

Introduction to R - Day 3

Lidiya Mishieva

03 February, 2026

Intro

- Day 1: basic syntax, classes, objects, functions
- Day 2: base package, tidy programming & tidyverse
- Day 3:
 - Part 1: Summarizing data / tables
 - Part 2: Reporting with RMarkdown / Quarto
- Day 4: Git & RStudio
- *Still missing: further operations on vectors. graphs, referencing in rmd/qmd, citations and bibliography*

Recap last week

- tidy paradigm
- working with data using base and tidyverse packages
- import and export, subsetting and filtering, rlong/wide format, variable transformations

Part 1

Summarizing data

- `table()`

```
1 table(penguins$species, penguins$island)
```

	Biscoe	Dream	Torgersen
Adelie	44	56	52
Chinstrap	0	68	0
Gentoo	124	0	0

- `prop.table()`

```
1 prop.table(table(penguins$species, penguins$island), margin = NULL) # margin=1 for rows, margin=2 for columns
```

	Biscoe	Dream	Torgersen
Adelie	0.1279070	0.1627907	0.1511628
Chinstrap	0.0000000	0.1976744	0.0000000
Gentoo	0.3604651	0.0000000	0.0000000

- `addmargins()`

```
1 addmargins(prop.table(table(penguins$species, penguins$island), margin = NULL)) # margin=1 for rows, margin=2 for columns
```

	Biscoe	Dream	Torgersen	Sum
Adelie	0.1279070	0.1627907	0.1511628	0.4418605
Chinstrap	0.0000000	0.1976744	0.0000000	0.1976744
Gentoo	0.3604651	0.0000000	0.0000000	0.3604651
Sum	0.4883721	0.3604651	0.1511628	1.0000000

Summarizing data

```
1 table(penguins$species, penguins$island, penguins$sex)
```

, , = female

	Biscoe	Dream	Torgersen
Adelie	22	27	24
Chinstrap	0	34	0
Gentoo	58	0	0

, , = male

	Biscoe	Dream	Torgersen
Adelie	22	28	23
Chinstrap	0	34	0
Gentoo	61	0	0

Summarizing data

- `summary()`

```
1 summary(penguins$bill_len, useNA="always")
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
32.10	39.23	44.45	43.92	48.50	59.60	2

- can be used as a general method to get a summary of different model outputs
- can be used to get a summary of an entire dataset

```
1 summary(penguins, useNA="always")
```

species	island	bill_len	bill_dep
Adelie :152	Biscoe :168	Min. :32.10	Min. :13.10
Chinstrap: 68	Dream :124	1st Qu.:39.23	1st Qu.:15.60
Gentoo :124	Torgersen: 52	Median :44.45	Median :17.30
		Mean :43.92	Mean :17.15
		3rd Qu.:48.50	3rd Qu.:18.70
		Max. :59.60	Max. :21.50
		NA's :2	NA's :2

flipper_len	body_mass	sex	year
Min. :172.0	Min. :2700	female:165	Min. :2007
1st Qu.:190.0	1st Qu.:3550	male :168	1st Qu.:2007
Median :197.0	Median :4050	NA's : 11	Median :2008
Mean :200.9	Mean :4202		Mean :2008
3rd Qu.:213.0	3rd Qu.:4750		3rd Qu.:2009
Max. :231.0	Max. :6300		Max. :2009
NA's :2	NA's :2		

Summarizing data

- dplyr package: `group_by()` and `summarize()`

```
1 penguins %>%
2   group_by(island) %>%
3   summarise(mean_bill_len = mean(bill_len, na.rm=TRUE))
```

```
# A tibble: 3 × 2
  island    mean_bill_len
  <fct>         <dbl>
1 Biscoe         45.3
2 Dream          44.2
3 Torgersen      39.0
```

```
1 penguins %>%
2   select(-c(year)) %>%
3   drop_na() %>%
4   group_by(island, sex) %>%
5   summarise(across(where(is.numeric), mean))
```

```
# A tibble: 6 × 6
# Groups:   island [3]
  island    sex    bill_len bill_dep flipper_len body_mass
  <fct>    <fct>    <dbl>    <dbl>    <dbl>    <dbl>
1 Biscoe  female     43.3     15.2     206.    4319.
2 Biscoe  male      47.1     16.6     213.    5105.
3 Dream   female     42.3     17.6     190.    3446.
4 Dream   male      46.1     19.1     196.    3987.
5 Torgersen female     37.6     17.6     188.    3396.
6 Torgersen male      40.6     19.4     195.    4035.
```


Summary tables for reports

- many different (not very practical) ways for producing summary tables
 - check out this guide:

<https://cran.r-project.org/web/packages/DescTools/vignettes/TablesInR.pdf>

- dataframe is a table and can be used always as such

