# Data Wrangling in R

http://sisbid.github.io/Data-Wrangling/
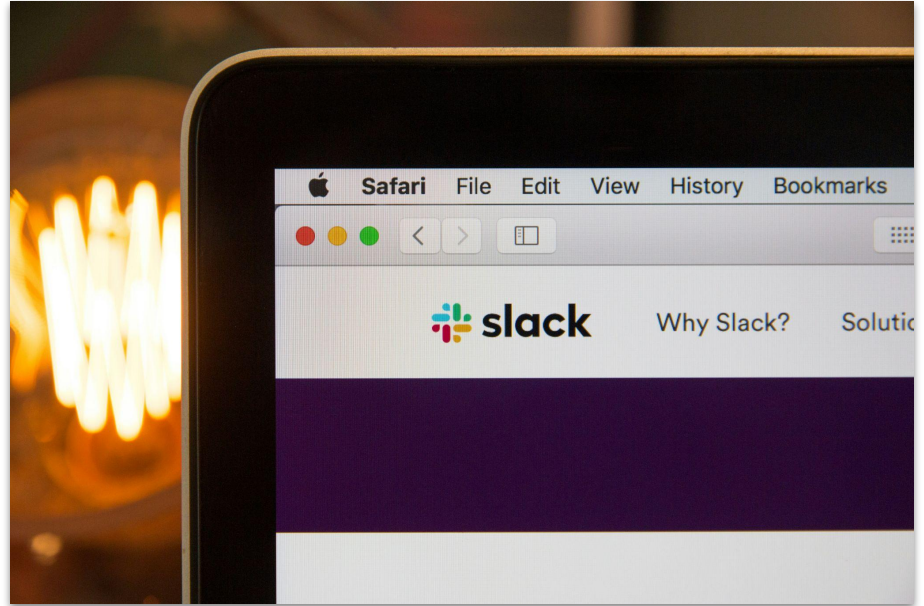
# Course Info

Course name    Data Wrangling in R

Instructors    <u>Carrie Wright</u> and <u>Ava Hoffman</u>

Course website    <u>http://sisbid.github.io/Data-Wrangling/</u>

Goals    Teach you how to get and clean data

Pre-reqs    Hopefully some R programming

Slack is helpful for:
- Getting help
- Talking with peers
- Accessing recordings

# What we will cover in this course:

1) Why data wrangling is important?
2) Importing data (and outporting)
3) Subsetting data
4) Summarizing data
5) Cleaning data
6) Reshaping data
7) Data merging and joining
8) Functional programming (efficiently applying functions)
9) Working with factors, strings, dates
10) Version control - Git/GitHub

But first, some jargon!

# Packages

A bundle or "package" of code (and or possibly data) that can be loaded together for easy repeated use or for sharing with others.

Packages are analogous to a software application like Microsoft Word on your computer.

# Function

A piece of code that allows you to do something in R. Packages often contain functions.

You can think of a function as <u>verb</u> in R.

A function might help you add numbers together, create a plot, or organize your data. More on that soon!

```
sum(1, 20234)

[1] 20235
```

# Argument

Something you pass to a function

```
round(0.627, digits = 2)

[1] 0.63
```

# An Object

Something that can be worked with or on in R - can be lots of different things!

You can think of an objects as a <u>noun</u> in R.

An object might be a data table, a plot, a function or more!

# Dataframes/DataTables/Spreadsheets

# Dataframes/Data tables

Rows = <u>samples</u> - individuals, locations, houses, viruses etc.

Columns = <u>variables</u> - aspects or features measured, color, count, type etc.

```
head(iris)
```

```
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1          5.1         3.5          1.4         0.2  setosa
2          4.9         3.0          1.4         0.2  setosa
3          4.7         3.2          1.3         0.2  setosa
4          4.6         3.1          1.5         0.2  setosa
5          5.0         3.6          1.4         0.2  setosa
6          5.4         3.9          1.7         0.4  setosa
```

# Need more help?

R jargon:

https://link.springer.com/content/pdf/bbm%3A978-1-4419-1318-0%2F1.pdf

# The Tidvyerse

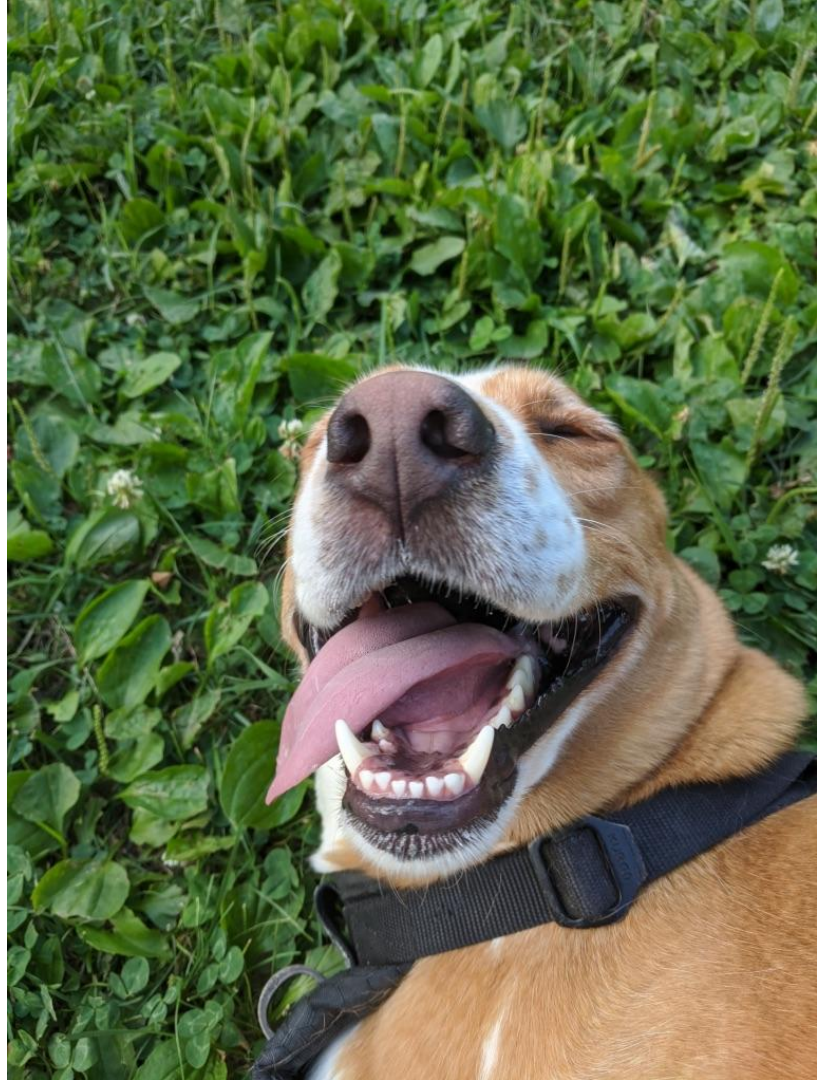Newer packages designed for data science that make R code more intuitive.

