

STAT 3027

2.	Vector
----	--------

Q	When are [brackets] used and when are parentheses used?
	<div style="display: flex; justify-content: space-around; align-items: center;"><div style="text-align: center;">↓ index</div><div style="text-align: center;">↓ compare</div></div>

② which methods are there to remove noise?

③ What are the difference between \leq and $=$?

18 $V \leftarrow C(2.71, 5, 3.14)$
combination.

length (v)

V

index	value
0	1
1	1
2	1
3	1
4	1
5	1
6	1
7	1
8	1
9	1
10	1
11	1
12	1
13	1
14	1
15	1
16	1
17	1
18	1
19	1
20	1
21	1
22	1
23	1
24	1
25	1
26	1
27	1
28	1
29	1
30	1
31	1
32	1
33	1
34	1
35	1
36	1
37	1
38	1
39	1
40	1
41	1
42	1
43	1
44	1
45	1
46	1
47	1
48	1
49	1
50	1
51	1
52	1
53	1
54	1
55	1
56	1
57	1
58	1
59	1
60	1
61	1
62	1
63	1
64	1
65	1
66	1
67	1
68	1
69	1
70	1
71	1
72	1
73	1
74	1
75	1
76	1
77	1
78	1
79	1
80	1
81	1
82	1
83	1
84	1
85	1
86	1
87	1
88	1
89	1
90	1
91	1
92	1
93	1
94	1
95	1
96	1
97	1
98	1
99	1

i

v[i]

1

271

5

"

3.14

words ← C("tree", "and", "chairman")

words [17]

* Basic Types, Specifying Constants

numeric: $\sqrt{\quad}$ real number

$$3.14e2 = 3.14 \times 10^2$$

chorale: r

$$u = 1 \quad \text{dimension} \quad u = 3/4$$

Escape sequence.

✓ (double gate) ∴

String = parts (sep. = " ", " ", " ")

cod (string)

$\text{cat } (\text{sep} = "$

logical : TRUE and FALSE

TEST

Sum(V>3)

=> 2

Vector (mode = "logical", length = 0)

creates a vector of the given mode & length

logical(0)

= 5

FALSE ... FALSE

numeric

= 5

0 0 0 0 0

char

change a vector's type

as.numeric()

as.character()

as.logical()

integer

complex

raw

as.numeric(V>3)

W = c("34", "12", "45")

W

"34" "12" "45"

Sum(W)

↳ error

W.numbers = as.numeric(W)

W.numbers

34 12 45

sum(w, numbers)

91

sum(as.numeric(w))

91

* Name attribute and a few functions

name(x) gets or set.

name(v) = c("e", "five", "pi")

v

e five pi

2.71 5.00 3.14

name(v) = NULL

remove names

== y = c(burger = 2.50, fries = 1.50)

y

burger

fries

2.5

1.5

names(y) = NULL

2.5 1.5

A Few Functions!

x <- c(1, 11, 111)

(3.14159 = pi) * 100

* Operators

%%/%, quotient

%% remainder

$$\text{sqrt}(\text{sum}((x - \text{mean}(x))^2) / (n-1))$$

seq (10, 15, by=2)

c(10, 12, 14)

seq (10, 15, length.out=3)

10, 12.5, 15.0

matching: %in%

1:3 %in% c(2,7)

FALSE TRUE FALSE

★ Indexing

4.

x[x<4]

negative \Rightarrow exclude.

indices = which(x < 4)

index

is a

x[indices]

x[x < 4]



which can be omitted,

x[(x %>% 1) == 0]

3. Vector (continued) and List

SORTING Functions

sort(x, decreasing = FALSE)

x = c(12, 11, 16, 11)

\Rightarrow 16, 11, 12, 11

$v = \text{rank}(x, \text{ties} = \text{"average"})$

第几大

... "first"

$v = \text{order}(x, \dots, \text{decreasing} = \text{FALSE})$

the index of smallest elements in x .

$x[v]$ is sorted.

