

<정태분 | 함수 총정리> (문자: chr 변수: factor) ok

데이터 파일 읽기 ← `read.csv("file.csv", (col.names=T) locale=locale('ko', encoding='euc-kr'))`: file를 데이터 파일 이름으로 불러와서 저장하기

★ 우선 `library(readr)` 먼저하기

`View(데이터)`: 데이터 보기

`str(데이터 또는 데이터 $변수)`: 데이터 또는 데이터의 변수에 대해 다양한 정보 보여줌

`summary(데이터)`: 범주형 변수에 대한 개발변에 대한 기본적인 정보 보여줌, NA가 어떤 변수에 있는지 알려줌

데이터 \$변수 ← `as.factor(데이터 $변수)`: 해당 변수의 종류를 범주형으로 변경

`summary(데이터)`: 데이터의 다양한 정보 보여줌 (색으로 할 때)

★ 우선 `library(dplyr)` 먼저하기

`table(데이터 $변수)`: 해당 변수의 빈도

`freq(데이터 $변수)`: 해당 변수의 빈도와 막대 그래프

★ 우선 `library(descr)` 먼저하기

`gplot(data=데이터, 변수)`: 해당 변수의 빈도를 막대 그래프로 표현

`gplot(data=데이터, 변수, fill=변수2)`: 변수1과 변수2의 조합된 빈도를 막대 그래프로 표현

★ 우선 `library(ggplot2)` 먼저하기 (색상 지정)

• 변수에는 " " 인덱스

`mean(데이터 $변수)`: 해당 변수의 평균

`var(데이터 $변수)`: 해당 변수의 분산

`sd(데이터 $변수)`: 해당 변수의 표준편차

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

`freq(ts, na(데이터 $변수))`: 데이터의 변수에서 평균치가 있는지 여부를 알고 싶으면 빈도 T로 표현

★ 우선 `library(descr)` 먼저하기 / 즉 NA=평균치가 몇 개인지

`describe(데이터 또는 데이터 $변수)`: 해당 데이터 또는 변수의 해당하는 정규분포의 수치

(`cskew-왜도`, `kurtosis-첨도` ...)

★ 우선 `library(psych)` 먼저하기

`hist(데이터 $변수, breaks=seq(M1, M2, k ...))`: 해당 변수를 M1/M2까지 간격을 k인 히스토그램

• 변수의 색도와 구상형 / M1, M2는 `summary(데이터 $변수)` 통해 색도, 최대(최소) 값

같다 ==
 같지 않다 !=
 크다 >
 작다 <
 이상 >=
 이하 <=
 그리고 &
 또는 | (shift + #)

<패키지명>

library(readr)
 library(dplyr)
 library(descr)
 library(ggplot2)
 library(psych)

table(데이터\$변수 %in% c("a", "b", "c")) : 데이터의 변수가 a, b, c 중에
 적어도 하나의 조건에 만족하는 것들의 빈도, (문자가 아니면 " " 필요)
 • 1 연산자 사용과 동일

데이터\$A <- B : B라는 값을 데이터의 새로운 변수 'A'에 저장
 ↳ 함수거나 어떤 데이터의 변수여도 OK

데이터\$A <- NA : 데이터의 A라는 변수 삭제

Weekdays(데이터\$변수) : 해당 변수(숫자)의 요일이 무엇인지 보여줌
 ↳ 숫자가 날짜(date) 여야 함

데이터\$변수 <- as.Date(데이터\$변수) : 해당 변수의 숫자를 날짜(date)로 변환
 ex) weekdays(as.Date("2020-01-01")) : 수요일

ex) weather\$요일 <- weekdays(weather\$일시) : weather의 변수 '일시'의
 요일을 위한 값들을 데이터 weather의 새로운 변수 '요일'에 지정

↑ R과 같은 환경 사용 가능

데이터\$변수 <- factor(데이터\$변수, levels=c("a", "b", "c", ...))
 'a, b, c' 순으로 변수를 배열하고, levels를 변수형으로 바꾸고 해당 데이터의 변수에 지정

