

Structuring your analysis with R Markdown



Geoffrey Arnold

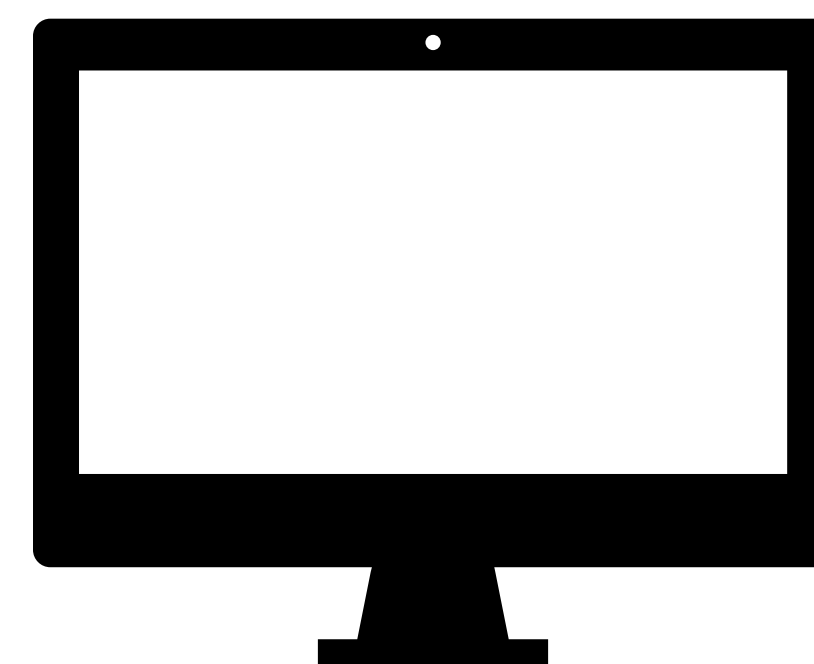
@geoffreylarnold 
geoffreylarnold 
geoffrey.arnold@pittsburghpa.gov 

Embedding R code



Code chunks

- Insert new chunks with
 - the Add Chunk button in the editor toolbar
 - typing the chunk delimiters ```{r}` and ````
 - the keyboard shortcut **Ctrl + Alt/Option + I**
- When you render your .Rmd file, R Markdown will run each code chunk and embed the results beneath the code chunk in your final report.



DEMO

02a-analysis.Rmd x

burritos x

← →

📄

📁

ABC ✓

🔍

🌐 Knit

⚙️

+

Insert

↑

↓

🏠 Run

🔄

☰

7

8

9

10

11

12

13

Load packages

{r load-packages}

library(tidyverse)

Code chunk

the data

data com

ps://www.kaggle.com/srcole/burritos-in-san-diego),

17

18

19

20

21

22

23

24

25

26

27

28

29

{r load-data}

burritos <- read_csv("../data/burritos_01022018.csv")

Mexican cuisine is often the best food option in southern California. And the burrito is the hallmark of the food: tasty, cheap, and filling. Appropriately, an effort was made to collect burritos across the county. At this time, the data set contains information on the following variables:

r `r nrow(burritos)` burritos fromd `r burritos %>%` restaurants.

ons of the San Diego burrito. * Volume * Tortilla

at quality * Non-meat filling quality * Meat-to-filling

Engine

Chunk label

Navigation

02A - Analysis

Load packages

Chunk 1: load-packages

The data

Chunk 2: load-data

31:1

The data

R Markdown



02a-analysis.Rmd x

burritos x

← →

📄

📁

ABC ✓

🔍

Knit

⚙️

+

Insert

↑

↓

🏠

Run

🔄

☰

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

31:1

Load packages

```{r load-packages}

library(tidyverse)

```

The data

The data come from [Kaggle.com](https://www.kaggle.com/srcole/burritos-in-san-diego),

```{r load-data}

burritos <- read\_csv("../data/burritos\_01022018.csv")

```

Mexican cuisine is often the best food option is southern California. And the burrito is the hallmark of delicious taco shop food: tasty, cheap, and filling. Appropriately, an effort was launched to critique burritos across the county o the lay burrito consumer. At this time, the data set r `r nrow(burritos)` burritos fromd `r burritos %>%` restaurants.

