

## 정제본 | 함수 총정리

음차 chr  
변수 factor

ok

데이터 파일 읽음  $\leftarrow$  read.csv("file.csv", (col.names=T), locale=locale("ko", encoding="euc-kr")) : file를 데이터 파일명으로 불러와서 저장하기

주요 함수 library(readr) 작성하기

view(데이터) : 데이터 보기

str(데이터) : 데이터 또는 데이터의 한 행에 대해 다양한 정보 보여줌

summary(데이터) : 해당 행의 모든 행에 대한 개략적인 정보 보여줌, NA가 어떤 행에 있는지 알려줌

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

glimPse(데이터) : 데이터의 다양한 정보 보여줌 (상세한 정보)

주요 함수 library(dplyr) 작성하기

table(데이터 \$ 변수) : 해당 변수의 빈도

freq(데이터 \$ 변수) : 해당 변수의 빈도와 막대 그래프

주요 함수 library(descr) 작성하기

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

plot(data=데이터, 변수1, freq=변수2) : 변수1과 변수2의 조합된 빈도를 막대 그래프로 표현

주요 함수 library(ggplot2) 작성하기 (생성권)

• 변수명은 " " 안씀

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 한 행에서 평균치(NA)를 제외하고 계산

freq(ts, na(데이터 \$ 변수)) : 데이터의 변수에서 결측치가 있는지 여부를 구하고 있으면 빈도 T로 표현

주요 함수 library(descr) 작성하기 / 즉 NA=결측치가 몇 개인지

describe(데이터 \$ 변수) : 해당 데이터 또는 변수의 해당하는 정량적 지표의 수치

(skewness, kurtosis, ...)

주요 함수 library(BSych) 작성하기

hist(데이터 \$ 변수, breaks=seq(M1, M2, K)) : 해당 변수를 M1~M2까지 구간을 나눠서 히스토그램

• 변수의 개수와 구간형 / M1, M2는 summary(데이터 \$ 변수) 통해 확인 가능

# <패러다임>

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

같다 ==  
달라 !=  
크다 >  
작다 <  
이상 >=  
아하 <=  
그리고 &  
또는 | (shift + #)

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

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

weekdays(데이터 \$ 변수) : 해당 변수(변수)의 요일이 무엇인지 보여줌  
↳ 첫도가 날짜 (date) 여야함  
데이터 \$ 변수 <- as.Date(데이터 \$ 변수) : 해당 변수의 첫도를 날짜 (date)로 변환  
ex) weekdays(as.Date("2020-01-01")) : 수요일

ex) weather \$ 요일 <- weekdays(weather \$ 날씨) : weather의 변수 '날씨'의  
요일을 원한 것들을 데이터 weather의 새로운 변수 '요일'에 저장

factor(데이터 \$ 변수, levels = c("a", "b", "c", ...))  
↳ a, b, c 순으로 변수를 배열하고 첫도를 변수명 바로 뒤에 해당 데이터의 변수에 지정  
factor 같은 함수 사용 가능

데이터 <- 데이터 %>% rename(변수명 = 기존 변수명) : 변수명 변경하기  
↳ \* 우선 library(dplyr) 먼저하기 / " " 안써도 됨  
데이터 \$ 변수 <- ifelse(데이터 \$ 변수 == "a", "b", 데이터 \$ 변수)  
↳ 만약 데이터의 변수 값이 a이면 b로 바꾸고 아니면 데이터의 원래 변수 값 그대로 두어야  
같은 변수가 아니어도 OK





데이터 ← 데이터 %>% relocate(변수, before(쪽은 1이라) = 변수2)

변수를 변수2 왼쪽과 오른쪽에 이동시키

- ①. before: 변수를 변수2 왼쪽에(앞)
- ②. after: 변수를 변수2 오른쪽에(뒤)

데이터 ← 데이터 %>% relocate(where (IS.character 쪽은 IS.factor))

: 문자형측도 쪽은 방형측도를 안알고 이동하기

데이터 ← 데이터 %>% relocate(where (IS.factor), before = where (IS.character))

: 방형측도 변수를 문자형측도 변수 앞으로 이동하기

데이터 %>% group\_by(변수) %>% summarise(mean or Var or sd ...)