데이터 <- 데이터 %>% rename(바꿀변수명 = 기존변수명) : 변수명 변경하기
 ↳ * 우선 library(dplyr) 먼저 하기 ↑ " " 안 써도 됨

데이터\$변수 <- ifelse(데이터\$변수 == "a", "b", 데이터\$변수)
 만약 데이터의 변수 값이 'a'이면 'b'로 바꾸고, 아니면 데이터의 원래 변수 값 그대로 두어라
 ↳ 같은 변수가 아니어도 OK

weather의 평균기압이 NA일 경우는

TS.NA(weather\$평균기압) 이지 weather\$평균기압 == NA는 응답

TS.NA(데이터\$변수): 해당변수가 결측치(NA)이다

데이터 전체의 행을 우선 (library(copyr) 복사하기)

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

⊕ var나 mean과 같은 기술통계량을 이용하기 위해서는 해당 사례들만

추출한 것들을 새로운 데이터에 assign해야 함

round(비율값, digits=k): 해당값에서 소수 k번째 자리까지

↳ 함수를 이용해 나가지 않고 바로 붙임

→ filter 함수 안에 다양한 여러 변수들을 기준으로 조건에 맞게 설정 가능

ex) exam %>% filter(class != 4, math >= 90 | history >= 95)

(4반이 아닌 학생 중 수학이 90점 이상이거나 역사 95점 이상)

데이터 %>% filter(변수 >= quantile(변수, probs = c(k))) : 데이터의 해당변수 값이

변수의 값 중에서 k인 값보다 이상인 사례만 추출 (즉 k가 0.9이면 0.9 이상인 값, 변수 상위 10%인 사례들만 추출 (변수의 값이 상위 1-k 인 것들만 추출))

데이터 %>% select(변수): 해당 변수들만 추출 (여러 변수를 ,로 연결해서 추출 가능)

데이터 %>% select(-변수): 해당 변수를 제외하고 추출

데이터 %>% select(contains("특정단어")) : 특정 단어가 포함된 변수 추출 (특정 단어를 면적으로 X 연산자 | 사용하기)

print(n=Inf): 밑에 행까지 다 보여주기

데이터 %>% arrange(변수): 변수에 대해서 오름차순 정렬

데이터 %>% arrange(변수): 변수에 대해서 내림차순 정렬 (desc) arrange(class, math)

↳ exam\$math의 형태가 배열이라 math, math 변수 그대로

summarise 쓸 때 변수명 지정해야 arrange 이용 가능

데이터 <- 데이터 %>% mutate(새변수 = []): []라는 과정을 거친

새변수를 데이터에 지정, mutate 함수로 동시에 여러 변수 생성 가능, 이 과정에서

주로 ifelse 함수가 쓰이는데 ifelse 안에 또 ifelse가 들어갈 있음

(변수가 1개로 구분되면 ifelse를 (if)로 사용)

반별로 오름차순 정렬한 것들에서 수학점수 내림차순해서 정렬하기

데이터 <- 데이터 %>% mutate(새변수 = case_when(기준변수 < k ~ " ", 기준변수 > k ~ " "))

ex) exam <- exam %>% mutate(test = case_when(total < 180 ~ "fail", total >= 180 ~ "pass"))

* 변수의 특정 값에 NA가 있다면 TRUE 조건이 아닌 total >= 180 ~ "pass" 줄이 새로 조건 이용해서

데이터 <- 데이터 %>% relocate(변수, before(혹은 after)=변수2)

: 변수를 변수2 왼쪽아 오른쪽에 이동시키기

* ① before: 변수를 변수2 왼쪽에(앞)

② after: 변수를 변수2 오른쪽에(뒤)

데이터 <- 데이터 %>% relocate(where(is.character 또는 is.factor))

: 문자형과도 같은 방위형과도 안 앞으로 이동하기

데이터 <- 데이터 %>% relocate(where(is.factor), before=where(is.character))

: 방위형과도 변수를 문자형과도 변수 앞으로 이동하기

데이터 %>% grab(b/(변수) %>% summarise(mean or Var or sd ...))

