Data Analysis with R in HPC

Weijia Xu, Ph.D Manager, Data Mining & Statistics Texas Advanced Computing Center xwj@tacc.utexas.edu Apr. 04, 2014



Outline

• Introduction of R basics

Data mining and analysis features in R

• Scaling up R with high performance computing resources



Introduction to R Basics



* Based on R tutorial by Lorenza Bordoli

R-project background

Origin and History

- initially written by Ross Ihaka and Robert Gentleman at Dep. of Statistics of U of Auckland, New Zealand during 1990s.
- International project since 1997
- Open source with GPL license
 - Free to anyone
 - In actively development
 - http://www.r-project.org/



What R does

R is a programming environment for statistical and data analysis computations.

Core Package

Statistical functions

plotting and graphics

Data handling and storage

predefined data reader

• textual, regular expressions

hashing

Data analysis functions

Programming support:

loops, branching, subroutines

Object Oriented

More additional developed packages.



Getting Started

- Download and install locally from – http://www.r-project.org/
- TACC
 - ssh stampede.tacc.utexas.edu
 - % module load R
 - % R



R command line interface on cluster

loginl\$ R

```
R version 2.15.1 (2012-06-22) -- "Roasted Marshmallows"
Copyright (C) 2012 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: x86 64-unknown-linux-gnu (64-bit)
```

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

```
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.
```

```
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
```





R Console interface on Mac

DE	0		R-bencl	nmark-25.R			
1	9		gcd2	;) (Q+ Help search		2
90							Merri
243	b <- solve(d	000	R Worksna	ce Browter			.0
245	cumulate <- c	A h A	R WOIKSPA	ce browser			
246	}	a p c				Q+ Help Search	
247	timing <- cumul-	Object	Type	1	Structure		
248	cat(c("Inverse	a	000		R Data Manager		
250	remove("a", "b"	b	(C		-	
251	if (R.Version()	c	Refresh List	Load Data	Q, Search		
252	times[, 2] <-		Data	Package	Description		
254	cat("		AirPassengers	datasets	Monthly Airline Passenger Numbers	19	
235	\n"))		BJsales	datasets	Sales Data with Leading Indicator		
256	cat(" TIT Pr		BJsales.lead (BJsa	datasets	Sales Data with Leading Indicator		
258	cat("		BOD	datasets	Biochemical Oxygen Demand		
259	if (R.Version()		CO2	datasets	Carbon Dioxide Uptake in Grass Plan	its	
261	# (1)		ChickWeight	datasets	Weight versus age of chicks on differ	ren	
262	cumulate <- 0;		DNase	datasets	Elisa assay of DNase		
264	a <- floor(Ru		EuStockMarkets	datasets	Daily Closing Prices of Major Europea	an	
265	invisible(gc(timing <- svs		Formaldehyde	datasets	Determination of Formaldehyde		Q.
267	b <- (phi^a		HairEyeColor	datasets	Hair and Eye Color of Statistics Stude	ent	
268	<pre>})[3] cumulate <- c.</pre>		Harman23.cor	datasets	Harman Example 2.3		
270	}		Harman74.cor	datasets	Harman Example 7.4		
271	times[1, 3] <- t	iming	Indometh	datasets	Pharmacokinetics of Indomethacin		
273	cat(c("3,500,000	Fibonacci numb	InsectSprays	datasets	Effectiveness of Insect Sprays		(3)
275	if (R.Version()\$	los "Win32"	JohnsonJohnson	datasets	Quarterly Earnings per Johnson & Joh	nns	1
277	# (2)	2000. 1	AirPassengers {data	sets}	R Documentatio	an a	
278	for (i in 1:runs	s) {					<u></u>
280	invisible(gc())) em.time({	Monthly	Y Airline 1	Passenger Numbers 1949-1960		
282	b <- rep(1:0	a, a); dim(b) <-	Description				
283 284	b <- 1 / (t(# Rem: this	b) + 0:(a-1)) is twice as fas					
285	5 # a <- 1:a; b <- 1 / outer(The classic Box & Jenkins airline data. Monthly totals of in		. Monthly totals of international airline passengers,				
286	<pre>})[3] cumulate <- cu</pre>	mulate + timina	1949 to 1900.				
288	}		Usage			1	
289	timing <- cumula	ite/runs					- C.
290	cat(c("Creation	of a 3000x3000	lirDsccondare				4
702	nemove("o" "h")						

RStudio: A better user interface of R

• RStudio is a open source graphical user environment for R users.

<u>https://www.rstudio.com/</u>

- RStudio allow users to
 - Interactive code development
 - Run R scripts
 - Exploring local file system
 - Viewing data file
 - Viewing graphical output from R

_ ...



