

Introduction: 统计工具分 supervised & unsupervised  $E(f(x)) = E(f(x) + \epsilon - f(x)) = [f(x)] - f(x) + \text{Var}(\epsilon)$  | 有监督有  $y_i$ , 无监督无  $y_i$  (cluster).  
Subset Selection Lasso Least Squares Generalized Additive Models Trees Bagging Boosting SVM  $\rightarrow$  Interpretability High to Low. | SVM 提前  $\leftarrow$  Flexibility Low  
Free Software: run for any purpose / study how the program work / adapt it to needs (Access to the source code) / redistribute copies / improve program, release it, benefit whole  
help. start("A") help("A")? A (精确) help. search() 返回相关 apropos() 返回相关obj 模糊 find("A") 返回包, example("A")  
**Working with R** R.version, string[getwd()] [setwd()] [dir("c:\\temp")] 显示目录下文件 [memory.size()] 内存使用 NA: 内存  
memory.limit() 显示或改变内存限制 [ls()] [object()] 显示工作空间的项目 [object.size(x)] 对象大小.  
object("package:base") base 中所有的 object [search()] 搜当前加载的包 [rm()] 清除 [rm(list=ls(all=TRUE))] 全部清除  
history(#) 显示最后 # 步命令 savehistory("myfile") 保存历史 命令命名为 myfile save.image("myfile") 保存工作空间.  
getwd() 返回 . [libpaths()] 库路径 [library()] 显示已安装包 install.packages(c("A","B")) 安装包 %% 取整 %/% 取余 119%7=0 判断整除.  
log() exp() sign() abs()  
log(x, base=a) log10() log2()  
log 显示, log 全部 option  
gamma(), beta() c() 向量  
R 中基本 unit 是 vector. length(x)  
2\*x+3 数乘, 5:1\*x 按顺序乘  
x+1:5 按顺序加, log(x) 逐个求  
xx[1:4] 取前4号.  
xx[xx>"a"] 取出 >"a" 的  
u<-xx>"a" 返回 FALSE, TRUE.  
xx[u] 取出 TRUE 位置值.  
xx[c(1,4)] 调用1和4号.  
xx[-c(2,3,5)] 调用除2,3,5号  
ones<-rep(1,10) 10个1的向量  
even<-seq(from=2,to=20,by=2)  
2到20步长2的向量.  
trend<-1981:2005 生成年份向量  
c(ones,even) 连接向量  
cumsum(x) 各个数累加值  
cumprod(x) 各个数累乘值  
min(z[1],z[2]) 求全数字最小  
pmin(z[1],z[2]) 对应位置最小  
返回向量.  
nrm(function(x) return(1/8))  
x: 自变量, 8: 起始值,  
nrm 求函数最小值  
返回 minimum, estimate 极值点  
gradient 极值处梯度  
iterations 迭代次数.  
D(expression(exp(x^2)), "x")  
求导, "x" 求导对象 deriv3  
integrate(function(x) x^2, 0, 1)  
积分, 0, 1 下限, 上限.  
x\*y 对应位置乘法.  
t(x)%\*%y=crossprod(x,y) (l)  
x%\*%ty=t(crossprod(x,y)) (r)  
t(x)转置, =outer(x,y)  
z<-matrix(c(c...), nrow=3,  
ncol=2, byrow=F)  
生成3行2列 F: 按列排  
dim(A) 维数(行,列)  
nrow(A) 行数, ncol(A) 列数  
A[1,] 取1st行, A[,1] 取1st列  
A[2,2] 取2行2列 (都返回向量)  
class(D) 向量返回内容类别, 非  
向量返回格式类别  
E<-A[2,2,drop=F] 强制返回向量  
A<-A[,2] 返回除第2列的阵  
det(A) A的行列式值.  
eigen(A) A的特征根.  
solve(A): 求A逆  
diag(2,3,4) 对角2,3行4列  
rep(c(1,2),c(3,3)) 重复3次2  
svd(A) 奇异值分解 qr(A) QR分解  
cbind(A,B) 左右联 bind 上下联  
runif(20) 长度20的均匀分布向量  
solve(matrix(c(1,-1,1),nrow=2),c(2,4))  
解  $\begin{pmatrix} 1 & -1 & 1 \\ x_1 & x_2 & x_3 \end{pmatrix}$   
Union(x,y) 并集, intersect(x,y) 交集  
setdiff(x,y) (x-y), %in% y 查在 y 内  
setequal(x,y) 检查 x 是否相同  
choose(n,k)=C<sub>n</sub><sup>k</sup>  
z<-3.5-8i (虚数)  
Re(z) 实部 Im(z) 虚部 Mod(z) 模长  
Conj(z) 共轭, is.complex 判是否  
most number are rounded to an  
accuracy of 53 binary digits.  
除用相同公式算否则2个浮点数  
不会相等, sqrt(2)\*sqrt(2) = 2  
sqrt(2)\*sqrt(2) = 2 但减其他正常  
x<0.3-0.2, y<0.1, x==y (F)  
identical(x,y) F all.equal(x,y) T  
命名: 区分大小写, 不以数/符号开头  
不能有空格, 中间随意  
① Data Structures.  
R可处理的数数据类型: numeric,  
character, logical, complex, raw  
R中数容器: scalars, vectors,  
matrices, arrays, data frame (list  
attributes(A) 返回A的各项属性  
A的类型  
向量只可装(数, 字符, 逻辑)  
同一个向量只可装一种类型.  
矩阵只可装(数, 字符, 逻辑)  