# How many people feel about data wrangling

# How we feel about data wrangling

# About us

Carrie

Ava

# About

Carrie Wright is a Senior Staff Scientist at the Fred Hutchinson Cancer Research Center. Dr. Wright's work is focused on innovating ways to make data science and computational biology more accessible. She is a member of the Open Case Studies team, the Genomic Data Science Community Network (GDSCN), and the Informatics Technology for Cancer Research (ITCR) Training Network (ITN). She also currently leads the ITCR Training and Outreach Working Group.

Previously, Dr. Wright was an Assistant Scientist in the Department of Biostatistics at the Johns Hopkins Bloomberg School of Public Health (JHSPH) and a member of the Johns Hopkins Data Science Lab (DaSL).

Prior to joining the JHSPH, Dr. Wright was a Postdoctoral Fellow at the Lieber Institute for Brain Development (LIBD), where her research focused on uncovering genetic mechanisms in psychiatric disease (with a particular emphasis on non-coding RNA) through the utilization of data science tools. At LIBD, Dr. Wright co-founded the LIBD rstats club, a community designed to encourage others to learn more about R programming and statistics. Dr. Wright has also served as an instructor for the Baltimore Underground Science Space and the Johns Hopkins Center for Talented Youth.

https://carriewright11.github.io/

# Open Case Studies

## What is the Open Case Studies (OCS) project?

The Open Case Studies project is an educational resource that educators can use in the classroom to teach students how to effectively derive knowledge from data in real-world challenges.



https://www.opencasestudies.org/

# AVA HOFFMAN

BALTIMORE, USA · AVAMARIEHOFFMAN @ GMAIL.COM

ABOUT

RESEARCH

RESOURCES

FUN STUFF

RESUME / CV

Hi! 👋

I'm a data scientist and ecologist. I'm currently exploring ways to make genomics research more accessible by connecting communities to cloud-based resources. I get really excited about things at the intersection of ecology and data, like population genetics and statistical modeling in nature. Lately, I've been especially interested in how plants evolve in man-made ecosystems (cities!) and how we can link those findings to public health outcomes.

I like coding, climbing things, everything DIY, and taking stuff apart to see how it works. I'm also working to get a community upcycling collective called Bed Roll Baltimore going!

https://www.avahoffman.com/

# Why this class

# rmarkdown

```
---
title: "My awesome website"
output:
    html_document:
        toc: true
        toc_float: true
        theme: cerulean
---
# This is Jeff's awesome website

![](https://media.giphy.com/media/d
rXGoW1iudhKw/giphy.gif)
```

# flexdashboard

```
---
title: "How does your BMI measure up?"
output: flexdashboard::flex_dashboard
runtime: shiny
---

Inputs {.sidebar}
-----------------------------------

```{r}
library(flexdashboard); library(NHANES); library(plotly);library(dplyr)
sliderInput("height", "Height in inches",0,100,72)
sliderInput("weight", "Weight in pounds",0,500,100)
sliderInput("age", "Age in years",0,120,50)

```

Column
-----------------------------------

### Chart 1

```{r}
nhanes = sample_n(NHANES,100)
renderPlotly({
  df = data.frame(bmi = c(nhanes$BMI,input$weight*0.45/(input$height*0.025)^2),
                  age = c(nhanes$Age,input$age),
                  who = c(rep("nhanes",100),"you"))
  ggplotly(ggplot(df) +
             geom_point(aes(x=age,y=bmi,color=who)) +
             scale_x_continuous(limits=c(0,90)) +
             scale_y_continuous(limits=c(0,60)) +
             theme_minimal()
           )
})
```
```

# httr

```r
library(httr)
library(dplyr)

username = 'janeeverydaydoe'

url_git = 'https://api.github.com/'


api_response =
GET(url = paste0(url_git, 'users/',
username, '/repos'))

content(api_response)[[1]]
```

```
$id
[1] 130377298
$node_id
[1] "MDEwOlJlcG9zaXRvcnkxMzAzNzcyOTg="
$name
[1] "first_project"
$full_name
[1] "JaneEverydayDoe/first_project"
$owner$gravatar_id
[1] ""
$owner$url
[1] "https://api.github.com/users/JaneEverydayDoe"

...
```

But also…

# Genomic signatures to guide the use of chemotherapeutics

Anil Potti[1,2], Holly K Dressman[1,3], Andrea Bild[1,3], Richard F Riedel[1,2], Gina Chan[4], Robyn Sayer[4], Janiel Cragun[4], Hope Cottrill[4], Michael J Kelley[2], Rebecca Petersen[5], David Harpole[5], Jeffrey Marks[5], Andrew Berchuck[1,6], Geoffrey S Ginsburg[1,2], Phillip Febbo[1,2,3], Johnathan Lancaster[4] & Joseph R Nevins[1,2,3]

**Using *in vitro* drug sensitivity data coupled with Affymetrix microarray data, we developed gene expression signatures that predict sensitivity to individual chemotherapeutic drugs. Each signature was validated with response data from an independent set of cell line studies. We further show that many of these signatures can accurately predict clinical response in individuals treated with these drugs. Notably, signatures developed to predict response to individual agents, when combined, could also predict response to multidrug regimens. Finally, we integrated the chemotherapy response signatures with signatures of oncogenic pathway deregulation to identify new therapeutic**

https://doi.org/10.1038/nm1491

Please watch!!!

Forensic Bioinformatics

When is Reproducibility an Ethical Issue? Genomics, Personalized Medicine, and Human Error

Keith A. Baggerly
Bioinformatics and Computational Biology
UT M. D. Anderson Cancer Center
kabagg@mdanderson.org

BIRS Workshop, Aug 14, 2013

http://www.birs.ca/events/2013/5-day-workshops/13w5083/videos/watch/201308141121-Baggerly.mp4

ℹ️ A Retraction to this article was published on 07 January 2011

ℹ️ A Corrigendum to this article was published on 01 August 2008

ℹ️ A Corrigendum to this article was published on 01 November 2007

NORTH CAROLINA

DURHAM COUNTY

DURHAM COUNTY
FILED
SEP 7 2011
4:03 O'CLOCK P M
CLERK OF SUPERIOR COURT

IN THE GENERAL COURT OF JUSTICE

SUPERIOR COURT DIVISION

1 CVS 4721

Richard Aiken, Jean K. Carroll,
Executrix of the Estate of Harold G.
Carroll, Jean K. Carroll, Individually,
Peggy Cox, as Administratrix of the Estate
of Paul F. Cox, Peggy Cox, Individually,
Helene L. Fligel, Jason Gannon, as
Personal Representative of the Estate of
Jennifer L. Gannon, John Haddock, as
Executor of the Estate of Karen Heath,
Walter Jacobs, as Executor of the Estate of
Juliet J. Jacobs, Walter Jacobs,
Individually, Polly Johnson, as Executor
of the Estate of Malcom W. Johnson, and
Polly Johnson, Individually,

               Plaintiffs

    vs.

