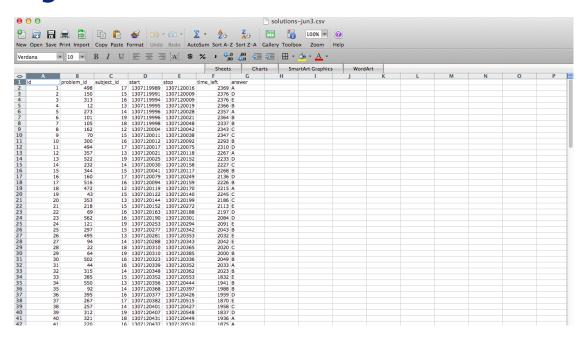


Data munging basics

Jeffrey Leek, Assistant Professor of Biostatistics Johns Hopkins Bloomberg School of Public Health

Recall Tidy Data



- 1. Each variable forms a column
- 2. Each observation forms a row
- 3. Each table/file stores data about one kind of observation (e.g. people/hospitals).

http://vita.had.co.nz/papers/tidy-data.pdf

Leek, Taub, and Pineda 2011 PLoS One

Where we would like to be

- Tidy data refers to the shape of the data
 - Variables in columns
 - Observations in rows
 - Tables holding elements of only one kind
- · Plus
 - Column names are easy to use and informative
 - Row names are easy to use and informative
 - Obvious mistakes in the data have been removed
 - Variable values are internally consistent
 - Appropriate transformed variables have been added

A partial list of munging operations

- · Fix variable names
- · Create new variables
- Merge data sets
- · Reshape data sets
- Deal with missing data
- Take transforms of variables
- · Check on and remove inconsistent values

These steps must be recorded

90% of your effort will often be spent here

A partial list of munging operations

- Fix variable names
- Create new variables
- Merge data sets
- Reshape data sets
- Deal with missing data
- Take transforms of variables
- · Check on and remove inconsistent values

Fixing character vectors - tolower(), toupper()

```
cameraData <- read.csv("./data/cameras.csv")
names(cameraData)

[1] "address"    "direction"    "street"    "crossStreet"

[5] "intersection" "Location.1"

tolower(names(cameraData))

[1] "address"    "direction"    "street"    "crossstreet"
[5] "intersection"    "location.1"</pre>
```

Fixing character vectors - strsplit()

- Good for automatically splitting variable names
- · Important parameters: x, split

```
splitNames = strsplit(names(cameraData),"\\.")
splitNames[[5]]
```

```
[1] "intersection"
```

```
splitNames[[6]]
```

```
[1] "Location" "1"
```

Data munging basics

Fixing character vectors - sapply()

- · Applies a function to each element in a vector or list
- · Important parameters: X,FUN

```
splitNames[[6]][1]

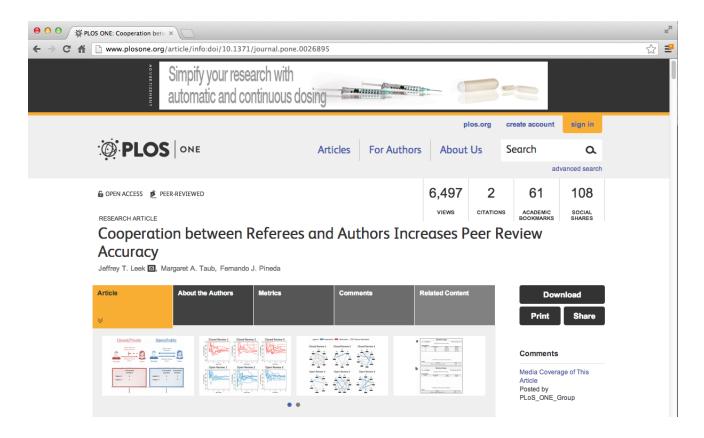
[1] "Location"

firstElement <- function(x){x[1]}
sapply(splitNames, firstElement)

[1] "address"    "direction"    "street"    "crossStreet"
[5] "intersection" "Location"</pre>
```

Peer review experiment data

· Data on submissions/reviews in an experiment



http://www.plosone.org/article/info:doi/10.1371/journal.pone.0026895

Peer review data

```
fileUrl1 <- "https://dl.dropboxusercontent.com/u/7710864/data/reviews-apr29.csv"
fileUrl2 <- "https://dl.dropboxusercontent.com/u/7710864/data/solutions-apr29.csv"
download.file(fileUrl1,destfile="./data/reviews.csv",method="curl")
download.file(fileUrl2,destfile="./data/solutions.csv",method="curl")
reviews <- read.csv("./data/reviews.csv"); solutions <- read.csv("./data/solutions.csv")
head(reviews,2)</pre>
```

```
head(solutions,2)
```

Fixing character vectors - sub(),gsub()

· Important parameters: pattern, replacement, x

Fixing character vectors - sub(),gsub()

```
testName <- "this_is_a_test"
sub("_","",testName)</pre>
```

```
[1] "thisis_a_test"
```

```
gsub("_","",testName)
```

```
[1] "thisisatest"
```

Quantitative variables in ranges - - cut()

· Important parameters: *x*,*breaks*

```
reviews$time_left[1:10]
```

```
[1] 1754 2306 2192 2089 2043 1999 2130 NA 1899 2024
```

```
timeRanges <- cut(reviews$time_left,seq(0,3600,by=600))
timeRanges[1:10]</pre>
```

```
[1] (1.2e+03,1.8e+03] (1.8e+03,2.4e+03] (1.8e+03,2.4e+03]
[4] (1.8e+03,2.4e+03] (1.8e+03,2.4e+03] (1.8e+03,2.4e+03]
[7] (1.8e+03,2.4e+03] <NA> (1.8e+03,2.4e+03]
[10] (1.8e+03,2.4e+03]
[10] (1.8e+03,2.4e+03] ... (3e+03,3.6e+03]
```