```
1 kableExtra::kable(penguins[1:4, ])
```

species	island	bill_len	bill_dep	flipper_len	body_mass	sex	year
Adelie	Torgersen	39.1	18.7	181	3750	male	2007
Adelie	Torgersen	39.5	17.4	186	3800	female	2007
Adelie	Torgersen	40.3	18.0	195	3250	female	2007
Adelie	Torgersen	NA	NA	NA	NA	NA	2007

Summary tables for reports

- programming and formatting usually with different packages
- `gtsummary` package for programming (unless it is a very specific table)
- package for formatting depends on the output file

Print Engine	Function	HTML	Word	PDF	RTF
<code>gt</code>	<code>as_gt()</code>	😊	😊	😊	⚠️
<code>flextable</code>	<code>as_flex_table()</code>	😊	😊	😊	😊
<code>huxtable</code>	<code>as_hux_table()</code>	😊	😊	😊	😊
<code>kableExtra</code>	<code>as_kable_extra()</code>	😊	🙅	😊	🙅
<code>kable</code>	<code>as_kable()</code>	😐	😐	😐	😐
<code>tibble</code>	<code>as_tibble()</code>	😞	😞	😞	😞

😊	Output fully supported
😐	Missing indentation, footnotes, spanning headers
😞	No formatted output
🙅	Output not supported
⚠️	Under development, missing indentation

Source: <https://www.danielsjoberg.com/gtsummary/articles/rmarkdown.html>
Introduction to R - Day 3

Tables with gtsummary

- summary tables from dataframes/tibbles
- regression tables from model outputs
- merging and stacking of tables
- some features for customization

Tables with gtsummary - descriptives

Variable	Treatment Received			p-value ²
	Indometh, N = 295 (49%) ¹	Placebo, N = 307 (51%) ¹	Overall, N = 602 (100%) ¹	
Age	44 (33, 54)	46 (36, 55)	45 (35, 54)	0.15
Missing	0	0	0	
Risk				0.2
Median (Q1, Q3)	2.50 (2.00, 3.00)	2.50 (1.50, 3.00)	2.50 (1.50, 3.00)	
Min, Max	1.00, 5.50	1.00, 4.50	1.00, 5.50	
Missing	0	0	0	
type				0.6
0_no SOD	47 (16%)	60 (20%)	107 (18%)	
1_type 1	38 (13%)	43 (14%)	81 (13%)	
2_type 2	139 (47%)	135 (44%)	274 (46%)	
3_type 3	71 (24%)	69 (22%)	140 (23%)	
Missing	0	0	0	

¹ Median (Q1, Q3); n (%)

² Two Sample t-test; Pearson's Chi-squared test

Tables with gtsummary - descriptives

```
1 # example datasets for different types of medical data
2 library(medicaldata)
3
4 # creating tables
5 library(gtsummary)
6
7
8 indo_rct %>%
9   # filter the columns of the dataset
10  select(c("rx", "age", "risk", "type")) %>%
11  # select(c(rx, age, risk, type)) %>% # also works here!
12
13  # recode the labels of the treatment variable rx
14  # factor is a character vector, that allows labeling and ordering
15  mutate(rx = factor(
16    rx,
17    levels = c("1_indomethacin", "0_placebo"),
18    labels = c("Indometh", "Placebo"))
19  ) %>%
20
21  # specify a summary table
22  tbl_summary(
23
24    # stratification by treatment
25    by = rx,
```

Tables with gtsummary - survival data

Variable	Survival				p-value ¹
	6 Month	12 Month	18 Month	24 Month	
Chemotherapy Treatment					0.2
Drug A	99% (97%, 100%)	91% (85%, 97%)	70% (62%, 80%)	47% (38%, 58%)	
Drug B	99% (97%, 100%)	86% (80%, 93%)	60% (51%, 70%)	41% (33%, 52%)	
¹ Log-rank test					

Tables with gtsummary - survival data