Structure summary, quantile

$\text{str}(\text{object})$

$\text{summary}(\text{object})$

quantile

$v = \text{quantile}(x, \text{prob} = c(0.25, 0.75))$

NA and NA, special values - $\text{is.na}(x)$

$\text{names}(x)$

$x[3] = \text{NA}$

$\text{sum}(x) = \text{NA}$

$\text{sum}(x, \text{na.rm} = \text{TRUE})$

$\text{is.null}()$

$\text{is.na}()$

$[1,]$

File input/output

$\text{scan}(\text{file} = \text{" "}, \text{what} = \text{numeric()})$

$\text{write}(x, \text{file} = \text{"data"})$

$[1](m)$

...

List: (if row = factor, not 'x'; row = y
not necessarily of the same type
can be different)

\$ operator = ...
y\$8 \Rightarrow null
list of [1] x
unlist(y, use.names = FALSE)

4. Data Framework

(R's fundamental data structure)

A data frame is \approx a list of vectors

Categorical Variables

Factor \Rightarrow a vector of categorical values
factor(x, levels, labels = levels)

table(...)

Manipulation Examples

m\$

m[, 1]

row column

m[1:3, 1:3]

dim(m)

n.rows = dim(m)[1]

n.cols = length(m)

head(m) = x > x\$head
 = del p "condition" = ...
 = m.p tail(m) ...
 rowname = (m) ...

m\$hp[30] = 25 ... = bottom ...

m = median(m\$hp)

★ sorted = m [order(m\$cy, m\$dis)]

write.table(x, file = "r12", as.csv = TRUE)

Formula

5. Graphics

1. Common parameters

formula, data

main, sub, xlab, ylab

xlim, ylim

pch

cex(1, 1, 1, 1, 1) = x

par(1, 1, 1, 1, 1) = x

"log" = 100

2. Numeric data

boxplot(x)

boxplot(x ~ g) # by group

boxplot(mtcars\$mpg, main="Gas mileage", ylab="mpg", ylim=c(0, 40))

boxplot(mpg ~ factor(cyl), data=mtcars, xlab="cyl", ylab="mpg")

stripchart(x, method="overplot")

stripchart(x ~ g)

overplot

jitter

> duplicate

hist(x, breaks="sturges", freq=FALSE)

density

★ cex.axis = 2

同时 cex.lab = 2

volume

plot(x, y)

}

4.17

points(x, y)

x = 1:5

y = 2 * x

plot(x, y)

? points

lines(x = c(1, 3, 5, 7, 9), y = c(8, 1, 4, 1, 8))

col = "red")

col = "red")

plot(density(x))

regression the plot points, which is a line

pairs (x)

curves (ex pr = $x \times \sin(1/x)$) from $-\pi/6$ to $\pi/6$, $n=200$

curve (ex pr = $x \times 1$, add = TRUE, col = "red")

Legends

legend(x, y, legend, col = position.val)

legend("top", legend = c(" $x \times \sin(1/x)$ ", "x"),
col = c("black", "red"), bty = "n")

expression(.)

? plot match

eng of fac(x, y)

stem (plot/20)

integrated data
bar plot (height, name s. org b = NULL)

Rgt Variable mosaicplot(x)

count = table()

barplot(counts, name, org = c(1, 2, 3))

spin Subplot labels

Multiple figures

matrix (data, nrow, ncol, byrow = FALSE)

layout(mat)

layout.show(3)

Write graphical output to a file

dev.off()

x = rnorm(100)

plot(x)

dev.off()

Test and Interval

One mean or the Difference of Two means

out = t.test(x, y = NULL, alternative = "two.sided")

mu = 0 conf.level = 0.95

\$p.value

\$statistic

\$conf.int

\$estimate

\$p.value

str(out)

out\$p.value

two means

$H_0: \mu_x - \mu_y = \mu_0 \Rightarrow \text{null}$ $\mu_0 = 0$ $\text{alt} = \text{two}$

F Test for Equality of Variances

out = var.test(x, y, ratio = 1, alternative = "two.sided", conf.level = .95)