**COMPLAINT**
**(JURY TRIAL DEMANDED)**

https://www.dukechronicle.com/article/2015/05/duke-lawsuit-involving-cancer-patients-linked-anil-potti-settled

Doesn't seem that important....

R Console

Typeset

```
> load("~/Documents/Work/workingpapers/openreview/data/processed-data-may11.rda")
> dim(dat)
[1] 730  15
> summary(glm(dat$correct ~ dat$study_type + dat$study_id,family="binomial"))

Call:
glm(formula = dat$correct ~ dat$study_type + dat$study_id, family = "binomial")

Deviance Residuals:
    Min       1Q    Median       3Q       Max
-1.6173  -1.4259    0.7941   0.9478    1.1431

Coefficients: (1 not defined because of singularities)
                        Estimate Std. Error z value Pr(>|z|)
(Intercept)               0.5675     0.1475   3.847  0.000120
dat$study_typenon-anon    0.4250     0.2182   1.948  0.051458
```
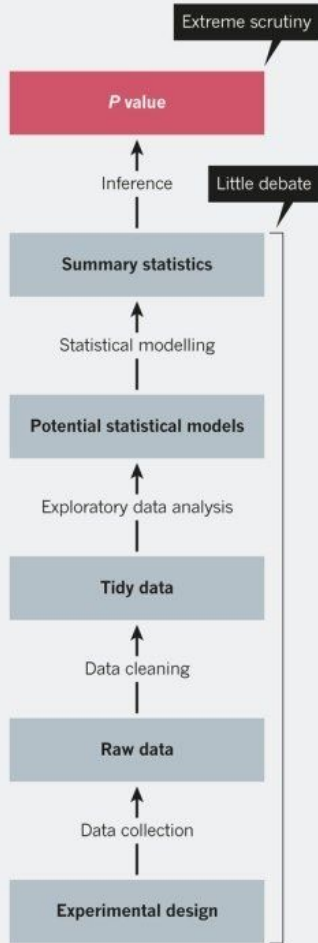
ON THE ONE
HAND...

**DATA PIPELINE**

The design and analysis of a successful study has many stages, all of which need policing.

Extreme scrutiny

*P* value

← Inference ← Little debate

Summary statistics

Statistical modelling

Potential statistical models

Exploratory data analysis

Tidy data

Data cleaning

Raw data

Data collection

Experimental design

- Most of the attention is on the last step
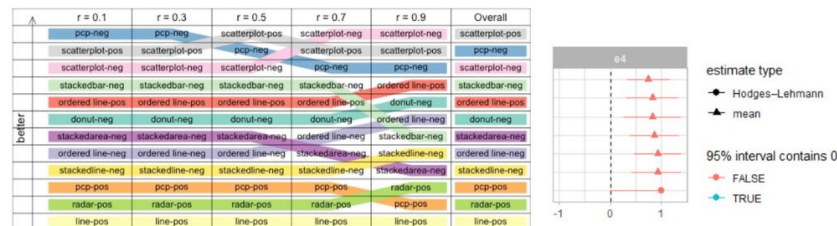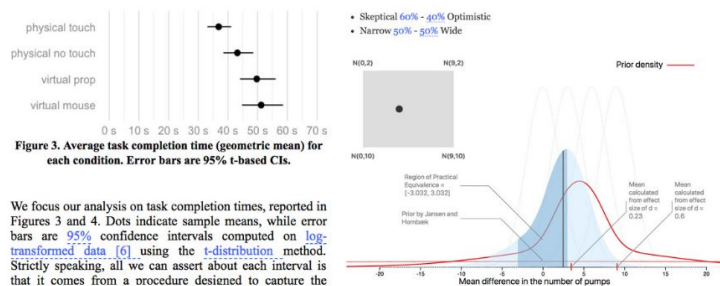- This course is about all the steps that come before
- They are *critical* for getting things rights

Leek and Peng Nature 2015

# The garden of forking paths: Why multiple comparisons can be a problem, even when there is no "fishing expedition" or "p-hacking" and the research hypothesis was posited ahead of time[*]

Andrew Gelman[†] and Eric Loken[‡]

14 Nov 2013

*"I thought of a labyrinth of labyrinths, of one sinuous spreading labyrinth that would encompass the past and the future ...I felt myself to be, for an unknown period of time, an abstract perceiver of the world."* — Borges (1941)

# Explorable Multiverse Analyses



Figure 3. Average task completion time (geometric mean) for each condition. Error bars are 95% t-based CIs.

We focus our analysis on task completion times, reported in Figures 3 and 4. Dots indicate sample means, while error bars are 95% confidence intervals computed on log-transformed data [6] using the t-distribution method. Strictly speaking, all we can assert about each interval is that it comes from a procedure designed to capture the

Figure 4. Perceptually-driven ranking of visualizations depending on the correlation sign (-neg / -pos), as a function of correlation value (r) and overall (right column).

Pierre Dragicevic (Inria), Yvonne Jansen (CNRS - Sorbonne Université), Abhraneel Sarma (University of Michigan)

Matthew Kay (University of Michigan), Fanny Chevalier (University of Toronto)

With **explorable multiverse analysis reports**, readers of research papers can explore alternative analysis options by interacting with the paper itself. This new approach to statistical reporting draws from two recent ideas: multiverse analysis, a philosophy of statistical reporting where paper authors report the outcomes of many different statistical analyses in order to show how fragile or robust their findings are; and explorable explanations, narratives that can be read as normal explanations but where the reader can also become active by dynamically changing some elements of the explanation.

