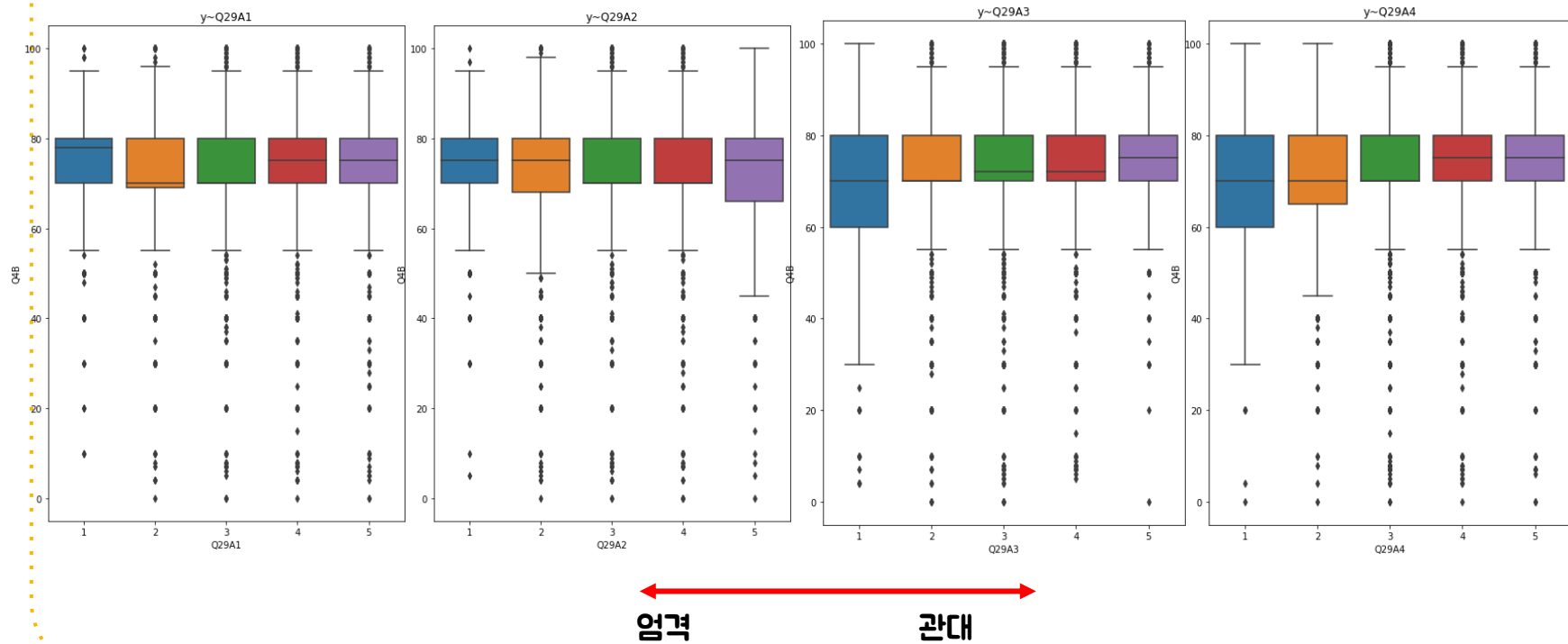
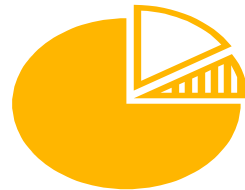
받을 의향

받을 의향 X

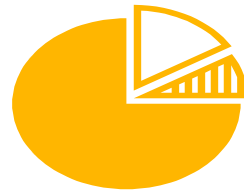
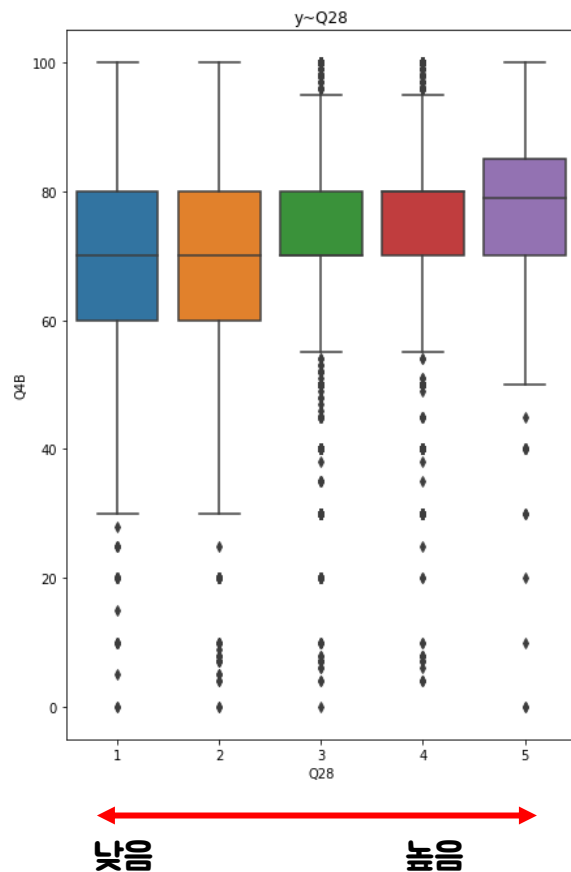
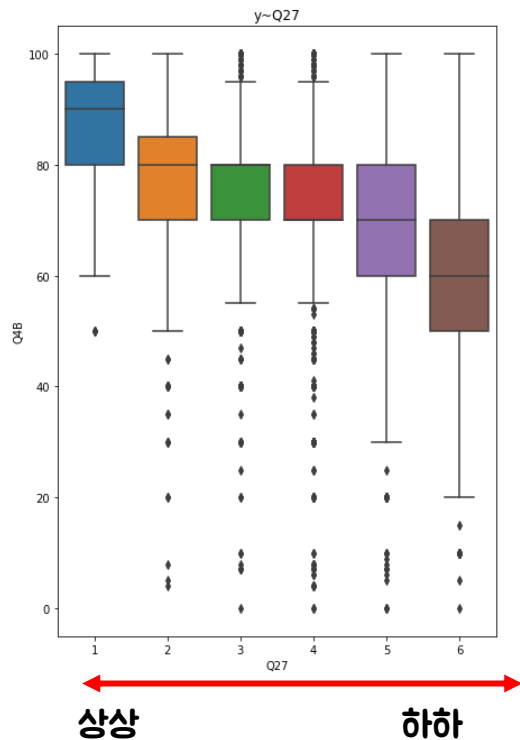
자원봉사활동 관련



