# 2. Manipulate

R$^3$ training
[noaa-iea.github.io/r3-train](noaa-iea.github.io/r3-train)

Ben Best
2021-06-21

# 1. Collaborate Recap

- Good job!
  - [github.com/noaa-iea/r3-train/**pulls**](github.com/noaa-iea/r3-train/pulls)
    - Open source project example: [icons/pull/58](icons/pull/58)
  - [github.com/noaa-iea/r3-train/**issues**](github.com/noaa-iea/r3-train/issues)
  - Merge pull request with Github Desktop; + resolve file conflicts with Atom
- ToDo for me:
  - Respond to [notes | r3-train - Google Docs](notes | r3-train - Google Docs)
  - Google Group
  - + Feedback menu
  - Github Discussions

# 2. Manipulate

# For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insights

nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-to-insights-is-janitor-work.html
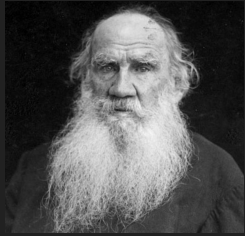
"Data scientists, according to interviews and expert estimates, spend from 50 to 80 percent of their time mired in the mundane labor of collecting and preparing data, before it can be explored for useful information." — NY Times (2014)

```
mdl      <- maxent(tbl_env, tbl_obs)
tbl_pred <- predict(mdl, tbl_envnew)
```

"Happy families are all alike; every unhappy family is unhappy in its own way."
— Leo Tolstoy



"Tidy datasets are all alike, but every messy dataset is messy in its own way."
— Hadley Wickham

5

A table is tidy if:

**A** **B** **C**

**&**

**A** **B** **C**

Each **variable** is in its own **column**

Each **observation**, or **case**, is in its own **row**

ie "long" vs "wide"

"Tidy datasets are all alike, but every messy dataset is messy in its own way."
— Hadley Wickham

# Tidy Manifesto

cran.r-project.org/web/packages/tidyverse/vignettes/manifesto.html

1. Reuse existing data structures

2. Compose simple functions with the pipe

3. Embrace functional programming

INPUT x

FUNCTION f:

OUTPUT f(x)

4. Design for humans



Fall into the Pit of Success

blog.codinghorror.com/
falling-into-the-pit-of-success

,,

# Tidyverse process   & R packages

**Import**
**readr**
**readxl**
DBI
jsonlite
googledrive

spatial:
sf
raster

**Tidy**
**tibble**
**dplyr**
**tidyr**
**purrr**

Program
devtools
usethis

Visualize
static:
ggplot2
ggraph
ggmap
interactive:
plotly
leaflet
dygraphs

Transform
**lubridate**
**stringr**
**scales**

Model
caret
broom
modelr

Communicate
rmarkdown
flexdashboard
bookdown
shiny

9

# Cheatsheets are

RStudio: Help > Cheatsheets > Data wrangling with dplyr

# Resources



[r4ds.had.co.nz](http://r4ds.had.co.nz)

[google.com](http://google.com)

[stackoverflow.com](http://stackoverflow.com)

# Lesson

[noaa-iea.github.io/r3-train/manipulate.html](noaa-iea.github.io/r3-train/manipulate.html)