```
1 library(survival)
2
3 # survival
4 tbl_survfit(
5   # compute survival curves
6   survfit(
7     Surv(ttdeath, death) ~ trt, trial
8     # default is kaplan-meier
9     # type=c("kaplan-meier", "fleming-harrington", "fh2")
10  ),
11
12  # specify timepoints for estimating survival probabilities
13  times = c(6, 12, 18, 24),
14
15  # change the header label
16  label_header = "**{time} Month**"
17  ) %>%
18
19  # render variable names in bold
20  bold_labels() %>%
21
22  # add log-rank test
23  add_p() %>%
24
25  # print the p-value in bold if below a threshold
```

Tables with gtsummary - regression tables

Tables with gtsummary - regression tables

Characteristic	OR	95% CI	p-value
Chemotherapy Treatment			
Drug A	—	—	
Drug B	1.13	0.60, 2.13	0.7
Age	1.02	1.00, 1.04	0.10
Grade			
I	—	—	
II	0.85	0.39, 1.85	0.7
III	1.01	0.47, 2.15	>0.9
Abbreviations: CI = Confidence Interval, OR = Odds Ratio			

Tables with gtsummary - regression tables

Characteristic	HR	95% CI	p-value
Chemotherapy Treatment			
Drug A	—	—	
Drug B	1.30	0.88, 1.92	0.2
Grade			
I	—	—	
II	1.21	0.73, 1.99	0.5
III	1.79	1.12, 2.86	0.014
Age	1.01	0.99, 1.02	0.3
Abbreviations: CI = Confidence Interval, HR = Hazard Ratio			

Tables with gtsummary - regression tables

Characteristic	Tumor Response			Time to Death		
	OR	95% CI	p-value	HR	95% CI	p-value
Chemotherapy Treatment						
Drug A	—	—		—	—	
Drug B	1.13	0.60, 2.13	0.7	1.30	0.88, 1.92	0.2
Age	1.02	1.00, 1.04	0.10	1.01	0.99, 1.02	0.3
Grade						
I	—	—		—	—	
II	0.85	0.39, 1.85	0.7	1.21	0.73, 1.99	0.5
III	1.01	0.47, 2.15	>0.9	1.79	1.12, 2.86	0.014
Abbreviations: CI = Confidence Interval, HR = Hazard Ratio, OR = Odds Ratio						

Tables with gtsummary - regression tables

```
1 # specify a logistic model
2 mod1 <- glm(response ~ trt + age + grade, trial, family = binomial)
3 # build a regression table
4 t1 <- tbl_regression(mod1, exponentiate = TRUE)
5
6 # specify a cox model
7 mod2 <- coxph(Surv(ttdeath, death) ~ trt + grade + age, trial)
8 # build a regression table
9 t2 <- tbl_regression(mod2, exponentiate = TRUE)
10
11 # merge tables
12 t3 <- tbl_merge(
13   tbls = list(t1, t2),
14   tab_spanner = c("**Tumor Response**", "**Time to Death**")
15 )
16
17 t1
18 t2
19 t3
```

Tables with gtsummary - regression tables