: 변수별로 그룹을 나누고 각 통계량을 조건에 맞게 제시

* summarise를 보여주는 결과값 앞에 변수명 지정 가능 ex) mean-math = mean(math)

mutate는 summarise를 이용하든 count = n() 빈도이용한 코드 사용 가능

변형 함수 코드

count/sum(count): 이미 기존에 변수로 그룹화하고 나누어진 사례들끼리의 전체 갯수에서 해당사례 비율

count/sum(데이터\$count): 모든 사례들(전체) 갯수에서 해당사례 비율

n-distinct(데이터\$변수): 변수의 정도와는 무관하게 중복되는 값이 아닌 고유값이 몇 개 존재하는가
(변수에 해당하는 사례가 중복되지 않고 몇 개나 있나?)

c(a:b): a부터 b까지의 연속된 숫자들의 배열

ex) c(1:7) 1 2 3 4 5 6 7

relocate(변수): 데이터상에서 변수가 제일 앞으로

ex) relocate(ID): 데이터상에서 'ID' 라는 변수가 제일 앞으로

데이터 <- left_join(데이터, 데이터2, by="변수"): 데이터와 데이터2를 데이터에
공동 변수를 이용해서 합침

데이터 <- left_join(데이터, 데이터2, by=c("변수1" = "변수2"))

: 공동 변수의 이름이 다를 때 변수1의 이름과 변수2의 이름이 같다고 지정해주기 데이터에 합치기

데이터 <- bind - rows (데이터1, 데이터2) : 데이터와 데이터2를 합쳐서 살펴보기

단, 변수명이 같아도 적어도 다른 2가지 방법 (처리가 동일해야 함) / 모든 변수가 동일
내용상 동일한 변수라도 변수명이 다른 다른 변수로 인식 / 모든 변수가 동일

- ⇒ 1) 변수명 일치하게 변경 : rename 함수
2) 새로 일체하게 변경 : as.factor 함수

데이터 <- 데이터 %>% distinct (변수, keep.all=T) : 해당 변수가 같은

사례들 중에 한개만 남기고 다 제거하기

ex) ID가 중복될 때 동일한 사례인지 확인

exam %>% group_by (ID) %>% summarise (count = n()) %>%
arrange (-count)

결과가 ID가 30인 것이 count가 2가 나온 것 동일했기 때문
distinct 써서 제거해야 함

데이터[a:b] : 데이터에서 a~b까지의 변수들

ex) exam[5:8] : exam에서 5번째 ~ 8번째의 변수들 (5열에서 8열까지)

데이터\$mean ± 2.57583 * 데이터\$sd : 이 기준을 넘어서면 이상치 처리

<+정형화>

① 데이터 프레임 만들기

데이터파일 읽기 <- read.csv ("file.csv")

② 보기 좋게 변수위치 조정 : 행순장 → 주자형 → 순자형

0 데이터 <- 데이터 %>% relocate 함수 이용

③ 이상치 검토 및 빈도 확인

descri <- describe (Htest[6:10])

descri <- descri %>% mutate (LL = mean - 2 * sd, UL = mean + 2 * sd)

④ 이상치를 제외한 데이터 프레임 생성

ex) ftest_new <- ftest %>% filter (빈도 ≤ k
UL 값)



1) 가설설정

$H_0: M_{Ho} = M_{Cs}$ (귀무가설)
 $H_a: M_{Ho} \neq M_{Cs}$ (대립가설)] 양측검정

2) 집단별로 데이터 프레임 만들기

ex) $Hest_Ho \leftarrow Hest_new \%>\% \text{filter}(\text{customer} == "Home\ office")$
 $Hest_Cs \leftarrow Hest_new \%>\% \text{filter}(\text{customer} == "Consumer")$

