STAT 3027

2. Vector

① When are [brackets] used and where are (parentheses) used?  'index      + continue

② Which methods are there to rename names?

③ What are the difference between <-  '='  and ==?

A   V <- c( 2.71. 5. 3.14)
                combination.

        length (V)

        V

        index    value
          i       V[i]
          1       2.71
          2        5
          3       3.14

        words <-   c("tree", "art", "chainsaw")

        words[i]


☆ Basic Types.    Specifying Constants
    numeric (real number)        3.14e2  = 314 × 10²
    character                              = 314

        Escape sequence.
                \" (double quate)

    String = paste (sep = " ", " ", " ")

        cat ( string )
        cat (sep = " ", " ", " " )
        logical : TRUE and FALSE

sum (V>3)
=> 2

Vector ( mode = "logical", length =0 )
creates a vector of the given mode & length
(logical 6)

= 5

FALSE .... FALSE

numeric        =5

0  0 0 0 0
           char

Change a vector's type
    a(.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 (w1))

91

☆ Names 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!                    ε
    x <- c  (e, 1, ...)

    sort (x, decreasing = ...)
⚡ Operators ...
    % / % : quotient
    % % remainder

$$sqrt(sum((x - mean(x))^2)/(n-1))$$

seq (10, 15, by=2)

c(10, 12, 14)

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

10, 12.5, 15.0

matching: %in%

3 %in% c(2,7)

FALSE    TRUE  FALSE

☆ Indexing

4.      X[v].

        negative ⟹ exclude.

        indice = which (x < 14)

            index              问 OPP

    X[indices]

    X[ x < 14 ]

☆  which  can be omitted.

    X[(x %% 2) == 0]


3.  Vector (continued) and List

    SORTING Functions

    sort (x, decreasing = FALSE)

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

            ⟹  14, 11, 12, 16

v = rank (x, ties.method = "average")
第几大

                "first"

v = order (x, ..., decreasing = FALSE)
    the index of smallest elements in x.
        x[v] is sorted.


Structure, summary, quantile
    str (object)

        summary (object)

        quantile

    v = quantile (x, probs = c (0.25, 0.75))


NULL and NA  special values -

names (x)

x[3] = NA
        sum (x)    = NA
        sum (x, na.rm = TRUE)


    is.null()    is.na()    — fn
                            [1,] n


File input/output
    scan (file = " ", what = numeric(3))

    write (x, file = "data")
                [1](n)

List:   (`a p...w = factor ...9..t .x.) dror = y`

not necessarily of the same type

can be ...joined

$ operator = ... w.s...l   ...) w..10 =.

y $ 8 = null ...f...b.i ot

l.st...2 .i [v]X

unlist (y) usenames=FALSE)

4. Data  Framework

(R's fundamental data structure)

A data frame is ($\approx$) a list of vector

Catergorical Variables

Factor  =). a vector of catergorical values

factor (x, levels, labels = levels)

table (...)

Manipulation Examples

m$ _

m[ ,1]

row  colum

m [1:3, 1:3]

dim (m)

^.rows = dim (m)[1]

n.cols = length (m)

head(m)

box(m)

tail(m)

rownames(m)

m$hp[30] = 2

M = median(m$hp)

A :order = m [order (m$cyl, m$d+disp))]

write.table(x, file=",1,2) 

read.csv( )

Formula.

## 5. Graphics
1. Common parameters
   formula , data.
   main, sub, xlab, ylab
   xlim, ylim,
   pch.
   cex( )
   par
   col =

2. Numeric data

boxplot (x)          boxplot ( x ~ g )

box plot ( mtcars $ mpg, main = "Gas mileage", y lab = "
                                        (n) lim + ylim = c (0, 40) )
                    by group                        y轴范围

boxplot ( mpg ~ factor ( cyl ), data =      , xlab.
    ylab =              );


Stripchart (x, method =" overplot")
strip chart (x ~ g)
        overplot
        jitter.        > duplicate
        stack

hist (x, breaks = "Sturges", freq = NULL )
                                    density

☆ cex. axis = 2
相对字号 cex. lab = 2.


plot (x, y)              图 4.1]
points (x, y)
    x = 1:5;    y = 2 * x, plot(x, y)
                                    ~ | x
    ? points
    lines ( x = c (1, 3, 5, 7, 9)
            y = c (8, 1, 4, 1, 8)
            col = "red")

plot (density (x))                    ...

rug(...)... the ... points. ... ... ...

pairs (x)

curves (expr = x * sin (1/x), from = -pi/6, to = pi/6, n=200)

curve (expr = x * 1, add = TRUE, col = "red")

... ... ... ...

Legends

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

legend ("top", legend = c ("x * sin(1/x)", "x"),
        col = c("black", "red"), lty C(1,1)

? legend ...

expression (...) ? ... ...

