

R Tutorial for Statistical Learning

Qidi Peng

1: Read, create and save data

Code	Comments	Results
<pre>a<-4; a=4; a</pre>	<p>Generate a value. Define “a” by 4.</p> <p>The same.</p> <p>Output the value of a (by pressing “enter” key)</p>	<pre>a [1] 4</pre>
<pre>c(1,2,3) 1:5 seq(from = 4.5, to = 2.5, by = -0.5) rep(3,4)</pre>	<p>Generate a vector manually</p> <p>Generate an increasing sequence of integers.</p> <p>Generate a decreasing sequence, starting from 4.5, ending at 2.5, with mesh -0.5 (by default, by=1).</p> <p>Generate a constant sequence with length 4</p>	<pre>[1] 1 2 3 [1] 1 2 3 4 5 [1] 4.5 4.0 3.5 3.0 2.5 [1] 3 3 3 3</pre>

<pre>matrix(1:8,nrow=2,ncol=4,byrow=TRUE)</pre>	<p>Generate a 2*4 matrix with data 1 to 8, by putting the digits one by one in row (by default).</p>	<pre>[,1] [,2] [,3] [,4] [1,] 1 2 3 4 [2,] 5 6 7 8</pre>
<pre>x=matrix(1:8,nrow=2,ncol=4,byrow=FALSE); x</pre>	<p>Generate a 2*4 matrix with data 1 to 8, by putting the digits one by one in column.</p>	<pre>x [,1] [,2] [,3] [,4] [1,] 1 3 5 7 [2,] 2 4 6 8</pre>
<pre>y=matrix(c("a", "b", 1, 2), nrow = 2); y</pre>	<p>Generate matrix where the coefficients are not values. Notice that here R regards "1" and "2" as symbols, not values.</p>	<pre>y [,1] [,2] [1,] "a" "1" [2,] "b" "2"</pre>
<pre>colnames(y) = c("char", "num"); y</pre>	<p>One assigns the column's names.</p>	<pre>y char num [1,] "a" "1" [2,] "b" "2"</pre>
<pre>rownames(y) = c("Sept.15", "Sept.16"); y</pre>	<p>One assigns the row's names.</p>	<pre>y char num Sept.15 "a" "1" Sept.16 "b" "2"</pre>
<pre>Age=c(23,21,22); Name=c("Jin","Kelly","Choi"); Gender=c("M","F","F"); Data=data.frame(Name,Age,Gender); Data=edit(Data)</pre>	<p>data.frame() functions are used to construct data set by typing in the values of each variable and assigning each vector its corresponding variable name. Manually import and change data in the data set "Data".</p>	<pre> Name Age Gender 1 Jin 23 M 2 Kelly 21 F 3 Choi 22 M</pre>

M454_Grades <pre>=read.table("C:/Users/pengq/Desktop/M454.txt");</pre> colnames(M454_Grades)=c("Name","Grades"); M454_Grades	Import data from a text file using <pre>read.table(directory/filename.txt)</pre> . One has to create this .tex fle first.	Names Grades 1 1 100 2 2 100 3 3 100 4 4 99 5 5 99 6 6 100 7 7 98 8 8 99 9 9 Absent 10 10 89
M454_Grades[1,]	Print the first row.	Name Grades 1 1 100
M454_Grades[,2]	Print the second column.	[1] 100 100 100 99 99 100 98 99
M454_Grades[9,2]	Print the element in Row 9, Column 2.	Absent 89 [1] Absent
read.table("website")	You can also read table from website.	
fpe = read.table("http://data.princeton.edu/wws509/data sets/effort.dat"); fpe; fpe=edit(fpe);	Examples: read data from website Read the table. A workable excel.	> fpe setting effort change Bolivia 46 0 1 Brazil 74 0 10 Chile 89 16 29 Colombia 77 16 25 CostaRica 84 21 29 Cuba 89 15 40 DominicanRep 68 14 21

		Ecuador 70 6 0 ElSalvador 60 13 13 Guatemala 55 9 4 Haiti 35 3 0
save(), sink(), write.table() save(fpe,file="data.Rda"); load("data.Rda")	Save data and variables Example: save the data.frame fpe. Load saved data.frame fpe.	
q()	Quit and save workspace (useful when too time consuming)	
help(write.table)	Use help to get the tutorial of functions	

