

R 기초

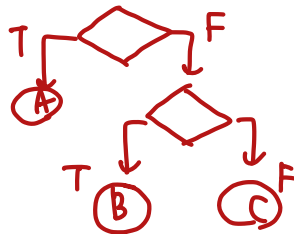
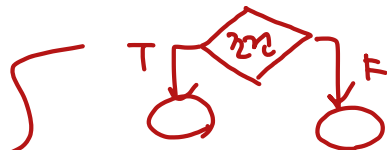
조건문과 반복문

조건문

특정조건만족 여부에 따라 별개의 작업을 진행

✓ R: if()와 ifelse()

- 조건이 하나: if
- 조건이 둘: if ~ else
- 조건이 여러 개: if A else if B (else if ...) ~ else C



✓ ifelse()는 조건결과가 벡터이면 각각의 결과에 대해 연산이 이루어짐

각연산개



반복문

SAS는 Do문

✧ 작업을 반복하기 위한 명령어

✓ R: for문, while문, repeat문

○ for: 반복수 지정

○ while: 특정 조건을 만족하는 동안 실행

○ repeat: 내부에서 중단되지 않은 동안 계속 실행

✓ 각 반복문의 구조 이해

📄 조건문 정리

✧ 특정 조건 만족여부에 따라 별개의 작업을 진행

✓ R: if()와 ifelse()

- 조건이 하나: if
- 조건이 둘: if ~ else
- 조건이 여러 개: if ~ else if ~ (else if ...) ~ else

✓ 주의점: if()와 ifelse()의 차이점

- 벡터 조건의 경우 if는 처음 결과만 ifelse는 각각의 결과에 대해 적용
- 논리연산에서 && vs &, || vs |와 유사

반복문 정리

✧ 작업을 반복하기 위한 명령어

✓ R: for문, while문, repeat문

- `for(x in vector) { }`
- `while(조건) { }`
- `repeat{ if (조건) break}`
- `{ }` 내에서 `break`를 통해 중단시킬 수 있음

R 기초

함수만들기와 기하평균추론

☞ 함수만들기

✧ 함수를 정의하는 방법

input값 ✓ 인수 인수의 설정과 인수의 적용범위

✓ 인수의 디폴트 지정방법

✓ 결과의 출력형태 지정



기하평균추론

복귀누적률, ...

◇ 기하평균: $\bar{x}_G = (x_1 \times \dots \times x_n)^{1/n} = \left(\prod_{i=1}^n x_i \right)^{1/n} = \exp\left(\frac{1}{n} \sum_{i=1}^n \log(x_i)\right)$

◇ 기하평균추론 **데이터가 반드시 양수여야함.**

✓ 분포가정: $Y_1 = \log(X_1), \dots, Y_n = \log(X_n) \sim iid N(\mu_Y, \sigma^2)$

✓ $\mu_G = \exp\left(E(\log(X))\right) = \exp(\mu_Y)$

✓ 기하평균 μ_G 의 신뢰구간: $(\exp(L), \exp(U))$

○ L : 로그 변환한 자료를 이용한 신뢰하한

○ U : 로그 변환한 자료를 이용한 신뢰상한

✓ 자료를 로그 변환 후 t.test() 적용하고 결과를 역변환

☞ 함수만들기

✧ 함수를 정의하는 방법: 함수이름 <- function(인수들)

✓ 인수의 설정과 인수의 적용범위

✓ 인수의 디폴트 지정방법: 인수 = 디폴트값

✓ 결과의 출력형태 지정

✓ 전역변수 지정: <<-

📋 기하평균추론

◇ 기하평균: $\bar{x}_G = (x_1 \times \cdots \times x_n)^{1/n} = \left(\prod_{i=1}^n x_i \right)^{1/n} = \exp\left(\frac{1}{n} \sum_{i=1}^n \log(x_i)\right)$

✓ $\min(x) \geq 0$ 이면 **`exp(mean(log(x)))`**

◇ 기하평균추론

✓ **자료에 음수가 있는지 확인**

✓ **자료를 로그 변환 후 `t.test()` 적용하고 결과를 역변환**

#####

조건문: if문

* if (condition) {statements1} else {statements2}

* else {statements2}는 생략가능하며 statements가 하나인 경우 {} 생략가능

* condition의 결과가 벡터이면 첫 번째 원소만 적용

조건이 여러개인 경우

* if (condition1) {

} else if (condition2) {

} else {}

벡터 조건문:

* ifelse(condition, statement1, statement2)

* condition의 결과가 벡터이면 개별 결과 적용

* 복습: 논리연산

x <- c(F,T); y <- T

x && y; x || y ⇒ F, T

x & y; x | y ⇒ FT, TT

x <- 1

if (x > 0) y <- "Yes" else y <- "No"

y ⇒ "Yes"

주의

if (x > 0) y <- "Yes"

else y <- "No" ⇒ Error : 다중문장으로 인식

y ⇒ "Yes"

x <- -1:1 (x: -1, 0, 1)

if (x > 0) y <- x+1 else y <- x-1

y ⇒ -2, -1, 0

Warning: the condition has length > 1 and only the first element will be used.

