# R for Bioinformatics
## Introduction, Programming, Data Analysis and Visualization
## Introduction to Data Analysis and R

Gang Chen
chengang@bgitechsolutions.com

November 30, 2013

# Outline

1. Data Analysis

2. Data Aanlysis and R

3. Hello R!

4. Development Environment

5. References

# Next

# Data Analysis

### Wikipedia

Analysis of data is a process of **inspecting**, **cleaning**, **transforming**, and **modeling** data with the goal of discovering useful information, suggesting conclusions, and supporting decision making.

## Data Analysis

1. Collecting $\rightarrow$
2. Cleaning $\rightarrow$
3. Transforming $\rightarrow$
4. Modeling $\rightarrow$
5. Visualizing $\rightarrow$
6. Knowledge

## Biological Data Analysis

1. Sequencing $\rightarrow$
2. QC $\rightarrow$
3. Alignment ...$\rightarrow$
4. GWAS, EWAS ...$\rightarrow$
5. Manhattan Plot, Q-Q plot ...$\rightarrow$
6. Paper?

# Next

# What is R?

### R

R is a **free** software environment for statistical computing and graphics.
----R-project.org

## Data Analysis

1. Collecting →
2. Cleaning →
3. Transforming →
4. Modeling →
5. Visualizing →
6. Knowledge

## R

- **Rcurl...**
- **gsub, unique...**
- **reshape...**
- **e1071...**
- **ggplot2...**
- **knitr...**

# Biological Data Analysis and R

## Biological Data Analysis

1. Sequencing $\rightarrow$
2. QC $\rightarrow$
3. Alignment ...$\rightarrow$
4. GWAS, EWAS ...$\rightarrow$
5. Manhattan Plot, Q-Q plot ...$\rightarrow$
6. Paper?

## R

- **Rsamtools, Affy...**
- **genomicRanges...**
- **reshape...**
- **e1071...**
- **ggbio...**
- **knitr, shiny...**

## History

- Version 0.16, This is the last alpha version developed primarily by Ihaka and Gentleman. Much of the basic functionality from the "White Book" (see S history) was implemented. The mailing lists commenced on April 1, 1997.

- Version 0.49, April 23, 1997, This is the oldest available source release, and compiles on a limited number of Unix-like platforms. CRAN is started on this date, with 3 mirrors that initially hosted 12 packages. Alpha versions of R for Microsoft Windows and Mac OS are made available shortly after this version.

- Version 0.60, December 5, 1997, R becomes an official part of the GNU Project. The code is hosted and maintained on CVS.

## History

- Version 1.0.0, February 29, 2000, Considered by its developers stable enough for production use.[28]
- Version 1.4.0, S4 methods are introduced and the first version for Mac OS X is made available soon after.
- Version 2.0.0, October 4, 2004, Introduced lazy loading, which enables fast loading of data with minimal expense of system memory.
- Version 2.1.0, Support for UTF-8 encoding, and the beginnings of internationalization and localization for different languages.
- Version 2.11.0, April 22, 2010, Support for Windows 64 bit systems.
- Version 2.13.0, April 14, 2011, Adding a new compiler function that allows speeding up functions by converting them to byte-code.

## History

- Version 2.14.0, October 31, 2011, Added mandatory namespaces for packages. Added a new parallel package.
- Version 2.15.0, March 30, 2012, New load balancing functions. Improved serialization speed for long vectors.
- Version 3.0.0, April 3, 2013, Support for numeric index values 231 and larger on 64 bit systems.

## R in China

- 2004, official documents are translated into Chinese
- 2006, some books on R in Bioinformatics
- 2008, the first R conference was hold at Renming University, Beijing.
- 2009 to 2013, China R Conference is hold at Beijing and Shanghai each year.
- 2012, popular R books are translated into Chinese.
- 2013, R in Action ggplot2 R in a nutshell ...are published in China.
- 2013, CUHK-R course is launched.

# Applications of R

## Applications

- Statistical analysis
- Data Mining
- Life Science
- Business Intelligence
- Data Visualization
- Social Network

- eCommerce
- Integrated Circuit
- Financial
- Media
- Consoluting
- ...

# Attendee of R community

## Pros and Cons

### Bo Cowgill, Google

``The best thing about R is that it was developed by statisticians. The worst thing about R is that ... it was developed by statisticians.''

# Next

# Hello R!

```
print("Hello R!")

## [1] "Hello R!"
```

# Hello Statistical Analysis

```
data(mtcars)
cor(mtcars$mpg, mtcars$wt)

## [1] -0.8677
```

# Hello Plot

```
data(mtcars)
plot(mtcars$mpg, mtcars$wt)
```

## Hello Plot

```
data(mtcars)
plot(mtcars$mpg, mtcars$wt, pch = 19)
```

# Hello Plot

```r
data(mtcars)
plot(mtcars$mpg, mtcars$wt, pch = 19, col = mtcars$gear)
```

## Hello Plot

```
data(mtcars)
plot(mtcars$mpg, mtcars$wt, pch = 19, col = mtcars$gear, xlab = "Mil
    ylab = "Weight")
```

## Hello Plot

```
data(mtcars)
plot(mtcars$mpg, mtcars$wt, pch = 19, col = mtcars$gear, xlab = "Mil
    ylab = "Weight", main = "Cars")
```