RStudio GUI





RStudio GUI

	Die Site State Mary Plate Service	in format made made stude	/Dropbox/Classe	s/GradSchool - RStudio			000	<u>×</u>
	Q - T - D - A	ou Bullect Bring Tools Helb					GradSchool — Classes	- -
	@_bridge.R × @_HW07.R × @_ProfilePi	Plots R × nyc ×		Workspace History			-0	
	2.12		168 observations of 7 variables	📑 📑 🔄 Import Dataset+ 🧃			Ø	
	Case Restaurant Pr	rice Food Decor Service East	8	Data				
	1 1. Daniella Ristorante 43	3 22 18 20 0		data	11 obs. of 7 variables			
Data Viewor	2 2 Tello's Ristorante 32	12 20 19 19 0	H	ms25	3x2 double matrix		0	
	3 3 Bricchino 54	14 21 13 18 0		ms35	3x2 double matrix			
	5 5 Da Umberto 54	4 24 19 21 0		nyc	168 obs. of 7 variables			Variahle/
	Le Madri 53	2 22 22 21 0		statesp	48 obs. of 9 variables	•		Variabici
	Le Zie 34	14 22 16 21 0		Values	0 7041 0000000007			Environment
	8 8 Pasticcio 34	4 20 18 21 1		le a	1			
	9 9 Belluno 39	19 22 19 22 1		Inc	1=[12]			Viewer
	10 10 Cinque Terre 44	4 21 17 19 1		la d	1-[12]			VIEWEI
	12 12 Marchi's 41	7 21 19 21 1		le ut	le[12]			
	13 13 Nicola Paone 52	2 21 19 20 1		ve:25160	outori c[2]			
	14 14 Notaro 35	15 19 17 19 1		1523100 mc2580	numeric[3]			
	15 15 Rossini's 41	7 20 18 21 1		052300	numeric[3]			
	16 16 Trattoria Alba 31	17 21 19 21 1		ns 2500	numeric[3]			
	17 17 Vila Berula 45	15 22 10 29 1		10.000	0.0			
	16 18 II POSDIO 5/	7 24 21 22 1 18 19 17 18 1		p val	0.35			
	20 20 San Giusto 51	1 22 20 22 1			1		2	
	21 21 Grifone \$4	4 28 20 28 1	~		*		2	
	Console -/Dropbox/Classes/GradSchoo	oll 🐵	-0-	Files Plots Packages Help	P		-0	
	R version 2 0 2 (2012.00.25)	*Erichen Cailing*	2	💽 New Folder 🔍 Delete 📝 Re	name 🛛 🎯 More -		C	
	Copyright (C) 2013 The R Foundat	tion for Statistical Computing	1	🗌 🏠 Home Dropbox Classes	GradSchool 642 Code 00_R-Code_al_files			•
	Platform: x05_64-redhat-linux-gn	nu (64-bit)		A Name		Size	Modified	
	R is free software and comes wit	th ABSOLUTELY NO WARRANTY.				1.6	and the state of the second	
	You are welcome to redistribute	it under certain conditions.		Anasset.R		1.9 KB	Jan 15, 2014, 9:19 MM	
	Type 'license()' or 'licence()'	for distribution details.		ancovplot_Soybean.R		802 bytes	Jan 15, 2014, 9:19 AM	
	Natural lang Install Packages	and the second se		ancovplots.R		1.5 KB	jan 15, 2014, 9:19 AM	
	B is a collabe install from:	(*) Configuring Repositories		A0VUsingR.txt		6 KB	Jan 15, 2014, 9:19 AM	
	Type 'contrabi Repository (CRAN)	-		D boxcox_Crabs.R		.2.4 KB	Jan 15, 2014, 9:19 AM	
	'citation[]' (and a second		D 2) boxcox_Crabs_V2.R		1.4 KB	Mar 4, 2014, 9:54 PM	
	Type 'demo() Packages (separate	multiple with space or comma):		Brand_AOV.R.bt		1.4 KB	jan 15, 2014, 9:19 AM	1
	'help.start()' snow			CapH calculations		848 bytes	lan 15, 2014, 9:19 AM	
	I SNOW D:	And the second se		0) crab cormiote R		781 bites	ian 15, 2014, 9-20 AM	
	[Workspace loa Snowball snowfall R/x86	_64-redhat-linux-gnu-library/3.0 ([+		() () Crabdata CUM D.P.		2 42	100 15, 2014, 0.10 MM	
	> y1 < c[1,2, snowFT					J ND	(an 15, 2014, 513 Her	
	> x1 < c(1,2, minstall dependen	ncies		Crabpower,trans.R		352 bytes	Jan 15, 2014, 9:20 AM	
Deelvere	> lm.w3 <- lm			P] FractionalFact.R		400 bytes	Jan 15, 2014, 9:20 AM	Eile
Package	> vif(lm.w3)	Install Cancel		PractionalFact_9-3.R		256 by	100 15 001 (0 30 W)	File
	> library[ca			C 2 friedman.R		618 b)		
wanagement	1,5,6)			GoldCondMeanPlot.R		753 bytes	Jan 15, 2014, 9;20 AM	vianagement
	> x2 < c('a', 'a', 'a', 'b', 'b	br, (br)		H04_Hmatrix.R		323 bytes	jan 15, 2014, 9:20 AM	
	> lm.w3 <- lm(yi-x1+x2)			2 HSU.R		2.1 KB	Jan 15, 2014, 9:20 AM	
	x1 x2			E & KruskalWallis_Crabs.R		1.3 KB	Jan 15, 2014, 9:21 AM	
	1.375 1.375		Ŷ	The Knickshifelie Physics		120 bites	150 15 2014 0-21 AM	
		8 🖸 🔊 🕸 😖 💿 👘	awaling bath=1.	R Buddy List Task	Coach +/hon 🔒 Estudio Sidec - 👔 Documente - Div	426 heter 8 🔏 🍯 🕷	100 15 2014 0-21 AM	



