

R programming

Philip J Cwynar

University of Pittsburgh
School of Information Sciences
and Intelligent Systems Program



(Big) Data Processing

Background

R is a programming language and software environment for statistical analysis, graphics representation and reporting.

R was created by Ross Ihaka and Robert Gentleman at the University of Auckland, New Zealand, and is currently developed by the R Development Core Team.

R is freely available under the GNU General Public License, and pre-compiled binary versions are provided for various operating systems like Linux, Windows and Mac.



(Big) Data Processing

Intro

The following important features of R:

R is a well-developed, simple and effective programming language which includes conditionals, loops, user defined recursive functions and input and output facilities.

R has an effective data handling and storage facility,

R provides a suite of operators for calculations on arrays, lists, vectors and matrices.

R provides a large, coherent and integrated collection of tools for data analysis.

R provides graphical facilities for data analysis and display either directly at the computer or printing at the papers.



(Big) Data Processing

R is world's most widely used statistics programming language.

It's the # 1 choice of data scientists and supported by a vibrant and talented community of contributors. R is taught in universities and deployed in mission critical business applications.

Download and install R here:

<https://www.r-project.org/>

<https://cran.cnr.berkeley.edu/>



(Big) Data Processing

R studio



Products Resources Pricing About Blog

Log In



Powerful IDE for R

RStudio is a powerful and productive user interface for R. It's free and open source, and works great on Windows, Mac and Linux.

Learn More



R Packages

Develop and expertly manage your own or someone else's R packages. Includes support for Jupyter notebooks, LaTeX, and more.

Learn More



Bring R to the web

Build a web-based user interface for your R code using Shiny. No web development experience required – just learn R and Shiny.

Learn More



(Big) Data Processing

R studio

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

<https://www.rstudio.com/products/rstudio/download/>



(Big) Data Processing

IDE Layout

R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for statistical computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (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.

```
> CODE GOES HERE!!!
```

(Big) Data Processing

Auto-Complete (use tab)

R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for statistical computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (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.

```
> mathtt
```

mathtt[base] mathtt[base] mathtt[base] mathtt[base] mathtt[base] mathtt[base]

(Big) Data Processing

Auto-Complete (even variables!)

R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for statistical computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (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.

```
> ggplot
```

ggplot[base] ggplot[base] ggplot[base] ggplot[base] ggplot[base]

(Big) Data Processing

Installing New Packages

R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for statistical computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (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.

```
> <-->
```

boot

(Big) Data Processing

Need Help with functions? Use ? or help()

R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for statistical computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (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.

```
> help(matrix)
?matrix
```

(Big) Data Processing

Loading a file

R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for statistical computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (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.

```
> load("start.RData")
```

Start.RData

(Big) Data Processing

R 101

DSL
Decision Systems Laboratory

(Big) Data Processing

In contrast to other Programming languages like C and java in R the variables are not declared as some data type.

The variables are assigned with R-Objects and the data type of the R-object becomes the data type of the variable just like in Python

There are many types of R-objects.

- The frequently used ones are
- Vectors
 - Lists
 - Matrices
 - Arrays
 - Factors
 - Data Frames

DSL
Decision Systems Laboratory

(Big) Data Processing

R - Data Types/Classes

| Data Type | Example | Command Line |
|-----------|----------------------------------|--|
| Logical | TRUE , FALSE | v <- TRUE print(class(v)) it produces following result: [1] "logical" |
| Numeric | 12.3, 5, 999 | v <- 23.5 print(class(v)) it produces following result: [1] "numeric" |
| Integer | 2L, 34L, 0L | v <- 2 print(class(v)) it produces following result: [1] "integer" |
| Complex | 3 + 2i | v <- 2+5i print(class(v)) it produces following result: [1] "complex" |
| Character | 'a', "good", "TRUE", "23.4" | v <- TRUE print(class(v)) it produces following result: [1] "character" |
| Raw | "Hello" is stored as 48 65 6c 6f | "Hello" 6c 6f v <- charToRaw("Hello") print(class(v)) it produces following result: [1] "raw" |

DSL
Decision Systems Laboratory

(Big) Data Processing

Variable Assignment

The variables can be assigned values using leftward, rightward and equal operator. The values of the variables can be printed using print() or cat() function. The cat() function combines multiple items into a continuous print output.