```
1 # table containing a set of univariable regressions
2 tbl_uvregression(
3   trial,
4   method = coxph,
5   y = Surv(ttdeath, death),
6   exponentiate = TRUE,
7   include = c("age", "grade", "response"),
8   pvalue_fun = label_style_pvalue(digits = 2)
9 )
```

Characteristic	N	HR	95% CI	p-value
Age	189	1.01	0.99, 1.02	0.33
Grade	200			
I		—	—	
II		1.28	0.80, 2.05	0.31
III		1.69	1.07, 2.66	0.024
Tumor Response	193	0.50	0.31, 0.78	0.003
Abbreviations: CI = Confidence Interval, HR = Hazard Ratio				

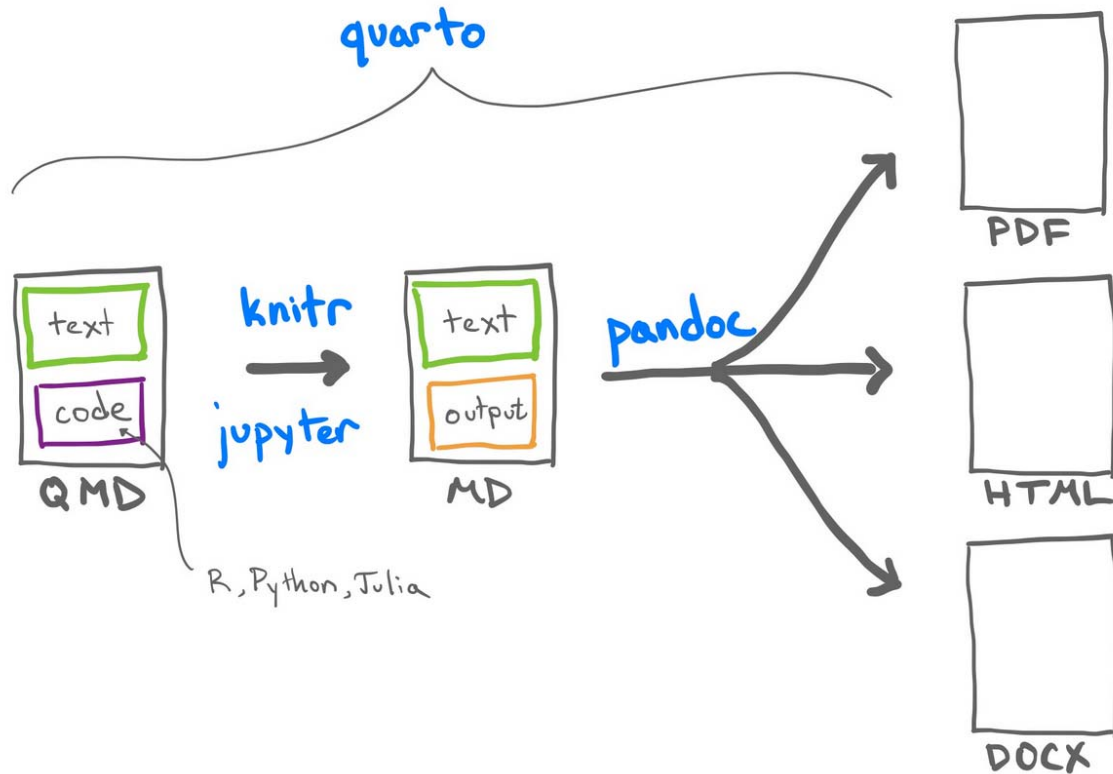
Part 2

Reporting with RStudio



Source: <https://allisonhorst.com>

A quarto document



Source: https://ubc-mds.github.io/DSCI_521_platforms-dsci_book/lectures/6-rmarkdown-quarto-slides-ghpages.html

Markup languages

- text-encoding systems that specify the structure & formatting of a document
- well-known examples: TeX, HTML, XML ...
- combine plain text & code used for formatting
- markdown (md) is a simplified markup language that uses minimal code

Markup languages - HTML example

```
<!DOCTYPE html>
<html>
<body>

<h1>My First Heading</h1>

<p>My first paragraph.</p>

</body>
</html>
```

My First Heading

My first paragraph.

Source: https://www.w3schools.com/html/tryit.asp?filename=tryhtml_basic_document

```
225 <li class="fragment"><p>Packages</p></li>
226 </ul>
227 </section>
228 <section id="tidy-paradigm---intuitive-readable-syntax" class="slide level2">
229 <h2>Tidy paradigm - intuitive [readable] syntax</h2>
230 <ul>
231 <li class="fragment"><p>Consider this story:</p>
232 <ol type="1">
233 <li class="fragment"><p><em>Little bunny Foo Foo</em></p></li>
234 <li class="fragment"><p><em>Went hopping through the forest</em></p></li>
235 <li class="fragment"><p><em>Scooping up the field mice</em></p></li>
236 <li class="fragment"><p><em>And bopping them on the head</em></p></li>
237 </ol></li>
238 <li class="fragment"><p>Translate into code:</p></li>
239 <li class="fragment"><div class="cell">
240 <div class="sourceCode cell-code" id="cb1"><pre class="sourceCode numberSource r number-lines code-wi
241 <span id="cb1-2"><a href=""></a>foo_foo <span class="ot">&lt;-</span> <span class="fu">little_bunny</
242 <span id="cb1-3"><a href=""></a><span class="co"># define the first operation</span></span>
243 <span id="cb1-4"><a href=""></a>foo_foo_1 <span class="ot">&lt;-</span> <span class="fu">hop</span></span>{f
244 <span id="cb1-5"><a href=""></a><span class="co"># define the second operation</span></span></span>
245 <span id="cb1-6"><a href=""></a>foo_foo_2 <span class="ot">&lt;-</span> <span class="fu">scoop</span></span>
246 <span id="cb1-7"><a href=""></a><span class="co"># define the thirs operation</span></span></span>
247 <span id="cb1-8"><a href=""></a>foo_foo_3 <span class="ot">&lt;-</span> <span class="fu">bop</span></span>{f
248 </div></li>
249 </ul>
```

HTML behind our Day 2 presentation

Markdown basics

Markdown basics - text

Markdown Syntax	Output
<code>*italics*, **bold**, ***bold italics***</code>	<i>italics</i> , bold , <i>bold italics</i>
<code>superscript^2^ / subscript~2~</code>	^{superscript²} / _{subscript₂}
<code>~~strikethrough~~</code>	strikethrough
<code>`verbatim code`</code>	<code>verbatim code</code>

Markdown basics - headings

Markdown syntax

Output

Heading 1

Heading 1

Heading 2

Heading 2

Heading 3

Heading 3

Heading 4

Heading 4

Markdown basics - footnotes

- Footnote reference

Here is a footnote reference^[1].

^[1]: Here is the footnote.

- Inline note

Here is an inline note.^[Inlines notes are easier to write, since you don't have to pick an identifier and move down to type the note.]

Markdown basics - inserting images