https://explorablemultiverse.github.io/

And so we data wrangle

# What you wished data looked like

# What it actually looks like

```
@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
GGGTGGGCATTTCCACTCGCAGTATGGGTTGCCGCACGACAGGCAGCGGTCAGCCTGCGCTTTGGCCTGGCCTTCGGAAA
+HWI-EAS121:4:100:1783:1611#0/1
a``^\__`_````^a``a`^a_^___]a_]\]`a_____`_^^`]X]_]XTV_\]]NX_XVX]]_TTTTG[VTHPN]VFDZ
@HWI-EAS121:4:100:1783:322#0/1
CGTTTATGTTTTTGAATATGTCTTATCTTAACGGTTATATTTTAGATGTTGGTCTTATTCTAACGGTCATATATTTTCTA
+HWI-EAS121:4:100:1783:322#0/1
abaa`^aaaaabbbaababbbbbb`bbbb_bbbbbbbb`bbbaV^_a``a``]``aT]a__V\]]_]^a`]a_abbaV__
@HWI-EAS121:4:100:1783:1394#0/1
GGGTCTTTATTGGTCTGGTGATCCCCCATATTCTCCGGTTGTGTGGTTTAACCGATCATCGCGCATTACTTCCCGGCTGC
+HWI-EAS121:4:100:1783:1394#0/1
```[aa\b^^[]aabbb][`a_abbb`a``bbbbbabaabaaaab_VZa_^___bab_X`[a\HV_[_]_[^_X\T_VQQ
@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_[]\ZTDUZZUSOPX]]POP\GS\WSHHD
@HWI-EAS121:4:100:1783:455#0/1
GGGTAATTCAGGGACAATGTAATGGCTGCACAAAAAAATACATCTTTCATGTTCCATTGCACCATTGACAAATACATATT
+HWI-EAS121:4:100:1783:455#0/1
abb_babbabaabbbbbbbbbbbbbbbba\`b`\abbbabbbbabbbbbbaabbbbb`bb`ab_O_bab_Q_bbabaa_a
```

# What it actually looks like

# What it actually looks like

```
------------------------------  ALLERGIES  ------------------------------     ------------------------------  MEDICATION HISTOR

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 ON
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
eaction:            URINARY RETENTION                                       Medication: IBUPROFEN 600MG TAB
llergy Type:        DRUG                                                    Instructions: TAKE ONE TABLET BY MOUTH FOUR TI
A Drug Class:       NON-OPIOID ANALGESICS                                   Status: Active
bserved/Historical: HISTORICAL                                             Refills Remaining: 3
omments:            gradually worsening difficulty emptying bladder         Last Filled On: 20 Aug 2010
                                                                           Initially Ordered On: 01 Jul 2010
```

Jenny Bryan @JennyBryan · Apr 20

I'm seeking TRUE, crazy spreadsheet stories. Happy to get the actual sheet or just a description of the crazy. Also: I can keep a secret.

Slide from Jenny Bryan
(https://github.com/jennybc/2016-06_spreadsheets/blob/master/2016-06_useR-stanford.pdf)

Slide from Jenny Bryan
(https://github.com/jennybc/2016-06_spreadsheets/blob/master/2016-06_useR-stanford.pdf)

**Desiree Narango**
@DLNarango

Today's updates on #otherpeoplesdata:

8:56 AM - 22 Oct 2018

1 Like

# Where you wish data was

# Where data actually is

# Where data actually is

# Raw & processed data

"Data are values of qualitative or quantitative variables, belonging to a set of items."

"Data are values of qualitative or quantitative variables, belonging to a set of items."

**Set of items**: Sometimes called the population; the set of objects you are interested in

"Data are values of qualitative or quantitative variables, belonging to a set of items."

**Variables:** A measurement or characteristic of an item

"Data are values of qualitative or quantitative variables, belonging to a set of items."

**Qualitative:** Country of origin, sex, treatment
**Quantitative:** Height, weight, blood pressure

# Data sharing

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

# Tidy Data

**Hadley Wickham**
RStudio

## Abstract

A huge amount of effort is spent cleaning data to get it ready for analysis, but there has been little research on how to make data cleaning as easy and effective as possible. This paper tackles a small, but important, component of data cleaning: data tidying. Tidy datasets are easy to manipulate, model and visualise, and have a specific structure: each variable is a column, each observation is a row, and each type of observational unit is a table. This framework makes it easy to tidy messy datasets because only a small set of tools are needed to deal with a wide range of un-tidy datasets. This structure also makes it easier to develop tidy tools for data analysis, tools that both input and output tidy datasets. The advantages of a consistent data structure and matching tools are demonstrated with a case study free from mundane data manipulation chores.

# A tidy data set

One variable per column
One observation per row
One table per "kind" of data with
Linking variables across tables

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 prodromal and schizophrenia. I'm not sure if this was ok. Is it appropriate to do a ttest
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
Study design quirks

# Recipe