Pa M

RStudio GUI





Web interface of Rstudio on Maverick

- Users can run an interactive web session with RStudio using maverick.tacc.utexas.edu
- The Job script template is in /share/doc/slurm/ job.RStudio

• On Maverick,

- Submitting the job with
 - Sbatch /share/doc/slurm/job.Rstudio
- After the job is running, a URL will be available for connection e.g.
 - <u>http://maverick.tacc.utexas.edu:12173</u>



• Then visit the URL and log in with your account credential

maverick.tacc.utexas.edu:12173/auth-sign-in		Google	Q 💽 🖡 👘
Rstudio			
	Sign in to RStudio Username:		
	Password:		



900	RStudio			EN.
RStudio +				
maverick.tacc.utexas.edu:12173		☆ ァ C 🛛 🚼 ▼ Google		Q 💽 - 🛃 🍙
File Edit Code View Plots Session Build	Debug Tools Help			xwj Sign Out
🔍 🔍 + 🚰 + 🔒 🗿 🗁 🖙 Ca so file/function				😰 Project: (None) 👻
Untitled1 ×	-0	Environment History		
🕞 🕞 Source on Save 🛛 💁 🖉 🗸	Run 🖘 Source ᠇ 📃	🚰 🔒 📑 To Console 🔤 To Source	e 🛛 🎸	Q
		Files Plots Packages Help Viewe New Folder Upload Delete Home Rprofile job.Rstudio	er Rename Size 232 B S.5 KB	More - C Modified Apr 3, 2014, 4:04 PM Apr 3, 2014, 4:03 PM
1.1 D (top Level) +	K Script +	Rstudio.out	1.5 KB	Apr 3, 2014, 9:13 PM
>				



Basic Math operations

R as a calculator
 +, -, /, *, ^, log, exp, ...

```
> (17*0.35)^(1/3)
[1] 1.812059
> log2(128)
[1] 7
> exp(1)
[1] 2.718282
> 3^-1
[1] 0.3333333
```



Variables

• Numeric



Character String



Logical

> c=(1+1==3) > c [1] FALSE



Assigning Values to Variables

•	"<-" or "="		> a=0	:(1, 2, 4, 7,	9)
	> a=4		> a		
	> a		[1] 1	. 2 4 7 9	
	[1] 4			$\geq 2\pi a a a n \langle \rangle$	
	> a<-40			1. 9	
	> a			2: 7	
	[1] 40			3: 4	
•	Assign mu	tiple values		4: 2	
	– Concaten	ate, c()		5: 1 6:	
	– From stdi	n, scan()		Read 5 items	
	– Series	> a=(1:6)		> a	
	• :	> a		[1] 9 7 4 2 1	
	• Sea()	[1] 1 2 3 4 5 6			
		> a=seq(1,6,0.5)			
		/ a [1] 1.0 1.5 2.0 2	.5 3.0 3.	5 4.0 4.5 5.0 5.	.5 6.0



NA: Missing Value

• Variables of each data type (numeric, character, logical) can also take the value NA: not available.

- NA is not the same as 0
- NA is not the same as ""
- NA is not the same as FALSE

•Any operations (calculations, comparisons) that involve NA may or may not produce NA:

> NA [1] NA	> NA TRUE [1] TRUE	
<pre>> 1+NA [1] NA [1] NA > log(NA) [1] NA</pre>	<pre>> NA FALSE [1] NA > NA & TRUE [1] NA > NA & FALSE</pre>	<pre>> max(c(1,2,3, NA)) [1] NA > max(c(1,2,3,NA), na.rm=T) [1] 3</pre>
	[1] FALSE	



Basic Data Structure

- Vector
 - an ordered collection of data of the same type
 - a single number is the special case of a vector with 1 element.
 - Usually accessed by index
- Matrix
 - A rectangular table of data of the same type

```
> a = c(1,2,3)
> a
[1] 1 2 3
> a[1]
[1] 1
> a[2]
[1] 2
> a*2
[1] 2 4 6
```



Basic Data Structure

• List

- an ordered collection of data of arbitrary types.
- name-value pair
- Accessible by name

```
> doe = list(name="john", age=28, married=F)
> doe$name
[1] "john"
> doe$age
[1] 28
> doe$married
[1] FALSE
> doe[1]
$name
[1] "john"
```



Basic Data Structure

• Hash Table

- In R, a hash table is the same as a workspace for variables, which is the same as an environment.
- Store Key-value pairs.
- Value can be accessed by key

```
> tab = new.env(hash=T)
> assign("A", list(id=682, description="six eight two"), env=tab)
> assign("B", list(id=77, description="seven seven"), env=tab)
> assign("C", list(id=77, mykey="the key is C"), env=tab)
> get("A", env=tab)
$id
[1] 682
$description
[1] "six eight two"
> get("C", env=tab)
$id
[1] 77
$mykey
[1] "the key is C"
```



Dataframes

R handles data in objects known as dataframes

 rows: data items;

- columns: values of the different attributes
 - Values in each column should be from the same type.