```
![*Source: https://allisonhorst.com](images/clipboard-4033369980.png){fig-align="left"}
```



Source: <https://allisonhorst.com>

Markdown basics

- There are many more formatting options such as formatting tables and diagrams, adding page breaks etc.
- Checkout the quarto website for details:
 - <https://quarto.org/docs/authoring/markdown-basics.html>

Maths using TeX

(1) inline maths

In simple linear regression, the model is $y = \beta_0 + \beta_1 x + \varepsilon$.

In simple linear regression, the model is $y = \beta_0 + \beta_1 x + \varepsilon$.

(2) standalone equation

In simple linear regression, the model is

```
$$  
\begin{equation}  
y = \beta_0 + \beta_1 x + \varepsilon.  
\end{equation}  
$$
```

In simple linear regression, the model is

$$y = \beta_0 + \beta_1 x + \varepsilon.$$

Cheatsheet: <https://tug.ctan.org/info/undergradmath/undergradmath.pdf>

Adding code to RMD / QMD

- Option 1: code chunk

```
```{r}
#| eval: false
#| include: false
#| output: false

round(
 mean(
 penguins$bill_len, na.rm = TRUE
), 2
)
```
```

- Option 2: inline code
 - Add output of the code within text

```
The mean in group XY was `r round(mean(penguins$bill_len, na.rm=TRUE), 2)`.
```

-
- The mean in group XY was 43.92.

YAML - Yet Another Markup Language

- In Quarto it is used for overall document configuration
- A YAML section is placed in the beginning of the document

```
1 ---
2 title: "Statistical Analyses - ADAPT HER2 IV Study"
3 subtitle: "Efficacy Analyses"
4 author: "Authors: Christine zu Eulenburg, Lidiya Mishieva, Ona Sauliene"
5 date: "`r Sys.Date()`"
6 output:
7   officedown::rdocx_document:
8     toc: yes
9     page_margins:
10       bottom: 1.27
11       top: 0
12       right: 2.54
13       left: 2.54
14       gutter: 0
15     reference_docx: template.docx
16 ---
```

- A code chunk can also contain a YAML section

```
{r chunk1, echo=TRUE = FALSE, message = FALSE, warning = FALSE}

summary(iris)

{r}
#| label: chunk1
#| echo: true
#| message: false
#| warning: false

summary(iris)
```

Code chunks - execution options

| Option | Description |
|----------------------------|------------------------------|
| <code>eval: true</code> | whether to evaluate the code |
| <code>echo: true</code> | whether to show the code |
| <code>include: true</code> | whether to include the code |
| <code>cache: true</code> | whether to cache the results |

Code chunks - figure control

| Option | Description |
|-----------------------------------|-----------------------------|
| <code>fig-cap: "My figure"</code> | figure caption |
| <code>fig-width: 8</code> | whether to show the code |
| <code>fig-height: 6</code> | width in inches |
| <code>fig-align: "center"</code> | height in inches |
| <code>fig-dpi: 300</code> | resolution |
| <code>out-width: "80%"</code> | output width |
| <code>layout-ncol: 2</code> | number of columns for plots |

Code chunks - output control

| Option | Description |
|-----------------------------|---------------------|
| <code>warning: false</code> | hide warnings |
| <code>message: false</code> | hide messages |
| <code>error: false</code> | hide error messages |
| <code>output: false</code> | hide all output |

Sourcing code

- **Option 1:** put all code into the code chunks or inline code
- **Option 2:** source your code from an external script
- You can do both (use results produced by the sourced code in the code chunks or inline code)

