# Achieving reproducibility

## Jeff Leek

@jtleek

www.jtleek.com

# A data sharing plan

1. The raw data.

2. A tidy data set

3. A code book describing each variable and its values in the tidy data set.

4. An explicit and exact recipe you used to go from 1 -> 2,3

# Raw data

```
@HWI-EAS121:4:100:1783:550#0/1
CGTTACGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACGGATCTCGTATGCGGTCTGCTGCGTGACAAGACAGGGG
+HWI-EAS121:4:100:1783:550#0/1
aaaaa`b_aa`aa`YaX]aZ`aZM^Z]YRa]YSG[[ZREQLHESDHNDDHNMEEDDMPENITKFLFEEDDDHEJQMEDDD
@HWI-EAS121:4:100:1783:1611#0/1
GGGTGGGCATTTCCACTCGCAGTATGGGTTGCCGCACGACAGGCAGCGGTCAGCCTGCGCTTTGGCCTGGCCTTC
+HWI-EAS121:4:100:1783:1611#0/1
a``^\__`_````^a`_a`^a_^__]a_]\]`a_____`_^^`]X]_]XTV_\]]NX_XVX]]_TTTG[
@HWI-EAS121:4:100:1783:322#0/1
CGTTTATGTTTTTGAATATGTCTTATCTTAACGGTTATATTTTAGATGTTGGTCTTATTCTAACGGTCA     TTTCTA
+HWI-EAS121:4:100:1783:322#0/1
abaa`^aaaaabbbaababbbbbb`bbbb_bbbbbbbb`bbbaV^_a``a``]``aT]a__V\]]_]^      bbaV__
@HWI-EAS121:4:100:1783:1394#0/1
GGGTCTTTATTGGTCTGGTGATCCCCCATATTCTCCGGTTGTGTGGTTTAACCGATCATCGCGCATTAC        CTGC
+HWI-EAS121:4:100:1783:1394#0/1
```[aa\b^^`[]aabbb][`a_abbb`a``bbbbbabaabaaaab_VZa_^___bab_X`[a\HV_[_]_[^_
@HWI-EAS121:4:100:1783:207#0/1
CCCTGGGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAACA
+HWI-EAS121:4:100:1783:207#0/1
abba`Xa\^\\`aa]ba__bba[a_O_a`aa`aa`a]^V]X_a^YS\R_\H_[]\Z
@HWI-EAS121:4:100:1783:455#0/1
GGGTAATTCAGGGACAATGTAATGGCTGCACAAAAAAATACATCTTTCATGTTCCA
+HWI-EAS121:4:100:1783:455#0/1
abb_babbabaabbbbbbbbbbbbbbba\`b`\abbbabbbbabbbbbbaabbbbb
```

**Processing**
**Computing**
**Summarizing**
**Deleting**

A tidy data set

One variable per column
One observation per row
One table per "kind" of variable

Reference: http://brianknaus.com/software/srtoolbox/s_4_1_sequence80.txt

anything doesn't make sense.

Files:

**1 Demographics**: tab 1 is schizophrenia patients, tab 2 is controls.
A. Cohort: M = Mannheim (Germany), C = Cologne (Germany), H= Hopkins. We had a few of our own patients so we included them too.
B. patient identification number
C. Age at time of CSF collection
D. Gender
E. BMI
F. Ethnicity (mostly Caucasian)
G. Diagnosis: DSM/ICD-10 diagnosis
H. Group: control, schizophrenia, or prodromal. I don't think we have enough power to run them as three groups so I combined prod... sure if this was ok. Is it appropriate to do a ttes... SZ?
I. Medication: mostly untreated
J. Education more or less than 13 years
K. current smoking status: yes or no

Variable names
Variable descriptions
Variable units

RStudio

geuvadis.Rmd

Knit HTML

Run    Chunks

```r
33  library(sva)
34  library(ffpe)
35  library(RColorBrewer)
36  library(corrplot)
37  library(limma)
38  trop = RSkittleBrewer('tropical')
39  ```
40
41
42  ## Load the data
43
44  You will need to download the GEUVADIS ballgown object from this site: https://github.com/a        zee
    /ballgown_code
45
46
47  ```{r loaddata,dependson="load"}
48  load("fpkm.rda")
49  pd = ballgown::pData(fpkm)
50  pd$dirname = as.character(pd$dirname)
51  ss = function(x, pattern, slot=1,...) sapply(strspli
52  pd$IndividualID = ss(pd$dirname, "_", 1)
53  tfpkm = expr(fpkm)$trans
54  ```
55
56  ## Subset to non-duplicates
57
58  You will need the GEUVADIS quality control information and population information available from these
```

1:1    (Top Level)    R Markdown

**R/Python Code
Input raw data -> output tidy
No parameters**

Step 1 - take the raw file, run version 3.1.2 of summarize software with parameters a=1, b=2, c=3
Step 2 - run the software separately for each sample
Step 3 - take column three of outputfile.out for each s
and that is the corresponding row in the output data

**Explicit instructions**
**Versions of software**
**Parameters included**

Step 1 - take the raw file, run version 3.1.2 of summarize software with parameters a=1, b=2, c=3
Step 2 - run the software separately for each sample
Step 3 - take column three of outputfile.out for each s
and that is the corresponding row in the output data

**Vague instructions**
**Missing versions**
**Skipped steps**

The Leek group guide to data sharing — Edit

⊙ **25** commits    ⌥ **1** branch    ⬙ **0** releases    ⬡ **8** contributors

⌥ branch: **master** ▾    **datasharing** / ⊡

Merge pull request **#9** from nikai3d/patch-1  ⋯

jtleek authored 6 days ago                    latest commit e53857faa4 ⬔

▤ README.md          fix typo                          6 days ago

▥ **README.md**

# How to share data with a statistician

This is a guide for anyone who needs to share data with a statistician. The target audiences I have in mind are:

- Scientific collaborators who need statisticians to analyze data for them
- Students or postdocs in scientific disciplines looking for consulting advice
- Junior statistics students whose job it is to collate/clean data sets

# Code

## Raw

## Literate programming