	Area	Slope	Vegetation	Soil.pH	Damp	Worm.density
Silwood.Bottom	5.1	2	Arable	5.2	FALSE	7
Gunness.Thicket	3.8	0	Scrub	4.2	FALSE	6
Oak.Mead	3.1	2	Grassland	3.9	FALSE	2
North.Gravel	3.3	1	Grassland	4.1	FALSE	1
South.Gravel	3.7	2	Grassland	4.0	FALSE	2
Pond.Field	4.1	0	Meadow	5.0	TRUE	6
Water.Meadow	3.9	0	Meadow	4.9	TRUE	8
Pound Hill	4.4	2	Arable	4.5	FALSE	5



Read Dataframes From File



- Read tab-delimited file directly.
- Variable name in header row cannot have space.
- To see the content of the dataframes (object) just type is name:
 - > worms



Selecting Data from Dataframes

• Subscripts within square brackets

- , means "all the rows" and
- ,1 means "all the columns"

• To select the first three column of the dataframe

\sim	LTO.	20200	م ۲	- 1 -	- 21
	wo	T 10	211	, ÷ •	

	Area	Slope	Vegetation
Silwood.Bottom	5.1	2	Arable
Gunness.Thicket	3.8	0	Scrub
Oak.Mead	3.1	2	Grassland
North.Gravel	3.3	1	Grassland
South.Gravel	3.7	2	Grassland
Pond.Field	4.1	0	Meadow
Water.Meadow	3.9	0	Meadow
Pound.Hill	4.4	2	Arable



Selecting Data from Dataframes

names()

- Get a list of variables attached to the input name

> names(worms)										
[1]	"Area"	"Slope"	"Vegetation"							
[4]	"Soil.pH"	"Damp"	"Worm.density"							

attach()

- Make the variables accessible by name:

> attach(worms)



Selecting Data from Dataframes

• Using logic expression while selecting:

	Area	Slope	Vegetation	Soil.pH	Damp	Worm.density				
Silwood.Bottom	5.1	2	Arable	5.2	FALSE	7				
Gunness.Thicke	t 3.8	0	Scrub	4.2	FALSE	6				
Oak.Mead	3.1	2	Grassland	3.9	FALSE	2				
North.Gravel	3.3	1	Grassland	4.1	FALSE	1				
South.Gravel	3.7	2	Grassland	4.0	FALSE	2				
Pond.Field	4.1	0	Meadow	5.0	TRUE	6				
Water.Meadow	3.9	0	Meadow	4.9	TRUE	8				
Pound.Hill	4.4	2	Arable	4.5	FALSE	5				
> worms[Area>4&Slope<1,]										
Ar	ea Slo	pe Ve	getation S	oil.pH	Damp V	Worm.density				
Pond.Field 4	.1	0	Meadow	5	TRUE	6				



Selecting Data From a Dataframe

More examples:

> worms[Damp]	,]							
	Area	Slope	Vegetation	Soil.pH	Damp	Worm.density		
Pond.Field	4.1	0	Meadow	5.0	TRUE	6		subset rows by a
Water.Meadow	3.9	0	Meadow	4.9	TRUE	8		logical vector
> worms\$Veget	tatior	1						
[1] Arable	Scru	ıb	Grassland 0	Grassland	d Gras	ssland Meadow	Meadow	subset a column
[8] Arable								
Levels: Arabi	le Gra	assland	d Meadow Sci	rub				
> worms\$Veget	tatior	l≕"Gra	assland"					comparison resulting
[1] FALSE FAI	LSE 1	TRUE 1	TRUE TRUE H	FALSE FAI	LSE FA	ALSE		in logical vector
> worms[worm	ms\$Veç	getatio	on=="Grassla	and",]				
	Area	Slope	Vegetation	Soil.pH	Damp	Worm.density		aubaat tha
Oak.Mead	3.1	2	Grassland	3.9	FALSE	2		Subset the
North.Gravel	3.3	1	Grassland	4.1	FALSE	1		selected rows
South.Gravel	3.7	2	Grassland	4.0	FALSE	2		



Sorting Data in Data frames

order()

State the Area for sorting order

State columns to be sorted

> worms[order(worms[,1]),1:6]

