

R Cheat Sheet for CS4375

Getting Started

<code>getwd()</code>	returns working dir
<code>setwd("dir")</code>	set wd to dir
<code>dir()</code>	show files in current dir
<code>ls()</code>	list objects in workspace
<code>rm("x")</code>	remove x from workspace
<code>?topic</code>	get help on topic
<code>?fun()</code>	get help on function
<code>library(x)</code>	load package x
<code>install.packages("x")</code>	install package x

Load Data

<code>data()</code>	list R built-in data sets
<code>data(iris)</code>	load built-in data set
<code>df <- read.table("x")</code>	read text file
<code>df <- read.csv("x")</code>	read csv file
<code>read.csv("x", header=TRUE)</code>	file has header
<code>read.csv("x", na.strings="NA")</code>	interpret NA

Also: `download.file("http://...", destfile="x")`

Avoid column names with blank spaces.

Common Data Structures

Data Frames

<code>is.data.frame(iris)</code>	returns t/f
<code>df <- as.data.frame(x)</code>	convert to data frame
<code>df <- data.frame(...)</code>	create a data frame

Columns may be of different data types.

Matrix

<code>m <- matrix(1:10,nrow=2,ncol=5)</code>	create matrix
<code>m <- is.matrix(x)</code>	returns t/f
<code>m <- as.matrix(x)</code>	convert to matrix

All columns must be of same data type.

Vector

<code>v <- 1:10</code>	create vector
<code>v <- seq(1,10)</code>	using seq
<code>v <- c("a","b","c")</code>	combine
<code>v <- rep(c(TRUE,FALSE),3)</code>	repeat

All items must be of same data type.

Common Data Types

<code>is.numeric(x)</code>	returns t/f
<code>as.numeric(x)</code>	converts
<code>df\$gender <- factor(df\$gender)</code>	make factor

is. and as. work for numeric, integer, character, etc.

Numeric: integer or real

Integer: integer only

Logical: TRUE, FALSE

Character: "1", "many"

Factor: qualitative data

Indexing Data

Indexing Vectors

<code>v[2]</code>	2nd element of v
<code>v[-2]</code>	all but 2nd element
<code>v[2:5]</code>	elements 2-5
<code>v[c(1:3,5)]</code>	elements 1,2,3,5
<code>v[-c(1:3,5)]</code>	all but el. 1,2,3,5
<code>v[v > 3 & v < 9]</code>	elements in range
<code>v[v %in% c(2,4)]</code>	elements in set

Indexing starts with 1, not 0!

Indexing Matrices

<code>m[2,3]</code>	el. at row 2, col 3
<code>m[2,]</code>	row 2
<code>m[,2]</code>	col 2
<code>m[,c(1,5)]</code>	cols 1, 5

Indexing Data Frames

<code>df\$x</code>	col "x"
<code>df\$x[df\$x>5]</code>	rows where col x > 5

Useful Data Functions

Vectors

<code>class(v)</code>	class
<code>attributes(v)</code>	attributes
<code>length(v)</code>	number of el
<code>mode(v)</code>	returns type
<code>min(v)</code>	minimum value
<code>max(v)</code>	maximum value
<code>range(v)</code>	min and max
<code>mean(v)</code>	average
<code>median(v)</code>	median
<code>sum(v)</code>	sum

Data Frames

<code>dim(df)</code>	dimensions
<code>nrow(df)</code>	num rows
<code>ncol(df)</code>	num cols
<code>names(df)</code>	col names
<code>summary(df)</code>	summary of df
<code>head(df)</code>	1st 6 rows
<code>tail(df, n=2)</code>	last 2 rows
<code>df <- na.omit(df)</code>	omit rows with na

Selecting

<code>which.max(v)</code>	index of max
<code>which.min(v)</code>	index of min
<code>which(v==5)</code>	indices where el is 5
<code>rev(v)</code>	reverse
<code>sort(v)</code>	sort
<code>unique(v)</code>	duplicates removed
<code>subset(df, select=c(2,4))</code>	subset cols 2,4
<code>subset(df, select=-x)</code>	subset all but x
<code>subset(df, y==0, select=-x)</code>	rows where y is 0 . . .