ons of the San Diego burrito. * Volume * Tortilla

at quality * Non-meat filling quality * Meat-to-filling

The data

Run all chunks up to this point

Run this chunk

02A - Analysis

Load packages

Chunk 1: load-packages

The data

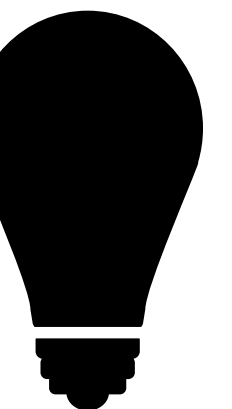
Chunk 2: load-data

R Markdown



Tips

- Avoid spaces in chunk labels, even though technically they are “allowed”, especially if you work with GitHub — *more on this later!*
- If you’re having a hard time coming up with a short label that describes what the chunk is doing, consider breaking it down into shorter chunks ➡ especially useful for troubleshooting!



What happens when there's an error in the R code in a chunk? What feedback/error does R give?

What about when there are multiple chunks with errors?



Inline code

Code results can be inserted directly into the text of a R Markdown file by enclosing the code with ``r``:

```
26 Mexican cuisine is often the best food option is southern California. And the
27 burrito is the hallmark of delicious taco shop food: tasty, cheap, and filling.
28 Appropriately, an effort was launched to critique burritos across the county
29 and make this data open to the lay burrito consumer. At this time, the data set
30 contains ratings from over `r nrow(burritos)` burritos fromd `r burritos %>%
count(Location) %>% nrow()` restaurants.
```

Mexican cuisine is often the best food option is southern California. And the burrito is the hallmark of delicious taco shop food: tasty, cheap, and filling. Appropriately, an effort was launched to critique burritos across the county and make this data open to the lay burrito consumer. At this time, the data set contains ratings from over 385 burritos fromd 102 restaurants.

Inline code styling

- R Markdown will always
 - display the results of inline code, but not the code
 - apply relevant text formatting to the results
- As a result, inline output is indistinguishable from the surrounding text
- Inline expressions do not take knitr options

bit.ly/2l1Jb66

bit.ly/2UwGGie

Your turn

- Open `sd-burritos.Rmd`, knit the document, and view the rendered file to get a sense of the data and the analysis.
- Add an inline R chunk to your document. This should be something that results in a numerical value or a character string.
- Add another inline R chunk and put some code in here that would result in a plot. What does a plot defined in an inline code chunk look like? Compare your answer with your neighbors'.
- **Stretch goal:** Update the date field in the YAML so that the date at the time of knitting the document is printed.



5_m 00_s



Chunk options

Safe to play with:

- `echo = FALSE`: code runs, code doesn't appear in rendered file, results do ➡ for when your audience doesn't need to see the code
- `include = FALSE`: code runs but neither code nor results appear in rendered file ➡ results can be used by other chunks
- `eval = FALSE`: code doesn't run but they appear in the rendered file ➡ when you want to show/teach code

Requires care:

- `message = FALSE`: hides messages ➡ especially useful/harmful for loading packages and data
- `warning = FALSE`: hides messages ➡ especially useful/harmful for functions that throw many warnings

Your turn

- Add labels to each chunk.
- Change the chunk options (`echo`, `eval`, `include`, `message`, `warning`) to explore what changes in the output. Then, decide on an appropriate option for each of the chunks. Compare your choices to your neighbors'.
- **Stretch goal:** What does the option `collapse` (set to `TRUE` or `FALSE`) do? What is the default setting for this option? Which code chunk(s) does changing this option affect?