Assignment using equal operator.
`var1 = "abcdefg"`

Assignment using leftward operator.
`var2 <- 12345`

Assignment using rightward operator.
`c(TRUE,1) -> var.3 print(var.1) cat ("var.1
is ", var.1 ,"\n") cat ("var.2 is ", var.2
, "\n") cat ("var.3 is ", var.3 ,"\n")`

DSL
Decision Systems Laboratory

(Big) Data Processing

Vectors

A vector is a sequence of data elements of the same basic type. Members in a vector are officially called components.

When you want to create vector with more than one element, you should use c() concatenate function which means to combine the elements into a vector.

```
# Create a vector.  
apple <- c('red','green',"yellow")  
>apple  
  
# Get the class of the vector.  
class(apple)  
When we execute above code, it produces following result:  
[1] "red" "green" "yellow"  
[1] "character"
```

DSL
Decision Systems Laboratory

(Big) Data Processing

Creating Vectors

The c() function can be used to create vectors of objects.

```
> x <- c(0.5, 0.6)      ## numeric  
> x <- c(TRUE, FALSE)   ## logical  
> x <- c(T, F)          ## logical  
> x <- c("a", "b", "c")  ## character  
> x <- 9:29              ## integer  
> x <- c(1+0i, 2+4i)    ## complex
```

Using the vector() function

```
> x <- vector("numeric", length = 10)  
> x  
[1] 0 0 0 0 0 0 0 0 0 0
```

DSL
Decision Systems Laboratory

(Big) Data Processing

Let's Make a Vector!

- Remember me? I'm a Python scalar.
someInt = 12
- someString = "hello Pittsburgh!"
- Remember me? I'm a Python list.
someIntList = [1,3,4,19]
- someStringList = ["apple", "banana", "kittens"]
- In R, I am handled like this:
> newVector = c(1,2,4,19)
> newVector
[1] 1 2 4 19
> newVector[2]
[1] 2 (Wait a second!) (2 == 2?) (Shouldn't 1 == 2?)



(Big) Data Processing

Mixing Objects

What about the following?

```
> y <- c(1.7, "a") ## character  
> y <- c(TRUE, 2) ## numeric  
> y <- c("a", TRUE) ## character
```

When different objects are mixed in a vector, coercion occurs so that every element in the vector is of the same class.



(Big) Data Processing

1 is the new 0

- "A single 0 in an index position returns an empty structure; x[0] returns named numeric(0)."
- So in Python it's 0.. In R, its 1
 - Got that? Good..

http://cran.r-project.org/doc/contrib/R_language.pdf



(Big) Data Processing

Matrices

All the elements of a matrix must be of the same type (numeric, logical, character, complex).

Matrices

Matrices are vectors with a dimension attribute. The dimension attribute is itself an integer vector of length 2 (nrow, ncol)

```
> m <- matrix(nrow = 2, ncol = 3)  
> m  
[,1] [,2] [,3]  
[1,] NA NA NA  
[2,] NA NA NA  
> dim(m)  
[1] 2 3  
> attributes(m)  
$dim  
[1] 2 3
```



(Big) Data Processing

Let's Add Some Labels!

```
>someMatrix = matrix(1:16, 4, 4, FALSE,  
  dimnames=list(c("A","B","C","D"),c("W","X","Y","Z")))  
> someMatrix  
      W     X     Y     Z  
A   1     5     9    13  
B   2     6    10    14  
C   3     7    11    15  
D   4     8    12    16  
> someMatrix["B","Y"]  
[1] 10
```

*Note: Use NULL if you don't want to append labels to either x or y axis



(Big) Data Processing

```
# Create a matrix. M = matrix( c('a','a','b','c','b','a'),  
  nrow=2, ncol=3, byrow = TRUE)
```

Matrix Example

```
> someMatrix = matrix(1:16,4,4)  
> someMatrix  
      [,1]     [,2]     [,3]     [,4]  
[1,]     1       5       9      13  
[2,]     2       6      10      14  
[3,]     3       7      11      15  
[4,]     4       8      12      16  
> someMatrix[10]  
[1] 10  
> someMatrix[2,3]  
[1] 10
```



(Big) Data Processing

Data Frames

Data frames are tabular data objects. Unlike a matrix in data frame each column can contain different modes of data. The first column can be numeric while the second column can be character and third column can be logical. It is a list of vectors of equal length.