↳ 첫번째값(-1)만 사용 → 이미 false로 판정

z <- ifelse(x > 0, x+1, x-1)

z ⇒ -2, -1, 2

↳ 개별적으로 비교

ifelse(x %% 2 == 1, "odd", "even")

홀짝 <- ifelse(x %% 2, "odd", "even")

홀짝 ⇒ "odd" "even" "odd"

문자에 대한 출력도 가능

#####

반복문

for 루프

* for (name in vector) statement1

* for (name in vector) {statements}

* vector는 문자도 가능

여러개의 명령어
묶어주기 :

계속연결되게

```
if (x > 0) {  
  y <- "Yes"  
  a <- "A"  
} else {  
  y <- "No"  
  a <- "a"  
}
```

```
## while 루프
## * while(condition) statement
## * while(condition) {statements}
## * condition이 참이면 statement(s)를 실행
## * 무한 루프에 빠질 가능성이 있음
```

```
## repeat 루프
## * repeat{statements}
```

비주

```
## break로 중단
## * 무한 루프에 빠질 가능성
```

```
# 1~10까지 합
n <- 10
csum <- 0
for (x in 1:n)
  csum <- csum+x
csum
```

*→ sum(1:10)
⇒ 55*

```
# 1~n까지 짝수의 합
n <- 1000000
x <- 1:n
evensum <- 0
for (value in x)
  if (value %% 2 == 0) evensum <- evensum+value
evensum
```

반복문은 속도가 느림
비교작업도

```
sum(x[x%%2 == 0])
sum(x[!x%%2])
```

```
## 문자처리 가능
item = c("C","B","D","A","F")
for (x in c("A","B"))
  cat(x, " is ",which(x==item),"-th item.\n")
```

*⇒ A is 4 -th item.
B is 2 -th item.*

```
for (x in c("A","B","K"))
  cat(x, " is ",which(x==item),"-th item.\n")
```

*⇒ A is 4 -th item.
B is 2 -th item.
K is -th item. → 없는 경우에도 출력됨.*

```
for (x in c("A","B","K"))
  if (x %in% item) ## %in% :vector 비교 가능(문자, 숫자 비교가능)
    cat(x, " is ",which(x==item),"-th item.\n")
```

```
## 반복횟수가 정해지지 않은 경우: while, repeat
## while(조건): 조건이 참이면 계속수행
## repeat: 최소한 한번은 수행
```

```
for (x in c("A","B","K"))
{
  if (x %in% item)
    cat ~
  else
    ~
}
```

⇒ 실행됨.

3가지 방법

→ 3가지 방법 x

```
# -1.0보다 작은 값이 나올때까지 표준정규난수를 발생
x <- numeric() # x <- NULL (빈집합)
for (i in 1:1000) #충분히 큰 수만큼 반복
{
  난수 <- rnorm(1)
  x <- c(x, 난수)
  if (난수 < -1.0)
    break
}
x
```

```
난수 <- rnorm(1)
x <- 난수
while(난수 >= -1.0)
{
  난수 <- rnorm(1)
  x <- c(x, 난수)
}
x
```

```
x <- numeric()
repeat{ #일반작업부터. (의downhile)
  난수 <- rnorm(1)
  x <- c(x, 난수)
  if (난수 < -1.0)
    break
}
x
```

seed정하기 : set.seed(7)

#####

함수만들기

#####

x <- rnorm(30) # 정규난수 → ? rnorm 설명보기

mean(x) # 평균 계산

x

x[21] <- NA # 결측값 추가

mean(x) NA (∵ na.rm=FALSE 설정돼있음 : default)

mean(x, na.rm=T) # 결측값 제거 후 평균 계산

⊖ (mean(x, trim=0.5, na.rm=T) # 50%절사평균 (trimmed mean)

median(x, na.rm=T) # 중앙값 → 상위 50%, 하위 50%를 trim했으므로 중앙값과 같음

* 기하평균 $\exp(-\infty) = 0$, exponential 지수 : skewed data (오른쪽 꼬리가 길다) 
x <- exp(rnorm(30))
exp(mean(log(x)))
ex) 감마분포 (x, 지수, ...)