10_m 00_s



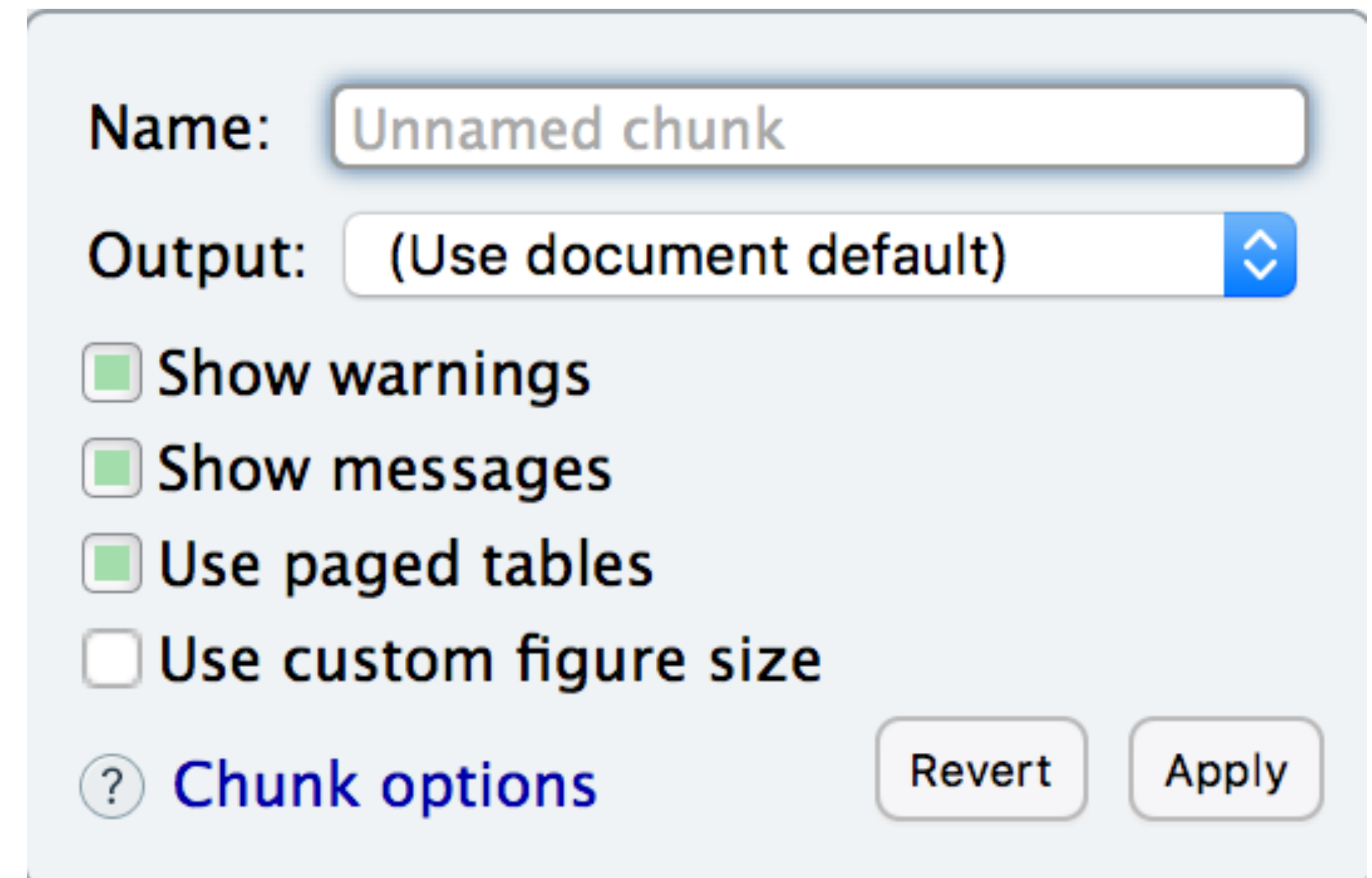
Chunk options for figures

- fig.height: Height in inches
- fig.width: Width in inches
- fig.align: "left", "right", or "center"
- fig.cap: Figure caption as character string

```
48
49 ```{r fig.height=5, fig.width=7, fig.align="right", fig.cap="Frequency distribution
  of reviewers"}
50 burritos_rev_count <- burritos %>%
51   mutate(Reviewer = fct_lump(Reviewer, n = 5)) %>%
52   count(Reviewer) %>%
53   mutate(Reviewer = fct_reorder(Reviewer, n, .desc = TRUE)) %>%
54   arrange(desc(n))
55 burritos_rev_count
56
57 ggplot(data = burritos_rev_count, mapping = aes(x = Reviewer, y = n)) +
58   geom_bar(stat = "identity") +
59   labs(title = "Distribution of reviewers", x = "", y = "")
60 ```
```

Setting chunk options via GUI

Some of the chunk options can be set via a handy GUI that you can access by clicking on the gear icon on a given chunk.



Name:

Output:

☒ Show warnings

☒ Show messages

☒ Use paged tables

☐ Use custom figure size

[? Chunk options](#)

```
5, fig.width=7, fig.align="right", fig.cap="Frequency distribution
t <- burritos %>%
  = fct_lump(Reviewer, n =
  %>%
  = fct_reorder(Reviewer, n, .desc = TRUE)) %>%
)
t

ritos_rev_count, mapping = aes(x = Reviewer, y = n)) +
  = "identity") +
  "Distribution of reviewers", x = "", y = "")
```

Options

So many more chunk options!