Data Frames are created using the `data.frame()` function.

```
# Create the data frame.
BMI <- data.frame( gender = c("Male", "Male","Female"),
height = c(152, 171.5, 165), weight = c(81,93, 78),
Age =c(42,38,26) )

print(BMI)
```

DSL Decision Sciences Laboratory (Big) Data Processing

DSL Decision Sciences Laboratory (Big) Data Processing

Read.Table()

- The majority of data is handled using a table (especially in ggplot)

USAGE: `read.table(file, header = FALSE, sep = "", quote = "\'", dec = ".", row.names, col.names, as.is = !stringsAsFactors, na.strings = "NA", colClasses = NA, nrow = -1, skip = 0, check.names = TRUE, fill = !blank.lines.skip, strip.white = FALSE, blank.lines.skip = TRUE, comment.char = "#", allowEscapes = FALSE, flush = FALSE, stringsAsFactors = default.stringsAsFactors(), fileEncoding = "R", encoding = "unknown", text)`

DSL Decision Sciences Laboratory (Big) Data Processing

Header

Made a slight change to "BeerAndWinePerCapita.txt"

| State | GallonsOfBeer | GallonsOfWine |
|-----------------|---------------|---------------|
| 1 NEVADA | 44 | 5.75 |
| 2 NEW HAMPSHIRE | 43.4 | 6.26 |
| 3 NORTH DAKOTA | 41.7 | 1.56 |
| 4 MONTANA | 41.5 | 3.06 |
| 5 SOUTH DAKOTA | 39 | 1.50 |
| 6 WISCONSIN | 38.2 | 2.63 |

DSL Decision Sciences Laboratory (Big) Data Processing

```
>mydata <- read.table("Retention.txt", header=TRUE, sep="\t")
>summary(mydata)

  spend      apret      top10      rejr
Min. :4125  Min. :18.75  Min. : 8.00  Min. : 0.00
1st Qu.: 7372  1st Qu.:45.37  1st Qu.:22.00  1st Qu.:19.17
Median :9265  Median :55.71  Median :30.00  Median :27.39
Mean :10975  Mean :56.72  Mean :38.46  Mean :30.65
3rd Qu.:12838 3rd Qu.:68.69  3rd Qu.:49.50  3rd Qu.:36.81
Max. :35863  Max. :95.25  Max. :98.00  Max. :84.07
  tspsc      pacc      strat      salar
Min. :48.12  Min. : 8.964  Min. : 7.20  Min. :38640
1st Qu.:61.11 1st Qu.:33.904  1st Qu.:13.40  1st Qu.:54650
Median :64.78  Median :40.850  Median :16.00  Median :61150
Mean :66.16  Mean :43.173  Mean :16.09  Mean :61358
3rd Qu.:70.45 3rd Qu.:51.773  3rd Qu.:18.57  3rd Qu.:67100
Max. :87.50  Max. :76.253  Max. :29.20  Max. :87900
```

DSL Decision Sciences Laboratory (Big) Data Processing

Factors

- Nominal (Categorical) –**
Two or more categories
Have no intrinsic ordering
Think: (Male/Female), (Blonde/Brunette/Red Hair), and (Pittsburgh/State College/Erie/etc.)
- Ordinal**
Similar to Nominal, but with order
Think: (Low/Medium/High), (Tall/Average/Short), and (Tall Latte/Grande Latte/Venti Latte no whip)
- Interval**
Same as Ordinal, but evenly spaced
Think: (Temperature), (Time), and (Measurements)

DSL Decision Sciences Laboratory (Big) Data Processing

ggplot2

DSL
Decision Sciences Laboratory

(Big) Data Processing

Some Basics of ggplot2

- Extended library of R
- Used to create a variety of visualizations without a lot of background knowledge
- Created in a layered fashion



DSL
Decision Sciences Laboratory

(Big) Data Processing

- Traditional plotting:** You are a painter
 - Manually place individual graphical elements
ggplot2: You employ a painter
 - Describe conceptually how data should be visualized

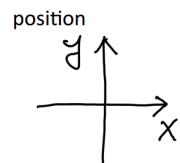
Source: <http://wilkelab.org/>

DSL
Decision Sciences Laboratory

(Big) Data Processing

Most confusing key concept: **aesthetic mapping**
Maps data values to visual elements of the plot

