# Data Analysis Pipelines

Data Analysis

# Pipelines are **automated** and **reproducible**

Pipelines are helpful when there is a **single question you will have to answer more than once**

# Appropriate Checks

- Input file in expected format?

- Necessary variables included in data?

- Observations coded as expected?

- etc...

# **Avoid hard-coding**
# whenever possible

let your code fill the values in for you!

# Good pipelines are **scalable** pipelines

Your pipeline should work in the future...when the dataset is likely larger
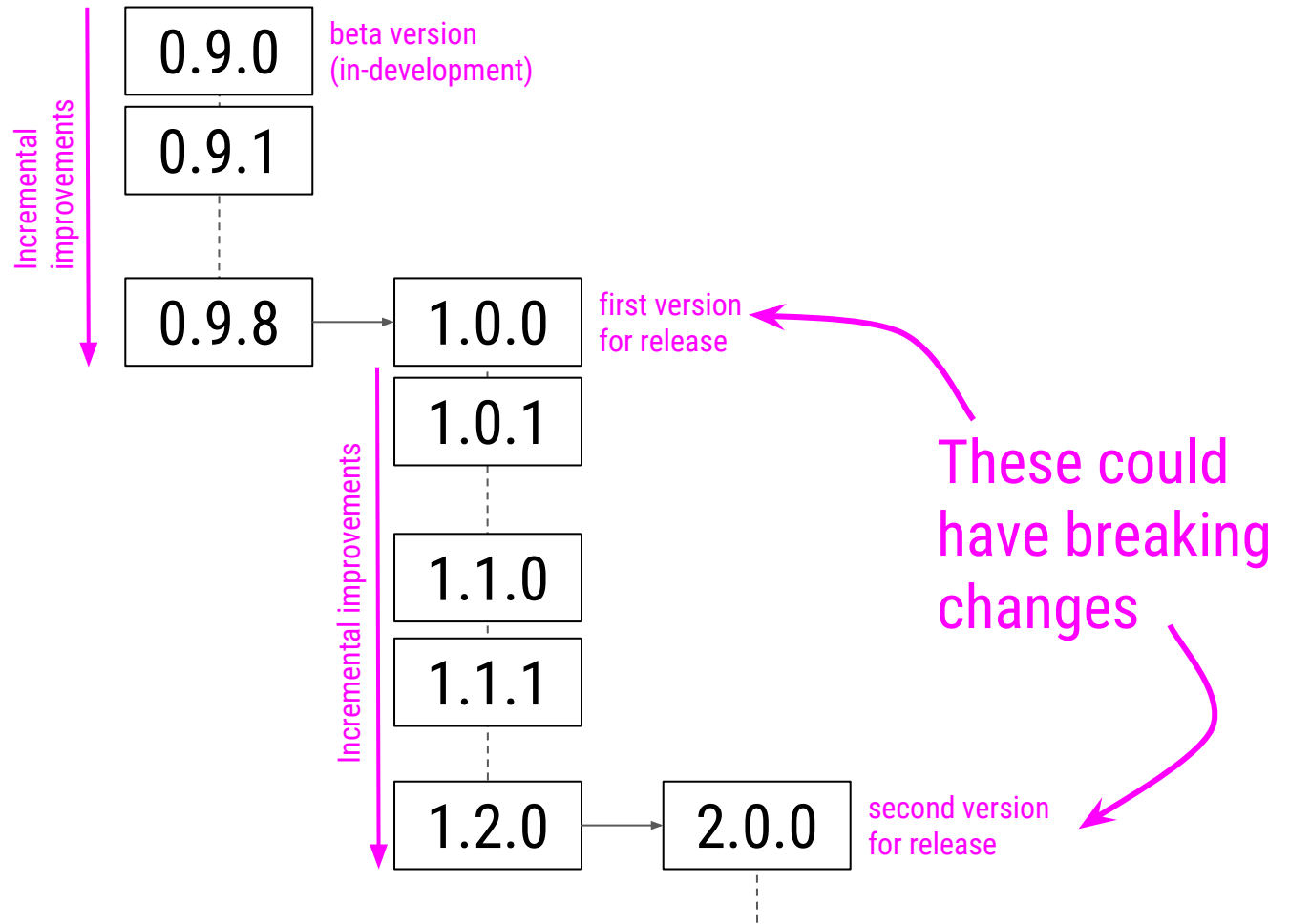
# Pipelines should be versioned