	Area	Slope	Vegetation	Soil.pH	Damp	Worm.density
Oak.Mead	3.1	2	Grassland	3.9	FALSE	2
North.Gravel	3.3	1	Grassland	4.1	FALSE	1
South.Gravel	3.7	2	Grassland	4.0	FALSE	2
Gunness.Thicket	3.8	0	Scrub	4.2	FALSE	6
Water.Meadow	3.9	0	Meadow	4.9	TRUE	8
Pond.Field	4.1	0	Meadow	5.0	TRUE	6
Pound.Hill	4.4	2	Arable	4.5	FALSE	5
Silwood.Bottom	5.1	2	Arable	5.2	FALSE	7



Sorting Data in Dataframes

More on sorting selected

sorted in descending order

<pre>> worms[rev(order(worms[,4])),c(4,6)]</pre>				
	Soil.pH Wo	orm.density		
Silwood.Bottom	5.2	7		
Pond.Field	5.0	6		
Water.Meadow	4.9	8		
Pound.Hill	4.5	5		
Gunness.Thicket	4.2	6		
North.Gravel	4.1	1		
South.Gravel	4.0	2		
Oak.Mead	3.9	2		



Flow Control



```
for(i in 1:10) {
    print(i*i)
}
```

```
i=1
while(i<=10) {
    print(i*i)
    i=i+sqrt(i)
}</pre>
```



Flow Control

• apply (arr, margin, fct)

 Applies the function fct along some dimensions of the vector/ matrix arr, according to margin, and returns a vector or array of the appropriate size.

> m				
	Soil.pH Worm.dens	sity		
Silwood.Bottom	5.2	7		
Pond.Field	5.0	6		
Water.Meadow	4.9	8		
Pound.Hill	4.5	5		
Gunness.Thicket	4.2	6		
North.Gravel	4.1	1		
South.Gravel	4.0	2		
Oak.Mead	3.9	2		
> apply(m, 1, su	um)			
Silwood.Bottom	Pond.Field	Water.Meadow	Pound.Hill	Gunness.Thicket
12.2	11.0	12.9	9.5	10.2
North.Gravel	South.Gravel	Oak.Mead		
5.1	6.0	5.9		
> apply(m, 2, su	um)			
Soil.pH Wor	cm.density			
35.8	37.0			



Flow Control

 lapply (list, fct) and sapply (list, fct) To each element of the list li, the function fct is applied. The result is a list where elements are the individual. 	
fct results.	<pre>> lapply(1:5, fct) [[1]] [1] 1 1 1</pre>
 Sapply, converting results into a vector or array of appropriate size 	[[2]] [1] 2 4 8
<pre>> fct = function(x) { return(c(x, x*x, x*x*x)) } > sapply(1:5, fct)</pre>	[[3]] [1] 3 9 27
[,1] [,2] [,3] [,4] [,5]	[[4]] [1] 4 16 64
[1,] 1 2 3 4 5 [2,] 1 4 9 16 25	[[5]]
[3,] 1 8 27 64 125	[1] 5 25 125



Create Statistical Summary

- Descriptive summary for numerical variables:
 - arithmetic mean;
 - maximum, minimum, median, 25 and 75 percentiles (first and third quartile);
- Levels of categorical variables are counted

> summary(worm	3)				
Area	Slope	Vegetation	Soil.pH	Damp	Worm.density
Min. :3.100	Min. :0.000	Arable :2	Min. :3.900	Mode :logical	Min. :1.000
1st Qu.:3.600	1st Qu.:0.000	Grassland:3	1st Qu.:4.075	FALSE:6	1st Qu.:2.000
Median :3.850	Median :1.500	Meadow :2	Median :4.350	TRUE :2	Median :5.500
Mean :3.925	Mean :1.125	Scrub :1	Mean :4.475	NA's :0	Mean :4.625
3rd Qu.:4.175	3rd Qu.:2.000		3rd Qu.:4.925		3rd Qu.:6.250
Max. :5.100	Max. :2.000		Max. :5.200		Max. :8.000



Create Plots

plot(...)
Create scatter plot.

> plot(Area, Soil.pH)

Automatically create a postscript file with default name





Other Common Plots

- Univariate:
 - histograms,
 - density curves,
 - Boxplots, quantile-quantile plots
- Bivariate:
 - scatter plots with trend lines,
 - side-by-side boxplots
- Several variables:
 - scatter plot matrices, lattice
 - 3-dimensional plots,
 - heatmap






Saving your work

- history(Inf)
 - To review the command lines entered during the sessions
- savehistory("history.txt")
 - Save the history of command lines to a text file
- loadhistory("history.txt")
 - read it back into R
- save(list=ls(),file="all.Rdata")
 - The session as a whole can be saved as a binary file.
- load("c:\\temp\\ all.Rdata")
 - Read back saved sessions.



Importing and exporting data

There are many ways to get data into R and out of R.

Most programs (e.g. Excel), as well as humans, know how to deal with rectangular tables in the form of tabdelimited text files.