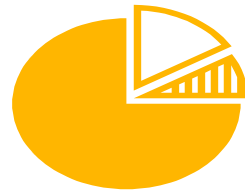
사회적 약자에 대한 태도



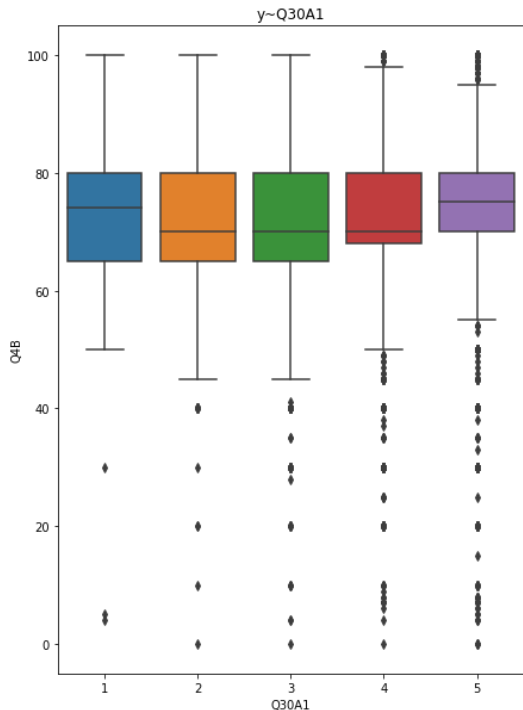
계층 인식



사회적 신뢰

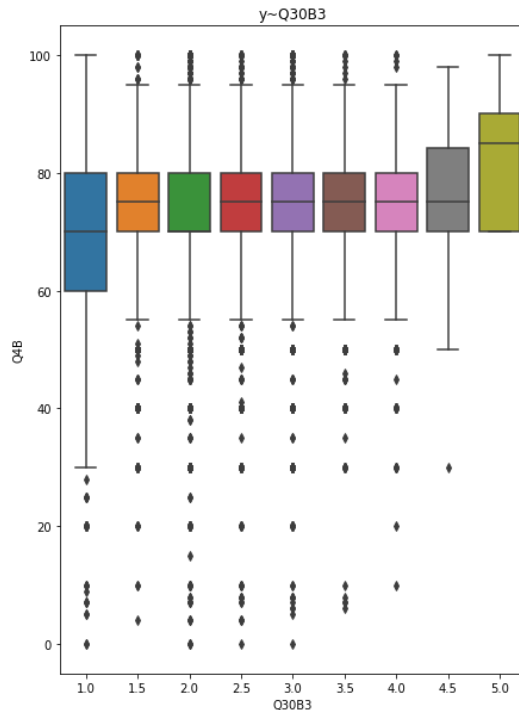


가족



신뢰 X

신뢰

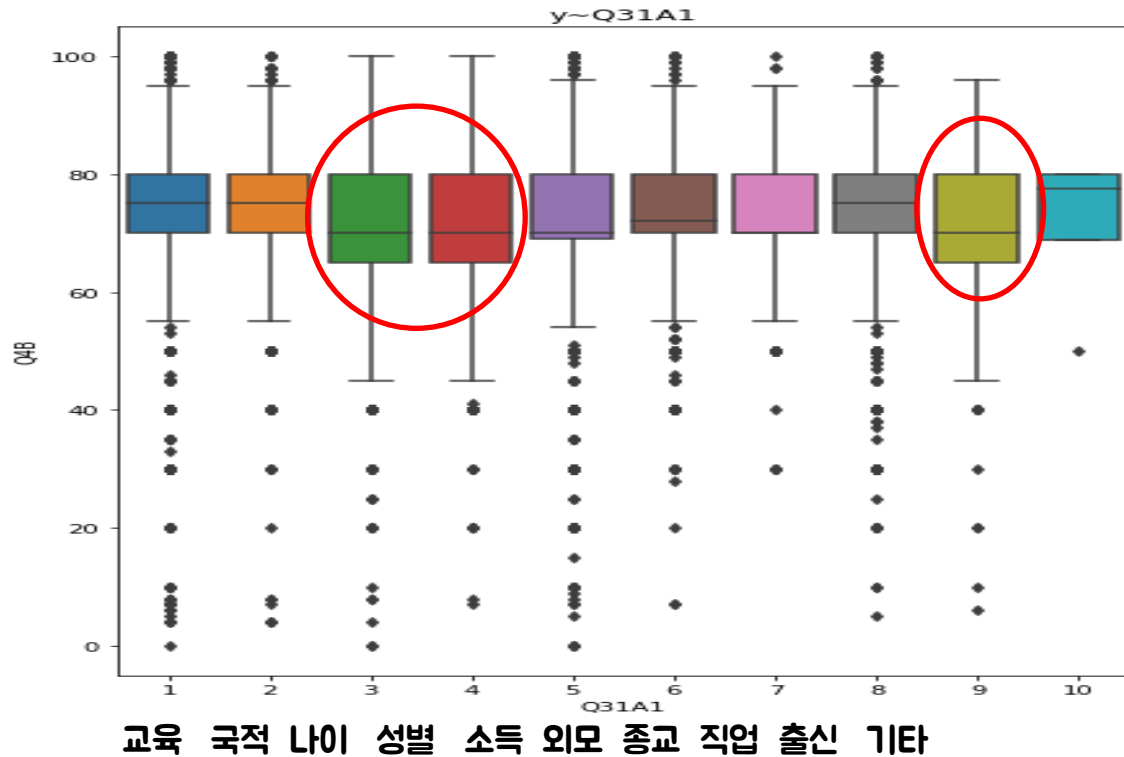
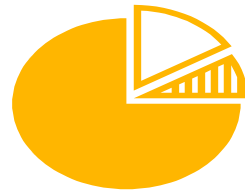


신뢰 X

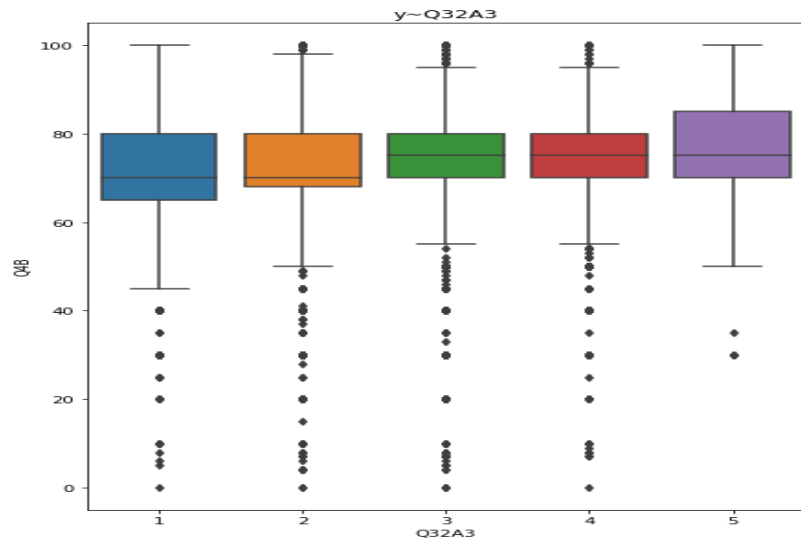
신뢰

외부 사람
= 처음 만난 + 다른
나라 사람

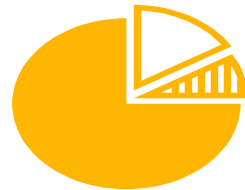
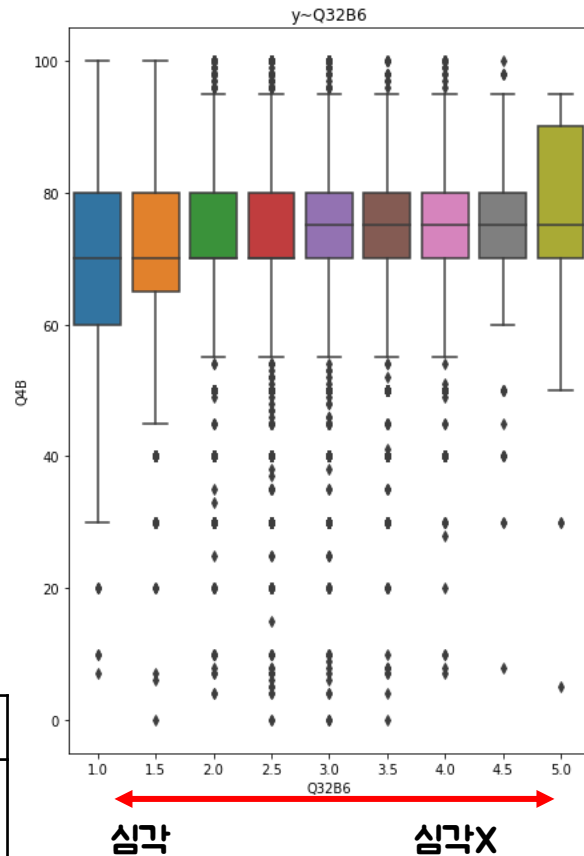
사회적 차별



