

Class 06

Jonathan Tannen

- Anatomy of a Project: Dan Lopez, Director of Software Engineering, OIT
- *10 minute break*
- Working Groups
- Better Engineering: Clean Functions

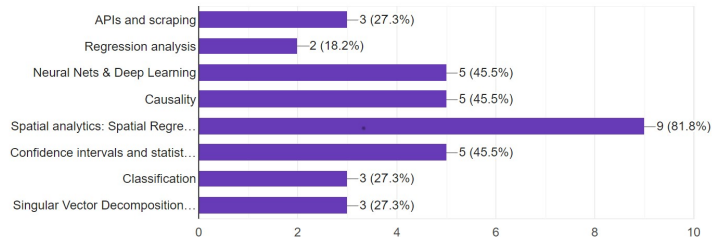
Project Proposal 1

Please remember to *upload to your github repo*.

Survey Results

What technical or GIS topics would you like covered in lecture?

11 responses



30 minutes

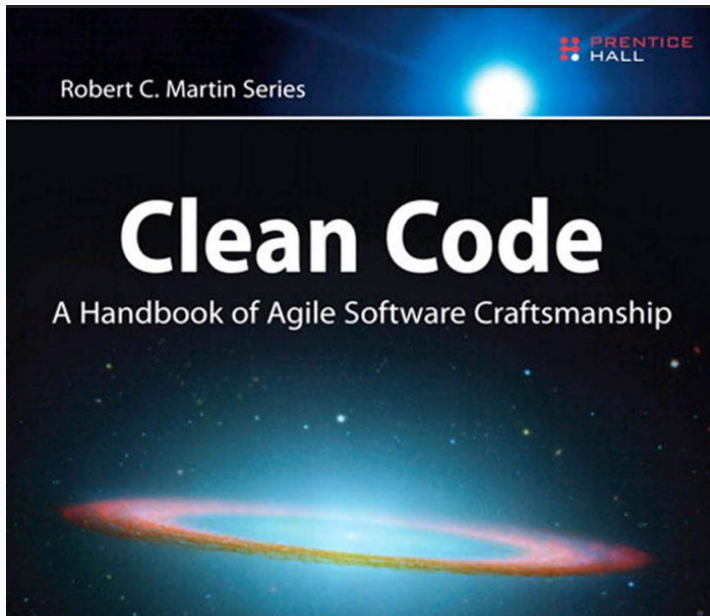
- A: GIS data wrangling in ArcMap
- B: GIS data wrangling in Python/R
- C: Building Interactive Tools / Markdown
- D: Using APIs / web tools
- E: Research question & statistical methods

Break into groups of 3.

Each person gets 10 minutes to discuss and actively debug.

Clean Functions

Inspired by Clean Code (Martin) and R for Data Science (Grolemund & Wickham).



Rule 1: Do not copy code more than once.

aka “Don’t Repeat Yourself” aka “Single Source of Truth”

```
df <- tibble::tibble(  
  a = rnorm(10),  
  b = rnorm(10),  
  c = rnorm(10),  
  d = rnorm(10)  
)  
  
df$a <- (df$a - min(df$a, na.rm = TRUE)) /  
  (max(df$a, na.rm = TRUE) - min(df$a, na.rm = TRUE))  
df$b <- (df$b - min(df$b, na.rm = TRUE)) /  
  (max(df$b, na.rm = TRUE) - min(df$a, na.rm = TRUE))  
df$c <- (df$c - min(df$c, na.rm = TRUE)) /  
  (max(df$c, na.rm = TRUE) - min(df$c, na.rm = TRUE))  
df$d <- (df$d - min(df$d, na.rm = TRUE)) /  
  (max(df$d, na.rm = TRUE) - min(df$d, na.rm = TRUE))
```



```
my_rescale <- function(x) {  
  rng <- range(x, na.rm = TRUE)  
  (x - rng[1]) / (rng[2] - rng[1])  
}
```

```
df$a <- my_rescale(df$a)  
df$b <- my_rescale(df$b)  
df$c <- my_rescale(df$c)  
df$d <- my_rescale(df$d)
```

```
my_rescale <- function(x) {  
  rng <- range(x, na.rm = TRUE)  
  (x - rng[1]) / (rng[2] - rng[1])  
}  
  
for(column in c("a","b","c","d")){  
  df[[column]] <- my_rescale(df[[column]])  
}
```