同矩阵用相同类型.  
matrix(c(c...), nrow=..., ncol=...,  
byrow=T/F, dimnames=list(row,col))  
y[1:2,c(3,4)] 取1-2行, 3和4列.  
z<-as.vector(y) 转成向量  
dim(z)<-c(4,5) 转成4x5矩阵  
str(z)=int[1:5,1:4] 1 2 3 4 ...  
z<-c(1:2,4) dim(z)=c("A1","A2")  
dim 2<-c("B1","B2","B3")  
dim 3<-c("C1","C2","C3","C4")  
z<-array(1:24,c(2,3,4),  
dimnames=list(dim1,dim2,dim3))  
A,B 嵌在 C1, C2, C3, C4 里  
aperm(A,c(2,3,1)) (调位)  
C>A 变成 C(3,4,2) 值  
数据框不同列可放不同类型  
A<-c(1,2,3,4) B<-c(a,b,c,d)  
cbind 变无名字, 1,2,3,4 变字符.  
data.frame(A,B) 都同时保留名字  
names(A) 显示框内名字  
A[1,2] 或 A["name1","name2"] 取  
不可 A[name1] 返回 data.frame.  
A\$name1 返回向量  
factor(数据, ordered=TRUE)  
有序因子 默认按字母表.  
factor(status, order=TRUE,  
levels=c("Poor","Improved","Exc"))  
指定排序, p<T<E.  
status2<-c("P","I","E")  
status3<-factor(status2, levels  
=c("P","I","E"), labels=c("C","B",  
"A")) 返回 [1] CBA [levels: CBA  
不在 labels 中的会被 miss.  
status3[4]<- "P" x 因子不能这样  
转回向量加好括弧加回  
! NOT & AND | OR  
A & B 仅 B 中 true 的去也.  
A | B 仅 B 中 true 的去也  
isTRUE(6>4) TRUE  
identical('twins','twins') FALSE  
xor(5==6,'FALSE') TRUE (或:)  
xor(1 is TRUE(TRUE), 6>1) TRUE  
xor(4>9, 8!=8.0) FALSE.  
8!=8.0 F, (T!=F)=!(6==7) (F)  
T & C(T,F,F) (T F F)  
T & C(C,T,F,F) (T T T)  
T | C(T,F,F) (T T T)  
5>8 | 6!=8 & 4>3.9 (F) 先 & &  
aaa<-gl(2,5) class(aaa) factor  
mode(aaa) numeric  
typeof(aaa) integer  
list 不要每列数据相同  
R中最 complex 的 data types  
lst<-list(name="Fred", wife  
="Mary", child.ages=c(9,7,9))  
lst[[2]] 取出 "Mary", lst[[3]][2]  
lst[[c("name")]] F: id = lst\$name  
z\$b<- "k" 增加1个属性  
z\$b<- NULL 删掉属性  
x\$a 模糊找, x[c] 精确找  
x[c, exact=false] 模糊找  
Formulas: f<-y~x  
NA 有分 integer... 类.  
NA N 属于 NA, NAN: not a number  
attach: 绑定到某数据上  
detach 解绑定  
with(mtcars, { }) 括号中绑定 mtcars  
{ } 中 <- 不保存, <-<- 保存  
Data Input.  
Mydata<-data.frame(age=numeric(),...)  
Mydata<-edit(Mydata) 窗口编辑  
x<-numeric(10) Data.entry(x): 输入10个  
a<-scan() 输入数不输字, 回车一次, 放空  
回车结束 (空行补齐, 扩展分割) -> "t"  
rt<-scan(file="~/rt.txt", sep="\n")  
回车分割 同行合并 默认空格分  
rt<-scan(file=file.choose, sep="t")  
scan(... what=double(), skip=1)  
双精度, 跳过前一行  
read.table(file, header=T/F, sep=" ",  
row.names=" ") header: 是否1st作names.  
colClasses 指定每列 class=c("integer"...)  
skip: 从第几行起 stringsAsFactors=FALSE  
nrow: 读几行. 禁止文字转因子  
文件结尾用 \n, 可读 URL  
read.delim(file) 默认 1st 作表头, 空格分  
url=URL("http:...") ds=readLine(url)  
read.xlsx("...xlsx", 1, header=T, as.data  
frame=T)  
read.dta("...dta") <- Stata  
read.table("clipboard") 从剪贴板读  
data() 读内存建 data,  
data(package=.package(all.available=T))  
读外包内置, try(data(package="M")看看中  
数据名字, history(int) 看全部历史(窗口)  
loadhistory(file) 保存历史  
pdf("~/1.pdf") data(mtcars)  
plot(mtcars\$wt~mtcars\$mpg)  
dev.off 将图存入 pdf  
write.table(data, file="...txt", sep=" ",  
quote=FALSE, append=FALSE, na="NA")  
p 为逗号.  
save(data, file="...Rdata")  
load("...Rdata")  
save(list=ls(all=TRUE), file="alt.Rdata")  
writeClipboard(data) 写入剪贴板  
file.exists("...txt") 确认文件存在  
data.frame 改名: names(A)[2]<-"B"  
A\$B<-factor(A\$B, levels=c(1,2), labels=c("male",  
"female")) 将数据 1,2 用 male, female 表示  
data\$sum<-data\$x1+data\$x2  
attach(data), data\$sum<-x1+x2 detach(data)  
data<-transform(data, sum=x1+x2) 三行用样