2. Create functions

Codes	Comments	Results
Summation=function(x,y){x+y};	Create a function Summation(x,y)=x+y.	[1] 8

Summation(3,5)		
Product=function(a,b,c){z=a*b+c;z}; Product(2,5,7)	Create a function Product(a,b,c)=ab+c.	[1] 17
fix(arrangement); function(n,k){factorial(n)/factorial(n-k)}; arrangement(3,3)	Create a function arrangement(n,k)=n!/(n-k)!	[1] 6
arrangement=edit(arrangement); fix(arrangement);	Edit the existing function. The same.	
fix(Taylor_expansion); function(n,x) {a=1;for (i in 1:n){a=a+x^i/factorial(i)};a}; Taylor_expansion(100,3) exp(3)	Use “for loops” to create Taylor expansion of exponential of x, of order n: $1+x+x^2/2!+\dots+x^n/n!$ Approximate e^3 with n=100 True value of e^3	[1] 20.08554 [1] 20.08554
fix(parity); function(x){for(i in 1:length(x)){if (x[i]/2==floor(x[i]/2)){print(x[i]); print(" is even")} else {print(x[i]);print("is odd")}};}; parity(1:8)	Use “if else” to create function, which tells the parity of any real valued vector.	[1] 1 [1] "is odd" [1] 2 [1] " is even" [1] 3 [1] "is odd" [1] 4 [1] " is even" [1] 5 [1] "is odd" [1] 6 [1] " is even" [1] 7 [1] "is odd" [1] 8 [1] " is even"
fix(sum_square); function(N){i=1;y=0; while(i<=N){ y=i*i+y; i=i+1 }}; sum_square(5)	Use “while loop” to generate the sum of squares: $1+2^2+3^2+\dots+n^2$. Determine $1+2^2+3^2+\dots+5^2$	 [1] 55
3+5, 4*8, 5*pi, exp(5), sin(pi), cos(pi), (c(1,2))^2, abs(-3), c(4,5)/3, complex(real=3,imaginary=5), c(1,2)%*%c(1,2)	Some useful functions.	

3. Plot functions

Codes	Comments	Results
<pre>fix(linear); function(x,a,b){a*x+b}; x=seq(0,1,0.01); plot(x, linear(x,2,-2), col="blue",typ="p",xlab="x",ylab="f(x)"); abline(h=-0.5,v=0.3,col="red"); legend("bottomright",title="Linear function", c("f", "line"),col=c(26,33),lty=c(3,1))</pre>	<p>Illustrate the trajectory of linear function $f(x)=ax+b$.</p> <p>Create function.</p> <p>Input values of x.</p> <p>Plot the corresponding values of $f(x)$</p> <p>Color =(blue, red, green...), line charts type=(l,p,o,b,s,h,...), xlab(name of xline), ylab(name of yline)</p> <p>Add lines without erasing ex plots.</p> <p>Add legend (these are underlines commands helping to make footnotes of the plot).</p>	
<pre>windows() quartz() split.screen(c(2,2)) screen(1) plot(1:5,1:5) screen(2) plot(1:5,-(1:5)) screen(3) plot(1:5,2:6) screen(4) plot(1:5,-(2:6))</pre>	<p>Open new windows for graphics.</p> <p>For windows system.</p> <p>For Mac system.</p> <p>Split screen into 2*2 small windows.</p> <p>Plot in the first window (topleft)</p> <p>Plot in the second window (topright)</p> <p>Plot in the third window (bottomleft)</p> <p>Plot in the fourth window (bottomright)</p>	[1] 1 2 3 4
<pre>pdf("name of your plot"); png(); jpeg();</pre>	<p>You can use “save” button to directly save your pictures or use these commands to save them.</p>	