가족 관련 위험

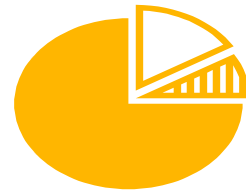
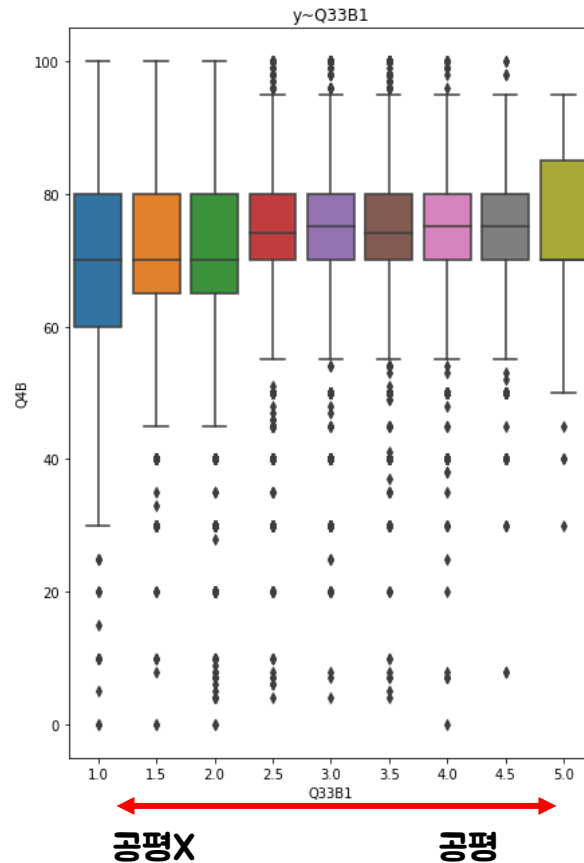


이혼 + 무자녀 또는 저출산	가족 형태
가족의 상부상조 기능 감소 + 재산 분배에 대한 갈등	가족 기능 (Q32B6)

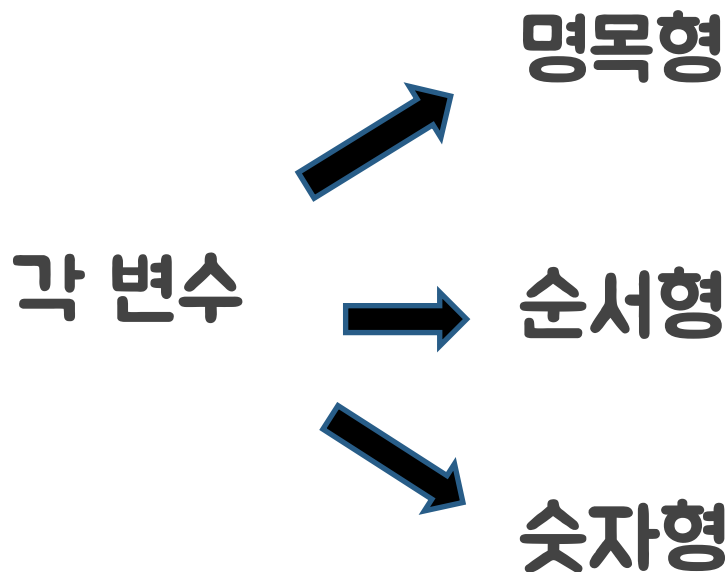
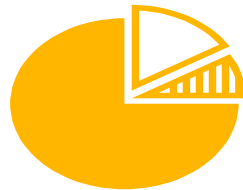


사회정의

수입과 소득 + 일자리/취업 기회	소득 (Q33B1)
대학교육의 기회	대학교육의 기회
수도권,지방의 발전 + 도시,농촌의 발전	지역 간 격차
사회복지 + 조세(세금) 정책	사회복지
남녀 평등 + 소수자의 권리	인권



각 변수 간의 독립성

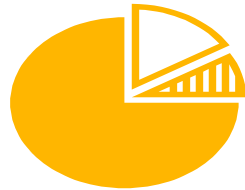


카이제곱 검정
~ Cramer's V

$$\phi_c = \sqrt{\frac{\chi^2}{N(k-1)}}$$

스피어만 상관계수

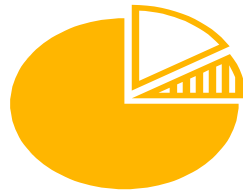
Cramer's V란?



$$\phi_c = \sqrt{\frac{\chi^2}{N(k-1)}}$$

2개 이상의 '범주'로 나눈 집단
간의 상관계수

Cramer's V란?



$$\phi_c = \sqrt{\frac{\chi^2}{N(k-1)}}$$

$$V = \sqrt{\frac{\chi^2}{n(q-1)}}$$

χ^2 : 카이제곱 공식에 의해 구함

n : 총사례수

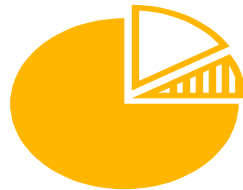
q : 줄 또는 칸의 유목수 중에 적은 숫자

특징1) 측정치들의 분포에 상관없이 적용

특징2) 항상 양수

=> 카이제곱검정의 효과크기

Cramer's V란?



$$\phi_c = \sqrt{\frac{\chi^2}{N(k-1)}}$$

χ^2 : 카이제곱 공식에 의해 구함

n : 총 사례 수

q : 줄 또는 칸의 수 중 적은 숫자

Cramer's V란?



EFFECT SIZE: CRAMER'S V

Contingency tables larger than the 2×2 table (i.e., that have more than two rows and/or two categories) typically use another measure of effect size, **Cramer's V**. Like phi values, **Cramer's V** calculates effect size values that range between 0.000 and 1.000.

The judgment for the magnitude of **Cramer's V** does not always use the guidelines we discussed for the other effect size measures (0.100 for small, 0.300 for medium, and 0.500 for large) that pertain to the other tests. This is because there are adjustments made to the formula caused by the shape of the table. Cohen (1988) provides an adjusted set of guidelines in a series of tables. For the 2×3 table in our example, we can use the same 0.100, 0.300, and 0.500 guidelines for judging the effect size for **Cramer's V** as we did for C. Larger tables will have reduced magnitude effect size criteria that determine small, medium, and large effects.⁴ (For example, if our table had been 3×4 , the effect size criteria for a large effect would be 0.354, not 0.500.)