What was achieved?

```
df <- tibble::tibble(  
  a = rnorm(10),  
  b = rnorm(10),  
  c = rnorm(10),  
  d = rnorm(10)  
)  
  
df$a <- (df$a - min(df$a, na.rm = TRUE)) /  
  (max(df$a, na.rm = TRUE) - min(df$a, na.rm = TRUE))  
df$b <- (df$b - min(df$b, na.rm = TRUE)) /  
  (max(df$b, na.rm = TRUE) - min(df$a, na.rm = TRUE))  
df$c <- (df$c - min(df$c, na.rm = TRUE)) /  
  (max(df$c, na.rm = TRUE) - min(df$c, na.rm = TRUE))  
df$d <- (df$d - min(df$d, na.rm = TRUE)) /  
  (max(df$d, na.rm = TRUE) - min(df$d, na.rm = TRUE))
```

What was achieved?

- Fewer mistakes, fewer places to change.
- Easier to understand.
- Clear SLA (Service Level Agreement)
- Encapsulation
- *Testable*.

```
df <- read.csv("C:/Users/Jonathan/papers/election_paper/data/data_2000.csv")
election_metadata <- read.csv("C:/Users/Jonathan/papers/election_paper/data/election_metadata.csv")
candidates <- read.csv("C:/Users/Jonathan/papers/election_paper/data/candidates.csv")

df$is_last_four_elections <- (ymd("2022-02-16") - ymd(df$date)) <= years(4)
df$age <- (ymd("2022-02-16") - ymd(df$date))
```

Do not copy code more than once: Globals

```
DIR <- "C:/Users/Jonathan/papers/election_paper/" # also could setwd().
DATA_DIR <- paste0(DIR, "data/")
DATE <- "2022-02-16"

df <- read.csv(paste0(DATA_DIR, "data_2000.csv"))
election_metadata <- read.csv(paste0(DATA_DIR, "elections.csv"))
candidates <- read.csv(paste0(DATA_DIR, "candidates.csv"))

df$is_last_four_elections <- (ymd(DATE) - ymd(df$date)) <= years(4)
df$age <- (ymd(DATE) - ymd(df$date))
```

Do not copy code more than once: Globals

```
DIR <- "C:/Users/Jonathan/papers/election_paper/" # also could setwd().
DATA_DIR <- paste0(DIR, "data/")
DATE <- "2022-02-16"

data_path <- function(file) paste0(DATA_DIR, file)

df <- read.csv(data_path("data_2000.csv"))
election_metadata <- read.csv(data_path("elections.csv"))
candidates <- read.csv(data_path("candidates.csv"))

df$age <- (ymd(DATE) - ymd(df$date))
df$is_last_four_elections <- df$age <= years(4)
```

- Use good names.
 - Full words
 - Searchable
- Functions should be small.
 - Usually not more than 20 lines.
- Functions should do one thing. (Single Responsibility Principle)

Small Functions

```
49 ▾ clean_and_save_raw_data <- function(){
50   data_regex <- "([0-9]+)_([a-z]+).csv"
51   all_files <- list.files(RAW_DATA_FOLDER, pattern = data_regex)
52
53   res <- list()
54 ▾   for(file in all_files){
55     year <- gsub(data_regex, "\\1", file)
56     election_type <- gsub(data_regex, "\\2", file)
57 ▾     if(year == "2020"){
58       res[[file]] <- clean_raw_data_2020(file)
59 ▾     } else{
60       res[[file]] <- clean_raw_data(file)
61 ▾     }
62     res[[f]]$year <- year
63     res[[f]]$election_type <- election_type
64 ▾   }
65
66   res <- bind_rows(res)
67
68   saveRDS(res, file = paste0(CLEANED_DATA_FOLDER, "res.RDS"))
69 ▾ }
70
71 ▾ if(FALSE){
72   clean_and_save_raw_data()
73 ▾ }
74
```

Unit Testing?

Two types of tests:

- Unit Tests: Test a single function.
- Integration Test: Test the full end-to-end systems.

- R: `testthat`
- Python: `unittest`

Good unit tests...

- Test base use cases
- Test edge cases
- Test settings (treatment of NAs, missing)