```
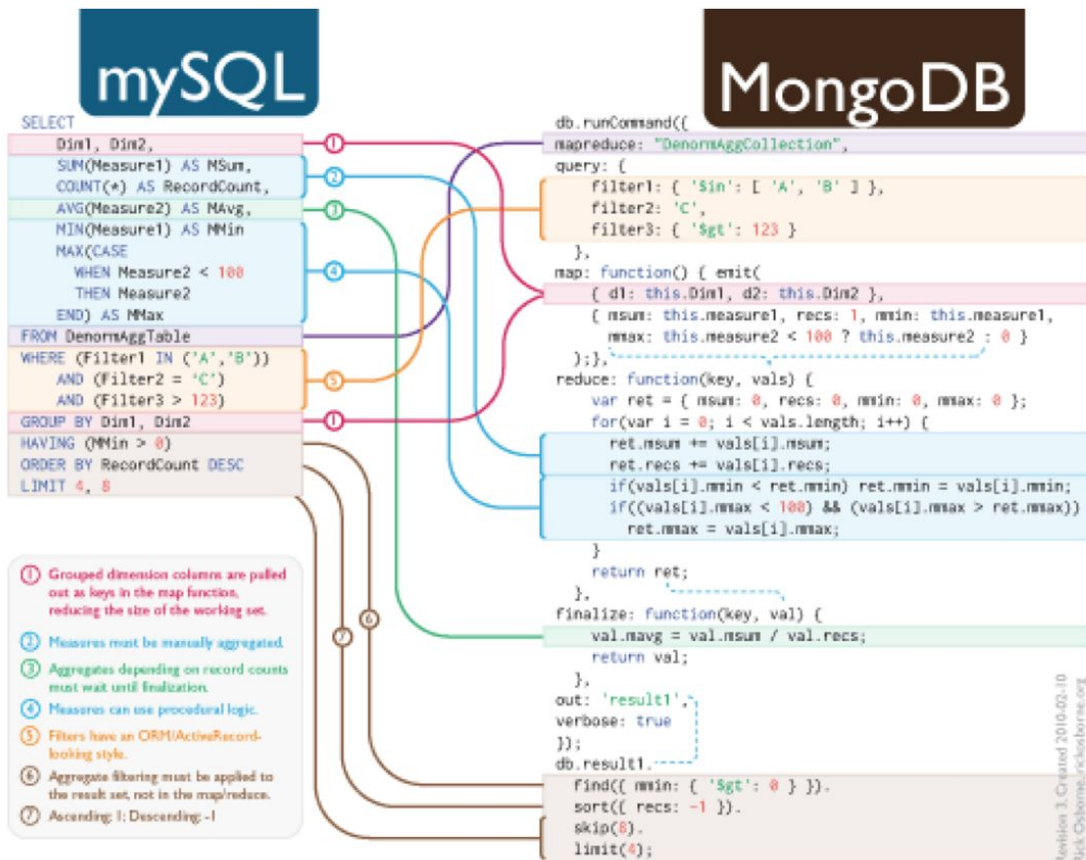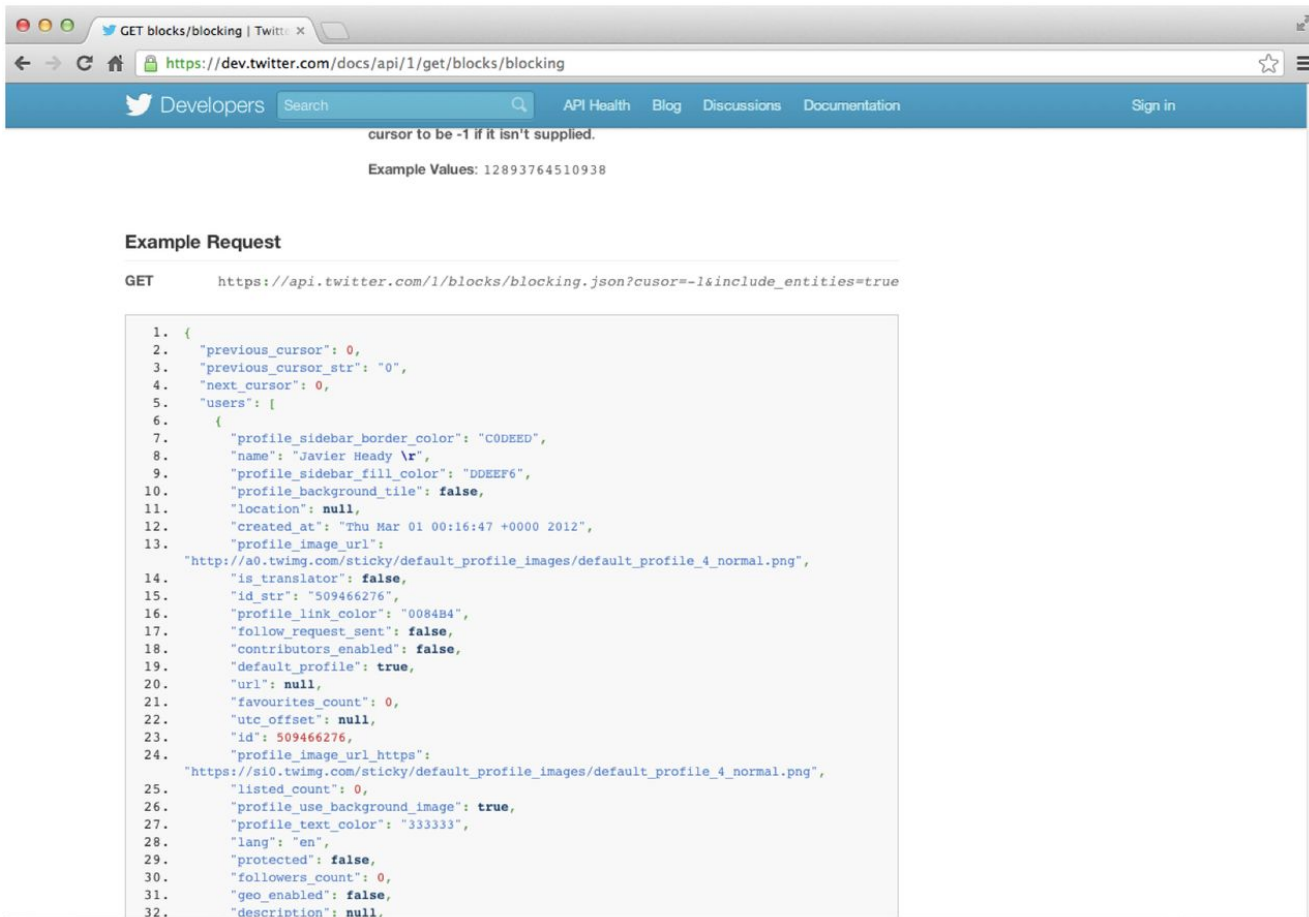33  library(sva)
34  library(ffpe)
35  library(RColorBrewer)
36  library(corrplot)
37  library(limma)
38  trop = RSkittleBrewer('tropical')
39  ```
40
41
42  ## Load the data
43
44  You will need to download the GEUVADIS ballgown object from this site: https://github.com/a      zee
    /ballgown_code
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(strsplit
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
```

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

RStudio toolbar: geuvadis.Rmd — ABC — ? — Knit HTML — Run — Chunks

1:1  (Top Level)  R Markdown

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**

| When.. | Be sure to… | So Do this… | Avoid this… | Why? |
|---|---|---|---|---|
| Naming variables (aka assigning column headers) | Use meaningful variable names | `AgeAtDiagnosis` | `ADx` | `ADx` is an unclear and uninformative abbreviation |
| Naming variables | Avoid spacing in column headers | `AgeAtDiagnosis` | `Age At Diagnosis` | Spacing in variable names makes the analyst's life more difficult |
| Naming variables | Use consistent capitalization | `AgeAtDiagnosis` | Using both `AgeAtDiagnosis` and `ageatdiagnosis` | Using consistent column names across tables/spreadsheets simplifies any merging the statistician may have to do. |
| Naming variables | Avoid using separators, but if it's necessary, use an underscore (`_`) | `IGF1` (or `IGF_1`) | `IGF.1`, `IGF-1`, `IGF/1`, `IGF,1` | Separators (commas, periods, hyphens, slashes, spaces etc.) often have different meanings in coding languages than they do in text. Avoiding them avoids error. |
| Coding variables | Avoid unnecessary spaces | 'male' | 'male ' | That extra space after 'male ' makes it different from 'male' without a space. |
| Coding variables | Be consistent! | 'male' | 'Male', 'male', and 'M', | In the eyes of the statistician, 'Male', 'male', and 'M' could be incorrectly perceived as three different values. |
| Coding variables | Be careful of spelling errors | 'male' | 'maale' | That extra 'a' makes these two different categories. |
| Coding date and time | Use ISO 8601 coding | 'YYYY-MM-DD' | 'MM/DD/YY' and `Month Day, Year` | Consistency simplifies the analyst's life, and YYYY-MM-DD will not be misconstrued if opened in Excel. |
| Coding missing data | Not leave any cells blank and use a consistent value | 'NA' | '0', '-9', red-highlighted blank cells, '. ', ' '-', … | Each cell should be filled with a consistent value. Pick a way to denote missingness (ideally 'NA') and stick with it. Avoid using numbers or punctuation to denote missing data. |
| Entering data | Stick to text and numbers | Convey all information with direct text/numerical entry | Using cell highlighting or font color to convey information | Your analyst may not use the same platform for analysis as you used for data entry, so avoiding font color and cell highlighting will minimize issues. |
| Generating an Excel file | Save the data in an appropriate format | Use one worksheet per table and save as CSV or text files | Multiple worksheets | Statisticians require this format to import your data onto other platforms. |
| Entering Data | Avoid entering unnecessary lines of text at the start | Start your first row with variable names | Adding lines of text | This violates the rules of tidy data and makes processing more difficult. Include this information in the "Code book" instead. |
| Opening files in Excel | Know and avoid its pitfalls | Consistently include one value per cell and be careful of date and time data. | Using macros, splitting cells, and merging cells | These formats are not amenable to data analysis on other platforms. |

Ellis SE, Leek JT. (2017) How to share data for collaboration. *PeerJ Preprints* 5:e3139v5 https://doi.org/10.7287/peerj.preprints.3139v5

# Rules for Tidy Spreadsheets

1. Be consistent

2. Choose good names for things

3. Write dates as YYYY-MM-DD

4. No empty cells

5. Put just one thing in a cell

6. Don't use font color or highlighting as data

7. Save the data as plain text files

# Organize thyself

"File organization and naming are powerful weapons against chaos."
- Jenny Bryan

Slide via Jenny Bryan:
http://www.slideshare.net/jenniferbryan5811/cm002-deep-thoughts

- ▼ 📁 code
  - ▶ 📁 final_code
  - ▶ 📁 raw_code
- ▼ 📁 data
  - ▶ 📁 raw_data
  - ▶ 📁 tidy_data
- ▶ 📁 figures
- ▼ 📁 products
  - ▶ 📁 writing

# 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
```

