

Problem Set 2

Data Analysis and Visualization using R

RStudy2020

Due by 2020년 11월 2일 월요일 오전 12시 (한국시간)

Information & Instructions

출제자: 박상훈

데이터 파일

첫 번째 문제는 여기에서 2012년 미국 총선 시계열 설문조사 자료를 사용.

두 번째 문제는 여기의 'Quality of Governance'의 시계열 자료를 사용.

제출

정해진 시간까지 정해진 문제에 대한 코드를 작성한 R 스크립트를 sp23@email.sc.edu 혹은 Dropbox의 4_QnA or Discussion 폴더에 업로드.

모든 코드는 다른 작업환경에서 열더라도 재생산가능하게 디렉토리 설정 등을 모두 고려한 결과물이어야 함.

Here starts the actual test

문제 1

ANES 2012 시계열 자료를 `anes`라고 하는 R의 객체로 저장하라.

A. 정당 일체감(`party ID`; `pid_x`) 변수를 민주당(매우 강한 민주당 지지, 약한 민주당 지지, 민주당 우호(leaning)) 일 경우 1, 공화당일 경우 0으로 나타내는 이항변수(binary variable)로 코딩하라. 이때, 무당파(independents)와 결측치를 나타내는 값들(응답거부(refused), 무응답(no answer), 기타(etc))은 제외하라.

B. 교육수준을 나타내는 변수(`dem_edu`)를 숫자형(numerical) 변수로 재코딩하되, 모든 결측치와 other 카테고리를 포함하지 않도록 하라. 결과적으로 숫자형 변수로 재코딩된 변수는 1부터 16까지의 값을 가져야만 한다.

C. 민주당 지지자들과 공화당 지지자들 간의 평균 교육 수준의 차이가 존재하는지 여부를 평균 차이 분석을 통해 검증하라. 두 집단의 차이가 유의미하다고 할 수 있는가? 두 정당 지지 집단 간의 교육 수준의 평균 차이가 실질적으로 중요한 차이라고 생각되는지에 대해 자신의 의견을 서술하라.

D. 새롭게 만든 교육 수준의 숫자형 변수의 측정 수준은 무엇인가(e.g. 명목형, 순위형, 연속형 등)? 이 측정수준이 C에서 요구한 평균 차이를 검증하는 데 있어서 문제의 소지가 있는지 혹은 없다고 생각하는지 자신의 의견을 서술하라.

E. 정당 일체감에 따른 평균 교육 수준의 차이가 존재하는지를 분석하기 위해 ANOVA를 시행하라. 이때, 정당일체감은 7점 척도의 변수를 사용하라(pid_x를 다시 사용하되, 결측치만 제외하고 일체감의 강도는 그대로 놔두어 분석하라).

F. ANOVA의 결과를 어떻게 해석할 수 있는가? 이때, ANOVA의 영가설은 무엇인가? 분석 결과는 ANOVA의 영가설을 기각할만한 충분한 근거를 제시하는가? 집단 간의 차이를 보여주는 데 있어서 ANOVA가 가지는 한계는 무엇인가?

문제 2.

Quality of Governance 시계열 자료를 qog라고 하는 R의 객체로 저장하라.

A. 가장 최근의 qog 데이터를 사용하여 각 지역(ht_region), 각 연도, 국가간 분쟁 경험 횟수 (UCDP-PRIO의 inter-state armed conflict 변수를 사용할 것)를 요약하여 보여주는 데이터셋, qog_agg을 만들어라. 하지만 단순히 분쟁 경험 횟수만을 측정하는 변수가 아니라 무력 분쟁 변수를 3개의 항목을 가진 분류형 변수로 재구성하라: (1) 0: 국가간 분쟁을 경험한 횟수가 없는 경우, (2) 1: 국가 간 분쟁을 한 번 경험한 경우, (3) 2: 2번 이상의 국가간 분쟁 경험이 있는 경우.

B. 새롭게 구축한 데이터셋의 분석 단위(level of analysis)는 무엇인가? 그리고 새롭게 구축한 무력분쟁 변수의 측정 수준은 무엇인가?

C. 지역과 무력분쟁 변수 간의 관계를 보여주는 교차표를 만들어라. 두 변수를 대상으로 카이스퀘어(χ^2) 검정을 수행하고, 각 변수의 관측치들이 서로 독립적이라는 카이스퀘어 검정을 충족시키는지 확인하라.

D. 카이스퀘어 검정 결과를 설명하라. 카이스퀘어 검정의 영가설은 무엇인가? c의 결과는 카이스퀘어 검정의 영가설을 기각할 수 있는가? 카이스퀘어 검정의 한계는 무엇인가?

E. 만약 지역별로 얼마나 많은 국가간 분쟁이 발생하였는지를 알고 싶다고 할 때, 우리가 재조작하여 만든 분쟁 변수가 가질 수 있는 문제점은 무엇인가?