```
32 {r source-code, include=FALSE, warning = FALSE, error=TRUE}
33
34 source("scripts/HER2-IV_EFFICACY_09.12.2025.R",
35        local = knitr::knit_global(),
36        encoding = "UTF-8")
37
38
```

Rendering tables

| Print Engine | Function | HTML | Word | PDF | RTF |
|--------------|-------------------------------|------|------|-----|-----|
| gt | <code>as_gt()</code> | 😊 | 😊 | 😊 | ⚠️ |
| flextable | <code>as_flex_table()</code> | 😊 | 😊 | 😊 | 😊 |
| huxtable | <code>as_hux_table()</code> | 😊 | 😊 | 😊 | 😊 |
| kableExtra | <code>as_kable_extra()</code> | 😊 | 🙅 | 😊 | 🙅 |
| kable | <code>as_kable()</code> | 😐 | 😐 | 😐 | 😐 |
| tibble | <code>as_tibble()</code> | 😞 | 😞 | 😞 | 😞 |

| | |
|----|--|
| 😊 | Output fully supported |
| 😐 | Missing indentation, footnotes, spanning headers |
| 😞 | No formatted output |
| 🙅 | Output not supported |
| ⚠️ | Under development, missing indentation |

Source: <https://www.danielsjoberg.com/gtsummary/articles/rmarkdown.html>

Rendering tables - flextable

- A package for creating and formatting tables
- Own class: flextable
- You can transform a gtsummary object to a flextable object
 - Add additional formatting
 - check out this webpage: <https://ardata-fr.github.io/flextable-book/index.html>
 - Else the format of gtsummary table is taken over to a flextable object
 - You can export it to .docx

| Name | Type | Value |
|------------|-----------------------------------|---|
| ft1 | list [c(label = 2.02594174985532, | List of length 7 |
| header | list [8] (S3: complex_tabpart) | List of length 8 |
| body | list [8] (S3: complex_tabpart) | List of length 8 |
| footer | list [8] (S3: complex_tabpart) | List of length 8 |
| col_keys | character [4] | 'label' 'estimate' 'conf.low' 'p.value' |
| caption | list [1] | List of length 1 |
| blanks | character [0] | |
| properties | list [8] | List of length 8 |

Flextable to docx - fitting to page

```
1 # =====
2 # Fit a flextable to a LANDSCAPE page
3 # =====
4
5 FitFlextableToPageLandscape <- function(
6
7   # select a flextable object
8   ft,
9
10  # usable page width in inches:
11  # A4 landscape width (29.7 cm)
12  # minus 5 cm margins, converted to inches
13  pgwidth = (29.7 - 5) / 2.54) {
14
15  # apply consistent formatting
16  ft_out <- ft %>%
17    fontsize(size = 9, part = "body") %>%      # body text size
18    fontsize(size = 10, part = "header") %>%    # header text size
19    autofit() %>%                                # initial auto column widths
20    height_all(0.25, part = "body") %>%         # compact row height
21    hrule(rule = "exact", part = "body")        # fixed row height (vs auto)
22
23  # rescale column widths to exactly fill the available page width
24  ft_out <- width(
25    ft_out,
```

Gtsummary to flextable to docx - workflow

```
1 # specify a logistic model
2 mod1 <- glm(response ~ trt + age + grade, trial, family = binomial)
3
4 # build a regression table
5 t1 <- tbl_regression(mod1, exponentiate = TRUE)
6
7 # transform to a flextable object and fit to page
8 FitFlextableToPageLandscape(as_flex_table(t1))
```

Debugging with RMD/QMD

- Generally more difficult to debug code in RMarkdown/Quarto
 - Errors are indicated in the output windows of the code chunks
 - When rendering the document, RMD/QMD will just stop in case of error
 - To debug, you need to run each code chunk separately (considering the order!)
- In a script, debugging is more straight forward, an an error will be shown in the console directly after the corresponding code line
- RMD/QMD will not render, if the sourced Rscript contains errors

Reporting with RStudio - example file

The end

R learners,



Source: <https://allisonhorst.com>

Move to the next session

Further operations on vectors

- sum, length, which

Graphs

- base R plots
- base plots into pane with `par(mfrow = c(2,2))` and `dev.off()`
- ggplot concept
- arranging plots with ggplot
- plots from models

Referencing with RMD/QMD

Citations and bibliography with RMD/ QMD