> x = read.delim("filename.txt")
also: read.table, read.csv

> write.table(x, file="x.txt", sep="\t")



Getting help

"?" Or "help" Details about a specific command whose name you know (input arguments, options, algorithm, results):

```
e.g.
                               t.test
                                                          package:ctest
                                                                                             R Documentation
    >? t.test
                               Student's t-Test
                               Description:
                                     Performs one and two sample t-tests on vectors of data.
    or
                               Usage:
                                    t.test(x, y = NULL, alternative = c("two.sided", "less", "greater"),
                                             mu = 0, paired = FALSE, var.equal = FALSE,
conf.level = 0.95, ...)
    >help(t.test)
                                    t.test(formula, data, subset, na.action, ...)
                               Arguments:
                                       x: a numeric vector of data values.
                                       y: an optional numeric vector data values.
                               alternative: a character string specifying the alternative hypothesis,
must be one of `"two.sided"' (default), `"greater"' or
`"less"'. You can specify just the initial letter.
                                      mu: a number indicating the true value of the mean (or difference
                                           in means if you are performing a two sample test).
                                 paired: a logical indicating whether you want a paired t-test.
                               var.equal: a logical variable indicating whether to treat the two
```



Data Mining with R



Data mining with R

- Many data mining methods are also supported in R core package or in R modules
 - Kmeans clustering:
 - Kmeans()
 - Decision tree:
 - rpart() in rpart library
 - Nearest Neighbour
 - Knn() in class library



. . .

Additional Libraries and Packages

- Libraries
 - Comes with Package installation (Core or others)
 - library() shows a list of current installed
 - library must be loaded before use e.g.
 - library(rpart)
- Packages
 - Developed code/libraries outside the core packages
 - Can be downloaded and installed separately
 - Install.package("name")
 - There are currently 2561 packages at <u>http://cran.r-project.org/web/packages/</u>
 - E.g. Rweka, interface to Weka.



Common Data Mining Methods

- Clustering Analysis
 - Grouping data object into different bucket.

- Classification
 - Assigning labels to each data object.
 - Requires training data.



Cluster Analysis

- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups
 - Inter-cluster distance: maximized
 - Intra-cluster distance: minimized





K-means Clustering

- Partitional clustering approach
- Each cluster is associated with a centroid (center point)
- Each point is assigned to the cluster with the closest centroid
- Number of clusters, K, must be specified
- The basic algorithm is very simple
 - 1: Select K points as the initial centroids.
 - 2: repeat
 - 3: Form K clusters by assigning all points to the closest centroid.
 - 4: Recompute the centroid of each cluster.
 - 5: **until** The centroids don't change



An Example of k-means Clustering





Examples are from Tan, Steinbach, Kumar Introduction to Data Mining



K-means clustering Example

login1% more kmeans.R x<-read.csv("../data/cluster.csv",header=F) fit<-kmeans(x, 2) plot(x,pch=19,xlab=expression(x[1]), ylab=expression(x[2])) points(fit\$centers,pch=19,col="blue",cex=2) points(x,col=fit\$cluster,pch=19)



> fitK-means clustering with 2 clusters of sizes 49, 51

Cluster means: V1 V2 1 0.99128291 1.078988 2 0.02169424 0.088660

Within cluster sum of squares by cluster: [1] 9.397754 7.489019

Available components: [1] "cluster" "centers" "withinss" "size" >







Classification Tasks





Support Vector Machine Classification

- A distance based classification method.
- The core idea is to find the best hyperplane to separate data from two classes.
- The class of a new object can be determined based on its distance from the hyperplane.



Binary Classification with Linear Separator

- Red and blue dots are representations of objects from two classes in the training data
- The line is a linear separator for the two classes
- The closets objects to the hyperplane is the support vectors.





SVM Classification Example install.packages("e1071") library(e1071) train<-read.csv("sonar_train.csv",header=FALSE)</pre> y<-as.factor(train[,61]) x<-train[,1:60] fit<-svm(x,y)</pre> 1-sum(y==predict(fit,x))/length(y))



SVM Classification Example

test<-read.csv("sonar_test.csv",header=FALSE) y_test<-as.factor(test[,61]) x_test<-test[,1:60]

1-sum(y_test==predict(fit,x_test))/length(y_test)



Scaling up R computation with high performance computing resources













What to do if the computation is too big for a single desktop

- A common user question:
 - I have an existing R solution for my research work.
 But the data is growing to big. Now my R program runs days to finish (/runs out of memory)
- Three strategies
 - Using automatically offloading with multicore/GPU/ MIC.
 - Break big computation with multiple job submission
 - Implement code using parallel packages.



Automatic offloading with latest hardware

- R is originally designed as for single thread execution.
 - Slow performance
 - Not scalable with large data
- R can be built and linked to library utilizes latest multiple core technology for automatic parallel execution for some operations, most commonly, linear algebra related computations.



Dynamic Library & R

 MKL provides BLAS/LAPACK routines that can "offload" to the Xeon Phi Coprocessor, reducing total time to solution





Automatic offloading with latest hardware

- Hardware supported:
 - Multiple cores on CPU
 - Intel Xeon Phi coprocessor (on Stampede)
 - GPGPU (on Stampede/Maverick)
- Libraries supporting automatic offloading
 - Intel Math Kernel Library (MKL)
 - Available on stampede and maverick for users
 - HiPlarB
 - Open source and freely available
 - http://www.hiplar.org/hiplar-b.html



R-2.5 benchmark performance with automatically hardware accelerration





- Advantage:
 - No code changes needed
 - User can run R solution as before without knowledge of the parallel execution.
- Limitations:
 - Only support limited computational operations.