3) 정규성 조건검토 (정규분포의 형태를 따는지)

summary, $htst$ 함수이용해서 파악하기

ex) $\text{summary}(Hest_ho\$sales)$

$htst(Hest_ho\$sales, breaks = seq(0, 9000, 50))$

=> 두 데이터 다 왼쪽으로 치우쳐진 히스토그램 모양이라 정규분포가 아님

$\text{shapiro.test}(\text{데이터\$변수})$: 데이터의 해당하는 변수에 대해 P-value 구하기

$P\text{-value} > 0.05$ 즉 P-value가 유의하지 않아 정규성 조건을 만족하는데

위 함수를 이용하면 $P\text{-value} < 0.05$ (α) 이기에 유의해서 정규성 조건에 부합X

=> 정규성 조건에 부합하게 정렬변수인 sales를 자연로그로 변환

ex) $Hest_ho \leftarrow Hest_ho \%>\% \text{mutate}(\lnsales = \log(sales))$

이후 $htst$, shapiro.test 이용하여 정규성 조건검토

이제 P-value 값이 여전히 예가까워서 ok

4) 등분산성 조건검토

var.test 함수이용해서 $P\text{-value} \leq 0.05$ 면 이분산 가정 $P\text{-value} \geq 0.05$ 면 등분산 가정

$\text{var.test}(Hest_ho\$lnsales, Hest_Cs\$lnsales)$

• $F = 1.3083$ 으로 예가까워서 표본분산이 비슷하고 $P\text{-value} = 0.5536 > 0.05$ α 이니 등분산

5) 독립표본 t-검정 실시 및 가설검정

$t.test(Hest_ho\$lnsales, Hest_Cs\$lnsales, alternative = "two.sided", var.equal = F)$

만약 4)과정에서 이분산이면 F

계산하면 $P\text{-value} = 0.2601$ 이라 0.05 (α) 보다 크고

$\text{mean of } X (Hest_ho\$lnsales) = 5.976144$

$\text{mean of } Y (Hest_Cs\$lnsales) = 5.915587$

이니 표본평균이 비슷하고 P-value도 유의하지 않고 $H_0 = M_{Ho} = M_{Cs}$ 의 귀무가설 채택

< 24분 > 0.05

one-way ANOVA - 일원분산분석

1 단계: 가설 수립

0 이상치 검토 anova-new 데이터 프레임 만들기

library(CPsych)

descr <- describe(anova | \$Prtce)

descr <- descr %>% mutate(UL = mean + 2 * sd, LL = mean - 2 * sd)

anova-new <- anova %>% filter(Prtce <= descr\$UL)

< 기말형 패키지 >

library(Creadr)

library(dplyr)

library(CPsych)

library(forcats)

library(CASIT.colae)

library(dunn.test)

library(Ccar)

library(CHH)

→ 한개변이하면 R 코드나 차분하기

0 이상치 검토 Prtce 만들기 → CHH와 High 포함

install.packages("forcats")

library(forcats)

anova-new <- anova-new %>% mutate(Prtor = fct - collapse(Prtor, "High" = c("CHH", "High")))

0 두가지 가설 수립

H0: 4가지 집단의 Prtce 평균이 같다

Ha: 적어도 한 집단의 Prtce 평균은 다른 집단과 다르다

2 단계: 집단간 데이터 프레임 생성하기

anova-H <- anova-new %>% filter(Prtor == "High")

* 각 Prtce의 변수 4개 간의 프레임 중첩해 생성 *

3 단계: 종속변 정형성 검토

summary(anova-H\$Prtce)

hist(anova-H\$Prtce, breaks = seq(0, 600, 10))

shapiro.test(anova-H\$Prtce)