<https://www.rstudio.com/resources/cheatsheets/>



R Markdown Reference Guide

Learn more about R Markdown at rmarkdown.rstudio.com

Learn more about Interactive Docs at shiny.rstudio.com/articles

Contents:

1. **Markdown Syntax**
2. Knitr chunk options
3. Pandoc options

Syntax

Plain text

End a line with two spaces
to start a new paragraph.

italics and *_italics_*

****bold**** and **__bold__**

superscript^{^2^}

~~~~strikethrough~~~~

[link] ([www.rstudio.com](https://www.rstudio.com))

### Becomes

Plain text

End a line with two spaces to start a new paragraph.

*italics* and *italics*

**bold** and **bold**

superscript<sup>2</sup>

~~strikethrough~~

[link](https://www.rstudio.com)





# Global options

- To set global options that apply to every chunk in your file, call `knitr::opts_chunk$set` in a code chunk
- Knitr will treat each option that you pass to `knitr::opts_chunk$set` as a global default that can be overwritten in individual chunk headers

Hide all code chunks

```
7  
8 ```{r}  
9 knitr::opts_chunk$set(echo=FALSE)  
10 ```  
11
```

# Your turn

- Add a new code chunk to `sd-burritos.Rmd` and set relevant options for that particular chunk. You create a plot, calculate summary statistics, or if you prefer, just do some basic calculation (without using the data).
- Remove the figure height and width options from individual chunks and set them as global options.



5<sub>m</sub> 00<sub>s</sub>



# Caching



There are only two hard things in Computer Science:  
cache invalidation and naming things.

*–Phil Karlton*



Martin Fowler, Two Hard Things. <https://martinfowler.com/bliki/TwoHardThings.html>.



# Caching

- If document rendering becomes time consuming due to long computations you can use caching to improve performance
- If `cache = TRUE` is set:
  - Cached chunks are skipped, but objects created in these chunks are (lazy-) loaded from previously saved databases (.rdb and .rdx) files
  - These files are saved when a chunk is evaluated for the first time, or when cached files are not found
  - Results of the code will still be included in the output even when cache is used, because knitr also caches the printed output of a code chunk as a character string

# Chunk options for caching

- `cache.path`: Directory to save cached results in (default = “cache/”)
- `dependson`: Chunk dependencies for caching (default = NULL)
- ...

# Your turn

- **Before you start:** Make sure all of your code chunks are labeled!
- Add the following to the chunk that creates the bar plot: `Sys.sleep(60)`
- Turn on caching for the code chunk that creates the bar plot by setting the relevant chunk option. Knit the document (this is going to take a while, a bit more than 60 seconds). Take a look at the folder where `sd-burritos.Rmd` lives. What else is new there?
- Knit the document again without making any changes to this particular code chunk. How long did knitting the document take the second time around?
- Add a code chunk before this one and slice and overwrite the data so that `burritos` is now just the first 50 observations in the original burritos dataset. Knit the document again. What is the problem?
- **Stretch goal:** How do you fix this?

10<sub>m</sub> 00<sub>s</sub>



# Other languages





# Other languages

knitr can execute code in many languages besides R. Some of the available language engines include:

- Python
- SQL
- Bash
- Rcpp
- Stan
- JavaScript
- CSS

```

1 ---
2 title: "Simple Language Demos"
3 output: html_document
4 ---
5
6 You can write code in languages other than R with R Markdown, e.g.
7
8 ## Bash
9
10 ```{bash}
11 ls *.Rmd
12 ```
13
14 ## Python
15
16 ```{python}
17 x = 'hello, python world!'
18 print(x.split(' '))
19 ```
20
21

```

Engine

# Simple Language Demos

You can write code in languages other than R with R Markdown, e.g.

## Bash

```
ls *.Rmd
```

```
## 1-example.Rmd
## 2-chunks.Rmd
## 3-inline.Rmd
## 4-languages.Rmd
```

## Python