⁴ We do not have the space in this book to develop this matter given its technical nature. However, SPSS reports the actual significance appropriate to any size of table, so in a practical sense, we do not need to delve more deeply into the subject at this point.

0.00–0.10

0.10–0.20

0.20–0.40

0.40–0.60

0.60–0.80

0.80–1.00

Interpretation of Φ in Chi-statistics or Cramér's V

Cramer's V is an effect size measure for two categorical fields are

Effect size is calculated in the

to determine which field has the

subtract 1 from the number of

multiply the result by the total

divide the chi-square value by

square test of independence

take the square root.

Cramer's V Coefficient (V)

Useful for comparing multiple χ^2 test statistics and is generalizable across contingency tables of varying sizes. It is not affected by sample size and therefore is very useful in situations where you suspect a statistically significant chi-square was the result of large sample size instead of any substantive relationship between the variables. It is interpreted as a measure of the relative (strength) of an association between two variables. The coefficient ranges from 0 to 1 (perfect association). In practice, you may find that a Cramer's V of .10 provides a good minimum threshold for suggesting there is a substantive relationship between two variables.

$$V = \sqrt{\frac{\chi^2}{n(q-1)}} \quad \text{where } q = \text{smaller \# of rows or columns}$$

Describing Strength of Association

1. Interpretation of effect size

Effect size	Interpretation
None	ES ≤ 0.2
Weak	0.2 < ES ≤ 0.3
Medium	0.3 < ES ≤ 0.5
Strong	ES > 0.5
Very strong	The result is strong. The fields are strongly associated.

Characterizations

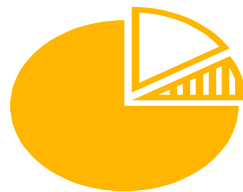
>.5	high association
.3 to .5	moderate association
.1 to .3	low association
0 to .1	little if any association

of-PH-in-Chi-statistics-or-Cramers-V_tbl2_311335682

g=PA98&lpg=PA98&dq=cramer%27s+v+interpretation+large+size&source=bl&ots=ytpqopG9Mko&sa=X&ved=2ahUKewjUgKaDINTpAhWIHqYKHUKDfMQ6AEwF3oECAGQAQ#v=onepage&q=cfalse

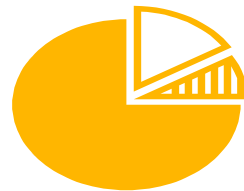
https://www.researchgate.net/figure/Interpretation-of-Phi-and-Cramers-V_tbl2_326885374

Cramer's V란?

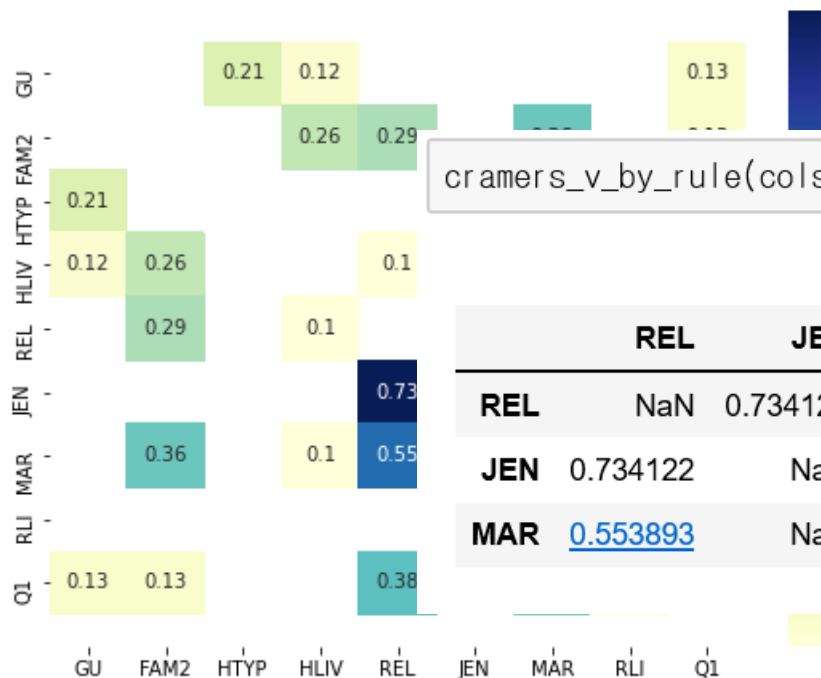


```
def direct_cramer_v_interpret(X, Y):  
    table = make_contingency_table(X, Y)  
    cramer = cramers_stat(table)  
    if cramer < 0.1:  
        print(X, ',', Y, ':', '연관성이 없다.')  
    elif cramer < 0.3:  
        print(X, ',', Y, ':', '약한 연관성이 있다.')  
    elif cramer < 0.5:  
        print(X, ',', Y, ':', '보통의 연관성이 있다.')  
    else:  
        print(X, ',', Y, ':', '강한 연관성이 있다.')
```

Cramer's v



Cramer V Correlation between Variables

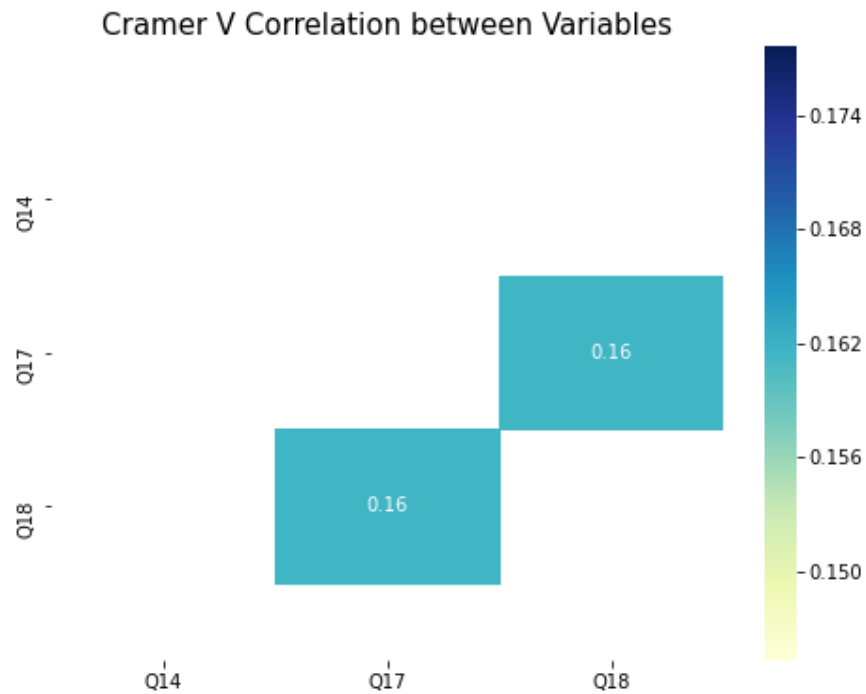
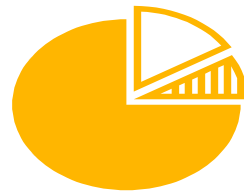


cramers_v_by_rule(cols)