: 변수별로 그룹을 나누고 기술통계량을 2번째 앞에 제시

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

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

변형 코드

count/sum(count): 이 변수에 대한 빈도 결과값이고 바뀌면 '사제들'의 전체 갯수에서 해당사제 비율

count/sum(데이터\$count): 모든사제들(변수) 갯수에서 해당사제비율

n\_distinct(데이터\$변수): 변수의 값으로는 무관하게 중복되는 값이 아닌 고유값이 몇개 중복되어 있는가 (변수값이 같은 사례가 중복되지 않고 몇개나 있는가?)

c(a:b): a부터 b까지의 연속된 숫자들의 배열

ex) c(1:7): 1~7

relocate(변수): 데이터상에서 변수가 제일앞으로

ex) relocate(ID): 데이터상에서 ID라는 변수가 제일앞으로

데이터 ← left\_join(데이터1, 데이터2, by="변수"), 데이터와 데이터2를 데이터1의 공통변수를 이용해서 합침

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

: 공통변수의 이름이 다를때 변수1의 이름과 변수2의 이름이 같다고 지정된 데이터 합치기



데이터 < bind - rows (데이터, 데이터2) : 데이터와 데이터2를 합쳐서 살펴볼 수  
 만. 변수명이 같아도 꼭 봐야 다른 소문자 (처리가 동일해야 함) / 모든 변수 동일  
 내용 동일할 때도 변수명이 다른 다른 변수일 수 / 모든 변수 동일  
 => 1) 변수명 일치하게 변경 : rename 함수  
 2) 식도 일치하게 변경 : as.factor 함수

데이터 <- 데이터 %>% distinct(변수, keep.all=T) : 해당 변수가 같은  
 시점을 제외 한 채만 남기고 다 제거하기

ex) ID가 중복한 일대 동일한 사례인지 확인

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

↓ 결과가 ID가 30일 것이 count가 2가 나옴 즉 동일한 것임  
 distinct 해서 제거함

데이터[a:b] : 데이터에서 a와 b까지의 변수들

ex) exam[5:8] : exam에서 5번째 ~ 8번째의 변수들 (5번째부터 끝까지)

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

## <+ 검증하기 >

### ① 데이터 프레임 확인

데이터파일 이름 <- read.csv("file.csv")

### ② 보기 좋게 병렬해서 표정 : 병렬성 → 지형 → 군지형

• 데이터 <- 데이터 %>% arrange 함수 이름

### ③ 이상치 정도 및 빈도 확인

descr <- describe(htest[6:10])

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

### ④ 이상치를 제외한 데이터 프레임 생성

ex) htest\_new <- htest %>% filter(변수 ≤ k  
 UL 값)



# 1) 가설 수립

$H_0: M_{Ho} = M_{Cs}$  (귀무가설)

$H_a: M_{Ho} \neq M_{Cs}$  (대립가설) 양측검정

# 2) 집합별 데이터 프레임 만들기

ex)  $Hest-Ho \leftarrow Hest-new \%>\% \text{filter}(customer == "Home office")$

$Hest-CS \leftarrow Hest-new \%>\% \text{filter}(customer == "Consumer")$

# 3) 정규성 조건 검토 (정규분포에 형태를 따는지)

summary, hist 함수를 통해서 관찰하기

ex) summary(Hest-ho\$ sales)

hist(Hest-ho\$ sales, breaks = seq(0, 9000, 50))

=> 두 데이터의 다윈쪽의 사분점의 히스토그램 모양이 대칭 정규분포가 아님

shapiro.test(데이터변수): 데이터의 해당변수 분포에 대해 P-value 검정하기

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

위 함수를 이용하면  $P\text{-value} < 0.05$  (x) 이기 때문에 정규성 조건에 부합X

=> 정규성 조건에 부합하지 않는 변수인 sales를 자연로그로 변환

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

이제 hist, Shapiro.test 이용하여 정규성 조건 검토

이때 P-value 값이 여전히 0.05 이하가 아님

# 4) 등분산성 조건 검토

var.test 함수를 이용해서  $P\text{-value} < 0.05$  면 이분산 + 검정  $P\text{-value} \geq 0.05$  면 등분산 + 검정

$\text{var.test}(Hest-ho\$lnsales, Hest-CS\$lnsales)$

$F = 1.3083$  이므로 이분산이 아니므로  $P\text{-value} = 0.5536 > 0.05$  (x) 이니 등분산

# 5) 독립표본 t-검정 실시 및 가설검정

$t\text{-test}(Hest-ho\$lnsales, Hest-CS\$lnsales, alternative = "two.sided")$

$\text{var.equal} = \text{FALSE}$

만약 4)와 정제해 이분산이면 F

계산하면  $P\text{-value} = 0.261$  이므로 0.05 (x) 보다 크고

mean of X (Hest-ho\$lnsales) = 5.976144

mean of Y (Hest-CS\$lnsales) = 5.915587

이러므로 등분산이 아니므로  $P\text{-value}$ 도 유의하지 않고  $H_0 = M_{Ho} = M_{Cs}$ 의 귀무가설 채택



24분 0.05

# 1. UNY ANOVA - 알뜰한 소비자

## 1 단계: 가설 설정

0 이상치 검토 anova-new 데이터 프레임 만들기

library(Psych)

desc<- describe(anova1\$PRICE)