? plot math

e.g. title(x, y) ... ...

... ... ... ... ...

... stem (plot()) ... ... o = un

... ... ... ...

categorical data ... ... ... ...

bar plot. (height, names. org = NULL)

... ...

mosaic plot (x)
varide

        count = table ( );

        barplot (counts, name, org = c(...))

                ... ...

Multiple figures $\qquad$ plot (x, ..., ...)

matrix (data, nraw, ncols, byraw =FALSE)...

layout (mat) $\qquad$ (x, ..., ...)

layout show (3) $\qquad$

Write graphical output to a file

dev.off (6)

x = rnorm (100)

pdf ("____.pdf")

plot.hops

dev.off (x)

## Test and Intervals

One mean or the Difference of Two mean

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

mu = 0, conf.level = 95) ... to test ' mu = 0

$ parameter : degree of freedom

$ statistic   student's test statistic ...

$ p-value ...

$ conf. int.

$ estimate

str(out) ...

out$p.value ...

two means

$H_0: \mu_x - \mu_y = \mu_0 \Rightarrow$ ...

## F Test for Equality of Variances

out = var.test(x, y, ratio $\sigma^2 = 1$, alternative

= "two sided", conf.level = q)

$H_0 \quad \dfrac{\sigma^2_x}{\sigma^2_y} = \text{ratio}.$

& parameter $\qquad (n_x - 1)\dfrac{s^2_x}{} + (n_y - 1)\dfrac{}{q}$

& statistic $\qquad f = \dfrac{s^2_x}{s^2_y}$

& p.value

& conf.int.

## chi-Squared Tests $\qquad df = m = f$

### Goodness of fit

- counts ··· (...)  probs = c(...)

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

& parameter

& statistic $\qquad (x) \ldots = A$

$(\quad) \times$

### Independence / Homogenity

& expected : expected counts under $H_0$

chisq.test();

## One Proportion or the Difference of Two Proportions

out = prop. test (x, n = p) alternative = "two sided"

cont. level

2-test     $H_0: p = p_0 = p$     [ ... ] of test ⌐     = .95

o iterior  1 - $ estimate $p = x/n$ ...

... nwt =

2-vetors x and n

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

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

= 7   ... 

## 7. Regression

### 1. Simple linear Regression

$y = mx + b$

cor(x, y)  /  correlation

lm(y~x, data)

$A = cor(x) \Rightarrow$  matrix A

$cor(x[,i], y[,])$

### 2. lm(y ~ x, data)

Call: ...

anova(m)

m$coefficients [1]

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

= abline (reg=m)

abline (a = mean (cor $ dist)

    horizontal

predict (node)

y.hat = predict (m, new data=dl)

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

plot ( m$ fitted values, m$residuals) → X   no pattern.

abline (0, 0)

    ε → indepnd
    & normal

2Qplot.

    N (µ= mean ( residuals) , σ = sd(residuals) )

    qqline (x)

e.g x = rnorm (n=100) ;  qqnorm (x);  qqline(x)
    w = rex p (100).
              random

    or use   plot (m)

multiple Linear Regression
    m = lm(y ~ X1 + X2 + ... , data= )

    confint (m,

8. Simulation

(1) repeatable

set. seed (seed)

set.seed(0); $a = rnorm(1)$

set. seed(0); $b = rnorm(1)$

(2) Repeat a calculation $N$ times.

replicate ( ... )

set.seed(0)

(3) ? distribution.

$\bar{X} \approx N(\mu, \frac{\sigma^2}{n})$

$t = replicate (N,$   $\{x = rnorm (n, mean =$

symbols  text(x=d, y=row, label= v3echar, cex=4)
identify ( XX, YY, n=1, plot= FALSE )
unique( , )%in%

replicate ( , _ )

↓ 替字

- 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 )
                          numeric

skip  if (c<7 { next }

3. Apply Function  lapply (X= , FUN= ) treating data frame as a list of vectors
   } sapply (X= , FUN= )  ⟹ vector, matrix, array
multiple argument  mapply ( FUN, x,y,z ) take the several vectors in. and applies FUN to all first elements
   apply  apply (X= , MARGIN, FUN, ...)  =1    rows    and then to
   tapply (x@= , Index@ FUN= )  =2    columns    all second elements

subset of vector
(4) tapply ( X = mtcars$mpg , INDEX = mtcars$cyl, FUN= )  分类!
    tapply (X=           INDEX = list (mtcars$cyl, mtcars$gear), FUN= )