Break Big Computations with multiple R jobs

- Running R in non-interactive session
- User can submit multiple R jobs with different command Line parameters
 - Similar to run R batch mode
 - Parameters is specified on the command line
 - Good for repeated runs of same computations or running script partially



Running R Session in Batch Mode

• R scripts

 Put the codes you would input when running interactively into a text file. e.g.

login1% more mtcars.Rdata(mtcars)# load built-in mtcars data tableattach(mtcars)# Attaching mtcars namesnames(mtcars)# show column namessummary(mtcars)# show statistical summary of all columns.detach()q()

`Batch mode

- ">R CMD BATCH /path/to/R_SCRIPT"
- Running R script stored in file "R_SCRIPT"
- By default the result is stored in R_SCRIPTOut



Running R Session in Batch Mode

login1% R CMD BATCH mtcars.R login1% more mtcars.Rout

```
R version 2.10.1 (2009-12-14)
Copyright (C) 2009 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
```

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

```
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.
```

```
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
```

[Previously saved workspace restored]

```
> data(mtcars) # load built-in mtcars data table
> attach(mtcars) # Attaching mtcars names
> names(mtcars) # show column names
[1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am" "gear"
[11] "carb"
```



Running R Session in Batch Mode

> names(mtcars) # show column names				
[1] "mpg" "cyl	" "disp" "hp"	"drat" "wt" '	'qsec" "vs" "a	m" "gear"
[11] "carb"				
> summary(mtcars) # show statistical summary of all columns.				
mpg	cyl	disp	hp	
Min. :10.40	Min. :4.000	Min. : 71.1	Min. : 52.0	
1st Qu.:15.43	1st Qu.:4.000	1st Qu.:120.8	1st Qu.: 96.5	
Median :19.20	Median :6.000	Median :196.3	Median :123.0	
Mean :20.09	Mean :6.188	Mean :230.7	Mean :146.7	
3rd Qu.:22.80	3rd Qu.:8.000	3rd Qu.:326.0	3rd Qu.:180.0	
Max. :33.90	Max. :8.000	Max. :472.0	Max. :335.0	
drat	wt	qsec	vs	
Min. :2.760	Min. :1.513	Min. :14.50	Min. :0.0000	
1st Qu.:3.080	1st Qu.:2.581	1st Qu.:16.89	1st Qu.:0.0000	
Median :3.695	Median :3.325	Median :17.71	Median :0.0000	
Mean :3.597	Mean :3.217	Mean :17.85	Mean :0.4375	
3rd Qu.:3.920	3rd Qu.:3.610	3rd Qu.:18.90	3rd Qu.:1.0000	
Max. :4.930	Max. :5.424	Max. :22.90	Max. :1.0000	
am	gear	carb		
Min. :0.0000	Min. :3.000	Min. :1.000		
1st Qu.:0.0000	1st Qu.:3.000	1st Qu.:2.000		
Median :0.0000	Median :4.000	Median :2.000		
Mean :0.4062	Mean :3.688	Mean :2.812		
3rd Qu.:1.0000	3rd Qu.:4.000	3rd Qu.:4.000		
Max. :1.0000	Max. :5.000	Max. :8.000		
> detach()				
> q()				
> proc.time()				
user system elapsed				
0.192 0.014	0.212			



```
login1% cat sample.R
          arg1 <-as.numeric(commandArgs()[4])
                                                     Parse arguments
          arg2 <- as.numeric(commandArgs()[5])</pre>
          paste("Input arguments are ", arg1, arg2, sep=" ")
interactive
                                                             -Do something
          paste("The sum is ", arg1+arg2, sep="")
          q()
          n
          login1% cat sample.R | R --slave --args 1 2
          [1] "Input arguments are 1 2"
          [1] "The sum is 3"
          login1% cat sample.R | R --slave --args 1231234 54532332
          [1] "Input arguments are 1231234 54532332"
          [1] "The sum is 55763566"
```



Like

mode

Running R Script with Parameters

- Enable more flexibility on computations of same R script.

q()