seq(from=0.04, by=0.01, length=11)  
seq(0.04, by=0.01, along=N)  
0.04 开始, N个, 步长0.01  
最后一个值超出范围则停止, 不超出.  
seq(1.4, 2.1, 0.3) : 1.4, 1.7, 2.0  
Sequence(c(4,3,4,4,4,5))  
生成 1:4, 1:3, 1:4, ...  
rep(4, 4, 2) 重复2次 1:4  
rep(1:4, each=2) 11223344  
rep(1:4, each=2, times=3)  
11223344 11223344 11223344  
rep(1:4, c(4,1,4,2))  
1x4, 2x1, 3x4, 4x2 次.  
rep( )  
重复值, 单个为数据, 向量对应  
gt(最大, 重复次数, 总长度)  
生成的是 levels, 长度不修再重  
Soft <- gl(3, 2, 6, labels=c("A", "B", "C"))  
Temp <- gl(2, 3, 6, labels=c("D", "E"))  
data.frame(Temp, Soft)  
修改变量:  
A\$B[A\$C > 5 & A\$E < 7.5] <- "X"  
B可以已有也可新建. 也可用:  
A <- within(A, {B <- NA; B[C>1]  
 <- "X"; B[F<9] <- "K"})  
可以保存新建量返回, with在  
{ }中新建的不保存(用<-保存)  
重命名变量: fix(A) (窗口)  
names(A)[2] <- "B"  
names(A)[1:2] <- c("A", "B")  
缺失值 NA, 不可能值 NaN  
A = NA 永远是 FALSE, 用 is.na(A)  
返回 T/F 的阵或 vector.  
X <- c(1, 2, NA, 3) Y <- X[!is.na(X)]  
print(y) [1] NA. 有 NA 算完全 NA.  
y <- sum(X, na.rm=TRUE) NA 转 0  
bad <- is.na(X) X[bad] 显示 NA  
X <- c(1, 2, NA) Y <- c("a", NA, "b")  
good <- complete.cases(X, Y)  
返回 T/F, X[good] 返回 Y[good] 返回  
na.omit() 删除含 NA 的行 (row)  
complete.cases(A) 对 A 的每一行查  
NULL: 不存在, 空. Z <- NULL 新建  
NULL 可添入东西, NA 不同 NULL 无 mode  
NULL 的 length 为 0, NA 的 为 1  
Dates -> Date class  
Times -> POSIXct or POSIXlt  
Dates 有从 1970-01-01 开始的天数  
Times 有从 1970-01-01 开始的秒数  
x <- as.Date("1970-01-01") unclass 为 0  
POSIXct 是一串数字, POSIXlt 是 list  
x <- Sys.time() 获取系统时间 (ct 型)  
date() 返回 "Wed Mar 09 23:45:16 wib"  
today <- Sys.Date() 只返回日期  
format(today, format="%m %d %Y")  
format(today, format="%Y")  
2 个 Date 型相减可得日期差  
difftime(today, dob, units="weeks")  
求差几周 myformat <- "%m/%d/%Y"  
as.Date(A\$B, myformat) 转格式  
可将 factor 转为 Date  
排定 A\$order(A\$B). 7 排 age 排序

