

학과관람의 시작 - 1912168 - 인포임

제출파일

중요도

1st  
2nd  
3rd

## < R 데이터 함수 총 정리 >

기초부터

데이터파일 읽기  $\leftarrow$  read.csv("파일명.csv", stringsAsFactors=F)

head(데이터, k): 데이터를 위에서 k행까지

tail(데이터, k): 데이터를 아래에서 k행까지

View(데이터): 데이터를 보임 (Data 창에서 바로 출력도 ok)

dim(데이터): 행과 열이 각각 몇 개씩 있는지 알기

str(데이터 or 데이터 \$ 변수): 해당 값의 특성을 보임 (변수 안자문자 알지)

데이터 \$ 변수  $\leftarrow$  NULL: 해당 변수 삭제

rm(데이터): Data 창에서 데이터 삭제

데이터 \$ 변수  $\leftarrow$  as.factor(데이터 \$ 변수): 해당 변수의 값을 범주형으로 변경

문자: chr  
범주: factor

summary(데이터): 해당 데이터의 기술통계량을 요약

mean(데이터 \$ 변수): 해당 변수의 평균

var(데이터 \$ 변수): 해당 변수의 분산 (분산이 클수록 편차가 큼)

sd(데이터 \$ 변수): 해당 변수의 표준편차

round(어떤 값, digits=N): 해당 값에서 소수 N번째 자리까지

함수를 이용한 값 찾기 ex) mean(exam \$ math)

hist(데이터 \$ 변수, breaks = seq(N1, N2, by=k)): 해당 변수를 N1 ~ N2까지 간격을 k인

히스토그램 그리기 \* 변수의 개수가 두 개일 / N1과 N2는 summary(데이터 \$ 변수) 통해 최대, 최대

table(데이터 \$ 변수): 해당 변수의 빈도

(크기가 유사한 범주의 빈도만 가능)

table(데이터 \$ 변수, 데이터 \$ 변수2): 두 변수를 조합한 빈도

install.packages("패키지"): 패키지를 설치함

library(패키지): 설치된 패키지를 불러옴

⊕ remove.packages("패키지 이름")을 통해 remove

\* R의 ggplot2 패키지 불러오기 library(ggplot2)

qplot(data=데이터 (변수)): 해당 변수의 빈도를 막대 그래프로 나타냄

my story 자이스토리

색상 구분  $\leftarrow$  문장형이므로 " " 안에서

qplot(data=데이터, 변수1, fill=변수2): 변수1과 변수2로 조합된 빈도를 막대 그래프로 표현

데이터 \$변수 ← as.Date(데이터 \$변수) : 해당변수의 형식을 날짜(Date)로 변경

Weekdays(데이터 \$변수) : 해당변수(날짜)의 요일이 무엇인지 알려줌

└ 이 변수의 형식도 날짜(Date)여야함

Weekdays(as.Date("2020-01-01")) : "2020-01-01"을 날짜(Date)

로 인식한 뒤 이 날짜의 요일을 알려줌

데이터 \$A ← B : B라는 값을 (항수로 구성) 데이터의 새로운 변수 'A'에 저장  
└ 새 변수  
└ 하나 어떤 데이터의 변수도 아 ex) 데이터 \$변수

mean/var/sd(데이터 \$변수, na.rm=T) : 데이터 변수에서 결측치(NA) 제외하고 계산

데이터 \$변수 ← factor(데이터 \$변수, levels=c("a", "b", "c" ... ))

: a, b, c 순으로 변수를 배열하고 형식을 범주형으로 바꾼 뒤 해당 데이터 변수에 지정

stat.desc(데이터 \$변수) : 해당변수의 여러 기술통계량 보여줌 % summary

★ └ 우선 'paste' 설치하고 불러오기

describe(데이터 \$변수) %>% stat.desc

└ 우선 'psych' 설치하고 불러오기

★ ★ 우선 "readxl" 패키지 설치하고 불러오기 ★ ★

데이터 파일 이름 ← read\_excel("파일명.xlsx", sheet=k) : 엑셀 파일에서 시트에서 k번째로 불러오기

write.csv(보낼때 저장할 이름, file="파일명.csv") : B를 A라는 이름으로 보내내기  
(CSV 형식) A B

★ ★ 우선 "writexl" 패키지 설치하고 불러오기 ★ ★

write\_xlsx(보낼때 저장할 이름, path="파일명.xlsx") : B를 A라는 이름으로 보내내기  
(Excel 형식) A B

★ ★ 예들 들어 welfare에 있는 education 변수를 welfare에 복사해서

1) welfare\$education ← welfare\$education

2) welfare\$education ← welfare\$education %>% mutate(education = welfare\$education)