A few examples of aesthetics



Source: <http://wilkelab.org/>

DSL
Decision Sciences Laboratory

(Big) Data Processing

Simple example: mean height and weight of boys/girls ages 10-20

| age (yrs) | height (cm) | weight (kg) | sex |
|-----------|-------------|-------------|-----|
| 10 | 138 | 32 | M |
| 15 | 170 | 56 | M |
| 20 | 177 | 71 | M |
| 10 | 138 | 33 | F |
| 15 | 162 | 52 | F |
| 20 | 163 | 53 | F |

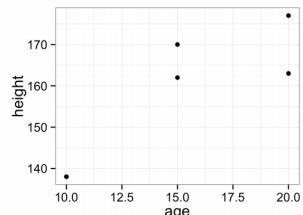
Source: <http://wilkelab.org/>

DSL
Decision Sciences Laboratory

(Big) Data Processing

Map age to x, height to y, visualize using points

`ggplot(data, aes(x=age, y=height)) + geom_point()`



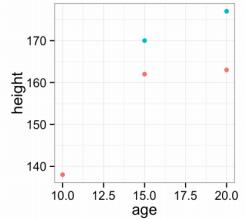
Source: <http://wilkelab.org/>

DSL
Decision Sciences Laboratory

(Big) Data Processing

Let's color the points by sex

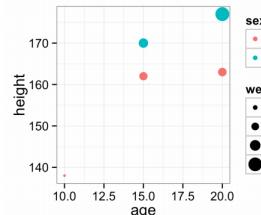
```
ggplot(data, aes(x=age, y=height,color=sex)) +geom_point()
```



Source: <http://wilkelab.org/>
DSL Decision Systems Laboratory
(Big) Data Processing

And change point size by weight

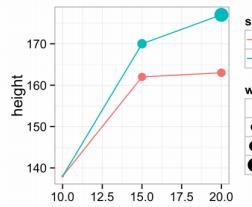
```
ggplot(data, aes(x=age, y=height,color=sex, size=weight)) +geom_point()
```



Source: <http://wilkelab.org/>
DSL Decision Systems Laboratory
(Big) Data Processing

The weight-to-size mapping should only be applied to points

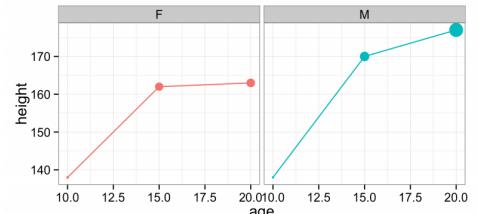
```
ggplot(data, aes(x=age, y=height,color=sex)) +  
  geom_point(aes(size=weight)) +  
  geom_line()
```



Source: <http://wilkelab.org/>
DSL Decision Systems Laboratory
(Big) Data Processing

We can also make side-by-side plots (called facets)

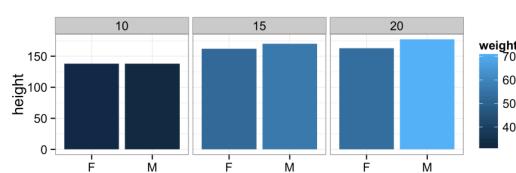
```
ggplot(data, aes(x=age, y=height,color=sex)) +  
  geom_point(aes(size=weight)) + geom_line() +  
  facet_wrap(~sex)
```



Source: <http://wilkelab.org/>
DSL Decision Systems Laboratory
(Big) Data Processing

Now let's facet by age, color by weight, and use bars to plot height

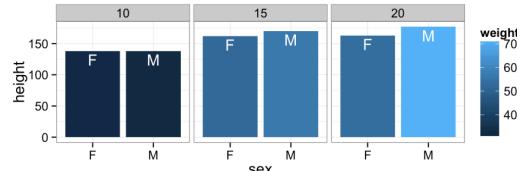
```
ggplot(data, aes(x=sex, y=height,fill=weight)) +  
  geom_bar(stat='identity') +facet_wrap(~age)
```



Source: <http://wilkelab.org/>
DSL Decision Systems Laboratory
(Big) Data Processing

Let's plot the sex also at the top of the bar

```
ggplot(data, aes(x=sex, y=height, fill=weight)) +  
  geom_bar(stat='identity') +  
  geom_text(aes(label=sex), vjust=1.3, color='white') +  
  facet_wrap(~age)
```