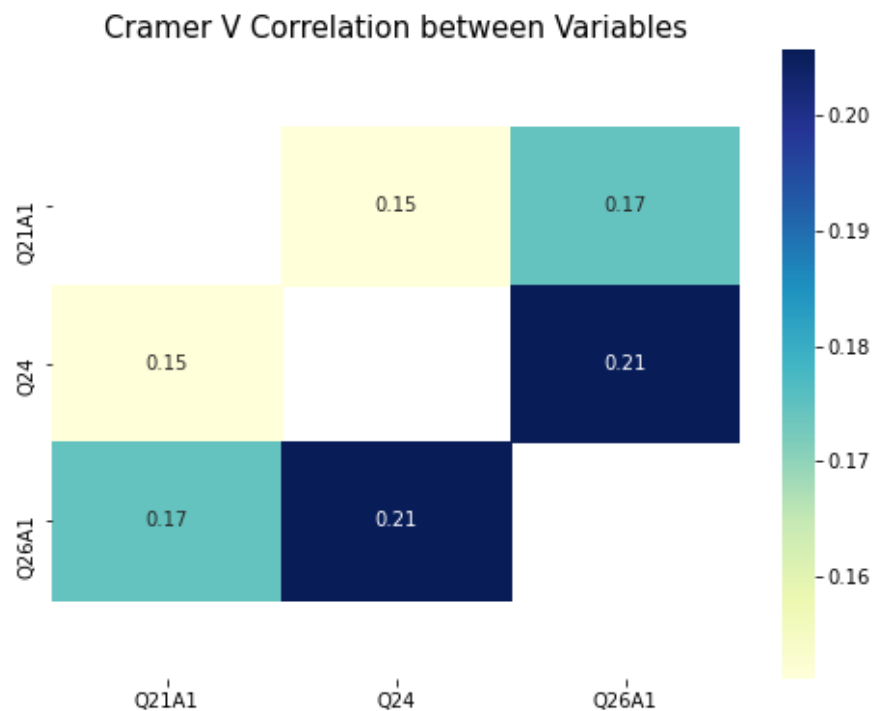
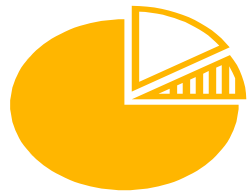
	REL	JEN	MAR
REL	NaN	0.734122	0.553893
JEN	0.734122	NaN	NaN
MAR	0.553893	NaN	NaN

	FAM2	REL	MAR	Q1
FAM2	NaN	NaN	0.363610	NaN
REL	NaN	NaN	NaN	0.381702
MAR	0.36361	NaN	NaN	0.344918
Q1	NaN	0.381702	0.344918	NaN

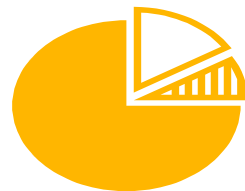
Cramer's v



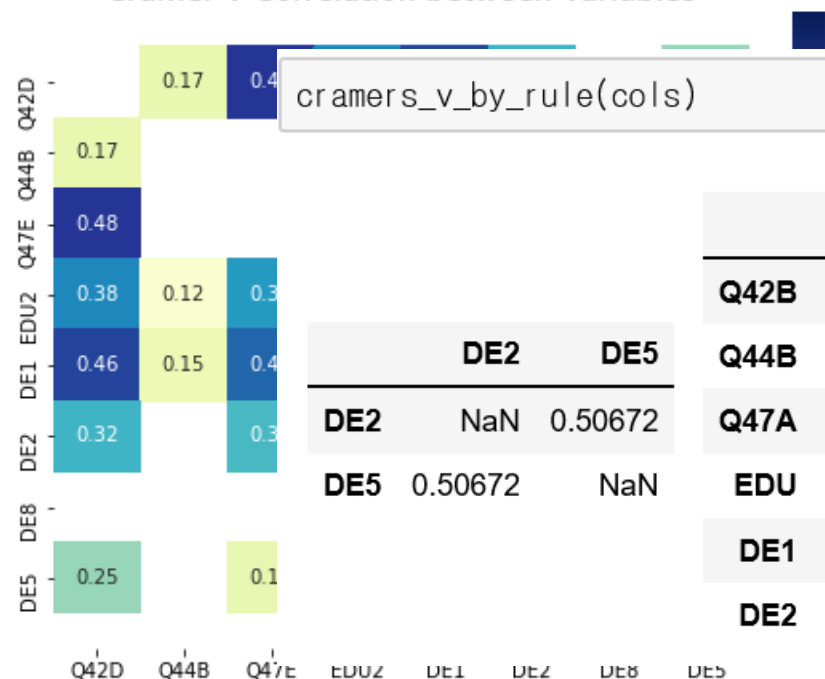
Cramer's v



Cramer's v



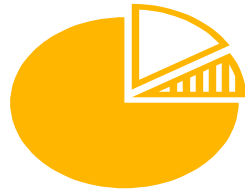
Cramer V Correlation between Variables



cramers_v_by_rule(cols)

	Q42B	Q44B	Q47A	EDU	DE1	DE2
Q42B	NaN	0.349537	NaN	NaN	0.326015	NaN
Q44B	0.349537	NaN	NaN	NaN	NaN	NaN
Q47A	NaN	NaN	NaN	NaN	0.316307	NaN
EDU	NaN	NaN	NaN	NaN	NaN	0.311649
DE1	0.326015	NaN	0.316307	NaN	NaN	NaN
DE2	NaN	NaN	NaN	0.311649	NaN	NaN

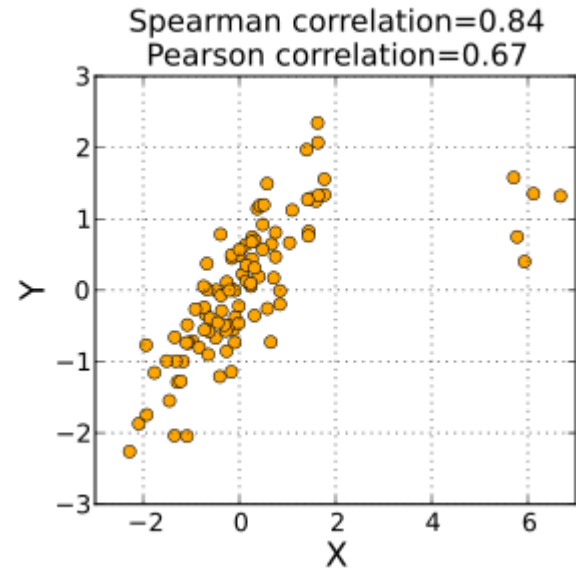
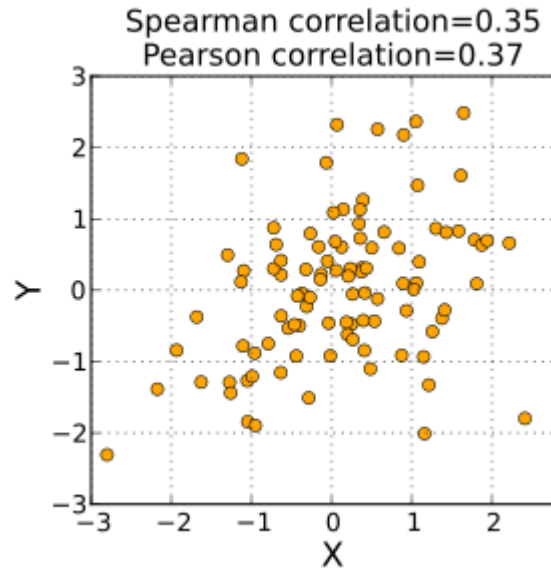
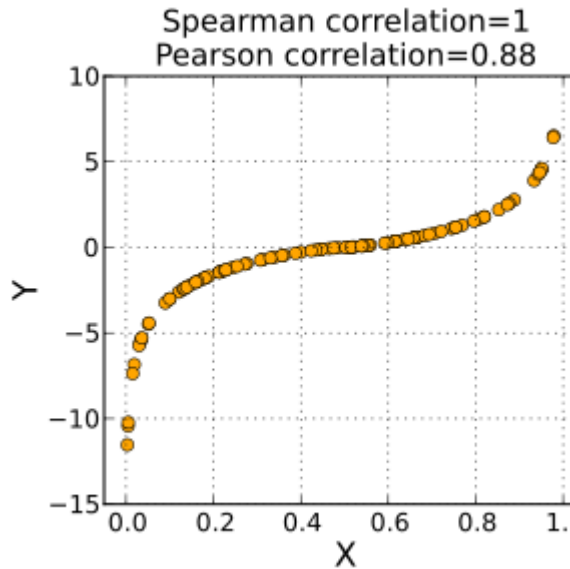
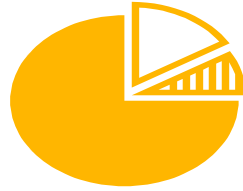
Spearman Rank Correlation