table (데이터 \$변수 == k) : 해당변수가 k인것 (문자형태면 " ")

table (데이터 \$변수 != k) : 해당변수가 k가 아닌 것

table (데이터 \$변수 > k) : 해당변수가 k보다 큰 것

< k : 해당변수가 k보다 작은 것

> k : 해당변수가 k보다 같거나 큰 것

≤ k : 해당변수가 k보다 같거나 작은 것

table (데이터 \$변수 + 비교연산자 & 데이터 \$변수 + 비교연산자) / (1)(shift + H)

~이, ~이면서 ~이거나, 또는

= / != / > < > < 다 가능

해당 조건을 만족하는 것

table (데이터 \$변수 %in% c("a", "b", "c...")) :

데이터 변수가 a이거나 b이거나 c인 것들의 것 (1 변수 사용과 동일)

→ %in% c(~) : ~값들을 다 통틀

★ 우선 "dplyr" 패키지 설치하고 불러들이 ★

1) 데이터 <- rename (데이터, 바꿀변수명 = 이전변수명)

여러 변수 동시에 바꿀 수 있음

2) 데이터 <- 데이터 %>% rename (바꿀변수명 = 이전변수명)

→ 변수가 문자형이어도 " " 안써도 OK

데이터 <- 데이터 %>% summarise (이 패키지에 들어있는 데이터만 가능)

데이터 <- 데이터 %>% summarise :: 데이터 : 데이터의 원본을 가져옴 (만 %>% 패키지에

안 들어 있는 데이터만 가능) ★ "summarise" 패키지 설치하고 불러옴 (ex) mpg, mtcars)

데이터 \$변수 <- ifelse (데이터 \$변수 == "a", "b", 데이터 \$변수) (ifelse 두 변수)

: 만약 데이터의 변수 값이 a이면 b로 바꾸고 아니면 그대로 데이터의 원래 변수명 그대로 두어라

table (is.na (데이터 \$변수)) : 데이터의 해당 변수의 NA (결측치)가 몇 개인지

table (! is.na (데이터 \$변수)) : 데이터의 해당 변수에서 NA (결측치)가 아닌게 몇 개인지

⊕ 예를 들어서 데이터 프레임 weather-new의 변수 평균값이 NA일 경우는

is.na (weather-new \$ 평균값) 이지

weather-new \$ 평균값 == NA는 답

%>% is.na (데이터 \$변수) : 해당 변수가 NA이다

→ filter 함수에 다양한 여러 변수들을 조건으로 조건에 맞게 설정가능

ex) filter(class != 4, math >= 90 | hstory >= 95)

4번이 아니면  
수학이 90이상  
영어 95이상

중요도 ☐  
재복습 2nd ☐  
3rd ☐

~~\*\*\*\*~~

" " 안써도됨

데이터 전체는 (우선 'dplyr' 패키지 설치 후 불러오기)

데이터 %>% filter(변수 == 사례) : 데이터 안에 있는 변수에서 해당 사례만 추출

④ var나 mean과 같은 기술통계량을 이용해 위해서는 해당 사례들만 추출한 것들을 새로운 데이터 파일에 asstion 해야함

데이터 %>% filter(변수 > quantile(변수, probs = c(k)))

: 데이터의 해당 변수의 값이 변수의 값 중에서 k인 값보다 이상인 사례만 추출

즉 k가 0.9이면 0.9 이상인 값, 변수가 상위 10% 인 사례들만 추출 (변수의 값이 상위 10% 인 것들만 추출)

데이터 %>% select(변수) : 해당 변수들만 추출 (여러 변수를 3로 변경해서 추출 가능)

데이터 %>% select(-변수) : 해당 변수를 제외하고 추출

데이터 %>% select(contains("특정단어")) : 특정 단어가 포함된 변수들

데이터 %>% arrange(변수) : 변수에 대한 오름차순 정렬

데이터 %>% arrange(-변수) : 변수에 대한 내림차순 정렬

④ head(k) : 정렬된 사례들 중에서 위에서 k까지

→ 이 변수는 무조건 변수의 상태만 가능 ex) exam\$math

이제라 math이나 -math 안됨

→ 기존 변수 활용 하 다른 데이터 이용 가능

데이터 <- 데이터 %>% mutate(새변수 = )

라는 과정을 거친 새변수를 데이터에 지정, mutate 함수로 동시에 여러 변수 생성 가능

(단순 값들만 보고 싶으면 데이터 파일에 asstion 안해도됨)

□ 과정에 ifelse 함수가 있는데 ifelse 안에 또 ifelse 가 들어갈 수 있음

데이터 %>% group-by(변수) %>% summarise(mean = mean(...))