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

anova.new<- anova1 %>% filter(PRICE <= desc\$UL)

<가설 검정 패키지 로드>

library(crandr)

library(car)

library(dplyr)

library(CHH)

library(Psych)

library(forcats)

library(ggfortify)

library(jorn.test)

→ 각 변수의 가설은 독립성이나 차이를 검정

## ANOVA를 prior 만들기 → categorical과 high 포함

install.packages("forcats")

library(forcats)

anova.new<- anova.new %>% mutate(prior=fct->discrete(prior))

"HIGH" = c("categorical", "HIGH"))

→ 연속변수

## 0 두 가지 가설 수립

H0: 4가지 브랜드의 PRICE의 평균이 같다

H1: 적어도 한 브랜드의 PRICE의 평균은 다른 브랜드와 다르다

## 2 단계: 간단한 데이터 프레임 생성하기

anova.H<- anova.new %>% filter(prior=="HIGH")

\* 각 prior의 변수 4개 간의 프레임들 생성하기

## 3 단계: 종속변수 가정성 검토

summary(anova.H\$PRICE)

hist(anova.H\$PRICE, breaks=seq(0, 600, 10))

shapiro.test(anova.H\$PRICE)

→ P-value가 0 이하가 아니라서 정규성 가정에 위반하지 않음

N의 개수가 많아서 2번째 연속성 가정에 위반하지 않고 가정하고 우도비교 검정

\* 모든 데이터 프레임에서 다 확인

## 4 단계: 등분산성 검토

install.packages("car")

library(car)

leveneTest(PRICE~prior, data=anova.new)

~~P(F) = 0.109~~ → 0.109가 0.05보다

커서 등분산 조건을 만족하므로 (가정) (가정)



종속변수

독립변수

5단계: ANOVA 검정 값 (독립성조기인 경우일 때)

$anova \leftarrow result \leftarrow aov(Payment \sim data = anova\_new)$

summary(anova\\_result)

P-value가 0.133이니  $\alpha=0.05$  작고

P-value가 유의하지 않고 귀무가설을 채택

즉 4가지의 관의 Price의 평균이 모두 같다.

관 귀무가설에 대한 5%수준의 독립성검정 (독립성조기인 경우일 때)

6단계: 사후검정 bonferroni test

install.packages("multcomp")

library(multcomp)

bonferroni.test(anova2\$result, "Payment", console=T)

anova2의 이름에 키를 입력

• P-value가 큰 관의 관하고 같은 group으로 나뉘는

관들의 P-value는 같다고 가정

•  $\text{관} \sim \text{expense P-value} > \text{관} \sim \text{expense P-value} = \text{관} \sim \text{expense P-value} \%$

관 data의 이름과 0.05이면 5단계에서 독립성검정 값이 유의하지 않음

oneway.test(expense ~ Payment, data = anova2\\_new)

• Payment가 0인 관의 관하고

독립성조기인 경우일 때 bonferroni.test (독립성조기인 경우일 때)

install.packages("dunn.test")

library(dunn.test)

dunn.test(anova2\\_new\$expense, anova2\\_new\$Payment, method="bonferroni")

데이터의 종속변수

데이터의 독립변수

• 관련관계와 회귀관계의 P-value가 1이다.  $\alpha/2$ 보다 작아서 유의하지 않고 P-value가 동일하다.

• 관련관계나 신용카드의 P-value가 0.133이니  $\alpha/2$ 보다 작아서 유의하고 P-value가 동일하지 않음에

신용카드가 더 크다 (관련관계 우측의 신용카드와 비교할 때  $-7.301456$ 이니)

회귀관계 = 관련관계

(회귀관계 - 신용카드  $-11.301456$ )

신용카드 > 관련관계

↑ 통계량 값

신용카드 > 회귀관계  $\%$

다행 P-value와  $\alpha/2 (=0.025)$  비교





문제: two-way ANOVA 수행 및 결과

1) two-anova-result ←  $anova(expense \sim gender \times OS, data = two\_anova\_new)$   
↳ 두 독립변수의 영향      전체 데이터

Summary(two-anova-result)

↳  $gender$ 의 P-value가  $1.38e-12$  ... 이므로  $ANOVA$  결과 유의하고 (5단계)  
 $gender \times OS$  P-value가  $0.000372$  ... 이므로  $ANOVA$  결과 유의하고 (1) ( $\alpha = 0.05$ 보다 작다)

( $H_0$ : 두 독립변수 간에 상호작용 효과가 없다) (대립가설 채택)

즉 IV1과 IV2의 조합에 따른 P-value를 모두 유의해야 하는 것이 조건 충족.