large release with
breaking changes

minimal,
non-breaking
features

# major.minor.patch

Includes new
non-breaking features
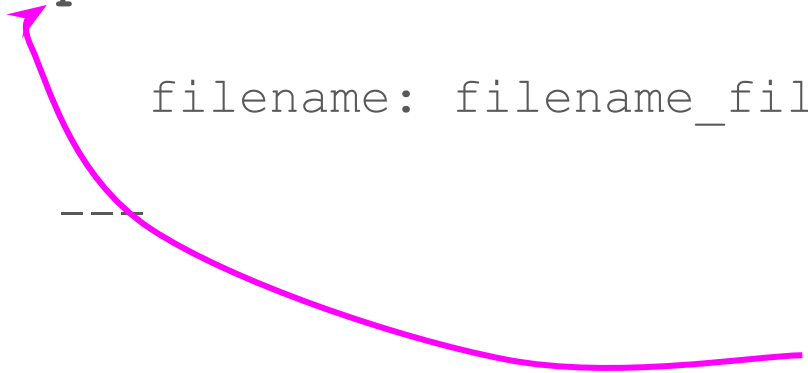
# The YAML will specify `params`

```
---

title: My Document

output: html_document

params:

    filename: filename_filedate.csv

---
```

Your YAML will include the parameters you'll use throughout your report

# `render()` will knit your parameterized report

```
rmarkdown::render("MyDocument.Rmd",
params = list(

   filename = "filename_filedate.csv")

)
```

To generate your document, you can use `render()` and specify `params` as a list.

In this example, **Sheet1** is the first week's data,
**Sheet2** includes the first two weeks' data

survey

| Sheet1 | Sheet2 |

| name | hrs_working | hrs_sleeping | hrs_fun | hrs_eating | hrs_socializing | hrs_other |
|---|---|---|---|---|---|---|
| Damon | 9 | 7 | 1 | 1 | 2 | 4 |
| Lilly | 7 | 8 | 2 | 1 | 1 | 5 |
| Will | 8 | 8 | 2 | 2 | 2 | 4 |
| Aisha | 8 | 6 | 2 | 1 | 4 | 3 |
| Hassan | 6 | 9 | 3 | 2 | 2 | 2 |
| Me | 10 | 8 | 2 | 2 | 1 | 1 |

**New R Markdown**

| | |
|---|---|
| 📄 Document | **Title:** Friend Survey |
| 📊 Presentation | **Author:** Jane Doe |
| Ⓡ Shiny | |
| ▦ From Template | |

**Default Output Format:**

🔘 **HTML**

Recommended format for authoring (you can switch to PDF or Word output anytime).

⚪ **PDF**

PDF output requires TeX (MiKTeX on Windows, MacTeX 2013+ on OS X, TeX Live 2013+ on Linux).

⚪ **Word**

Previewing Word documents requires an installation of MS Word (or Libre/Open Office on Linux).

[ OK ]  [ Cancel ]

We'll start with a typical R Markdown document.

```
friend_survey.Rmd

1   ---
2   title: "Friend Survey"
3   author: "Jane Doe"
4   output: html_document
5   date: '`r format(Sys.Date(), "%B %d, %Y")`'
6   params:
7     file_url: "ask"
8     worksheet: 1
9   ---
10
```

`params` are specified within your YAML

install and load
necessary R packages

```r setup, include=FALSE}
## install packages (if needed)
list.of.packages <- c("ggplot2", "googlesheets", "dplyr", "reshape2")
new.packages <- list.of.packages[!(list.of.packages %in% installed.packages()[,"Package"])]
if(length(new.packages)){
  install.packages(new.packages)}

## load packages
library(googlesheets)
library(ggplot2)
library(dplyr)
library(reshape2)
```

# Read in the Google Sheet

```r data, include=FALSE
## read Google Sheet in
survey <- gs_url(params$file_url)

df <- survey %>%
  gs_read(ws = params$worksheet)
```

# Write checks for your data and analysis

```{r checks, echo=FALSE}
columns <- c("name","hrs_working",  "hrs_sleeping", "hrs_fun", "hrs_eating",
"hrs_socializing", "hrs_other")
if(!sum(colnames(df) %in% columns)==length(columns)){
  stop("The input columns are unexpected - check to make sure the Google Sheet you
specified is the correct URL.")
}
```

## messages, warnings, and stop

**`message()`** - prints a message, code continues to run

**`warning()`** - prints a warning, code continues to run

**`stop()`** - stops code from running, prints error message

# Write checks for your data and analysis

```{r checks, echo=FALSE}
columns <- c("name","hrs_working",  "hrs_sleeping", "hrs_fun", "hrs_eating",
"hrs_socializing", "hrs_other")
if(!sum(colnames(df) %in% columns)==length(columns)){
  stop("The input columns are unexpected - check to make sure the Google Sheet you
specified is the correct URL.")
}
```