→ P-value 가 0 이하가 아니라서 정형성 조건에 부합하지 않지만

n의 크기가 많아서 조건에 만족할 것이라고 가정하고 추후 분석 진행

* 모든 데이터 프레임에서 다 확인 *

4 단계: 등분산성 검토

install.packages("car")

library(Ccar)

leveneTest(Prtce ~ Prtor, data = anova-new)

종속변 독립변

전체 데이터 프레임

~~P(F > F) = 0.107~~이라 d인 0.05보다

커서 등분산 조건을 만족해 버림 (유의성) 없음



5단계: ANOVA 설정 및 (독립성조건만족할때) / 양 P-value와 α 값 비교 (F검정)

`anova <- result <- aov(PITce ~ PTO, data = anova2_new)`
`summary(anova <- result)`
 P-value가 0.193이니 α 인 0.05가 크고
 P-value가 유의하지 않고 귀무가설을 채택
 즉 4가지의 관천의 PITce의 평균이 모두 같다.

만약 귀무가설이 아니라 실험 2처럼 대립가설 채택이면 사후검정 실시 (독립성조건에서 대립가설 채택일때)

6단계: 사후검정 Duncan test 실시
`install.packages("duncan")`
`library(duncan)`
`duncan.test(anova2_result, "Payment", console=T)`
 aov함수 이용 예이더 독립변수

- 평균이 큰 관천의 관천과 같은 group으로 나뉘는 관천끼리의 평균은 같다고 가정

• 신용카드 expense 평균 > 계좌이체 expense 평균 = 관천별 expense 평균 %

만약 α 가 0.05이면 5단계에서 독립성조건이 이분산조건을 실시해야함

`oneway.test(expense ~ payment, data = anova2_new)`
 • P-value가 0에 가까워서 대립가설 채택

이분산조건 사후검정 Duncan test 실시 (이분산조건에서 대립가설 채택할때)

`install.packages("dunn.test")`
`library(dunn.test)`
`dunn.test(anova2_new$expense, anova2_new$payment, method="bonferroni")`
 데이터\$종속변수 데이터\$독립변수

- 관편검제와 계좌이체의 P-value가 0이다. $\alpha/2$ 보다 작아서 유의하지 않고 P-평균이 동일하다.
- 관편검제와 신용카드의 P-value가 0에 가까워서 $\alpha/2$ 보다 작아서 유의하고 P-평균 동일하지 않음에 신용카드가 더 크다 (관편검제 수정된 신용카드와 대응할때 -11.301456 이니)

계좌이체 = 관편검제 (계좌이체 - 신용카드 = -11.301456)
 신용카드 > 관편검제 L + 통계량 값
 신용카드 > 계좌이체 %

대하는 P-value와 $\alpha/2$ (=0.025)와 비교

Two-way ANOVA 이론 공부

1단계: 가설설정

가설설정

H_0 : 두 독립변수 간에 상호작용 효과가 없다.

H_a : 두 독립변수 간에 상호작용 효과가 있다.

이상치 검토 및 제거

library(psych)

descr <- describe(two-anova\$expense)

descr <- descr %>% mutate(VL = mean + 3 * sd + 1, LL = mean - 3 * sd)

two-anova-new <- two-anova %>% filter(expense <= descr\$VL)

2단계: 서브데이터 프레임 만들기

two-anova-male <- two-anova-new %>% filter(gender == "Male")

* 각 gender의 다른 Female도 같이 만들어주기 *

3단계: 정규성 검토

summary(two-anova-male\$expense)

hist(two-anova-male\$expense, breaks = seq(0, 2000, 40))

shapiro.test(two-anova-male\$expense)

→ P-value가 0.001238이라서 정규성 가설의

반박 증거가 충분하므로 기각 (n) 2개로 진행

4단계: 등분산성 검토

library(car)

leveneTest(expense ~ gender, data = two-anova-new)

→ P-value가 0.001238이라서 0.05보다 작아서 유의함

등분산 가설을 기각하지 못함

5단계: 이분산 가정 one way anova 시행

oneway.test(expense ~ gender, data = two-anova-new)