$H_0: \frac{\sigma_x^2}{\sigma_y^2} = \text{ratio}$

$\frac{\sigma_x^2}{\sigma_y^2} = 1$ $2 - 19 - 19 = 0$

\$ parameter

$$(n_x - 1)s_x^2 + (n_y - 1)s_y^2$$

\$ statistic

$$f = \frac{s_x^2}{s_y^2}$$

\$ p.value

\$ conf.int

chi-Squared Tests

$$df = x + y - 1$$

Goodness of fit

counts = c(...) probs = c(...) (6, 2) 10

out = chisq.test(x = counts, p = probs)

\$ parameter

\$ statistic χ^2 (x) 100 = A

$$[(1, 1), (1, 2), (2, 1), (2, 2)] \times 100$$

Independence / Homogeneity

$$(n_{ij} - (n_{i.} \cdot n_{.j}) / n) \leq$$

\$ expected .. expected counts under H_0

chisq.test();

$$(n) \text{ rows}$$

17 statistic as χ^2

One Proportion or the Difference of Two Proportions

old = prop.test(x, n, p0, alternative = "two.sided",

2-test $H_0: p = p_0 = p$ conf.level = .95

alternative 1 - estimate $\hat{p} = \sum x_i / n$

2-vectors x and n

$$H_0: p_1 - p_2 = 0 \quad (p_1 = p_2)$$

$$\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

7. Regression

1. Simple Linear Regression

$$y = mx + b$$

cor(x, y) / correlation

lm(y ~ x, data)

$$A = \text{cor}(x) \Rightarrow \text{matrix } A$$

$$\text{cor}(x[,i], y[,j])$$

2. lm(y ~ x, data)

Call: lm(y ~ x, data)

anova(m)

coef(m)

abline (a, b) $y = a + bx$
 = abline (reg=m)

abline (a = mean (cov \$ dist)

horizontal

(1) horiz = n (10) horiz.to

(1) horiz = d (10) horiz.to

predict (model)

$y_{hat} = \text{predict}(\text{model}, \text{newdata} = \text{data})$

points (x = d \$ speed, y = y.hat, pch = 19, cex = 3)

plot (m \$ fitted.values, m \$ residuals)

abline (0, 0)

= , x = 1, y = x

(1) horiz = n (10) horiz.to

$\epsilon \rightarrow$ independent & normal

QQplot.

$N(\mu = \text{mean}(\text{residuals}), \sigma = \text{sd}(\text{residuals}))$

qqline (x)

e.g. $x = \text{rnorm}(n=100)$; $\text{qqnorm}(x)$; $\text{qqline}(x)$
 $w = \text{rexp}(100)$

or use plot (m)

multiple Linear Regression

$m = \text{lm}(y \sim x_1 + x_2 + \dots, \text{data} =)$

confint (m,

8. Simulation

①

repeatable

set.seed(seed)

set.seed(0); a = norm(1)

set.seed(0); b = norm(1)

homosced

(when together)

② Repeat calculation \wedge time

replicate(N, expr = a+b) = list

set.seed(0)

③ ? distribution

$\bar{X} \approx \sqrt{\frac{1}{N}} \left(\sum_{i=1}^N x_i \right)$

replicate(N, expr)

t = replicate(N, {x = norm(1, mean =

(x) = rep(0, N) ; (x) = rep(0, N)

(x) = rep(0, N) ; (x) = rep(0, N)

(x) = rep(0, N) ; (x) = rep(0, N)

symply text (x=ch, y=row, label= usechar, cex=4)

identity (xx, yy, n=1, plot=FALSE)

unique () %in%

replicate (,)

function
parity = ifelse ((X%%2) == 0, "even", "odd")

2. loop: while () { }

repeat { }

if () { break }

Loop one or more times: repeat { expression

if (condition) {

break

}

input: xx = scan (what = character, n=1, quiet=TRUE)

skip: if (c-z) next }

3. Apply Function lapply (X=, FUN=) treating data frame as a list of vectors

* sapply (X=, FUN=) ⇒ vector, matrix, array

multiple argument apply [FUN, x, y] take the several vectors in and applies FUN to all first elements

apply (X=, MARGIN, FUN, ...) ⇒ 1 rows

tapply (x=, Index=, FUN=...) ⇒ 2 columns

subset of vector

↳ tapply (X= matrix, INDEX = not car & cyl, FUN=) 分类!