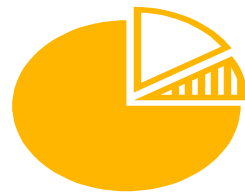
$$r_s = \rho_{rg_X, rg_Y} = \frac{\text{cov}(rg_X, rg_Y)}{\sigma_{rg_X} \sigma_{rg_Y}}$$

스피어만 상관 계수는 두 변수의 순위 사이의 통계적 의존성을 측정하는 비모수적인 척도이다. 이는 두 변수의 관계가 **단조 함수**를 사용하여 얼마나 잘 설명될 수 있는지를 평가한다. 스피어만 상관 계수는 순위가 매겨진 변수 간의 피어슨 상관 계수 로 정의된다.

Spearman Rank Correlation

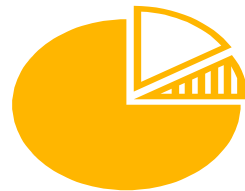


Spearman Rank Correlation

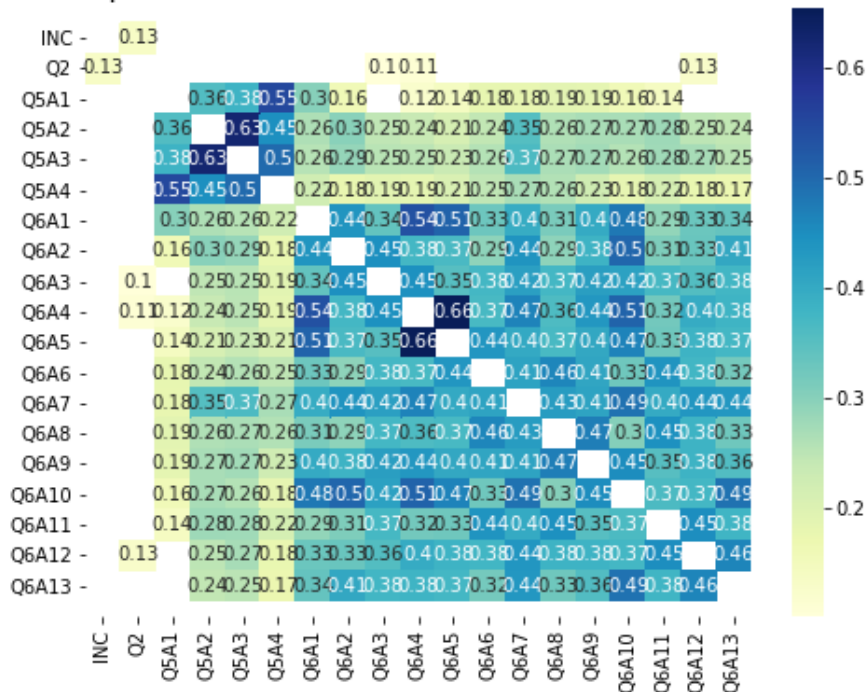


```
def spearmanr_interpret(X, Y):  
    spearman_rho = calculate_spearmanr(X, Y)  
    if abs(spearman_rho) < 0.1:  
        print(X, ',', Y, ':', '연관성이 없다.')  
    elif abs(spearman_rho) < 0.3:  
        print(X, ',', Y, ':', '약한 연관성이 있다.')  
    elif abs(spearman_rho) < 0.5:  
        print(X, ',', Y, ':', '보통의 연관성이 있다.')  
    else:  
        print(X, ',', Y, ':', '강한 연관성이 있다.')
```

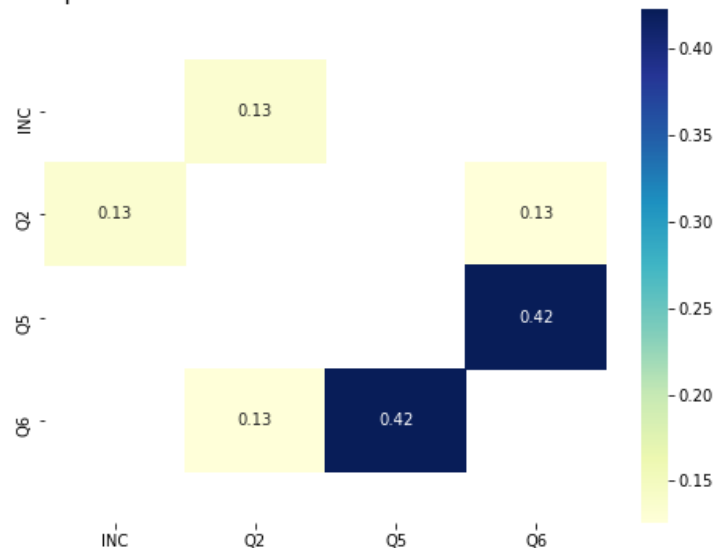
Spearman 상관계수



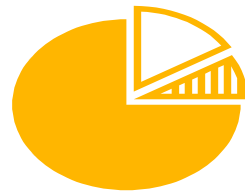
Spearman Rank Correlation between Variables



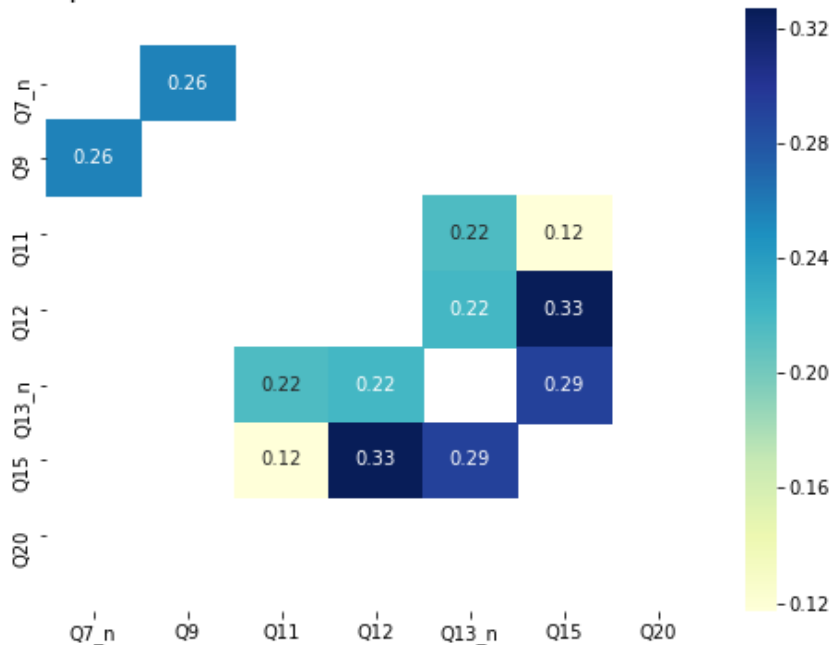
Spearman Rank Correlation between Variables