attach(A) B <- A\$order(gender, -age)  
detach(A) 先排列性别, 再按年龄大小排  
myvars <- names(A)[!%in% c("B", "C")]  
返回 T/F. 判断 A 各变量是否是 B 或 C  
A[!myvars] 或 A[c(-1,-5)] 剔除  
选择变量: A[which(A\$X=="M" & A\$Y>20)]  
which 返回 c(...) 代入 A 取数据.  
A <- Subset(B, C) > 35 | C < 10, select  
= c(X, Y) 选出 B 中 C>35 或 C<10 的  
条目的 X, Y 两变量  
A <- Subset(B, k="M" & C>30, select  
= X:Y. 选出 X 到 Y 之间全部数据  
A[sample(1:nrow(A), 3, replace=F)]  
抽 A 中的样本 (范围, 次数, 是否放回)  
sample 返回 1 个向量.  
加载: aggregate(price ~ cut, dia, mean)  
cut price 返回以 cut 分类的 dia  
Fair 39 数据中 price 的 mean.  
aggregate(price ~ cut + color, dia, mean)  
按 cut, color 交叉分类求 price 均值  
aggregate(cbind(price, car) ~ cut, dia, mean)  
以 cut 分类返回 price, car 的 mean  
merge(A, B) A: 

a	b	c	d
m	n	h	x
o	p	q	r
q	r	x	x

 B: 

a	b	c	d
m	n	h	x
o	p	q	r
q	r	x	x

  
merge(A, B, all=T) ->  
merge(x=A, y=B, by="x")  
= c("m", "n"), by="y"  
= c("m", "n") 按 y 的 m 和 n 来合并  
(ply 包) join(X=A, Y=B, by=c("m", "n"))  
rbind 上下连 dataframe 时不用使 A B 列  
名顺序对应 会自动匹配 但像一样  
重命名 matrix 的行名为 rownames.  
colnames(A) <- c("...", "...")  
Creshape 包: melt(A, id=c("ID", "Time"))  
A: 