Source: <http://wilkelab.org/>
DSL Decision Systems Laboratory
(Big) Data Processing

All the geom's with all their options are described on the ggplot2 web page

<http://docs.ggplot2.org/current/>



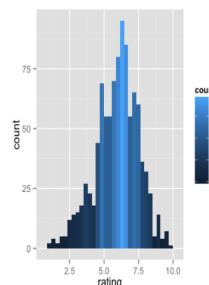
(Big) Data Processing

More Examples



(Big) Data Processing

Histogram



http://docs.ggplot2.org/current/geoms_histogram.html



(Big) Data Processing

Facet Grid

ggplot2 Quick Reference: **facet**

The faceting approach supported by ggplot2 partitions a plot into a matrix of panels. Each panel shows a different subset of the data. There are two faceting approaches:

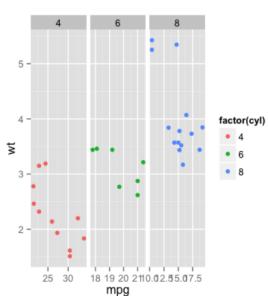
facet_wrap(~cell) - univariate: create a 1-d strip of panels, based on one factor, and wrap the strip into a 2-d matrix

facet_grid(row~col) - (usually) bivariate: create a 2-d matrix of panels, based on two factors



(Big) Data Processing

Facet Grid



http://docs.ggplot2.org/current/facet_grid.html

[http://www.cookbook-r.com/Graphs/Facets_\(ggplot2\)/](http://www.cookbook-r.com/Graphs/Facets_(ggplot2)/)



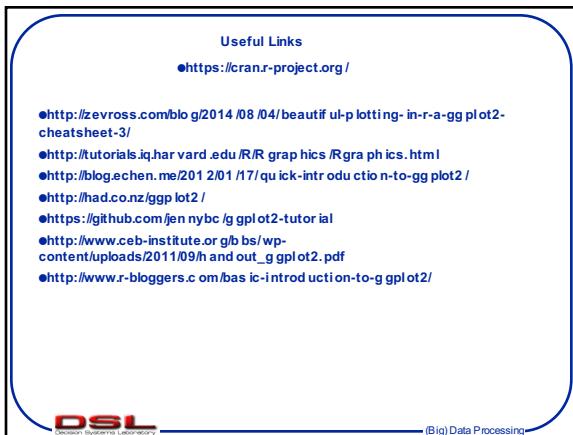
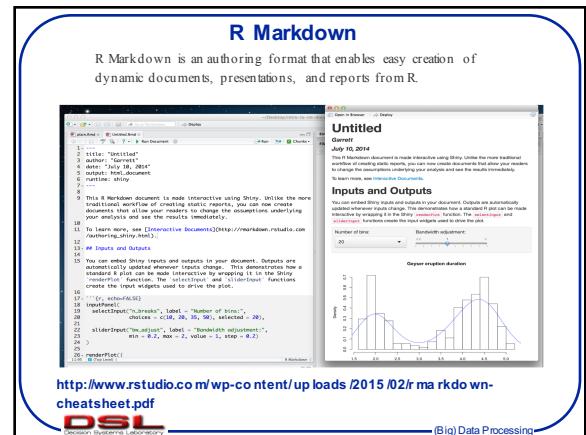
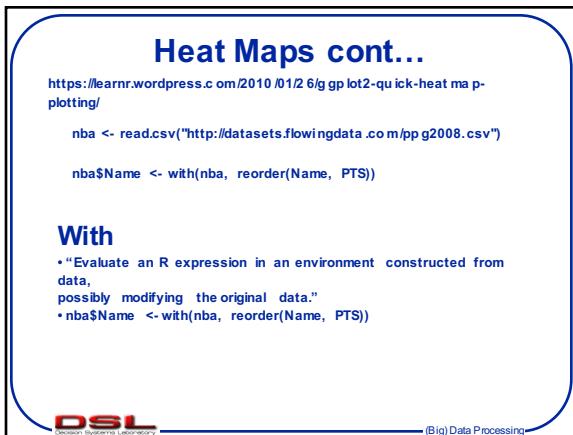
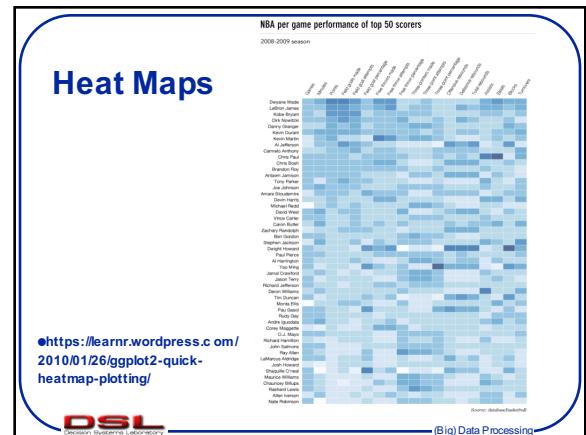
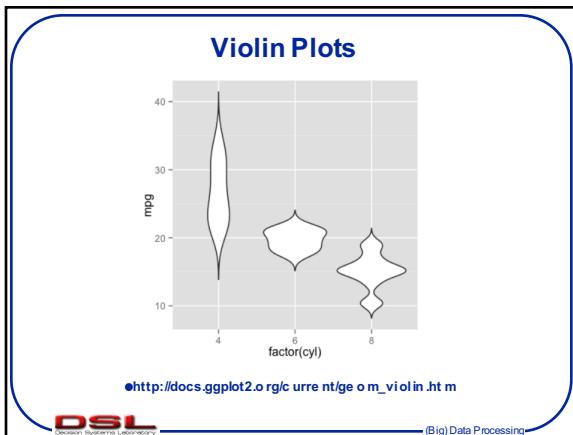
(Big) Data Processing