2) two-anova-new\$gender ← factor(two-anova-new\$gender, levels = c("Male", "Female"))

↳ (one-way ANOVA에서 표본평균이 작은 순서대로 집단을 먼저 출력해줌  
따라서 작은 Male이 먼저 출력하게 순서 바꾸기)

install.packages("HH")

library(HH)

InteractionPlot(expense ~ gender \* OS, data = two-anova-new)

- ↳ OS와 무관하게 Male과 Female 간에 expense의 평균 차이가 있다.
- ↳ OS가 IOS일 때 Male과 Female 간에 expense의 평균 차이가 더 크다 (가설이 맞다)
- ↳ OS가 Indroid일 때 Male과 Female 간에 expense의 평균 차이가 작다 (가설이 맞다)
- ↳ 상호작용의 정도 차이

문제: 후속분석 | 집단을 세개 이상에서 비교하기 위해 one-way ANOVA

집단 네개 MA, MT, FA, FT (성별, 여장 / 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 함수에서 새로운 변수 genderos만 만들어 비교하기 위한 것  
또한 one-way ANOVA 문제 2번 문제

원본 데이터와 이원분할 데이터에서 각각 카운트, 카우치에서 어떤 차이가 있는 것은  
P-value와  $\alpha$  (0.05) 간에 비교하는 것임. 이후 결과에 차이가 있는지는  
사후분석에서  $post\_hoc$  values VS  $\alpha$  이거나 P-value VS  $\alpha/2$  사이를 비교하면 된다.



이와 독립검정

이분변치에 대한 가설검정

## ① 데이터베이스 불러오기

library(readr)

proptest &lt;- read\_csv("proptest.csv")

## ② 빈도 기준 교차표 만들기

table(proptest\$customer == 20 &amp; proptest\$trans == "Yes") = 48

\* customer와 trans의 모든 변수의 조합에 대해 하기\*

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

prop &lt;- matrix(c(48, 49, 30, 73), nrow=2, ncol=2, byrow=T)

[행이 2개] [열이 2개] [행과 열의 조합이 4개]

rownames(prop) &lt;- c(20, 30)

colnames(prop) &lt;- c("Yes", "No")

→ 각 행과 열의 이름을 대어 지정해 주기

## ③ 비율 기준 교차표 만들기

prop.table(prop, margin=1)

[prop]은 빈도표에 있는 값을 비율로 바꾸는 것에 해당. 값의 범위가 0~1로 나와

## ④ 가설검정

 $H_0: p_{20} - p_{30} = 0$  $H_1: p_{20} - p_{30} \neq 0$  $p_{20} - p_{30}$ 의 점추정치를 2개 이상의 조건인  $|OP \times 11(HP)| \geq 5$ 에 확인을 위한 해야 함•  $97 \times 0.1495 = 14.51575$ •  $97 \times 0.1505 = 14.60575$ •  $103 \times 0.291 = 29.973$ •  $103 \times 0.709 = 73.027$ 모든 범위가 5보다 크니  $p_{20} - p_{30}$ 은 점추정치를 2한다

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

P-value의 값이 0.005031 인데 < 0.05 보다 작으니 즉 유의해마 대립가설 채택  
 (양측검정)  $P_{20} > P_{30}$  (대립가설 채택)  $P_{prop} 0.14948434 > P_{H0} 0.2912621$

### 다양도 지수 측정 방법

① 데이터프레임 중 변수의 인덱스 파악, 현재 데이터 check

library(readr)

telecom <- read\_csv("telecom.csv")

table(telecom\$past) = skt: 120 / kt: 100 / LG: 80

table(telecom\$current) = skt: 149 / kt: 85 / LG: 66

② 첫 번째 가설검정: 과거사형에 3개의 사망형률들이 동일한지

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

각 변수를 넣기

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

③ 두 번째 가설검정: 현재사형에 3개의 사망형률들이 동일한지

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

→ P-value = 6.8e-09 이므로 < 0.05 보다 작게 나왔으니 마찬가지로 현재의 이동통신 3사의 사망형률들은 서로 같지 않다.