ID	Time	Var	value
1	1	X1	3
1	2	X1	6
2	1	X1	2
2	2	X1	4
1	1	X2	4
1	2	X2	5
2	1	X2	6
2	2	X2	4

  
-> 

ID	Time	Var	value
1	1	X1	3
1	2	X1	6
2	1	X1	2
2	2	X1	4
1	1	X2	4
1	2	X2	5
2	1	X2	6
2	2	X2	4

  
Case (md, ID + Variable  
- Time) 把 time 值当分类  
Cast (md, Time + Variable,  
mean)  
country, developed GDProw GDProw GDProw  
China 0 5000 5500 5010  
USA 0 6000 6500 6010  
Japan 1 7000 7500 7001  
long <- reshape(Data, idvar="country",  
varying=list(names(Data)[3:5]),  
v.names="GDP", timevar="year",  
times=c(2000, 2005, 2010), direction="long")  
long <- reshape(long\$country, )  
China.2000 China 0 2000 5000  
China.2005 China 0 2005 5500  
China.2010 China 0 2010 5010  
Japan.2000 Japan 1 2000 7000  
Japan.2005 Japan 1 2005 7500  
Japan.2010 Japan 1 2010 7001  
USA.2000 USA 0 2000 6000  
USA.2005 USA 0 2005 6500  
USA.2010 USA 0 2010 6010  
wide <- reshape(long, v.names="GDP",  
idvar="country", timevar="year",  
direction="wide") rownames(wide) <- NULL  
回到 md 但 GDP.2000 带点.

paste("A", c("B", "C", "D"), c("E", "F", "G"))  
"A B E" "A C F" "A D G"  
paste("A", "B", "C", sep="/") "A/B/C"  
x="a", y="b" paste("A", x, y) "A a b s"  
sprintf("A %s B %s", x, y) "A a B b s"  
sprintf("A %s B %s C %s", c("a", "b", "c"),  
c("1", "2", "3"), c("4", "5", "6")) "A a B b C c"  
AM B3 C Y AQB F C Y AM B7 C Y  
(string 包) str-split(string=A\$B,  
pattern="-") 去掉 A 中所有的 "-" (List)  
Reduce(rbind, A) 将 A 中所有行合并成  
str\_sub(string=A\$B, start=1, end=3)  
取前 3 个字符.  
A[str\_sub(string=A\$B, start=4, end=4)==1,]  
抽出 A 中 B 的第 4 个字符为 1 的条目  
Str\_detect(string=A\$B, pattern="k") 找出带有 k  
的条目返回位置向量, 位置大小写  
str\_detect(string=A\$B, "k") 也可用, 包括  
str\_detect(string=A\$B, ignore.case=T)  
str\_split(string=A, pattern="H", n=2) 看到  
H 分割, 分 2 段, H 后去掉, 以第 1 个为准.  
str\_trim(A) 去掉开头和结尾空格和 N.  
str\_extract(string=A, pattern="k") 提取 k 元  
k 的条目返回 NA, 有 k 则提取出.  
[0-9] 代表 0-9 数字, [0-9][0-9][0-9][0-9] 9 位数字  
[0-9]{4} 代表 4 位数字, \\d{4} 数字  
\\d{1,3}, 1-2-3 个数字组合  
"\\d{4}" 开头是 4 位数字, "\\d{4}\$" 结尾是 4 位数字,  
"\\d{4}\$" 开头结尾都是 4 位数字, str\_replace 反之  
str\_replace\_all(string=A, pattern="\\d",  
replacement="x") 全部替换成 x  
unlist(str\_split(Sys.time(), "")) [2]  
for(A) 命令 或 for(A) { } 可换行  
A[i in 1:4] (i in seq\_along(x)) x: 1 的长度  
(letter in X) -> 用 X 中的字母替换 A 中  
for(i in 1:100) { iTotal <- iTotal + i  
cat("Sum of 1-100:", iTotal, "\n", sep="")  
szSymbols <- c("MSFT", "GOOG", "AAPL", ...)  
for(SymbolName in szSymbols) {  
cat(SymbolName, "\n", sep="")  
}  
x <- matrix(1:6, 2, 3)  
for(C in seq\_len(nrow(x))) {  
for(J in seq\_len(ncol(x))) {  
print(X[C,J]) } }  
count <- 0 while(count < 10) { print(count)  
count <- count + 1 }  
i < 11; Total <- 0 while(i < 100) { iTotal <-  
iTotal + i; i < i + 1 }  
z < 5 while(z >= 3 && z <= 10) { print(z)  
coin <- rbinom(1, 1, 0.5) # 伯努利分布 p=0.5  
if(coin == 1) { z <- z + 1 } else { z <- z - 1 }  
i < 1; iTotal <- 0 repeat { iTotal <- iTotal + i  
i < i + 1; if(i <= 100) { next } else break }  
n < 3; if(a == 1) { print("a=1") } else { print  
("a!=1") } } else 合符  
grade <- as.factor(c("grade1", "grade2"))  
if(!is.factor(grade)) { grade <- as.factor(g  
rade) } else { print("Grade already is a  
factor") } } 判断是否是因子, 不是就转  
a <- 4 if(a == 1) { print("a=1") } else if  
(a == 2) { print("a=2") } else { print("Not  
1 & 2") } } 多条件