4 Matrix  m = matrix ( data = vector, nrow= , ncol= , byrow = TRUE )    extra
          dimnames = list (c( colname ), c( rowname ))        probs=c(
   cbind (matrix, )          rbind (matrix, )
   m[ row(m) == col(m) ]  main diagonal
   m[ row (m) + col (m) ==r-c ]th diagonal through (r,c)
   m    t    +    reverse.
   A * B element-wise product    A%*%B  matrix product
   solve (a=A, b=b)    Ax=b  A% x % x  matrix -vector product

5. pattern
   grep (pattern, x, ignore.case =FALSE, value = FALSE)  indice of elements of chara
   grep(pattern =         ...  TRUE ⟹ values

   sub (pattern= , replacement = , x= )  a copy of x after replacing 1st occurence
   of pattern with replacement
   gsub (): all    ( isTRUE ( ) )

\d  0123456789    \w  \s  \b  不与叵配怖

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.  {^} exactly n times / {^,} n or more times  * ({0,1}, + {1,})
          {n,m} n—m times inclusive   ? ⇒{0,1} or "optional"
e.g. grep (pattern = "\\d{4}$", x=a)   4 digits, end of lines
   grep (pattern = "□ \\d{4}$", x=a)  # space, 4 digits, end of line
   grep (pattern = " \\d{4,5}$", x=a) # space, 4 or 5 digits, end of lie
      \N refers to what the Nth enclosed expression matched
— sub (pattern = "(\\w+), (\\w+ ) + (\\d+)(\\w+). x , replacement = "\\2, \\1, \\4, \\3", x= )
 | means or   grep (pattern = "Joe | Jack", x=a)   grep (pattern = "[(0)a)", x=a )
  \\
    ⇒ ? regex
 —Splitting strings   strsplit (x=a, split=",")  strsplit (x=a, split=" +" =) 折分所有 space
    ⇒ write(x= csv, file= "  ")   d= read. csv(" ", header=FALSE, col names =c(  ))
6. ggplot ① qplot (data= ,   x= ___ y= ___ , geom= "boxplot.")
    ⊕ ggplot 1. gg.1 <- ggplot (data = trees, aes (x= __ , y= 1)
       gg1 ⇐ + geom_ boxplot()+ geom_jitter()  + violin()   + coord_flip() (flip
2. +xlab ("  ")  +ylab("  ") . + labs (title=" ") + theme (plot.title = element.text (size=22,
3. ggplot1   aes=(x. y. fill= x or y)   + guildes(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 ( ))   shape =factor()
+ facet_grid( . ~ )   + scale_x_sqrt (limits =c (0,50))
7. web scrap   TEAM <- "http:___ "  lines <- readLines (TEAM)
tem.lines <- grep (pattern ="  ", x=lines, value = TRUE)  team.names<-sub (x=tem.lines, pattern=" (*)
replace ="\\1")   link<-paste0( __, __, __ )  table <-read HTMLTable (link)
table [[3]] [ , ]    [ ] as. numeric (as. character ( ))
Vector sort( ) order   x[order (x)]   na.rm=TRUE   list[[ ]]
write table (x, file ="  ", ...)   table = read.table (...)
csv   .csv
    stopifnot ( function( ) == "  " )
turn a vector into a character ⇒ paste (V, sep=" ", collapse=" ")
    ⇒   = data.frame ( , )

## Web scraping

```r
library("XML" ); TEAM <- "http://www.nfl.com/teams" ; lines <- readLines(TEAM)
team.lines <- grep(pattern="statistics", x = lines, value=TRUE)
```
 • grab just those team names and the abbr from "\t\t\t\t\t\t\t\t<option
value=\"/teams/baltimoreravens/statistics?team=BAL\">Statistics</option>"
```r
team.names <- sub(x = team.lines, pattern=".*teams/(.*)/statistics.*", replace = "\\1")
team.abbreviation <- sub(x = team.lines, pattern = ".*team=(.*)\">.*", replace=" \\1" )
```
 • specifically i want to get information from:
http://www.nfl.com/teams/baltimoreravens/statistics?season=2014&team=BAL&seasonType=REG#
```r
link <- paste0(TEAM, "/", team.names[1], "/statistics?season=2014&team=", team.abbreviation[1],
"&seasonType=REG#" ); tables <- readHTMLTable(link)
```
 • for loop for all the teams
```r
rushing <- NULL; receiving <- NULL
for (i in 1:length(team.names)) {
    link <- paste0(TEAM, "/", team.names[i], "/statistics?season=2014&team=", team.abbreviation[i],
"&seasonType=REG#")
    tables <- readHTMLTable(link); print(link)
    rushing[i] <- as.numeric(as.character(tables[[3]][2,3])); receiving[i] <- as.numeric(as.character(tables[[4]][2,3])) }
```