tapply (X=, INDEX = list (mtcar, cyl, notcar & gear), FUN=)

4. Matrix m = matrix (data = vector, nrow =, ncol =, byrow = TRUE)

dimnames = list (C (), C (rowname)) probs = c (

cbind (matrix,) rbind (matrix,)

m [row (m), :col (m)] main diagonal

m [row (m) + col (m) == r - 1] anti-diagonal through (1, C)

m + reverse

A * B element-wise product A %*% B matrix product

solve (a = A, b = b) A %>% A %x% x ... matrix-vector product

5. pattern

grep (pattern, x, ignore.case = FALSE, value = FALSE) ⇒ index of elements of char

grep (pattern =, TRUE) ⇒ values

sub (pattern, replacement = x) ⇒ a copy of x after replacing 1st occurrence

of pattern with replacement

gsub () all

(TRUE ())

`gsub(pattern = "[aeiou]", replacement = "", x=a)` # Strip vowels
`gsub(pattern = "[^aeiou]", replacement = "", x=a)` # Strip non vowels
`^` matches the beginning of a line (`$` matches the end) e.g.
`grep(pattern = "^r", x=a)` `grep(pattern = "r", ignore.case = TRUE, x=a)`
`\\>` matches the beginning of a word (`\\<` end)

e.g. `grep(pattern = "e\\>", x=a)`
 repetition: `{n}` exactly n times / `{n,}` n or more times * `{0,1}`, `{1,}`
`{n,m}` n-m times inclusive `{0,1}` or "optional"

e.g. `grep(pattern = "\\d{4}$", x=a)` 4 digits, end of lines
`grep(pattern = "\\s\\d{4}$", x=a)` # space, 4 digits, end of line
`grep(pattern = "\\s\\d{4,5}$", x=a)` # space, 4 or 5 digits, end of line
`N` refers to what the Nth enclosed expression matched

`sub(pattern = "(\\w+)(\\w+)(\\w+)", replacement = "\\2, \\1, \\4, \\3", x=a)`
`|` means or `grep(pattern = "Joe | Jack", x=a)` `grep(pattern = "J(o|a)", x=a)`
`\\` \Rightarrow ? regex

Splitting strings `strsplit(x=a, split=",")` `strsplit(x=a, split=" ")` # split on space
`write(x=csv, file=" ")` `d=read.csv(" ")` header=FALSE, col.names=c()

6. ggplot `ggplot(data=, x=, y=, geom="boxplot")`
`ggplot1` `gg1 <- ggplot(data=trees, aes(x=, y=))`
`gg1 + geom_boxplot() + geom_jitter() + violin() + coord_flip() (flip)`

2. `+xlab(" ") + ylab(" ") + labs(title=" ")` + theme(plot.title = element_text(size=22, just="right"))
 3. `ggplot1` `aes(x, y, fill=x or y)` + `guides(fill=FALSE)`

4. reorder species by mean `ggplot(data, aes(x=reorder(x, y, fun=mean), y=,))`
 5. outlier shape `+geom_boxplot(outlier.shape=21)` 6. change to log scale `+scale_y_log10(labels=function(l) {paste(l, "in.")})`
 7. color brewer `-scale_fill_brewer(palette="Set1")`

`geom_density()` `geom_point()` `geom_smooth()` `geom_line()` `aes(color=factor(1))` `shape=factor(1)`
`+facet_grid(~)` `+scale_x_sqrt(limits=c(0,50))`

7. web scrap `TEAM <- "http:—" lines <- readLines(TEAM)`
`ten.lines <- grep(pattern=" ", x=lines, value=TRUE)` `team.names <- sub(x=ten.lines, pattern="") (*)`
`replace="=1")` `link <- paste0(, , ,)` `table <- readHTMLTable(link)`
`table[[3]][1,]` \rightarrow as.numeric(as.character())
Vector `sort()` `order` `x[order(x)]` `na.rm=TRUE` `list[]`
`write.table(x, file=" ", ...)` `table = read.table(...)`
`stopifnot(function() == " ")`
 turn a vector into a character \Rightarrow `paste(v, sep=" ", collapse=" ")`
 \Rightarrow `=data.frame(,)`