# Clean your data!

```r
```{r clean, include=FALSE}
## check for data entry errors
## remove samples where total hours != 24h
df_filtered <- df %>%
  select(2:ncol(df)) %>%
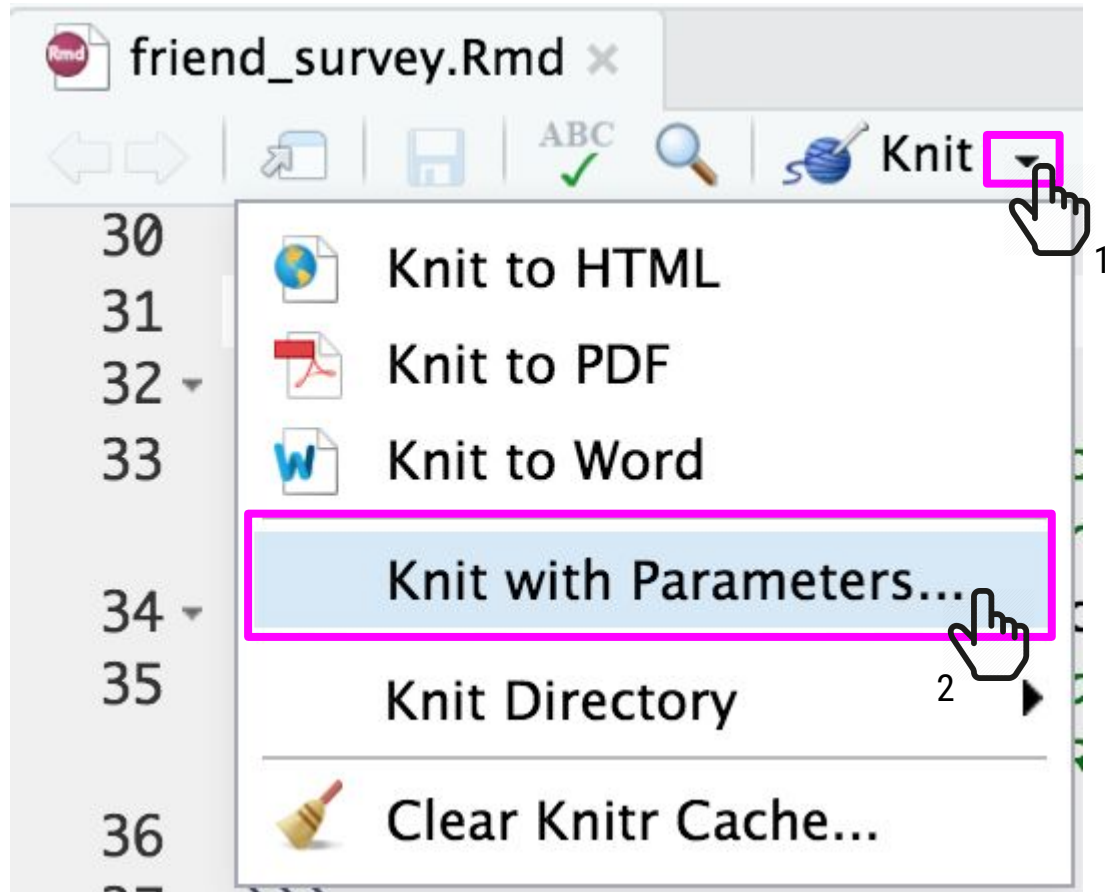  filter(rowSums(.)==24)
```
```

In this analysis, data for `r nrow(df)` individuals were read in; however, only `r nrow(df_filtered)` were included for analysis. Indivdiuals whose total number of hours was not equal to 24 were removed from analysis (N = `r nrow(df)-nrow(df_filtered)`).

Avoid hard-coding!

# Generate plot of interest

```r
```{r analyze, echo=FALSE, message=FALSE }
## generate plot
df_filtered %>%
  melt() %>%
  ggplot(aes(x=variable, y=value)) +
  geom_boxplot()
```
```