X <- matrix(1:6, 2, 3) if else (X > 0, sqrt(X), NA)  
若 x=0 开方, 否则 NA 可开方除 0 个  
n < 1 | switch(n, print("option1"), print("op  
tion2"), print("options")) 输出 option1  
ccc <- c("b", "aa", "a", "A", "bb")  
for(ch in ccc) cat(ch, " ", switch(EXPR = ch,  
a=1, b=2, 3), "\n") EXPR 固定开头, ch 被判断  
元素, 若符合到 a 则输出 1, b 则输出 2, 3  
用户函数 my <- function(a1, a2, ...) { ... }  
函数可作为参数传递到另一个函数, 可嵌套  
colmean <- function(y) { h <- ncol(y); mean <-  
numeric(h); for(i in 1:h) { means[i] <- mean  
(y[, i]) } } means } 有 NA 出 NA  
colmean <- function(y, removeNA=TRUE) {  
h <- ncol(y); means <- numeric(h); for(i in 1:h) {  
means[i] <- mean(y[, i], na.rm=removeNA) } }  
means } 可总是去除 NA.  
median(x) 中位数.  
parboth <- function(a, b) { c <- pmax(a, b); d <- pmin(a, b)  
answer <- list(median(c), median(d)); names(  
answer)[1:2] <- "..."; names(answer)[2] <- "..."  
return(answer) } 也可 return 结果  
twosam <- function(y1, y2) { n1 <- length(y1); n2 <-  
length(y2); yb1 <- mean(y1); yb2 <- mean(y2); S1 <-  
Var(y1); S2 <- Var(y2); SC <- ((n1-1)\*S1 + (n2-1)\*S2) /  
(n1+n2-2); tstat <- (yb1-yb2) / sqrt(SC\*(1/n1+1/n2))  
tstat } T 检验 sd: 标准差可转 na.rm  
args() 返回函数参数.  
lm(data=mydata, y~x, model=F, l=100)  
= lm(y~x, mydata, l=100, model=F)  
... 写在前面不可按位置匹配