Running R Script with Parameters

login1% cat mtcars2.R | R --slave --args mtcars [1] "mpg" "cvl" "disp" "hp" "drat" "wt" "gsec" "vs" "am" "gear" [11] "carb" disp hp mpg cvl Min. :10.40 Min. :4.000 Min. : 71.1 Min. : 52.0 1st Ou.:15.43 1st Ou.:4.000 1st Ou.:120.8 1st Ou.: 96.5 Median :19.20 Median :6.000 Median :196.3 Median :123.0 Mean :20.09 Mean :6.188 Mean :230.7 Mean :146.7 3rd Qu.: 22.80 3rd Qu.: 8.000 3rd Qu.: 326.0 3rd Qu.: 180.0 Max. :33.90 Max. :8.000 Max. :472.0 Max. :335.0 drat wt asec VS Min. :2.760 Min. :1.513 Min. :14.50 Min. :0.0000 1st Ou.: 3.080 1st Ou.: 2.581 1st Ou.: 16.89 1st Ou.: 0.0000 Median :3.695 Median :3.325 Median :17.71 Median :0.0000 Mean :3.597 Mean :3.217 Mean :17.85 Mean :0.4375 3rd Ou.: 3.920 3rd Ou.: 3.610 3rd Ou.: 18.90 3rd Ou.: 1.0000 Max. :4.930 Max. :5.424 Max. :22.90 Max. :1.0000 carb am gear Min. :0.0000 Min. :3.000 Min. :1.000 1st Ou.:0.0000 1st Ou.:3.000 1st Ou.:2.000 Median :0.0000 Median :4.000 Median :2.000 Mean :0.4062 Mean :3.688 Mean :2.812 3rd Ou.:1.0000 3rd Ou.:4.000 3rd Ou.:4.000 Max. :1.0000 Max. :5.000 Max. :8.000 login1% cat mtcars2.R | R --slave --args mpg [1] "mpg" "cyl" "disp" "hp" "drat" "wt" "gsec" "vs" "am" "gear" [11] "carb" Min. 1st Ou. Median Mean 3rd Ou. Max. 10.40 15.42 19.20 20.09 22.80 33.90 login1% cat mtcars2.R | R --slave --args cyl [1] "mpg" "cvl" "disp" "hp" "drat" "wt" "gsec" "vs" "am" "gear" [11] "carb" Min. 1st Ou. Median Mean 3rd Ou. Max. 4.000 4.000 6.000 6.188 8.000 8.000



Text Analysis of HathiTrust Corpus ('tm' package, ~1M books)



Guangchen Ruan, Hui Zhang, et al. <u>http://www.hathitrust.org/htrc</u>


Advantages

- Users only need to develop job submission scripts
- Each job can use existing R code
- Good for repeated analysis with different data set or many independent analysis tasks over large data set.
- Limitations
 - A "data-parallel" solution that may not suitable for simulation based workflow



Running R with parallel packages

- There are many parallel packages available to enable parallelism with R
- Two most common approaches included with R distribution
 - Multicore
 - Snow/Rmpi







Multicore

- Utilizes multiple processing core within the same node.
- Replace several common functions with parallel implementations
- No need of significant changes on the existing coding process control.
- Scalability is limited by the number of core and memory available within single node



Multicore -- mcapply

- Iapply → mcapply
 - lapply(1:30, rnorm)
 - mclapply(1:30, rnorm)
- mc.cores
 - The maximum number of cores to use
- mc.preschedule
 - TURE, computation is first divided by the number of cores.
 - FALSE, one job is spawned for each value sequentially



Multicore –parallel and collect

 parallel(expr, name, mc.set.seed = FALSE, silent = FALSE)

- Starts a parallel process for evaluating expr,

- collect(jobs, wait = TRUE, timeout = 0, intermediate = FALSE)
 - Collects the result from the parallel process.

p <- parallel(1:10)
q <- parallel(1:20)
collect(list(p, q)) # wait for jobs to finish and collect all results</pre>



Snow

- Developed Based on Rmpi package,
- Simplify the process to initialize parallel process over cluster.

```
cl <- makeCluster(4, type='SOCK')</pre>
```

```
birthday <- function(n) {
    ntests <- 1000
    pop <- 1:365
    anydup <- function(i)
    any(duplicated(
        sample(pop, n,replace=TRUE)))
    sum(sapply(seq(ntests), anydup)) / ntests}</pre>
```

x <- foreach(j=1:100) %dopar% birthday (j)

stopCluster(cl)

Ref: http://www.rinfinance.com/RinFinance2009/presentations/UIC-Lewis%204-25-09.pdf



Snow

- Provide similar MPI functions on snow cluster:
 - clusterSplit, clusterCall, ClusterEvalQ, clusterApply,

```
clusterApply(cl, 1:2, get("+"), 3)
clusterEvalQ(cl, library(boot))
x<-1
clusterExport(cl, "x")
clusterCall(cl, function(y) x + y, 2)</pre>
```



Snow

• Provide parallel version of common functions:

- parLapply, parApply, parSapply
- Similar to mcapply from mutlicore
- Need to setup the snow cluster first

cl <- makeCluster(4, type='SOCK')
parSapply(cl, 1:20, get("+"), 3)</pre>



- Advantage
 - Do whatever you want with them
 - Get the best performance
- Limitations
 - Need code development
 - In some case, the analysis workflow may need be changed.



Further references

• R

- M. Crawley, Statistics An Introduction using R, Wiley
- J. Verzani, SimpleR Using R for Introductory Statistics http://cran.r-project.org/doc/contrib/Verzani-SimpleR.pdf
- Programming manual:
 - http://cran.r-project.org/manuals.html

- Using R for data mining
 - Data Mining with R: Learning with case studies, Luis Togo
- Contact Info
 - Weijia Xu <u>xwj@tacc.utexas.edu</u>