4: Generating random variables

Code	Comments	Results
runif(1);	Generate a uniform random variable over [0,1].	[1] 0.6197008
runif(4);	Generate a sequence of 4 independent uniform random variables over [0,1].	[1] 0.23239163 0.87730174 0.06837577 0.54853920
runif(1,3,4)	Generate a uniform random variable over [3,4].	[1] 3.892476

dunif(0.5)	The density function of the uniform distribution over [0,1], at 0.5.	[1] 1
punif(0.5)	The probability distribution function of the uniform distribution over [0,1], at 0.5.	[1] 0.5
floor(runif(1,1,7))	Generate a uniform random number over {1,2,3,4,5,6} (consider a die of 6 faces).	[1] 2
sample(0:1, 10, replace=T)	Toss a coin independently 10 times.	[1] 1 1 0 1 0 0 0 1 0 1
sample(1:6, 10, replace=T)	Roll a die independently 10 times.	[1] 3 1 3 5 4 1 4 5 2 4
sample(1:9, 3, replace=FALSE)	Randomly and equally likely pick 3 numbers from numbers 1 to 9, without replacement.	[1] 7 4 6
sample(1:3,1,replace=T,c(1/4,1/2,1/4))	Generate a non-uniform random number over {1,2,3}, with probability distribution (1/4,1/2,1/4).	[1] 2
sample(10)	Permute randomly and equally likely numbers 1 to 10.	[1] 5 3 4 6 7 1 2 9 8 10
rbinom(1,1,0.5)	Sample of Bernoulli(0.5).	[1] 1
rbinom(1,3,0.5)	Sample of Binomial distribution Binomial(3,0.5).	[1] 1
rexp(1,0.5)	Sample of exponential distribution Exponential(0.5).	[1] 2.365289
rnorm(1)	Sample of standard normal N(0,1).	[1] 0.9059011
qnorm(0.975)	Quantile of standard normal at level 0.975.	[1] 1.959964
qt(0.975,18)	Quantile of 18 degrees of freedom student t-distribution at level 0.975.	[1] 2.100922
qf(0.95,10,8)	Quantile of (10,8) degrees of freedom F-distribution at level 0.95.	[1] 3.347163
mean(); var(); median(); sum(); cumsum(); length();	Average of vector; Sample variance of a vector (unbiased) Median of vector; Summation of all elements of a vector All partial sums of a sequence (this gives a vector of the same size as the sequence) Cardinal (number of elements, dimension) of vector	[1] 1 2 3

<pre>X=1:6; X[X<=3]; var(), sd(), sort(), sort(1:5,decreasing=TRUE), quantile(), sqrt(),median()</pre>	<p>Extract all elements of X (in order) which are less or equal to 3.</p> <p>Help yourself to find the meaning of these operators for vectors.</p>	
------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------	--

5: Linear regression

Code	Comments	Results
<pre>install.packages("ISLR"); library(ISLR); library(MASS); fix(Boston); names(Boston); ?Boston; lm.fit=lm(medv~lstat,data=Boston); lm.fit; names(lm.fit); coef(lm.fit); confint(lm.fit,level=0.95); summary(lm.fit); predict(lm.fit,data.frame(lstat=c(5,10,15)),interval="prediction",level=0.95); plot(Boston\$lstat,Boston\$medv,pch="+",col="blue"); abline(lm.fit,lwd=3,col="red");</pre>	<p>Install the data set package ISLR, which includes the data we use in this class..</p> <p>Load ISLR.</p> <p>Load MASS, which contains some exotic statistical functions which are not contained in the base functions of R For more information on this package, search it on website.</p> <p>Open data.frame "Boston" in MASS.</p> <p>Check variable names of the columns.</p> <p>Check all information about the data set Boston.</p> <p>Simple linear regression with $Y = \text{medv}$, $X = \text{lstat}$.</p> <p>Watch the estimates!</p> <p>Statistical properties contained in the linear regression.</p> <p>Least squares coefficient estimates.</p> <p>0.95 confidence interval of the real parameters β_0 and β_1.</p> <p>More detailed information on this linear regression.</p> <p>0.95 prediction interval of the response variable $Y = \text{fit}$, with predictor $X = 5, 10, 15$.</p> <p>Plot $(X, Y) = (\text{lstat}, \text{medv})$;</p> <p>Add the regression line without erasing the ex-plots.</p>	<p>Call:</p> <pre>lm(formula = medv ~ lstat, data = Boston)</pre> <p>Coefficients: (Intercept)</p> <p>lstat</p> <p>34.55</p> <p>-0.95</p>