```
x = 'hello, python world!'
print(x.split(' '))
```

```
## ['hello,', 'python', 'world!']
```



# Output options

# Output options

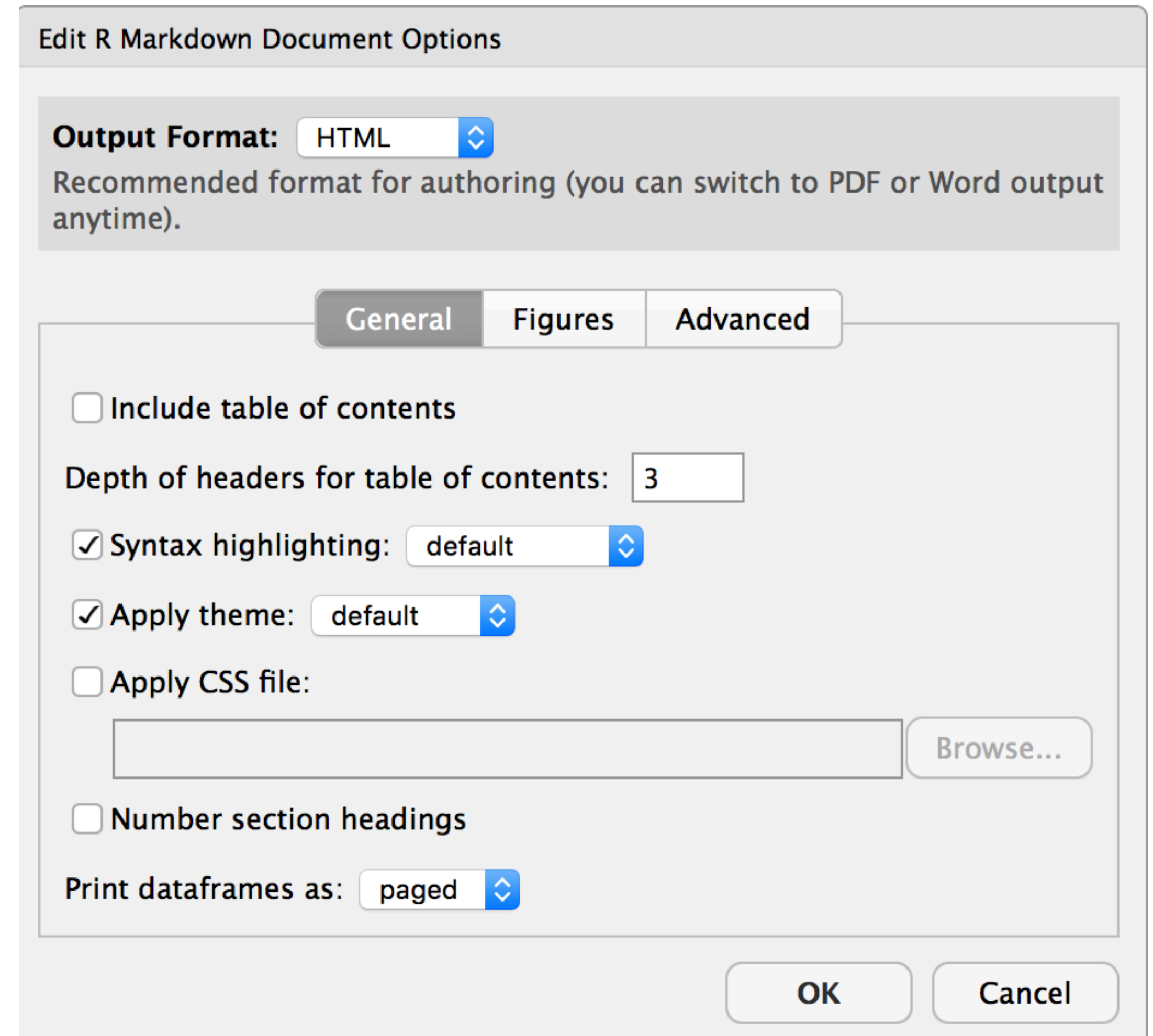
- Output options are defined in the YAML
- They are similar to setting global knitr options
- To learn which arguments a format takes, read the format's help page in R, e.g. `?html_document`