# Hello Plot

```r
data(mtcars)
plot(mtcars$mpg, mtcars$wt, pch = 19, col = mtcars$gear, xlab = "Mile Per Gallon",
     ylab = "Weight", main = "Cars", cex = 2)
lines(loess.smooth(mtcars$mpg, mtcars$wt), col = rgb(1, 0, 0, 0.5), lwd = 10)
```

# Hello Plot

```r
data(mtcars)
plot(mtcars$mpg, mtcars$wt, pch = 19, col = mtcars$gear, xlab = "Mile Per Gallon",
     ylab = "Weight", main = "Cars", cex = 2)
lines(loess.smooth(mtcars$mpg, mtcars$wt), col = rgb(1, 0, 0, 0.5), lwd = 10)
text(30, 5, paste("Cor", round(cor(mtcars$mpg, mtcars$wt), 4), sep = ":"), cex = 2)
```

# Next

# Download and Installation

### Download

# CRAN

### Installation

- R: Linux(apt, yum), Mac OS, Windows
- Rtools: Windows
- packages: CRAN, devtools, github, local file

## R Commands

### R Commands

`R CMD command args`

``command'':
   INSTALL  Install add-on packages
   REMOVE   Remove add-on packages
   BATCH    Run R in batch mode

# R Command Options

## R Command Options

| | |
|---|---|
| -h, --help | Print usage message and exit |
| --version | Print version information and exit |
| --save | Do save workspace at the end of the session |
| --no-save | Don't save it |
| --restore | Do restore previously saved objects |
| --no-restore | Don't restore anything |
| --vanilla | Combine --no-save, --no-restore, --no-site-file, --no-init-file and --no-environ |
| -f file, --file=file | Take input from ``file'' |
| -e expression | Use `exression' as input |

# Editors and IDEs

### Editors

- R terminal
- Rgui
- VIM + Vim-R-plugin
- Emacs + ESS
- Notepad++ + NppToR
- ...

# R Terminal and Rgui

# R Terminal and Rgui

## R

- Ctrl + R : run
- Tab: auto complete
- arrow up and down: history

## R and Texteditor

- copy and paste

- `source("source.R")`

## source

```
sourceDir <- function(path, trace = TRUE, ...) {
    for (nm in list.files(path, pattern = "[.][RrSsQq]$")) {
        if (trace)
            cat(nm, ":")
        source(file.path(path, nm), ...)
        if (trace)
            cat("\n")
    }
}
```

# VIM + Vim-R-plugin

# Notepad++ + NppToR

# Emacs + ESS

## Emacs + ESS

### What is ESS?

Emacs Speaks Statistics (ESS) is an add-on package for emacs text editors such as **GNU Emacs** and XEmacs. It is designed to support editing of scripts and interaction with various statistical analysis programs such as **R**, S-Plus, SAS, Stata and JAGS.

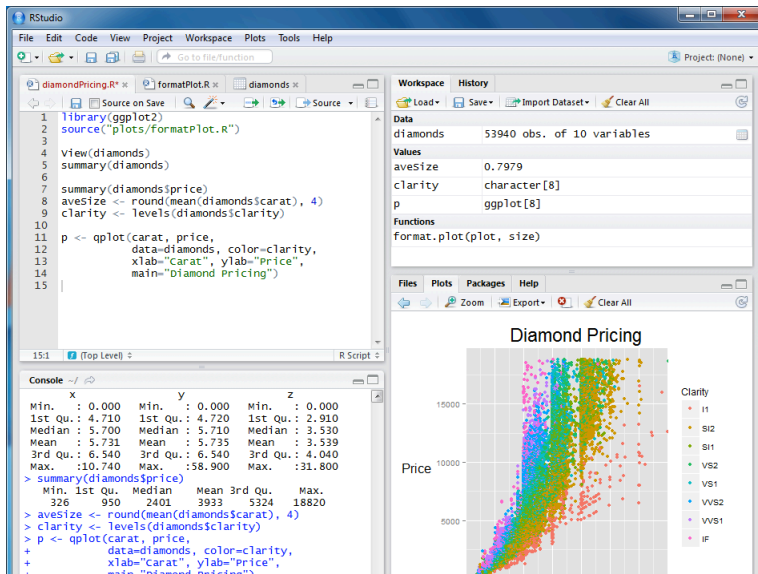### ESS Website

http://ess.r-project.org/

## IDEs

### IDEs

- RStudio: local and cloud-based
- TinnR
- StatET: eclipse for R
- ...

# RStudio

# Next

1. Data Analysis

2. Data Aanlysis and R

3. Hello R!

4. Development Environment

5. References

## Books

- R in action (also in Chinese)
- Introduction to R (also in Chinese)
- R for beginner (also in Chinese)
- R in a Nutshell (Chinese version is in press)
- The art of R programming (also in Chinese)
- ggplot2. Elegant Graphics for Data Analysis (also in Chinese)

## Websites

- R-project and CRAN
- COS.name (Chinese)
- Quick-R
- http://had.co.nz/, Hadley Wickham
- Twitter, github, RForge
- Google

## Websites

- R-project and CRAN
- COS.name (Chinese)
- Quick-R
- http://had.co.nz/, Hadley Wickham
- Twitter, github, RForge
- Google Baidu?

## Journals

- The R Journal
- Journal of Statistical Software
- BMC Bioinformatics: Software
- Bioinformatics: Application Note