<pre>lm.fit2=lm(medv~lstat+age,data=Boston); summary(lm.fit2); lm.fit3=lm(medv~.,data=Boston); summary(lm.fit3).</pre>	<p>Multiple linear regression (here are 2 predictors $X=(lstat,age)$).</p> <p>Statistical information on the regression.</p> <p>Linear regression, $Y=medv$, X contains all other variables contained in the data set Boston.</p>	
<pre>lm.fit4=lm(medv~lstat+l(lstat^2),data=Boston); summary(lm.fit4);</pre>	<p>Non linear regression (here one assumes $Y=\beta_0+\beta_1X+\beta_2X^2+\epsilon$).</p>	

6: Classification

Code	Comments	Results
<pre># Logistic regression summary(Smarket); cor(Smarket); cor(Smarket[, -9]); Smarket\$Direction; glm.fit=glm(Direction~Lag1+Lag2+Lag3+Lag4+Lag5 +Volume,data=Smarket,family=binomial); coef(glm.fit); glm.probs=predict(glm.fit,type="response"); contrasts(Smarket\$Direction); plot(glm.fit); glm.pred=rep("down",1250);</pre>	<p>Consider data set Smarket (stock market data) in the package ISLR.</p> <p>Produce correlation matrix of all variables contained in Smarket. Error detected when meeting qualitative data.</p> <p>Exclude the last 9 columns (variables) of Smarket. The variable "Direction" in Smarket is qualitative.</p> <p>Multiple logistic regression with response "Direction" and predictors Lag1 to Volume. Observe that here the p-values are relatively large, so there is no clear evidence of a real association between Direction and those variables.</p> <p>Estimates of the coefficients.</p> <p>Predict the probability $P(Y=1 X)$ for each training data X. To see the meaning of event $Y=1$, use contrasts().</p> <p>Convert numbers to dummy variable.</p> <p>The glm function also provides figure representations. It's slow. Please wait and watch.</p> <p>Create a vector of class prediction based on whether</p>	<p>Error in <code>cor(Smarket) : 'x' must be numeric.</code></p>

<pre>glm.pred[glm.probas>0.5]="up"; table(glm.pred,Smarket\$Direction); mean(glm==Smarket\$Direction)</pre>	<p>the predicted probability of a market increase is greater than 0.5.</p> <p>Produce a confusion matrix in order to determine how many observations were correctly or incorrectly classified. The diagonal elements of the confusion matrix indicate correct predictions, while the off-diagonals represent incorrect predictions.</p> <p>Compute the fraction of correct predictions.</p>	
<pre># Discriminant analysis lda.fit=lda(Direction~Lag1+Lag2,data=Smarket,subset=Smarket\$train); lda.fit; plot(lda.fit); lda.pred=predict(lda.fit,Smarket); lda.class=lda.pred\$class; qda.fit=qda(Direction~Lag1+Lag2,data=Smarket,subset=Smarket\$train);</pre>	<p>Linear discriminant analysis function from package MASS.</p> <p>The result shows 49.2% of the training observations correspond to days during which the market went down.</p> <p>Illustration It produces the linear discriminant, obtained by computing $\beta_1 \text{Lag1} + \beta_2 \text{Lag2}$ for each of the training observations.</p> <p>The prediction function returns 3 variables: class, probability of predictor belonging to this class; and linear discriminant.</p> <p>Check the classes for each training observation.</p> <p>Quadratic discriminant analysis by using the same training data set.</p>	
<pre># KNN library(class); train=(Smarket\$Year<2005); train.X=cbind(Smarket\$Lag1,Smarket\$Lag2)[train,]; test.X=cbind(Smarket\$Lag1,Smarket\$Lag2)[!train,]; train.Direction=Smarket\$Direction[train]; set.seed(1); knn.pred=knn(train.X,test.X,train.Direction,k=3);</pre>	<p>The K-nearest neighbors function is contained in the package “Class”.</p> <p>Consider the historical data before 2005.</p> <p>Input 1: a matrix of predictors X associated with the training data;</p> <p>Input 2: a matrix of predictors X associated with the data we want to make predictions;</p> <p>Input 3: a vector of class labels Y.</p> <p>Input 4: the value of K (here K=3).</p> <p>The result is random, we fix it by using set.seed.</p>	