⊖ (product(x)^(1/30)) 기하평균 데이터의 개수가 많은 경우 Inf가 나와버림
처리할 수 있는 수의 범위가 있으므로 좋지 않은 방법

평균편차와 평균이 바뀌게 → log(x)
($\frac{\sigma}{\mu}$ 또는 $\frac{s}{\bar{x}}$ ⇒ 변동계수 CV)

지수
곱
* 함수만들기:
* "함수이름 <- function(인수들)" 형태로 정의
* 인수, default 지정
* 인수의 조건 만족여부 확인
- 기하평균: 자료는 0보다 커야함

minus1 <- function(x) x <- x-1 # 함수의 인수는 내부에서만 적용(지역변수)

minus1 <- function(x) {
 x <- x-1
 return(x)
}

x <- 1:5 ⇒ minus1(x) 실행해도 지는 바뀌지 않음

minus5 <- function(y) x-y # 함수의 인수가 아닌 변수는 앞에서 정의된 값 적용

minus5(5) ⇒ -4 → -2 → -1 0 → x-y 계산해서 출력하라는 의미

→ 결과반환
전역변수 값은 아직 바뀌지 않음.
출력할게 없으므로 return값이 x만 써도 됨.

minusA <- function(y) a+y # 변수 a는 앞에서 정의되지 않음
minusA(5) # error: 객체 a를 찾을 수 없습니다

geomean <- function(x)
exp(mean(log(x))) 기하평균 함수

minus1 <- function(x) {
 x <- x-1
 }
⇒ x를 전역변수로

geomean(x)

x <- exp(rnorm(20))

x[21] <- NA

geomean(x) → NA

점검할 것

geomean <- function(x) (① 양수인지 (원래는 0도 안되지만 R에서 $\exp(-\infty)=0$ 으로 계산)
② NA가 아닌지

{
 x <- x[!is.na(x)]

if (min(x) >= 0.0) 최솟값이 0 이상이면 모든 데이터가 0 이상이다.

return(c(length(x), exp(mean(log(x)))))

else

cat("자료 중 0보다 작은 값이 있어 기하평균을 계산할 수 없습니다.")

ex) $\text{geommean}(1) \Rightarrow \frac{29.00000}{1.08097}$
 다른결과값이 실패상태나 같이 실패상태로 출력됨.

```
geommean <- function(x)
{
  x <- x[!is.na(x)]
  if (min(x) >= 0.0)
  {
    Result <- list()
    Result$n <- length(x)
    Result$Geomean <- exp(mean(log(x)))
    Result
  }
  else
    cat("자료 중 0보다 작은 값이 있어 기하평균을 계산할 수 없습니다.")
}
```

$\text{geommean}(1) \Rightarrow$ \$n
 [1] 29
 \$Geomean
 [1] 1.08097

결과 <- geommean(x)
 (리스트와 똑같이 사용가능)
 결과\$Geomean 또는 결과[[2]]

7. t-test (Var.equal=FALSE가 default)

```
geommean.test <- function(x, mu=1.0, alternative="two.sided", conf.level=0.95)
{
  if (mu <= 0.0 | min(x, na.rm=T) <= 0.0)
    cat("모수 설정이 잘못되거나 자료에 0보다 작거나 같은 값이 있습니다.")
  else
  {
    n <- sum(!is.na(x))
    if (n < length(x))
      x <- x[!is.na(x)]
    logx <- log(x)
    logmu <- log(mu)
    result <- t.test(logx, mu=logmu, alternative=alternative, conf.level=conf.level)
    result$null.value <- exp(result$null.value)
    result$conf.int <- exp(result$conf.int)
    result$estimate <- exp(result$estimate)
    return(result)
  }
}
```

geommean.test(1)

```
> geommean.test(x)

One Sample t-test

data:  logx
t = 0.34103, df = 28, p-value = 0.7356
alternative hypothesis: true mean is not equal to 1
95 percent confidence interval:
 0.6771886 1.7255096
sample estimates:
mean of x
 1.08097
> |
```

```
> t.test(log(x))

One Sample t-test

data:  log(x)
t = 0.34103, df = 28, p-value = 0.7356
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 -0.3898054  0.5455224
sample estimates:
mean of x
0.07785851
```

```
> exp(-0.3898054)
[1] 0.6771886
> |
```