# Organizing a Data Analysis

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

# Data analysis files

- Data

  - Raw data

  - Processed data

- Figures

  - Exploratory figures

  - Final figures

- R code

  - Raw scripts

  - Final scripts

  - R Markdown files (optional)

- Text

  - Readme files

  - Text of analysis

# Raw Data

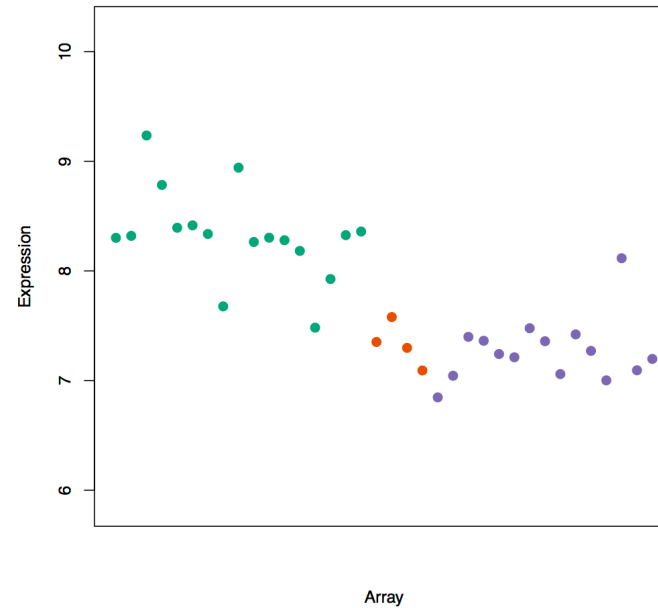

- Should be stored in your analysis folder

- 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.

- They do not need to be "pretty"

# Final Figures



- Usually a small subset of the original figures

- Axes/colors set to make the figure clear

- Possibly multiple panels

# Raw scripts



- May be less commented (but comments help you!)

- May be multiple versions

- May include analyses that are later discarded

# Final scripts



- Clearly commented

  - Small comments liberally - what, when, why, how

  - Bigger commented blocks for whole sections

- Include processing details

- Only analyses that appear in the final write-up

8/12

# 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:
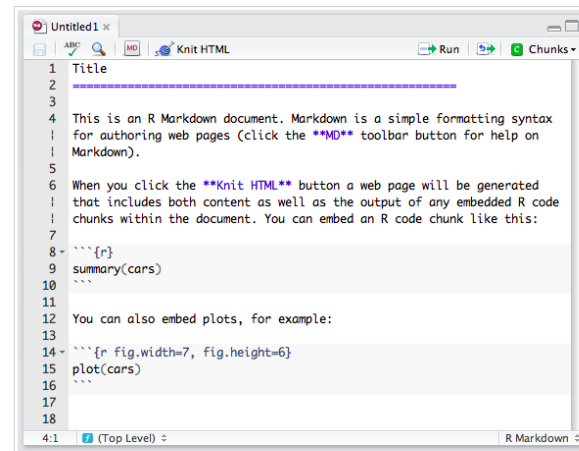
```
Untitled1 ×                                          □ □
  ABC  🔍    MD    Knit HTML                    Run  ⟲ ▶  C Chunks ▾
1  Title
2  ==========================================
3
4  This is an R Markdown document. Markdown is a simple formatting syntax
   for authoring web pages (click the **MD** 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    (Top Level) ⇕                                 R Markdown ⇕
```
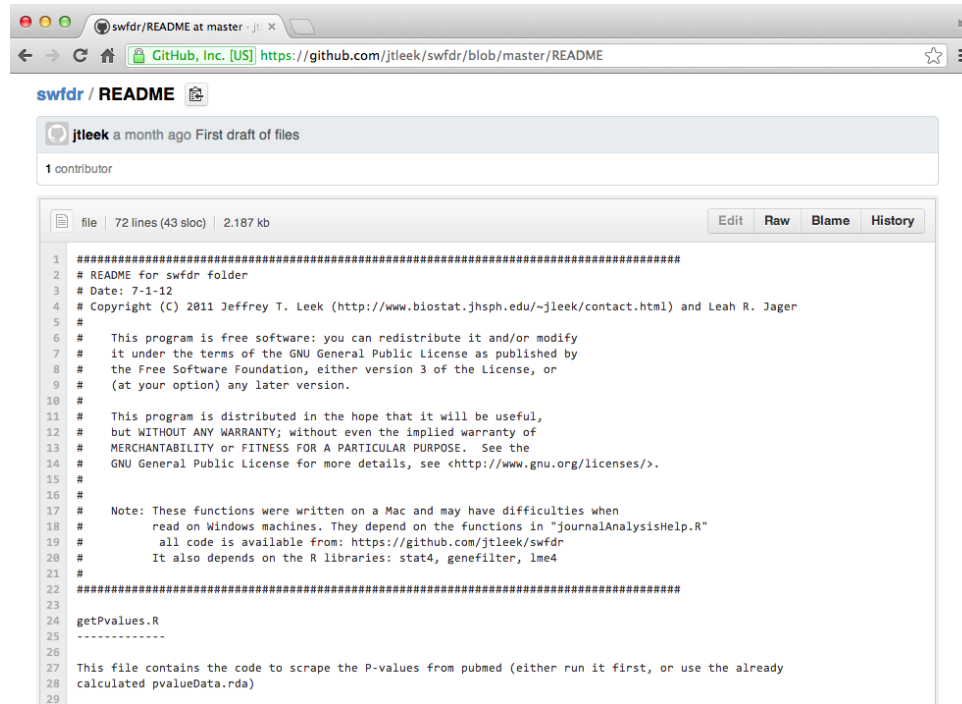
Note that the toolbar provides some useful tools for working with R Markdown:

- **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.

· R markdown files can be used to generate reproducible reports

· Text and R code are integrated
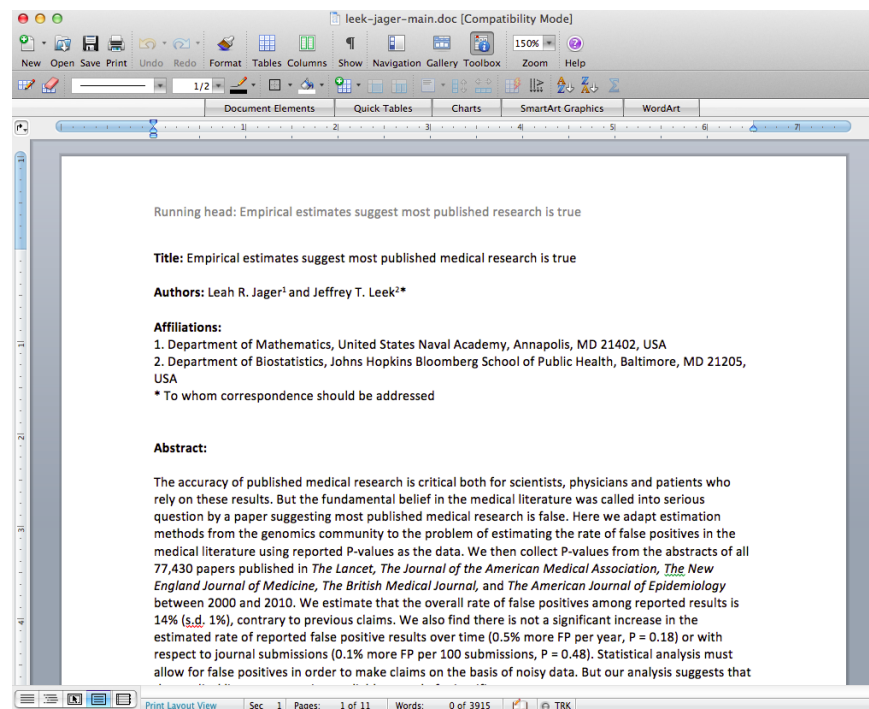
· Very easy to create in Rstudio

# Readme files



- Not necessary if you use R markdown

- Should contain step-by-step instructions for analysis

- Here is an example https://github.com/jtleek/swfdr/blob/master/README

# Text of the document



- It should include a title, introduction (motivation), methods (statistics you used), results (including measures of uncertainty), and conclusions (including potential problems)

- It should tell a story

- *It should not include every analysis you performed*

- References should be included for statistical methods

11/12

# 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