→ P-value가 0.001238이라서 귀무가설을 기각하고

대립가설 채택 ~~→~~ gender에 따라 집단을 두개로 구분했을 때 ~~중간~~ expense의

평균이 차이가 있다 → $m_{female} > m_{male}$ (553 > 359) < 1

two-anova-new %>% group-by(gender) %>% summarise(mean = mean(expense))



7단계: two-way ANOVA 수행 및 그래프

1) two_anova_result <- aov(expense ~ gender * OS, data = two_anova_new)

summary(two_anova_result)

→ gender의 P-value가 1.38e-12... 이므로 0보다 작고 유의하고 (5단계)
 gender * OS P-value가 0.0003172... 이므로 0보다 작고 유의하기 (α=0.05보다 작고)
 (H_a: 두 독립변수 간에 상호작용 효과가 있다) (패러미터 추정)

→ 즉 IV1과 IV2의 곱에 따른 P-value가 0보다 유의해야 하는 것이 조건

2) two_anova_new\$gender <- factor(two_anova_new\$gender, levels = c("Male", "Female"))

→ one-way ANOVA에서 표본평균이 작은 순서대로 집단을 먼저 출력해줌
 그래서 작은 Male이 먼저 출력되게 순서 바꾸기

install.packages("HH")

library(HH)

Interaction2wt(expense ~ gender * OS, data = two_anova_new)

→ OS와 무관하게 Male과 Female 간에 expense의 평균값이 차이가 있다.
 → OS가 IOS일 때 Male에서 Female로 갈 때 expense의 평균값이 더 커진다 (각을 더 곱해서)
 → OS가 Indroid일 때 Male에서 Female로 갈 때 expense의 평균값이 작아진다 (각을 더 곱해서)
 → 상호작용의 정도 차이

8단계: 추가 분석 집단 내 비교를 위해선 세분화한 one-way ANOVA

*** 집단 내 비교 MA, MT, FA, FT (별명, 여성 / Indroid, IOS)

two_anova_new <- two_anova_new %>% mutate(genderOS = ifelse(gender == "Male" & OS == "Android", "MA", ifelse(gender == "Male" & OS == "IOS", "MT", ifelse(gender == "Female" & OS == "Android", "FA", "FT"))))

*** mutate와 ifelse 함수를 써서 새로운 변수 genderOS만 만들어서 비교의 집단을 만든다
 → 이후 one-way ANOVA 1단계 ~ 6단계 수행

원본 변수와 야원변수에서 패러미터, 키가 같을 때 어떤게 채택인지 확인하는 것은 P-value와 α(표준) 간에 비교하는 것이고, 이후 얼마만큼 차이가 있는지는 사후분석에서 오차 P-value VS α이거나 P-value VS α/2 사이를 비교하면 된다.

비교와 독립성검정
비율사이비대항가설검정

① 데이터제일 불러오기

library(readr)
 propTest <- read_csv("propTest.csv")

② 빈도 기준 교차표안출기

table(propTest\$customer == 20 & propTest\$trans == "Yes") = 48

* customer과 trans의 모든 범주의 조합끼리 대응해주기

	Yes	No
20대	48	49 (97)
30대	30	73 (103)

prop <- matrix(c(48, 49, 30, 73), nrow=2, ncol=2, byrow=T)
 행의갯수 열의갯수 행을 기준으로 4칸

rownames(prop) <- c(20, 30)
 colnames(prop) <- c("Yes", "No")
 ↳ 각 행과 열의 이름을 대서 지정해주기

③ 비율 기준 교차표안출기

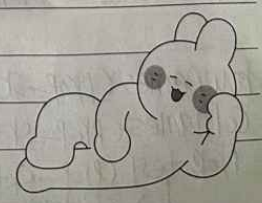
prop.table(prop, margin=1)
 ↳ prop이라는 빈도표에있있는것을 비율로 바꾸기 위해 행이라 값의 행이 1로 해야함