# 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

# Raw scripts



```r
1   library(chron)
2   library(affy)
3   library(oligoClasses)
4   celfiles <- list.celfiles("~/Projects/batchreview/",listGzipped=T)
5   dts <- sapply(celfiles,celfileDate)
6
7   ll <- strsplit(dts,"-")
8
9   yy <- as.numeric(lapply(ll,function(x){x[1]}))
10  mm <- as.numeric(lapply(ll,function(x){x[2]}))
11  dd <- as.numeric(lapply(ll,function(x){x[3]}))
12
13  jul <- julian(mm,dd,yy)
14
15  # Identify the arrays corresponding to CEU parents
16  ceuparents <-scan("~/Documents/Work/workingpapers/CHEUNG/CEU_parents.txt",what="character")
17  tmp <- list.files("~/Documents/Work/workingpapers/CHEUNG/CEU_data")
18
19  rep <- rep(c(0,1),each=100)
20  for(i in 1:length(ceuparents)){
21
22  }
23
24
25  tmp <- tmp[9:272]
26  array <- as.character(sapply(strsplit(tmp,"_"),function(x){x[1]}))
27  sample <- as.character(sapply(strsplit(tmp,c("_")),function(x){x[2]}))
28  sample <- as.character(sapply(strsplit(sample,c("\\.")),function(x){x[1]}))
29  rp <- as.character(sapply(strsplit(tmp,"_"),function(x){x[3]}))
30  rp <-   as.character(sapply(strsplit(rp,c("\\.")),function(x){x[1]}))
31
32
33  ceufiles <- array[sample %in% ceuparents]
34
35
```

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

- May be multiple versions

- May include analyses that are later discarded

# Final scripts



```r
1  f.pvalue <- function(dat,mod,mod0){
2      # This is a function for performing
3      # parametric f-tests on the data matrix
4      # dat comparing the null model mod0
5      # to the alternative model mod.
6      n <- dim(dat)[2]
7      m <- dim(dat)[1]
8      df1 <- dim(mod)[2]
9      df0 <- dim(mod0)[2]
10     p <- rep(0,m)
11     Id <- diag(n)
12
13     resid <- dat %*% (Id - mod %*% solve(t(mod) %*% mod) %*% t(mod))
14     resid0 <- dat %*% (Id - mod0 %*% solve(t(mod0) %*% mod0) %*% t(mod0))
15
16     rss1 <- resid^2 %*% rep(1,n)
17     rss0 <- resid0^2 %*% rep(1,n)
18
19     fstats <- ((rss0 - rss1)/(df1-df0))/(rss1/(n-df1))
20     p <-  1-pf(fstats,df1=(df1-df0),df2=(n-df1))
21     return(p)
22 }
23
24 setwd("cheung/")
25 # Load data and create group variable
26 dat <- read.table("full.data")
27
28 jpt.names <- scan("JPT.cname.txt",what="character")
29 chb.names <- scan("CHB.cname.txt",what="character")
30 ceu.names <- scan("CEU_parents.txt",what="character")
31 nceu <- length(ceu.names)
32 njpt <- length(jpt.names)
33 nchb <- length(chb.names)
34
```

- 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

# This is the README file for my_first_project

Last updated: 02-Mar-2018

The folders in this project are:

- *data* - is the folder where you can find all the collected data.
- *figures* - is where you can find all the plots, data pictures, and other images.
- *code* - is where you can find code files for collecting, cleaning up, or analyzing data.
- *products* - is where you can find reports, presentations, or products

Data on crime is obtained from International Crime Data collected between 2015-2018 and is publicly available. Data on happiness is collected from the Survey of International Happiness.

Contributors:

- Jane Everyday Doe, jane.everyday.doe@gmail.com
- John Everyday Doe, john.everyday.doe@gmail.com

Cite: Doe, J, and Doe, J, Sample Analysis Using Sample Data, Working Paper, 2018

# Just no

key principles of file naming for data science projects:

- Machine readable
- Human readable
- Be nicely ordered

Source: Jenny Bryan

| Bad Naming | Good Naming |
|---|---|
| 2013 my report.md | 2013_my_report.md |
| malik's_report.md | maliks_report.md |
| 01_zoë_report.md | 01_zoe_report.md |
| AdamHooverReport.md | adam-hoover-report.md |
| executivereportpepsiv1.md | executive_report_pepsi_v1.md |

2018_jan_sales_cust001_prod001.md
2017_mar_sales_cust001_prod001.md
2016_may_sales_cust001_prod008.md
2017_jan_sales_cust120_prod007.md
2015_oct_sales_cust034_prod001.md
2015_oct_sales_cust034_prod002.md

| Year | Month | Type | Customer ID | Product ID |
|------|-------|------|-------------|------------|
| 2018 | jan | sales | 001 | 001 |
| 2017 | mar | sales | 001 | 001 |
| 2016 | may | sales | 001 | 008 |
| 2017 | jan | sales | 120 | 007 |
| 2015 | oct | sales | 034 | 001 |
| 2015 | oct | sales | 034 | 002 |

# Which one is better?

analysis.R
or
2017-exploratory_analysis_crime.R?

# Which one is better?

05-21-2017-analysis-cust001.R
or
2017-05-21-analysis-cust001.R?