④ 세 번째 가설검정: 과거사형률들과 현재사형률들이 모두 동일한지 여부

chisq.test(c(120, 100, 80), p=c(0.4, 1/3, 4/15))

과거 빈도수

과제정해설문제에  
과제시점 현재빈도

ex)  $\frac{100}{300}$  (KT동선 4 비율)

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

⑤ 세 번째 가설검정에 대한 세부보기

prop\_skt <- matrix(c(c(120, 149, 151), 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 사망형률이 더  
작은 통신사도 많았음



이제 kt와 LG+도 똑같은 방식으로 하려면 되는데 Prop-kt는 2P-value=0.259  
 즉 0.105보다 크기 때문에 유효하지 않고 kt의 시정값들은 2개의 현재값과 변화가 없음  
 Prop-LG+도 2P-value가 0.216이라 0.105보다 크기 때문에 유효하지 않고  
 마찬가지로 LG+의 시정값들은 2개의 현재 값과 변화가 없음.

## 동원장점

### 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 바꿔줌          ↳ 소문자 변경함수

2) tennis$surface <- tolower(tennis$surface): tennis의 변수의 속성값을 소문자로  

   ↳ 독립변수 surface와 마찬가지로 종속변수 result에도 해주어야함

3) surface의 속성값 clay와 clay(t)를 하나로 hard와 hard(t)를 하나로  

   ↳ library(formats)
   tennis$surface <- fct_collapse(surface, "clay" = c("clay", "clay(t)"))

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

↳ 변수의 속성값 hard, hard(t)도 같은 방식으로 한거임 / 두 방법 다 가능 (a,b 중에 골라쓰셈)
```

만약 surface의 속성값 중에서 정독해야 할 것이라면 filter 이용해서 제외하고 새로운 데이터 만들기

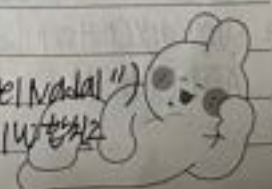
```
tennis_new <- tennis %>% filter(!is.na(surface))
```

### 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 뽑아

w, wr, ww를 뽑아



step  
t.test  
p-value

③ 독립성검정

$X_{tabs} \sim \text{surface} + \text{result}; \text{tennis} - \text{nadal}$   
 $\text{prop.table}(X_{tabs} \sim \text{surface} + \text{result}, \text{tennis} - \text{nadal}), \text{margin} = 1)$   
 $\text{chisq.test}(X_{tabs} \sim \text{surface} + \text{result}, \text{tennis} - \text{nadal})$   
 - p-value가 2.41e-09 이므로 사실상 0이고 surface과 result는 독립하지 않다고 하는 대립가설에서 즉, surface가 result에 영향을 미침, chisq 테스트에서 영향이 더 높다  
 -  $\text{prop.table}$ 에서 확인

대응표본 t-검정 (종속변수 없이 두 종속변수끼리의 평균차이)

① step 1: 가설설정

$H_0: \mu_m - \mu_{w(RM)} = 0$   
 $H_a: \mu_m - \mu_{w(RM)} \neq 0$

② step 2: 파일 불러오기 및 사이변수 만들기

$\text{library(readr)}$   
 $\text{Pttest} \leftarrow \text{read\_csv}("Pttest.csv", \text{locale} = \text{locale}("ko", \text{encoding} = "euc-kr"))$   
 - 데이터에 한글이 있으면 이항부호를 사용

$\text{library(dplyr)}$

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

③ step 3: d에 대한 정규성검토

$\text{shapiro.test}(Pttest\$d)$   
 - p-value가 0.806422라서 0에 가까우나 유의해서 정규성검토에 만족하지 않음  
 - 한계치인 한계값의 한계에서 중심극한정리에서 두 종속변수가 정규분포를 따기에 표준편차의 차이로 정규분포를 연다고 생각하고 넘기기

④ summary(Pttest\$d)

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



step 4: 대응표본 t-검정을 통한 가설검정

f. test (Phtest \$ morning, Phtest \$ weekend, alternative = "two.sided",  
p.value = T)

[대응표본 t-검정이다]

- Phtest에서 morning의 수 < weekend의 수였고  
2. p-value가  $2.2e-16$  보다 작으니 0에 가깝고 <인 0.05보다 작으니까 유의하고.  
두 행동별 수의 평균 차이가 0이 아니라는 대립가설 채택. 즉 (한달동안 주말과 평일 방문객 수)  
차이 발생 우연률) 이런 차이는 통계적으로 유의한 차이이다.