# Organizing a Data Analysis

Roger D. Peng, Associate Professor of Biostatistics

May 18, 2016

# Data analysis files

- Data
- Raw data
- Processed data
- Figures
- Exploratory figures
- Final figures
- R code
- Raw / unused scripts
- Final scripts
- R Markdown files
- Text
- README files
- Text of analysis / report

# Raw Data



```
----------------------------  ALLERGIES  ----------------------------        ----------------------------  MEDICATION HISTORY  ----------------------------

ast Updated: 01 Dec 2011 @ 0851                                              Last Updated: 11 Apr 2011 @ 1737

                                                                            Medication: AMLODIPINE BESYLATE 10MG TAB
llergy Name:        TRIMETHOPRIM                                            Instructions: TAKE ONE TABLET BY MOUTH TAKE ONE-HALF TABLET FOR
ocation:            DAYT29                                                  GRAPEFRUIT JUICE--
ate Entered:        09 Mar 2011                                             Status: Active
eaction:                                                                    Refills Remaining: 3
llergy Type:        DRUG                                                    Last Filled On: 20 Aug 2010
A Drug Class:       ANTI-INFECTIVES,OTHER                                   Initially Ordered On: 13 Aug 2010
bserved/Historical: HISTORICAL                                             Quantity: 45
omments:            The reaction to this allergy was MILD (NO SQUELAE)      Days Supply: 90
                                                                            Pharmacy: DAYTON
llergy Name:        TRAMADOL                                               Prescription Number: 2718953
ocation:            DAYT29
ate Entered:        09 Mar 2011                                            Medication: IBUPROFEN 600MG TAB
eaction:            URINARY RETENTION                                      Instructions: TAKE ONE TABLET BY MOUTH FOUR TIMES A DAY WITH FOOD
llergy Type:        DRUG                                                   Status: Active
A Drug Class:       NON-OPIOID ANALGESICS                                  Refills Remaining: 3
bserved/Historical: HISTORICAL                                            Last Filled On: 20 Aug 2010
omments:            gradually worsening difficulty emptying bladder        Initially Ordered On: 01 Jul 2010
                                                                            Quantity: 360
```
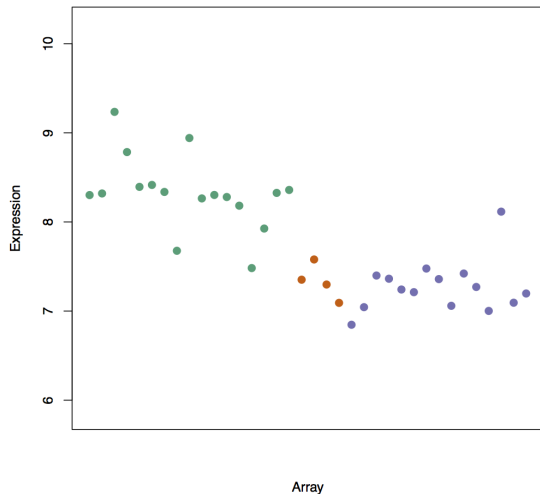
- ▶ Should be stored in your analysis folder
- ▶ If accessed from the web, include url, description, and date accessed in README
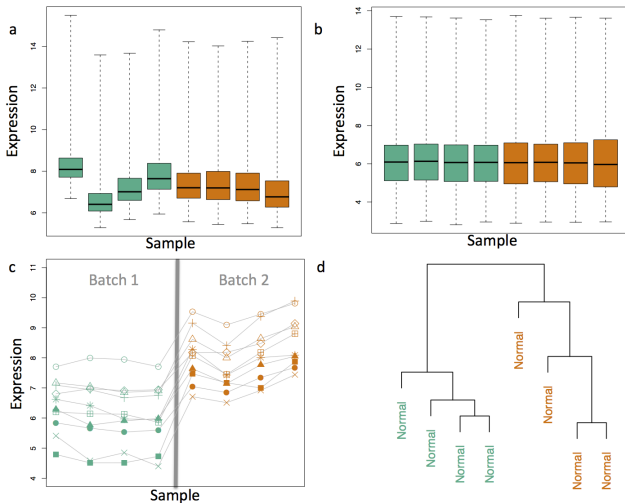
# Processed data



Processed data should be named so it is easy to see which script generated the data. * The processing script - processed data mapping should occur in the README * Processed data should be tidy

# Exploratory figures



- Figures made during the course of your analysis, not necessarily part of your final report.

# Final Figures



- Usually a small subset of the original figures
- Axes/colors set to make the figure clear

# Raw scripts

```r
source("regmodel.R")

dp <- ddm[, c("group", "pm25_0", "pm25_1", "symfree0", "symfree1")]
dp$p_id <- row.names(dp)

fitx0 <- lm(pm25_1 ~ pm25_0 + age + no2_0 + pm10_0, data = subset(ddm, group ==
fitx1 <- lm(pm25_1 ~ ns(pm25_0, 2) + age + no2_0 + pm10_0, data = subset(ddm, gro

fity0 <- glm(cbind(symfree1, 14-symfree1) ~ symfree0 + age + factor(gender), data
fity1 <- glm(cbind(symfree1, 14-symfree1) ~ symfree0 + age + factor(gender), data

y10 <- predict(fity0, subset(ddm, group == 1), type = "response") * 14
y01 <- predict(fity1, subset(ddm, group == 0), type = "response") * 14
p10 <- predict(fitx0, subset(ddm, group == 1))
p01 <- predict(fitx1, subset(ddm, group == 0))

yy <- data.frame(p_id = as.integer(c(names(y10), names(y01))),
                 symfree00 = c(y10, y01))
pp <- data.frame(p_id = as.integer(c(names(p10), names(p01))),
                 pm25_00 = c(p10, p01))

m <- merge(dp, yy, by = "p_id")
mm <- merge(m, pp, by = "p_id")
```