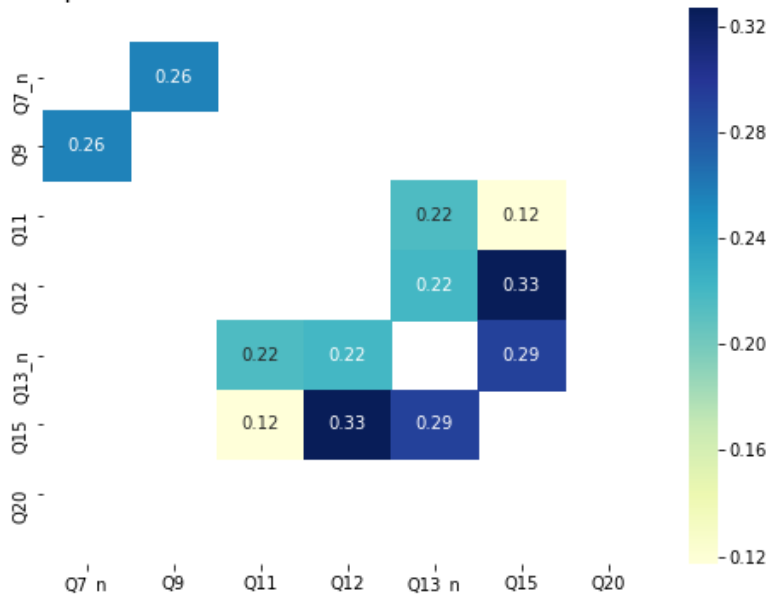
Spearman 상관계수



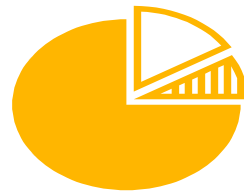
Spearman Rank Correlation between Variables



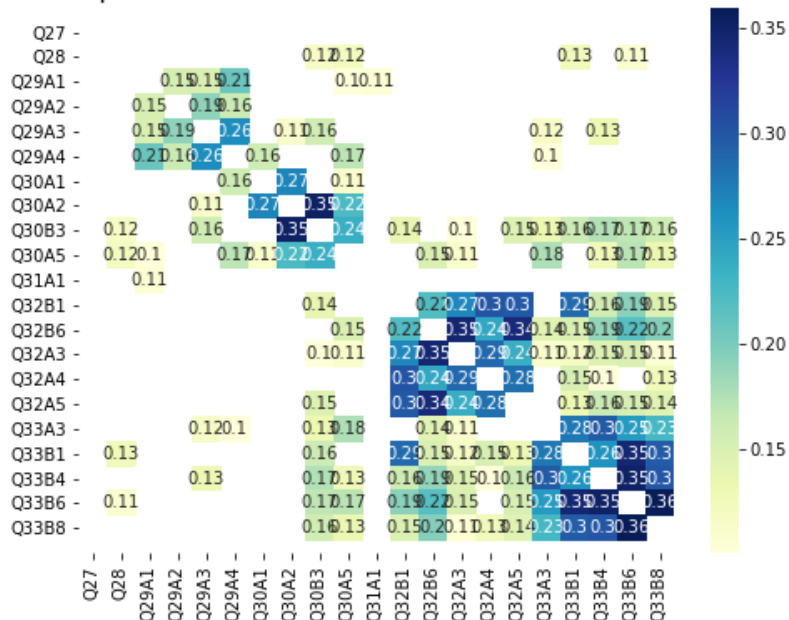
Spearman Rank Correlation between Variables



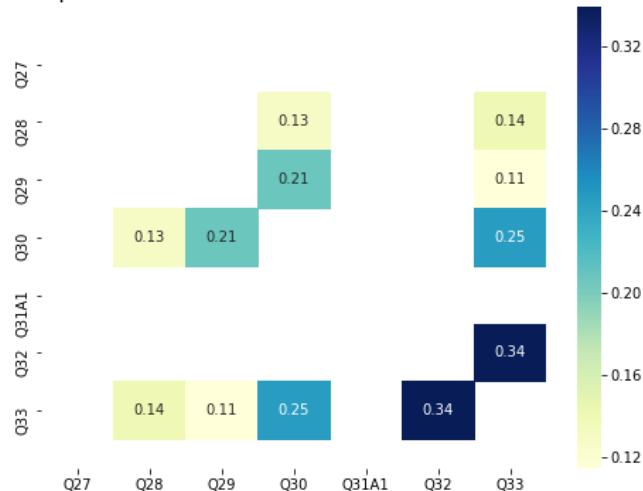
Spearman 상관계수



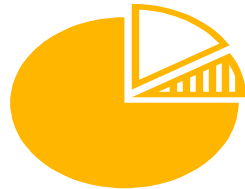
Spearman Rank Correlation between Variables



Spearman Rank Correlation between Variables



Spearman 상관계수



```
#Q5
data = data.assign(Q5 = np.mean(data.loc[:, 'Q5A1': 'Q5A4'], axis=1)) # A1~A4 평균값
data.drop(['Q5A1', 'Q5A2', 'Q5A3', 'Q5A4'], axis=1) #기존 변수를 제거 #inplace=True는

#Q6.
data = data.assign(Q6 = np.mean(data.loc[:, 'Q6A1': 'Q6A13'], axis=1)) # A1~A13평균값
data.drop(['Q6A1', 'Q6A2', 'Q6A3', 'Q6A4', 'Q6A5', 'Q6A6', 'Q6A7', 'Q6A8', 'Q6A9', 'Q6A10', 'Q6A11', 'Q6A12', 'Q6A13'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

#Q7,8
data = data.assign(Q7_n = np.mean(data.loc[:, 'Q7': 'Q8'], axis=1)) # Q7~Q8평균값으로
data.drop(['Q7', 'Q8'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

#Q13, Q13A
data = data.assign(Q13_n = np.mean(data.loc[:, 'Q13': 'Q13A'], axis=1)) # Q13~Q13A평균값
data.drop(['Q13', 'Q13A'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

#Q29
data = data.assign(Q29 = np.mean(data.loc[:, 'Q29A1': 'Q29A4'], axis=1)) # Q13~Q13A평균값
data.drop(['Q29A1', 'Q29A2', 'Q29A3', 'Q29A4'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

#Q30
data = data.assign(Q30 = np.mean(data.loc[:, 'Q30A1': 'Q30A5'], axis=1)) # Q13~Q13A평균값
data.drop(['Q30A1', 'Q30A2', 'Q30A3', 'Q30A4', 'Q30A5'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

#Q32
data = data.assign(Q32 = np.mean(data.loc[:, 'Q32A1': 'Q32A7'], axis=1)) # Q13~Q13A평균값
data.drop(['Q32A1', 'Q32A2', 'Q32A3', 'Q32A4', 'Q32A5', 'Q32A6', 'Q32A7'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

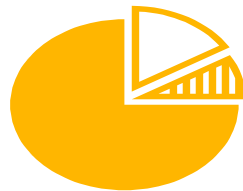
#Q33
data = data.assign(Q33 = np.mean(data.loc[:, 'Q33A1': 'Q33A9'], axis=1)) # Q13~Q13A평균값
data.drop(['Q33A1', 'Q33A2', 'Q33A3', 'Q33A4', 'Q33A5', 'Q33A6', 'Q33A7', 'Q33A8', 'Q33A9'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음

#Q45
data = data.assign(Q45 = np.mean(data.loc[:, 'Q45A1': 'Q45A3'], axis=1)) # A1~A3평균값
data.drop(['Q45A1', 'Q45A2', 'Q45A3'], axis=1) #기존 변수를 제거 #inplace=True는 아직 안 넣음
```