```r
f.pvalue <- function(dat,mod,mod0){
  # This is a function for performing
  # parametric f-tests on the data matrix
  # dat comparing the null model mod0
  # to the alternative model mod.
  n <- dim(dat)[2]
  m <- dim(dat)[1]
  df1 <- dim(mod)[2]
  df0 <- dim(mod0)[2]
  p <- rep(0,m)
  Id <- diag(n)

  resid <- dat %*% (Id - mod %*% solve(t(mod) %*% mod) %*% t(mod))
  resid0 <- dat %*% (Id - mod0 %*% solve(t(mod0) %*% mod0) %*% t(mod0))

  rss1 <- resid^2 %*% rep(1,n)
  rss0 <- resid0^2 %*% rep(1,n)

  fstats <- ((rss0 - rss1)/(df1-df0))/(rss1/(n-df1))
  p <-  1-pf(fstats,df1=(df1-df0),df2=(n-df1))
  return(p)
}

setwd("cheung/")
# Load data and create group variable
dat <- read.table("full.data")

jpt.names <- scan("JPT.cname.txt",what="character")
chb.names <- scan("CHB.cname.txt",what="character")
ceu.names <- scan("CEU_parents.txt",what="character")
nceu <- length(ceu.names)
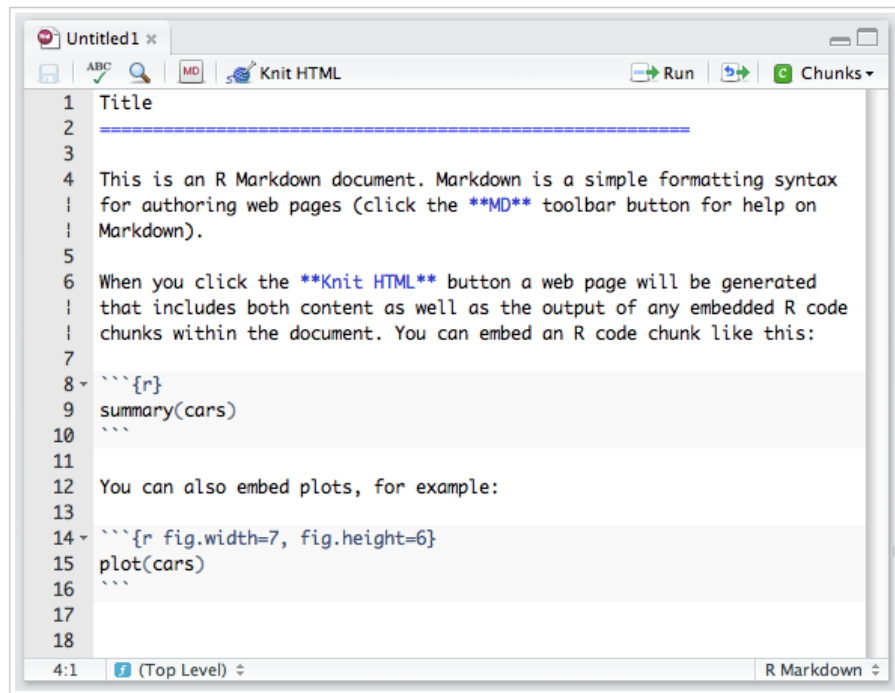njpt <- length(jpt.names)
nchb <- length(chb.names)
```

# R Markdown Documents

To work with R Markdown (.Rmd) files in RStudio you first need to ensure that the knitr package (version 0.5 or later) in installed.

To create a new R Markdown file, go to **File | New |** and select **R Markdown**. A new file is create with a default template to get you oriented:

```
Untitled1

  ABC  🔍   MD   ✏ Knit HTML                      ➡ Run  ↩  C Chunks▾

 1  Title
 2  ============================================================
 3
 4  This is an R Markdown document. Markdown is a simple formatting syntax
 ⋮  for authoring web pages (click the **MD** toolbar button for help on
 ⋮  Markdown).
 5
 6  When you click the **Knit HTML** button a web page will be generated
 ⋮  that includes both content as well as the output of any embedded R code
 ⋮  chunks within the document. You can embed an R code chunk like this:
 7
 8▾ ```{r}
 9  summary(cars)
10  ```
11
12  You can also embed plots, for example:
13
14▾ ```{r fig.width=7, fig.height=6}
15  plot(cars)
16  ```
17
18

 4:1   f (Top Level) ↕                                        R Markdown ↕
```

Note that the toolbar provides some useful tools for working with R Markdown:

- **Quick Reference** — Click the **MD** toolbar button to open a quick reference guide for Markdown.
- **Knit HTML** — Click to knit the current document to HTML, see the **Knitting to HTML** section below for more details.
- **Run** — Run the current line or selection of lines in the console. This allows running R code inside a code chunk similar to a normal R source file.
- **Chunks** — The chunks menu provides assistance with inserting, running, and chunk navigation. See the **Chunk Menu and Options** section below for more details.