: 변수별로 그룹을 나누고 기술통계량을 조건에 맞게 제시

• summarise를 보여주는 프로그래밍 언어 변수명 지정 가능 ex) mean-math = mean(math)

~~\*\*\*\*~~ mutate든 summarise를 이용하든 count = n() 변수 이용한 코드

count / sum(count) : 해당 변수들까지 나뉘어진 사례들 개리의 전체 개수에서 해당 사례 비율

my story 자이스트리

count / sum(데이터 \$ count) : 모든 사례들 (전체) 개수에서 해당 사례 비율

print(n=k) : k를까지 다 보여줌 (앞까지 다 보고 싶으면 k를 100으로 지정)



데이터 <— 데이터 %>% relocate (변수, before = 변수2) : 변수를 변수2 왼쪽에

데이터 <— 데이터 %>% relocate (변수, after = 변수2) : 변수를 변수2 오른쪽에

데이터 <— left\_join (데이터1, 데이터2, by = "변수") : 데이터와 데이터2를  
공통변수를 이용해 합침 (데이터에 저장)

데이터프레임명 <— data.frame (변수1 = c(변수1의 사례내용), 변수2 = c(변수2의 사례내용))  
: 변수1과 변수2를 갖는 새로운 데이터 프레임 생성 (변수가 문자형이면 " ")

~~\*\*\*~~ 사례들이 단일 숫자인 : 이용해서 표 (x) 1:7 (단 1~7명씩만)

→ 공통변수의 적도와 동일한 경우 str() 이용해서 as.factor()로 적도를 범주형으로 일치

데이터 <— left\_join (데이터1, 데이터2, by = c("변수1" = "변수2"))

: 공통변수가 없을 때 변수와 변수가 같다고 저장한 뒤 데이터를 같이 합치기

데이터 <— bind\_rows (데이터1, 데이터2) : 데이터와 데이터2를 합쳐서 사례추가

데이터 <— 데이터 %>% distinct (변수, keep\_all = T) :

해당변수가 같은 사례들중에 한개만 남기고 다 제외하기 → 변수 최대 3개가 적당 (값이 같은 변수만)  
→ (exam에 3개만 알아서 자른 사례도)

~~\*\*\*~~ mtlwest 3층은 우선 'ggplot2' 패키지 설치 후 불러오기

데이터 <— ggplot2 :: 데이터 이용해서 파일 불러오기 후 실행

데이터 <— 데이터 %>% drop\_na() : 데이터에서 다양한 변수들중에서

그 값이 단 하나라도 NA가 있다면 그 사례는 삭제 → 우선 'tidyr' 패키지 설치 후 불러오기

데이터 %>% select (A : last\_col()) : 데이터에서 A부터 마지막 변수까지 추출

~~\*\*\*~~ sav 형태의 파일은 'foreign' 패키지 설치 후 불러오기

데이터 <— read.sav (fille = "처음 저장할 때 파일 이름", to.data.frame = T)

: sav 파일 불러서 데이터 프레임 만들기 → 파일명.sav 형태

ggplot (데이터, aes (변수1, 변수2, fill = 변수2)) + geom\_bar (stat = "identity")

: 변수1과 변수2를 각각 x축 y축이 되고 '값'들을 나타내는 막대 그래프 만들기

'ggplot2' 패키지 불러오기

\*\*\* 원 Library (ggplot2) 하기 \*\*\*

ggplot (데이터, aes (reorder (변수1, 변수2), 변수2, fill = 변수2))  
+ geom\_bar (stat = "identity") + theme (axis.text.x =  
element\_text (size = 7.5, angle = 50))

: 변수1과 변수2를 표현하는 막대 그래프 만들기

새로운데이터 <- 기존데이터 : 기존데이터의 복사본인 새로운데이터 만들기

A <- C : (1:30) : 1부터 30까지의 수라는 이름의 열벡터 생성

Colors () : 색깔을 나타내줌 (그래프나 텍스트 마이닝 이용)