④ 가설검정
 $H_0: p_{20} - p_{30} = 0$
 $H_A: p_{20} - p_{30} \neq 0$

$p_{20} - p_{30}$ 이 정규분포를 근사한 조건인 $np \& n(1-p) \geq 5$ 인지 확인을 우선해야함

- $97 \times 0.495 = 48.015 > 5$
- $97 \times 0.505 = 48.985 > 5$
- $103 \times 0.29 = 29.87 > 5$
- $103 \times 0.71 = 73.027 > 5$

모든 범위가 5보다크니 $p_{20} - p_{30}$ 은 정규분포를 근사



Prop.test(Prop, alternative="two.sided", correct=T)
 P-value의 값이 0.005031 인데 <인 0.05보다 작으니 즉 유의해따 대립가설 채택
 (양방향임) $P_{20} > P_{30}$ (대립가설 채택) $Prop1: 0.14948454 > Prop2: 0.291262$

다항도표의 적합성검정

① 데이터프레임 불러오기 및 각각 과거, 현재에서 빈도수 check

library(readr)

telecom <- read_csv("telecom.csv")

table(telecom\$Past) = skt: 120 / kt: 100 / Lt: 80

table(telecom\$Current) = skt: 149 / kt: 85 / Lt: 66

② 첫번째 가설검정: 과거시점에 3개의 시장 점유율이 동일하지

chisq.test(c(120, 100, 80))

각 빈도를 넣어

P-value = 0.01832로 <인 0.05보다 작게 나왔으니 대립가설 채택 즉
 과거의 이동통신 3사의 시장 점유율은 서로 같지 않다.

③ 두번째 가설검정: 현재시점에 3개의 시장 점유율이 동일하지

chisq.test(c(149, 85, 66))

P-value = 6.06e-09 이므로 <인 0.05보다 작게 나왔으니
 마찬가지로 현재의 이동통신 3사의 시장 점유율은 서로 같지 않다.

④ 세번째 가설검정: 과거 시장 점유율과 현재 시장 점유율의 유무동일한지 여부

chisq.test(c(149, 85, 66), P=c(0.4, 1/3, 4/15))

현재 빈도수

과거시점 해당통신사 비율
 과거시점 전체 개수

OK) $\frac{100}{300}$ (KT 통신사 비율)

P-value = 0.002868 이라 <인 0.05보다 작게 나왔으니 대립가설 채택하고
 과거에 대해서 이동통신 3사의 시장 점유율은 변화가 있다.

⑤ 세번째 가설검정에 대한 사후분석

Prop_skt <- matFX(c(120, 149, 150), nrow=2, ncol=2, byrow=T)

skt가 아닌 나머지 통신사들의 과거, 현재의 개수

rownames(Prop_skt) <- c("Past", "Current")

colnames(Prop_skt) <- c("skt", "not skt")

Prop.test(Prop_skt, alternative="two.sided", correct=T)

과거보다 시장 점유율이 더
 많은 KT, SK, LG

이제 1과 LG+도 똑같은 방식으로 해주면 되는데 Prop-kt는 2p-value=0.259
 즉 α 인 0.05보다 크기 때문에 유의하지 않고 kt의 시정점효용은 과거와 현재 비교해 변화가 없음
 Prop-LG+도 2p-value가 0.216이라 α 인 0.05보다 크기 때문에 유의하지 않고
 마찬가지로 LG+의 시정점효용은 과거와 현재 비교해 변화가 없음.