Quantitative variables in ranges - - cut()

class(timeRanges)

[1] "factor"

table(timeRanges,useNA="ifany")

Quantitative variables in ranges - cut2() {Hmisc}

You may need to run install.packages("Hmisc")

```
library("Hmisc")
timeRanges<- cut2(reviews$time_left,g=6)
table(timeRanges,useNA="ifany")</pre>
```

Quantitative variables in ranges - cut2() {Hmisc}

```
table(timeRanges,useNA="ifany")
```

Adding an extra variable

```
timeRanges <- cut2(reviews$time_left,g=6)
reviews$timeRanges <- timeRanges
head(reviews,2)</pre>
```

Merging data - merge()

- · Merges data frames
- Important parameters: x,y,by,by.x,by.y,all

```
names(reviews)
```

```
[1] "id" "solution_id" "reviewer_id" "start"
[5] "stop" "time_left" "accept" "timeRanges"
```

```
names(solutions)
```

```
[1] "id" "problem_id" "subject_id" "start" "stop"
[6] "time_left" "answer"
```

Merging data - merge()

```
mergedData <- merge(reviews, solutions, all=TRUE)
head(mergedData)</pre>
```

```
stop time left solution id reviewer id accept
  id
          start
  1 1304095119 1304095169
                                 2343
                                                NA
                                                             NA
                                                                    NA
  1 1304095698 1304095758
                                 1754
                                                 3
                                                             27
                                                                     1
  2 1304095119 1304095183
                                 2329
                                                NA
                                                             NA
                                                                    NA
   2 1304095188 1304095206
                                 2306
                                                             22
                                                                     1
                                                 4
  3 1304095127 1304095146
                                 2366
                                                NA
                                                             NA
                                                                    NA
  3 1304095276 1304095320
                                 2192
                                                 5
                                                             28
                                                                     1
   timeRanges problem id subject id answer
                                  29
1
         <NA>
                      156
                                           В
2 [1496, 1909)
                       NA
                                  NA
                                        <NA>
3
         <NA>
                      269
                                  25
                                           C
4 [1909,2306]
                                  NA
                                        <NA>
                       NA
         <NA>
                       34
                                  22
                                          C
6 [1909,2306]
                       NA
                                  NA
                                        <NA>
```

Merging data - merge()

```
mergedData2 <- merge(reviews, solutions, by.x = "solution_id", by.y = "id", all=TRUE)
head(mergedData2[,1:6],3)</pre>
```

```
reviews[1,1:6]
```

Data munging basics

Sorting values - sort()

· Important parameters: x, decreasing

mergedData2\$reviewer_id[1:10]

8/28/13

[1] 26 29 27 22 28 22 29 23 25 29

sort(mergedData2\$reviewer_id)[1:10]

[1] 22 22 22 22 22 22 22 22 22 22

Ordering values - order()

· Important parameters: list of variables to order, na.last, decreasing

```
mergedData2$reviewer_id[1:10]
```

```
[1] 26 29 27 22 28 22 29 23 25 29
```

```
order(mergedData2$reviewer id)[1:10]
```

```
[1] 4 6 14 22 23 24 27 32 37 39
```

Ordering values - order()

mergedData2\$reviewer id[order(mergedData2\$reviewer id)]

```
        [1]
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        22
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
        23
```

Reordering a data frame

```
head(mergedData2[,1:6],3)
```

```
sortedData <- mergedData2[order(mergedData2$reviewer_id),]
head(sortedData[,1:6],3)</pre>
```

```
      solution_id id reviewer_id
      start.x
      stop.x time_left.x

      4
      4
      2
      22 1304095188 1304095206
      2306

      6
      6
      16
      22 1304095303 1304095471
      2041

      14
      14
      12
      22 1304095280 1304095301
      2211
```

Reordering by multiple variables

```
head(mergedData2[,1:6],3)
```

```
sortedData <- mergedData2[order(mergedData2$reviewer_id,mergedData2$id),]
head(sortedData[,1:6],3)</pre>
```

```
      solution_id
      id
      reviewer_id
      start.x
      stop.x
      time_left.x

      4
      4
      2
      22
      1304095188
      1304095206
      2306

      14
      14
      12
      22
      1304095280
      1304095301
      2211

      6
      6
      16
      22
      1304095303
      1304095471
      2041
```

Reshaping data - example

```
misShaped <- as.data.frame(matrix(c(NA,5,1,4,2,3),byrow=TRUE,nrow=3))
names(misShaped) <- c("treatmentA","treatmentB")
misShaped$people <- c("John","Jane","Mary")
misShaped</pre>
```

http://vita.had.co.nz/papers/tidy-data.pdf

8/28/13

Reshaping data - melt() {reshape2}

You may need to run install.packages("reshape2") before using the melt() function

· Important parameters: id.vars, measure.vars, variable.name

```
library("reshape2")
melt(misShaped,id.vars="people",variable.name="treatment",value.name="value")
```

```
people treatment value

1 John treatmentA NA

2 Jane treatmentA 1

3 Mary treatmentA 2

4 John treatmentB 5

5 Jane treatmentB 4

6 Mary treatmentB 3
```

More resources

- Tidy data and tidy tools
- Andrew Jaffe's Data Cleaning Lecture
- Hadley Wickham on regular expressions
- · Long, painful experience :-)