# Structure of a filename

processed_pvalue_data_from_pubmed_oct24.rda

# What did I do to this data

processed_pvalue_data_from_pubmed_oct24.rda

# What kind of data is this?

processed_pvalue_data_from_pubmed_oct24.rda

# Where did it come from?

processed_pvalue_data_from_pubmed_oct24.rda

# When did I get it?

processed_pvalue_data_from_pubmed_oct24.rda

# Underscores/slashes not dots/whitespace

processed_pvalue_data_from_pubmed_oct24.rda

# Consistency is the main rule

processed_pvalue_data_from_pubmed_oct24.rda
raw_pvalue_data_from_pubmed_oct24.rda

Your closest collaborator is you six months ago, but you don't reply to emails

- Karl Broman

http://kbroman.org/Tools4RR/assets/lectures/06_org_eda.pdf

Step 1: slow down and document.
Step 2: have sympathy for your future self.
Step 3: have a system.

- Karl Broman

http://kbroman.org/Tools4RR/assets/lectures/06_org_eda.pdf

# R + RStudio

# The R Project for Statistical Computing

**Download**

CRAN

**R Project**

About R
Contributors
What's New?
Mailing Lists
Bug Tracking
Conferences
Search

**R Foundation**

Foundation
Board
Members
Donors
Donate

## Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To **download R**, please choose your preferred CRAN mirror.

If you have questions about R like how to download and install the software, or what the license terms are, please read our answers to frequently asked questions before you send an email.

## News

- **The R Journal Volume 7/1** is available.

- **R version 3.2.1 (World-Famous Astronaut)** has been released on 2015-06-18.

- **R version 3.1.3 (Smooth Sidewalk)** has been released on 2015-03-09.

- **useR! 2015**, will take place at the University of Aalborg, Denmark, June 30 - July 3, 2015.

- **useR! 2014**, took place at the University of California, Los Angeles, USA June 30 - July 3, 2014.

www.r-project.org

https://www.rstudio.com/

RStudio

□ □ □                                          RStudio                                          ⬒

[New] ▾  [Open] ▾  ▣  ▣  🖨   [→ Go to file/function          ]                    🔲 Project: (None) ▾

┌─ 🔲 Untitled1 ✕ ─────────────────────────────────┐  ┌─ Workspace  History ──────────── ▭ ▢ ┐
│                                          ▭ ▢ │  │                                            │
│ 🔲 [← →] [💾] ⬚ Source on Save [🔍] [✦▾]  [↗ Run] [⟳] [□ Source ▾] [▥] │  ┌─ Files  Plots  Packages  Help ──────── ▭ ▢ ┐
│  1 │                                       │  │ [← →] | [🏠] | [🖨] [↗] [⟳]    [🔍           ] │
│    │                                       │  │                                            │
│    │                                       │  │ Markdown Quick Reference ▾  [Find in Topic ] │
│    │                                       │  │ 3. Item 3                                   │
│    │                                       │  │    * Item 3a                                │
│    │                                       │  │    * Item 3b                                │
│    │                                       │  │                                            │
│    │                                       │  │ *Manual Line Breaks*                        │
│    │                                       │  │ End a line with two or more spaces:         │
│    │                                       │  │ ┌──────────────────────────────────────┐   │
│    │                                       │  │ │ Roses are red,                       │   │
│    │                                       │  │ │ Violets are blue.                    │   │
│ 1:1 │ 🔲 (Top Level) ⇕                R Script ⇕ │  │ └──────────────────────────────────────┘   │
└──────────────────────────────────────────┘  │ *Links*                                     │
┌─ Console ~/Dropbox/Jeff/teaching/2013/modules/ALL UNUSED CONTENT/toolBox/ ↗ ▭ ▢ ┐  │ Use a plain http address or add a link to a phrase: │
│                                              │  │ ┌──────────────────────────────────────┐   │
│ > |                                          │  │ │ http://example.com                   │   │
│                                              │  │ │                                      │   │
│                                              │  │ │ [linked phrase](http://example.com)  │   │
│                                              │  │ └──────────────────────────────────────┘   │
│                                              │  │ *Images*                                    │
│                                              │  │ Images on the web or local files in the same directory: │
│                                              │  │ ┌──────────────────────────────────────┐   │
│                                              │  │ │ ![alt text](http://example.com/logo.png) │ │
│                                              │  │ │                                      │   │
│                                              │  │ │ ![alt text](figures/img.png)         │   │
│                                              │  │ └──────────────────────────────────────┘   │
│                                              │  │ *Blockquotes*                               │
│                                              │  │ ┌──────────────────────────────────────┐   │
│                                              │  │ │ A friend once said:                  │   │
│                                              │  │ │ > It's always better to give         │   │
│                                              │  │ │ > than to receive.                   │   │
│                                              │  │ └──────────────────────────────────────┘   │
│                                              │  │ *R Code Blocks*                             │
│                                              │  │ R code will be evaluated and printed        │
│                                              │  │ ┌──────────────────────────────────────┐   │
│                                              │  │ │ ```{r}                               │   │
│                                              │  │ │ summary(cars$dist)                   │   │
└──────────────────────────────────────────┘  └──────────────────────────────────────────┘

https://www.rstudio.com/

# Welcome to AnVIL

The NHGRI AnVIL (Genomic Data Science Analysis, Visualization, and Informatics Lab-space) is a project powered by Terra for biomedical researchers to **access data**, **run analysis tools**, and **collaborate**.

Find how-to's, documentation, video tutorials, and discussion forums ↗

### View Workspaces

Workspaces connect your data to popular analysis tools powered by the cloud. Use Workspaces to share data, code, and results easily and securely.

### View Examples

Browse our gallery of showcase Workspaces to see how science gets done.

### Browse Data

Access data from a rich ecosystem of data portals.

POWERED BY Terra   BETA

https://anvil.terra.bio/#

https://anvil.terra.bio

# What is AnVIL???

https://anvilproject.org/

# AnVIL is "renting computers"



Standard Computing
- You buy a laptop one time
- You get that one laptop
- You pay little per use

Cloud computing
- You use any web browser
- You rent the computers
- You pay per hour/gigabyte/etc.

# It can feel a little weird





Purchased car
- You buy the car
- You fill up at a station
- You pay less per mile

ZipCar
- You don't buy the car
- You pay by the mile
- You may pay more per mile

# AnVIL: Data + Sharing + Platforms

# AnVIL Dataset Catalog

| Search | | Consortium | Cohorts | Diseases | Cohorts | Data Types | Cohorts | Consent Code | Cohorts | Access | Cohorts |
|---|---|---|---|---|---|---|---|---|---|---|---|
| e.g. disease, study name, dbGaP I | | ☐ 1000 Genomes | 1 | ☐ Alzheimer's disease | 3 | ☐ Exome | 138 | ☐ DS | 4 | ☐ Consortium Access | 99 |
| | | ☐ CCDG | 196 | ☐ asthma | 1 | ☐ RNAseq | 1 | ☐ DS-ASD | 18 | ☐ Controlled Access | 145 |
| | | ☐ CMG | 36 | ☐ atherosclerosis | 1 | ☐ Whole Genome | 108 | ☐ DS-ASD-IRB | 5 | ☐ Open Access | 2 |
| | | ☐ Convergent Neuro | 2 | ☐ atrial fibrillation | 13 | | | ☐ DS-ASD-IRB-COL | 1 | | |
| | | + 4 more | | + 17 more | | | | + 69 more | | | |

No selected terms.