- ► May be less commented (but comments help you!)
- ► May be multiple versions
- ► May include analyses that are later discarded

# Final scripts

```
49  ###############################################################
50  ## Main 'pgibbs()' function
51
52
53  pgibbs <- function(gibbsState,
54                     maxit = 80000,
55                     verbose = TRUE,
56                     dbfile = "statepgibbs",
57                     deleteCache = FALSE,
58                     singleAgeCat = TRUE,
59                     sigmaE = NULL,
60                     delta = NULL) {
61      library(MASS)
62
63      ## Setup database of results
64      if(file.exists(dbfile)) {
65          if(deleteCache) {
66              message("removing existing cache file")
67              file.remove(dbfile)
68          }
69          else
70              stop(sprintf("cache file '%s' already exists", dbfile))
71      }
```
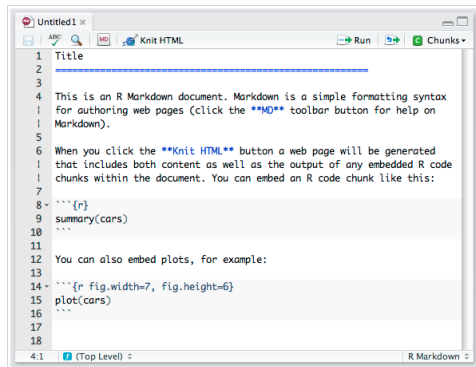
- ▶ Clearly commented
- ▶ Small comments liberally - what, when, why, how
- ▶ Bigger commented blocks for whole sections

# R markdown files

## 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:
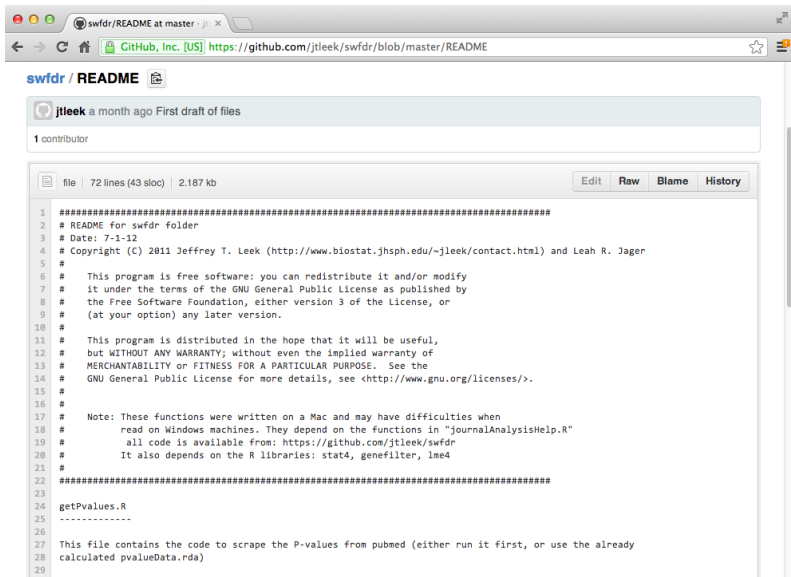


Note that the toolbar provides some useful tools for working with R Markdown:

- **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.

▶ R markdown files can be used to generate reproducible reports
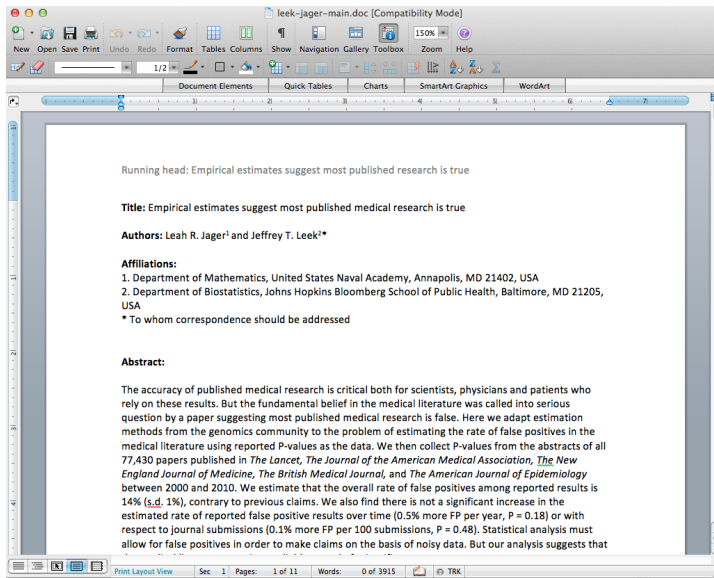▶ Text and R code are integrated

# Readme files



- Not necessary if you use R markdown
- Should contain step-by-step instructions for analysis

# Text of the document



- It should include a title, introduction (motivation), methods (statistics you used), results (including measures of

# Further resources

- Information about a non-reproducible study that led to cancer patients being mistreated: The Duke Saga Starter Set
- Reproducible research and Biostatistics
- Managing a statistical analysis project guidelines and best practices
- Project template - a pre-organized set of files for data analysis