## Apply
 • lapply (" list apply" ) applies function FUN to each element of vector or list X, returning a list of the same length
```r
lapply(X=mtcars, FUN=mean); sapply(mtcars, mean); sapply(mtcars, mean, simplify=FALSE, USE.NAMES =
FALSE);
```
 • sapply - simplified apply, return numbers
 • mapply(FUN, ...)(" multiple arguments apply" ) takes the several vectors in ... and applies FUN to all first elements,
then to all second elements, etc.
```r
mapply(sum, x, y, z); apply(X=m, MARGIN=1, FUN=sum) # keep dimension 1 (rows)
```
 • apply- applies the specified function over "margins" (MARGIN=1 for row, MARGIN=2 for column)
 • tapply - applied the function over a subset indicated with INDEX
```r
tapply(X=mtcars$mpg, INDEX=mtcars$cyl, FUN=mean); tapply(X=mtcars$mpg, INDEX=list(mtcars$cyl,
mtcars$gear), FUN=mean)
tapply(X=mtcars$mpg, INDEX=mtcars$cyl, FUN=quantile, probs=c(.25, .75))
```

## Matrix
 • main diagonal: m[row(m) == col(m)] • diagonal through (r, c): m[row(m) - col(m) == r - c] • reverse diagonal through
(r, c): m[row(m) + col(m) == r + c]
 • A * B is an element-wise product; A %*% B is the usual matrix product
 • solve(a=A, b=b) gives the solution x to the system of linear equations, A*x = *b

## Pattern
grep(pattern, x, ignore.case=FALSE, value=FALSE) 不分大小写, 显示满足条件的index。value=TRUE: 显示满足条件

的x中的结果。

sub(pattern = "e", replacement = "E", x = a)替换第一次出现的e; gsub(pattern = "e", replacement = "E", x = a)替换所有的e

- \\w, \\s, \\d: words, space, digit　　　　・\\W, \\S, \\D 大写表示不是word, space, digits;
- [aeiou]，其中任何一个　　　　・[^aeiou]不包括所有的　　　　・grep(pattern = "Joe|Jack", x = a)或者
- ^a以a开头　　・pattern = " e\\>" 以e结尾
- . \ | ( ) [ { ^ $ * + ? 遇到这些前面要加\\ (double backslash)

grep(pattern = " \\d{4,5}$", x = a)　# space, 4 or 5 digits, end-of-line
sub(pattern = "(\\w+),(\\w+) +(\\d+) (\\w+).*", replacement = "\\2,\\1,\\4,\\3", x=a) 满足pattern条件的内容按照2，1，4，3的顺序排列
sub(pattern=".*<a href=(.*)>.*", replacement="\\1", x=link) 满足pattern条件的只去处括号里的内容
- Splitting Strings for character vector x

  strsplit(x=a, split=" ," ); strsplit(x=a, split=" +" );　strsplit(x=a, split=",|( +)")

ggplot : library(ggplot2 )
- box plot或其他图　qplot(data = trees, x = species, y = dbh, geom="boxplot" ) ・ggplot(data = trees, aes(x =, y =)) + geom_boxplot()
- 加点，加线：ggplot(data = anscombe, aes(x = x1, y=y1)) + geom_point() + geom_smooth(method=lm) + geom_line()
- 不同颜色分组画图：ggplot(data = mtcars, aes(x=mpg, y=qsec, color=factor(cyl), shape=factor(gear)))
- 自己设定颜色：+ scale_color_manual(values = c("pink", " black" ))
- 加标题,调字大小：+ ylab("Diameter") + xlab("Species") + labs(title="Some trees")+ theme(plot.title = element_text(size=22, vjust=2, hjust=1))
- 重新排序，填色，不要注释：ggplot(data = trees, aes(x= reorder(species, dbh,fun=mean), y = dbh, fill = species))+ guides(fill= FALSE)
- 把y变成log形式，填色种类：+ scale_y_log10(labels = function(l) {paste(l, " in." )} + scale_fill_brewer(palette = "Set1")
- 画多张图on same plot :  + facet_grid(.~gear) ;  + facet_grid(carb~gear) ; + facet_wrap(carb~gear)
- scale_x_sqrt(limits=c(0, 50))
- error bar : se <- sqrt(diag(vcov(mod1))) ;  gg2 <- ggplot(data = sim, aes(x = Species, y = pred, fill=Species)) + geom_bar(stat="identity" )
    - □ 标准差：gg2 + geom_errorbar(aes(ymax = pred + se, ymin = pred - se), width=.25)
    - □ 置信区间：gg2 + geom_errorbar(aes(ymax = pred + qnorm(.975)*se, ymin = pred - qnorm(.975)*se), width=.25)
    - □ 画两端不出头的range：把 geom_errorbar 换成geom_pointrange；粗粗的range：换成geom_crossbar
    - □ 在bar上加字母：+ geom_text(aes(x = Species, y = pred + se + 1 任意高度, label = group要加的字母))