Download TSV ⬇   Copy URL ▢

## Search Summary

| Consortium | Cohorts | Samples | Subjects | Size (TB) |
|---|---|---|---|---|
| 1000 Genomes | 1 | 3,202 | 3,202 | 72.98 |
| CCDG | 196 | 250,770 | 243,226 | 2,381.24 |
| CMG | 36 | 11,424 | 10,063 | 73.61 |
| Convergent Neuro | 2 | 304 | 304 | 5.32 |
| GTEx (v8) | 1 | 17,382 | 979 | 182.14 |

https://anvilproject.org/data

# AnVIL Data Dashboard

Extensive unrestricted and protected data sets already available within AnVIL

- 246 cohorts (CCDG, CMG, GTEx,1000G, eMerge)
- 285k subjects
- 3Pb and rapidly growing

- Open access (e.g. 1000G), dbGaP authenticated (e.g. GTex) and consortium authenticated (e.g. CCDG) options available



https://anvilproject.org/data

# AnVIL Analysis Platforms



**Jupyter**

+ Code, text and plots in one document

+ Supports coding in Python or R

- Least scalable, not a complete IDE

**Galaxy**

+ Graphical interface for thousands of tools and workflows

+ Highly accessible and reproducible

- Tools must be preconfigured to use

**R Studio®**

+ Feature rich IDE for programming in R

+ Rich statistics & ML and visualizations

- Limited support for other programming languages

**{wdl}**

+ Extremely scalable and flexible

- Most technically demanding

- Unpredictable and potentially large costs

# AnVIL Analysis Platforms



R Studio®

+ Feature rich IDE for programming in R

+ Rich statistics & ML and visualizations

- Limited support for other programming languages



THIS ONE

Also great!

https://rstudio.cloud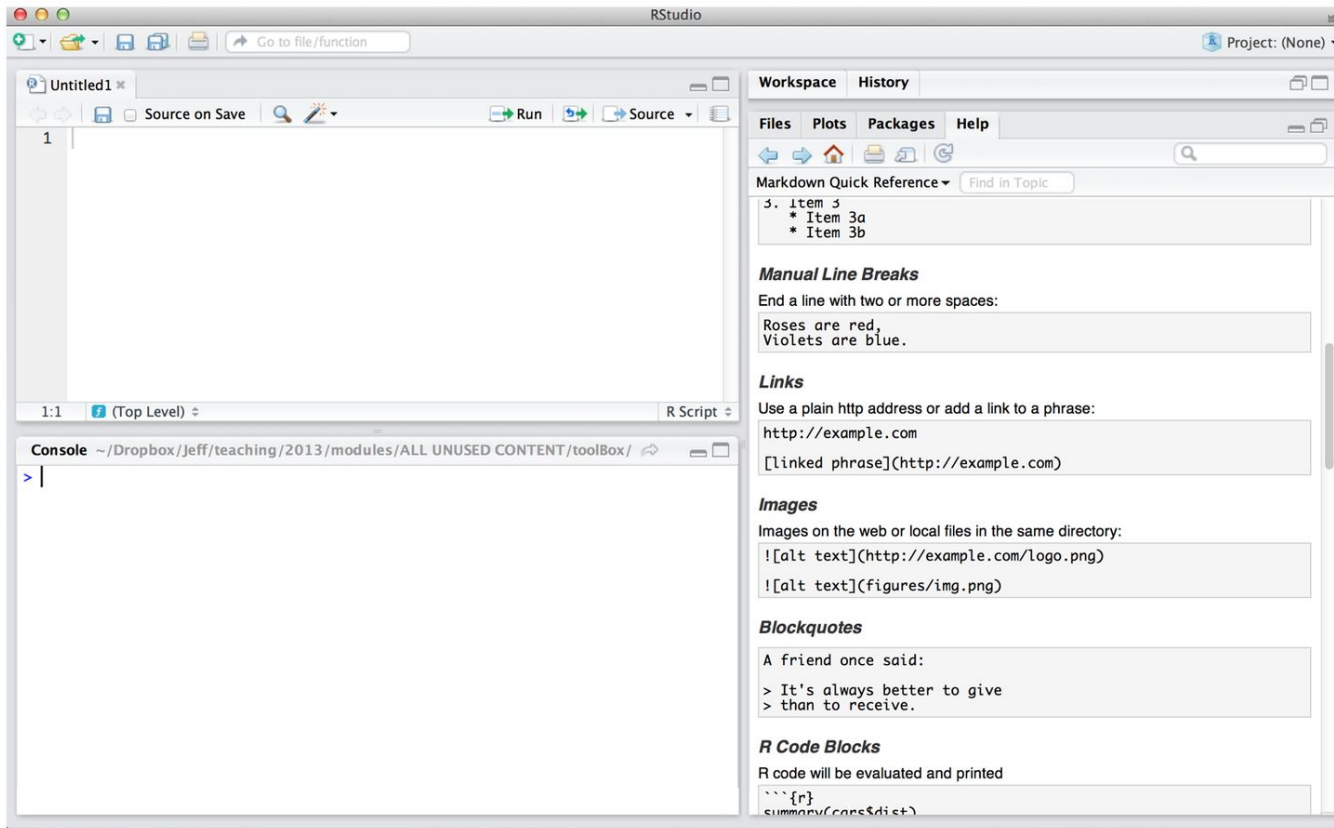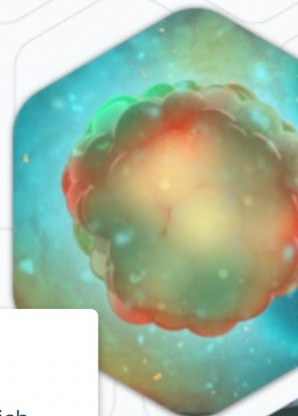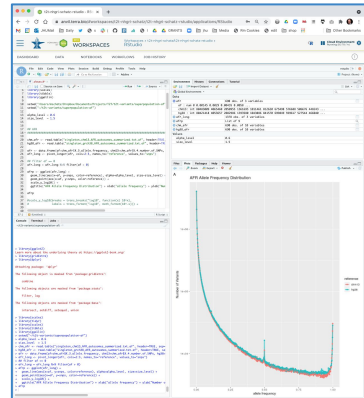