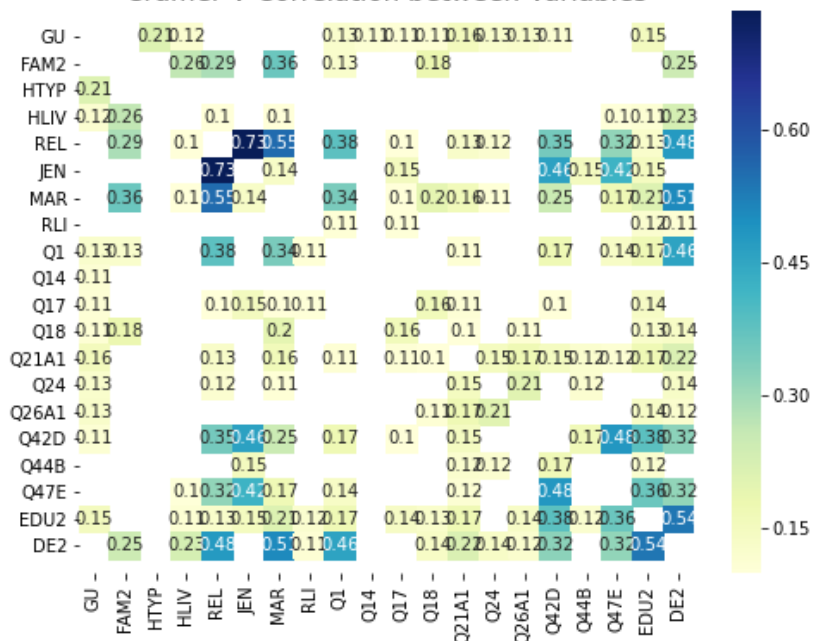
-
-
-
-
-

-
-
-
-

명목형 변수 Total

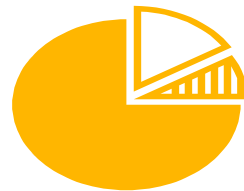


Cramer V Correlation between Variables



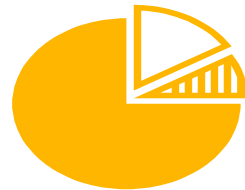
	REL	JEN	MAR	EDU2	DE2
REL	NaN	0.734122	0.553893	NaN	NaN
JEN	0.734122	NaN	NaN	NaN	NaN
MAR	0.553893	NaN	NaN	NaN	0.506720
EDU2	NaN	NaN	NaN	NaN	0.538172
DE2	NaN	NaN	0.506720	0.538172	NaN

전체 Total / Strong Corr



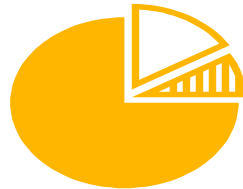
	Q42C3	REL	JEN	MAR	Q42D	EDU2	DE2
Q42C3	NaN	NaN	NaN	NaN	0.689181	NaN	NaN
REL	NaN	NaN	0.734122	0.553893	NaN	NaN	NaN
JEN	NaN	0.734122	NaN	NaN	NaN	NaN	NaN
MAR	NaN	0.553893	NaN	NaN	NaN	NaN	0.506720
Q42D	0.689181	NaN	NaN	NaN	NaN	NaN	NaN
EDU2	NaN	NaN	NaN	NaN	NaN	NaN	0.538172
DE2	NaN	NaN	NaN	0.506720	NaN	0.538172	NaN

전체 Total / Moderate Corr



	Q12	Q23A	Q42C3	FAM2	REL	JEN	MAR	Q1	Q24	Q42D	Q47E	EDU2	DE2
Q12	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	0.317194	NaN
Q23A	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	0.312008	NaN	NaN	NaN	NaN
Q42C3	NaN	NaN	NaN	NaN	0.306921	0.402452	NaN	NaN	NaN	NaN	0.461518	0.389835	NaN
FAM2	NaN	NaN	NaN	NaN	NaN	NaN	0.363610	NaN	NaN	NaN	NaN	NaN	NaN
REL	NaN	NaN	0.306921	NaN	NaN	NaN	NaN	0.381702	NaN	0.350983	0.316603	NaN	0.476977
JEN	NaN	NaN	0.402452	NaN	NaN	NaN	NaN	NaN	NaN	0.463026	0.422675	NaN	NaN
MAR	NaN	NaN	NaN	0.36361	NaN	NaN	NaN	0.344918	NaN	NaN	NaN	NaN	NaN
Q1	NaN	NaN	NaN	NaN	0.381702	NaN	0.344918	NaN	NaN	NaN	NaN	NaN	0.456939
Q24	NaN	0.312008	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Q42D	NaN	NaN	NaN	NaN	0.350983	0.463026	NaN	NaN	NaN	NaN	0.481647	0.384876	0.324049
Q47E	NaN	NaN	0.461518	NaN	0.316603	0.422675	NaN	NaN	NaN	0.481647	NaN	0.362904	0.316871
EDU2	0.317194	NaN	0.389835	NaN	NaN	NaN	NaN	NaN	NaN	0.384876	0.362904	NaN	NaN
DE2	NaN	NaN	NaN	NaN	0.476977	NaN	NaN	0.456939	NaN	0.324049	0.316871	NaN	NaN

앞으로...



1) Feature Extraction
PCA < FAMD

2) Analysis
SVM
XGboost
Group Lasso
Random Forest
*stacking

〈Feature Extraction 관련 생각〉

생각1. 잠재변수를 3개 또는 5개로 설정하고 싶다.

=> FA를 했을 때 3개 또는 5개로 나오면 좋겠다.

WHY?

=> 1) 3개: 행복의 구성요소 3개(LS, PA, NA)

=> 2) 5개: 성격 5요인(외향성, 개방성, 신경증, 성실성, 원만성)

생각1-1.

위의 생각에 따라, 미리 각 요인에 해당하는 값들에 적합할
항목들을 나눠보고

그것들의 밸런스를 위해 추가적인 변수를 제작해보자. ex) 부정적
감정을 측정하기 위해 어떤 것들 ratio로 제시해본다던가

Q&A

감사합니다^^