```
html_notebook(toc = FALSE, toc_depth = 3, toc_float = FALSE,  
  number_sections = FALSE, fig_width = 7, fig_height = 5,  
  fig_retina = 2, fig_caption = TRUE, code_folding = "show",  
  smart = TRUE, theme = "default", highlight = "textmate",  
  mathjax = "default", extra_dependencies = NULL, css = NULL,  
  includes = NULL, md_extensions = NULL, pandoc_args = NULL,  
  output_source = NULL, self_contained = TRUE, ...)
```




# Setting output options via GUI

Some of the output options for some of the output types can be set via a handy GUI that you can access by clicking on the gear icon on the toolbar



The screenshot shows the 'Edit R Markdown Document Options' dialog box with the 'General' tab selected. The 'Output Format' is set to 'HTML'. Below it, a note states: 'Recommended format for authoring (you can switch to PDF or Word output anytime)'. The 'General' tab contains several settings: 'Include table of contents' is unchecked; 'Depth of headers for table of contents' is set to 3; 'Syntax highlighting' is checked and set to 'default'; 'Apply theme' is checked and set to 'default'; 'Apply CSS file' is unchecked with an empty text field and a 'Browse...' button; 'Number section headings' is unchecked; and 'Print dataframes as' is set to 'paged'. At the bottom right are 'OK' and 'Cancel' buttons.


Edit R Markdown Document Options


**Output Format:** HTML   
Recommended format for authoring (you can switch to PDF or Word output anytime).

General Figures Advanced

☐ Include table of contents

Depth of headers for table of contents: 3


☒ Syntax highlighting: default 

☒ Apply theme: default 

☐ Apply CSS file:

Browse...

☐ Number section headings

Print dataframes as: paged 

OK Cancel

# Your turn

- Add a floating table of contents to the html document output.
- Also set all figures to be 4 x 7.



3<sub>m</sub> 00<sub>s</sub>





# Output formats



# Documents

- Most commonly used output is `html_document` — HTML document w/ Bootstrap CSS
- Other document options are as follows:
  - `html_notebook` - Interactive R Notebooks
  - `pdf_document` - PDF document (via LaTeX template)
  - `word_document` - Microsoft Word document (docx)
  - `odt_document` - OpenDocument Text document
  - `rtf_document` - Rich Text Format document
  - `md_document` - Markdown document (various flavors)

```
1 ---
2 title: "San Diego Burritos"
3 author: "Mine Çetinkaya-Rundel"
4 date: "2018-01-23"
5 output:
6   html_document:
7     highlight: pygments
8     theme: cosmo
9 ---
```

# Your turn

Convert `html_document` output to `pdf_document`. What changed?



# Presentations

- `ioslides_presentation` - HTML presentation with ioslides
- `revealjs::revealjs_presentation` - HTML presentation with reveal.js
- `slidy_presentation` - HTML presentation with W3C Slidy
- `beamer_presentation` - PDF presentation with LaTeX Beamer
- `xaringan::moon_reader` - remark.js slides

# Other output formats

- `flexdashboard::flex_dashboard` - Interactive dashboards
- `tufte::tufte_html` - HTML handouts in the style of Edward Tufte
- `html_vignette` - R package vignette (HTML)
- `github_document` - GitHub Flavored Markdown document





# Your turn

- **Before you start:** Delete any files and folders created during the caching exercise. Turn off caching, and remove the `Sys.sleep(60)` command.
- Change output to `github_document`. Knit the document.
- What changed, other than cosmetic changes in the output? Discuss with your neighbors.
- **Stretch goal (if you are a git/GitHub user):** Push this document, and any other necessary files, so that it can be previewed on Github (with figures and output).



3<sub>m</sub> 00<sub>s</sub>



# Your turn

- Decide which you want to work on: `xaringan` slides or `tufte_html`
- Make sure the package for the one you choose (`xaringan` or `tufte`) is installed and loaded
- Go to New File ➡ R Markdown... ➡ From Template, and then choose either Ninja Presentation (`xaringan`) or Tufte Handout (`tufte`)
- Convert `sd-burritos.Rmd` into one of these format

## San Diego Burritos

Geoffrey Arnold

2019-04-05

```
library(tidyverse)
```

```
## — Attaching packages —
```

```
## ✓ ggplot2 3.1.0   ✓ purrr  0.3.2
## ✓ tibble  2.1.1   ✓ dplyr  0.8.0.1
## ✓ tidyr   0.8.3   ✓ stringr 1.4.0
## ✓ readr   1.3.1   ✓ forcats 0.4.0
```

```
## — Conflicts —
```

```
## * dplyr::filter() masks stats::filter()
## * dplyr::lag()     masks stats::lag()
```

### The data

Kaggle: *SD Burritos*

The data come from [Kaggle.com](https://www.kaggle.com/geoffreyarnold/sd-burritos):

Mexican cuisine is often the best food option in southern California. And the burrito is the hallmark of delicious taco shop food: tasty, cheap, and filling. Appropriately, an effort was launched to critique burritos across the county and make this data open to the lay burrito consumer.

```
burritos <- read_csv("../data/burritos_01022018.csv")
```

## San Diego Burritos

Geoffrey Arnold

2019-04-05

10<sub>m</sub> 00<sub>s</sub>