Data Wrangling

<code>with(data, expr)</code>	eval expr on data
<code>by(data, indices, func)</code>	eval by factor level
<code>length(which(df\$x==1))</code>	how many xs are 1
<code>sapply(df, func)</code>	apply func to cols of df
<code>lapply(df, func)</code>	apply func to cols of df
<code>tapply(v1, v2, func)</code>	apply func to v1, grouped by v2

Examples:

```
tapply(iris$Sepal.Length, iris$Species, mean)
length(which(iris$Species=="setosa"))
sapply(iris[,1:3], mean) # returns vector
lapply(iris[,1:3], mean) # returns list
```

Check Correlation

<code>cor(x,y)</code>	correlation
<code>pairs(df)</code>	correlation matrix

Basic Graphs

<code>hist(v)</code>	histogram
<code>plot(x,y)</code>	plot
<code>, xlab="x label"</code>	optional
<code>, ylab="y label"</code>	optional
<code>, main="heading"</code>	optional
<code>, pch=19</code>	point symbol
<code>, cex=1.5</code>	point size
<code>, color="red"</code>	point color
<code>par(mfrow=c(1,2))</code>	1x2 layout
<code>cdplot(y~x)</code>	cond. density

Plot creates scatterplots for numeric data
and box plots for qualitative data

Example:

```
plot(Petal.Length, Petal.Width, pch=21,
     bg=c("red","green3","blue")
     [unclass(Species)], main="Iris Data")
```

Basic Coding

<code>if(cond) expr</code>	if
<code>if(cond) expr1 else expr2</code>	if else
<code>ifelse(test, yes, no)</code>	if-else
<code>for(var in seq) expr</code>	for
<code>while(cond) expr</code>	while
<code>function (args) expr</code>	function
<code># comments</code>	

Loops are often avoided by vector operations.

Example: `5 * v` # multiply each el by 5

Operand Surprises

<code><-</code>	assignment
<code>&</code>	vector and
<code> </code>	vector or
<code>&&</code>	logical and
<code> </code>	logical or

Split Train and Test

```
      by row number
train <- df[1:500]
test  <- df[501:700]
      randomly
set.seed(1234)
i <- sample(nrow(df), floor(nrow(df)*0.8), replace=F)
train <- df[i,]
test  <- df[-i,]
```

Across Algorithms

df <- na.omit(df)	omit rows with na
df\$x <- scale(df\$x)	scale
df\$new <- ifelse(df\$x>n,1,0)	real to binary
I(x^2)	inhibit
contrasts(x)	factor coding
summary(model1)	model summary
anova(model1, model2)	compare models

Common Metrics

accuracy	mean(pred==test\$target)
mse	mean(residuals ²)

Linear Regression

lm1 <- lm(y~x, data=mydata)	returns model
lm1 <- lm(y~., data=mydata)	all predictors
lm1 <- lm(y~x+z, data=mydata)	predictors x and z
lm1 <- lm(y~.-x, data=mydata)	all but x
lm1 <- lm(y~.x*z, data=mydata)	interaction
lm1 <- lm(log(y)~.x+z, data=mydata)	transformation
coef(lm1)	coefficients
confint(lm1)	confidence intervals
summary(lm1)	model summary
predict(lm1, newdata)	predict on new data
plot(lm1)	residual plots

Logistic Regression

```
glm1 <- glm(y~x, family=binomial)
probs <- predict(glm1, newdata=test, type='response')
pred <- ifelse(probs>0.5,1,0)
table(pred=pred, actual=df$response)
mean(pred==df$response)
```

Naive Bayes

```
library(e1071)
nb1 <- naiveBayes(df[, -1], df[, 1], data=train)
pred <- predict(nb1, newdata=test[, -1], type="class")
pred <- predict(nb1, newdata=test[, -1], type="raw")
```

Clustering: kNN

```
      knn regression
library(caret)
fit <- knnreg(df[, -1], df[, 1], k=3)
pred <- predict(fit, test[, -1])
cor(pred, test[, 1])
      knn classify
library(class)
pred <- knn(train=x, test=y, cl=x$Labels, k=3)
mean(pred==x$Labels)
```

Clustering: k-Means

```
km.out <- kmeans(data, centers=5, nstart=25)
fit$centers
```

Clustering: Hierarchical

```
d <- dist(x.scaled)
fit <- hclust(d, method="average")
plot(fit, han=-1, cex=.8)
```

SVM

```
library(e1071)
svm1 <- svm(y~., data=train, kernel="linear", cost=10)
svm1 <- svm(y~., data=train,
            kernel="radial", cost=10, gamma=1)
pred <- predict(svm1, newdata=test[, -1], type="class")
pred <- predict(svm1, newdata=test[, -1], type="raw")
```

Decision Tree

```
      regression tree
library(tree)
tree1 <- tree(y~., df, subset=train)
yhat <- predict(tree1, newdata=test)
mean((yhat-df$response)^2)      MSE
plot(tree1)
text(tree1, pretty=0)
cv.t <- cv.tree(tree1)
plot(cv.t$size, cv.t$dev, type='b')
      classification tree
tree2 <- tree(y~., df)
```