독립성검정

① 데이터 불러오기

library(readr)

tennis <- read_csv("tennis.csv")

② 데이터 전처리

1) names(tennis) <- tolower(names(tennis)): tennis 데이터의 변수들의 이름을 다 소문자로
 [변수들만 바꾸는 함수] [소문자 변경 함수]

2) tennis\$surface <- tolower(tennis\$surface): tennis의 변수의 속성값을 소문자로
 * 독립변수 surface와 마찬가지로 종속변수 result에도 해주기 *

3) surface의 속성값 clay와 clay(t)를 합치고 hard와 hard(t)를 합치기

library(forcats)

tennis\$surface <- fct_collapse(tennis\$surface, "clay" = c("clay", "clay(t)"))

library(forcats)

library(dplyr)

tennis <- tennis %>% mutate(surface = fct_collapse(surface, "clay" = c("clay", "clay(t)")))

* 변수의 속성값 hard, hard(t)도 같은 방식으로 합치기 / 두 방법 다 가능 (a, b 중에 골라 써요)

만약 surface의 속성값 중에서 결측치가 존재하면 filter 이용해서 제거하고 새로운 데이터 만들기

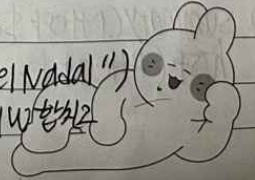
tennis_new <- tennis %>% filter(!is.na(surface))

4) player가 Rafael Nadal 인 것만 추출해서 데이터 만들기

tennis_nadal <- tennis_new %>% filter(player == "Rafael Nadal")

5) 4, 3)의 코드를 이용해서 Nadal의 result의 속성값 중에서 l, tr, w를 합치기

w, wr, ww를 합치기



③ 독립성검정

$X_{\text{tabs}}(\sim \text{surface} + \text{result}, \text{tennis} = \text{nadal})$
 $\text{prop.table}(X_{\text{tabs}}(\sim \text{surface} + \text{result}, \text{tennis} = \text{nadal}), \text{margin} = 1)$
 $\text{chisq.test}(X_{\text{tabs}}(\sim \text{surface} + \text{result}, \text{tennis} = \text{nadal}))$
 - P-value가 2.47e-09 이므로 사설 0이고 surface와 result는 독립하지 않다고 하는 대립가설 채택 즉, surface가 result에 영향을 미침, chi2코드로서 응용이 더 좋다
 - prop.table에서 확인

대응표본 t-검정 (독립변수 없이 두 종속변수끼리의 평균의 차이)

① step 1: 가설설정

$H_0: M_m - M_w(\text{md}) = 0$
 $H_a: M_m - M_w \neq 0$

② step 2: 파일 불러오기 및 사이변수 만들기

library(readr)
 $\text{pttest} \leftarrow \text{read_csv}("pttest.csv", \text{locale} = \text{locale}("ko", \text{encoding} = "euc-kr"))$
 - 데이터에 한글이 있으면 이항나코드를 사용

library(dplyr)

$\text{pttest} \leftarrow \text{pttest} \%>\% \text{mutate}(\text{d} = \text{morning} - \text{weekend})$
 - 두 종속변수의 차이에 대한 변수(d) 생성

③ step 3: d에 대한 정규성검토

$\text{shapiro.test}(\text{pttest} \$ \text{d})$
 - P-value가 9.70e-12라서 0에 가까워서 유의해서 정규성조건에 만족하지 않겠지만 앞의 양의 값에서 중심극한정리에서 두 종속변수가 정규분포를 따기에 표본평균의 차이로 중요성을 없다고 생각하고 넘기기

④ summary(pttest \$d)

$\text{hist}(\text{pttest} \$ \text{d}, \text{breaks} = \text{seq}(-15, 8, 1))$

step 4: 대응표본 t검정을 통한 가설검정

• test (pttest \$ morning, pttest \$ weekend, alternative = "two.sided",
p.value = 0.05)

↳ 대응표본 t검정이다

• pttest에서 morning의 수 < weekend의 수였고

2. p-value가 $2.2e-16$ 보다 작으니 0에 가깝고 α 인 0.05보다 작으니까 유의하고.

두 행동편의의 평균 차이가 0이라는 대립가설 채택. 즉 (한달동안 주말배동우운평균 > 새벽배동우운평균) 이런 차이는 통계적으로 유의한 차이이다.