Violin Plots

The violin plot is similar to [box plots](#), except that they also show the [probability density](#) of the data at different values (in the simplest case this could be a [histogram](#)). Typically violin plots will include a marker for the median of the data and a box indicating the interquartile range, as in standard box plots. Overlaid on this box plot is a [kernel density estimation](#).



(Big) Data Processing



Kaggle Titanic Assignment

DSL
Decision Sciences Laboratory

(Big) Data Processing

| File Name | Available Formats |
|------------------|-----------------------|
| train | CSV (59.76 kb) |
| gendermodel | CSV (3.18 kb) |
| genderclassmodel | CSV (3.18 kb) |
| test | CSV (27.96 kb) |
| gendermodel | PY (3.59 kb) |
| genderclassmodel | PY (5.43 kb) |
| myfirstforest | PY (3.99 kb) |

See, fork, and run a random forest benchmark model through Kaggle Scripts

DSL
Decision Sciences Laboratory

(Big) Data Processing

Titanic Data Graphics

- (1) Go to kaggle.com > Titanic: Machine Learning from Disaster and download (train.csv)
- (2) Generate Descriptive Statistics
- (3) You will create a plot (for each, so 5 plots in total) in ggplot2 using:
 - a. Whisker-plot
 - b. Histogram
 - c. Facet grid
 - d. Violin plot
- (4) Write up a 500 – 1500 word document talking about the assignment using the graphics to describe the passengers of Titanic.

DSL
Decision Sciences Laboratory

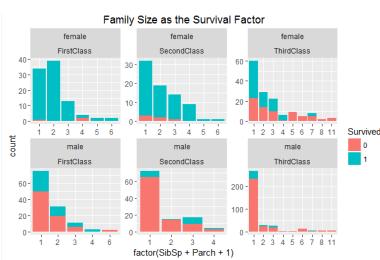
(Big) Data Processing

Titanic Examples

DSL
Decision Sciences Laboratory

(Big) Data Processing

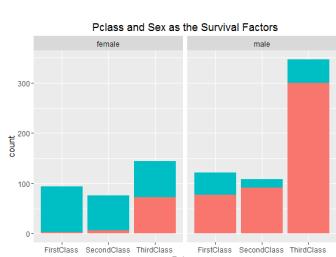
```
ggplot(tmal, aes(factor(factor(SibSp + Parch + 1))) + geom_bar(aes(fill = Survived)) + facet_wrap(~Sex*Pclass, nrow = 2, scales = "free") + ggtitle("Family size as the survival factor"))
```



Source: Kaggle scripts

DSL
Decision Sciences Laboratory

```
ggplot(tmal, aes(Pclass)) + geom_bar(aes(fill = survived)) + facet_grid(~sex) + ggtitle("Pclass and Sex as the Survival Factors")
```



Source: Kaggle scripts

DSL
Decision Sciences Laboratory

(Big) Data Processing