## Knit with Parameters

**file_url**

> ask

**worksheet**

> 1

Cancel    **Knit**

**Knit with Parameters** ✕

**file_url**

oogle.com/spreadsheets/d/1MpGE4YHB14qBgrg3lqa1eq_Mb7L8TMOwTZtVzqQ0mmA
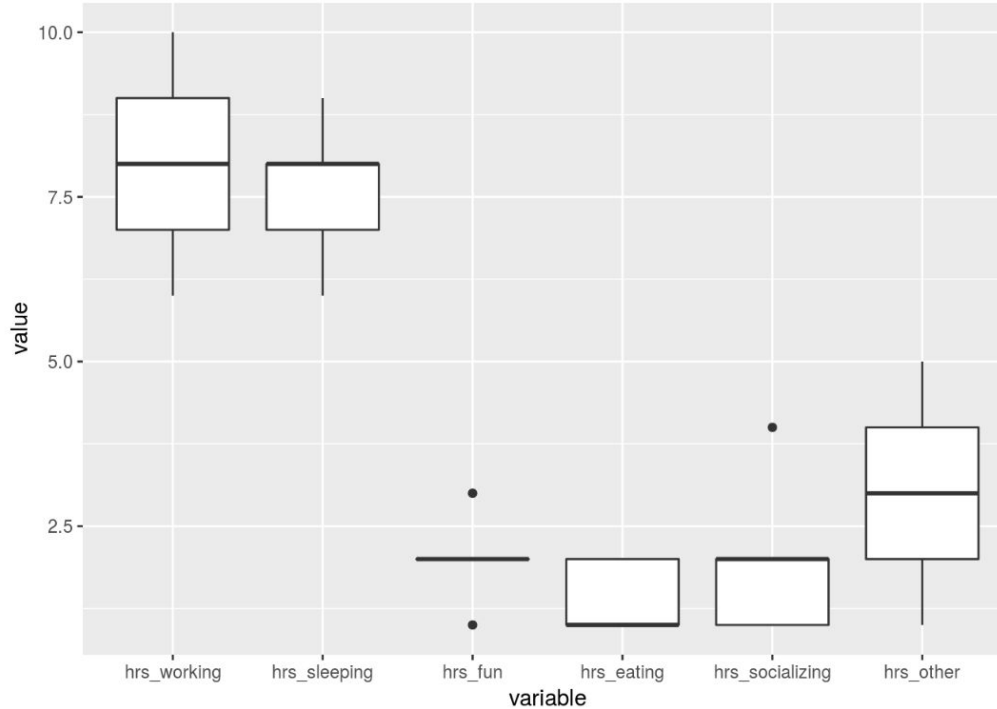
**worksheet**

1

Cancel    **Knit**

# Friend Survey

*Jane Doe*

*August 06, 2018*

In this analysis, data for 6 individuals were read in; however, only 5 were included for analysis. Individuals whose total number of hours was not equal to 24 were removed from analysis (N = 1).

**Knit with Parameters**

**file_url**

https://docs.google.com/spreadsheets/d/1MpGE4YHB14qBgrg3lqa1eq_Mb7L8TMOw⁻
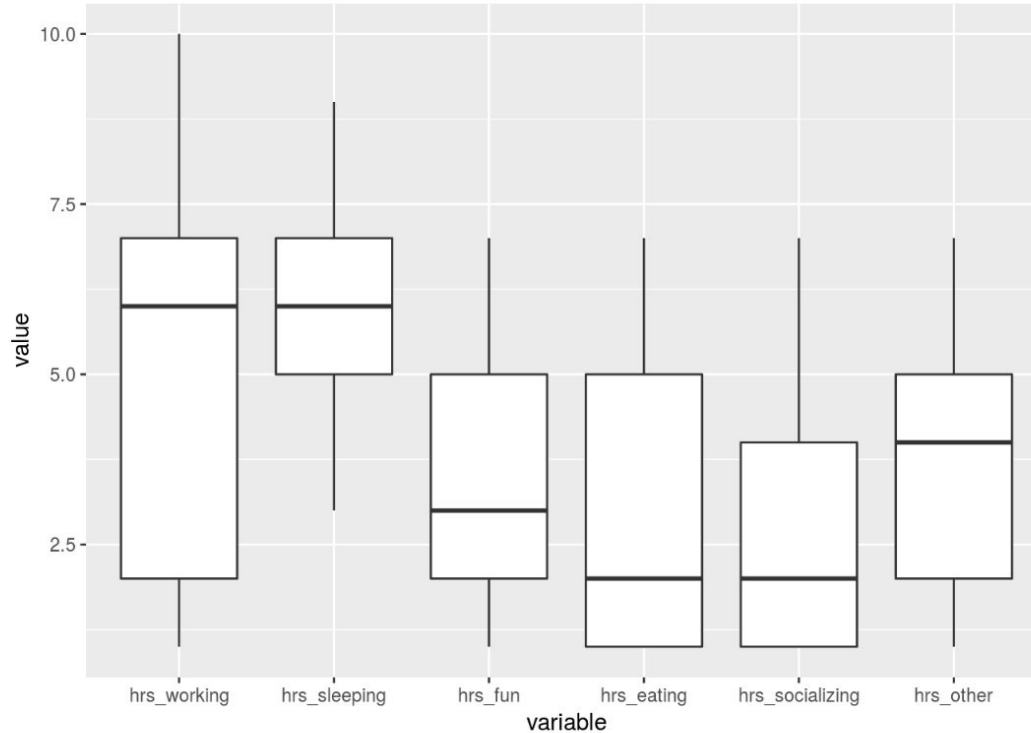
**worksheet**

2

Cancel    Knit

# Friend Survey

*Jane Doe*

*August 06, 2018*

In this analysis, data for 28 individuals were read in; however, only 25 were included for analysis. Indivdiuals whose total number of hours was not equal to 24 were removed from analysis (N = 3).



Values are updated using the code in the R Markdown document!

The same pipeline with updated data automatically updates the report!