PHUSE EU
Connect 2025

OS17: End to End R: Where to Start?

Luke Thompson
Statistical
Programmer
Astellas
11th October 2025

Agenda

- Background
 - Useful R resources
- Development of an R Environment

SDTM, ADaM & TLG **Programming**

Ш Utility macros Conclusion

The speaker is a paid employee of Astellas. This presentation is intended for informational purposes only and does not replace independent professional judgment. This presentation is not intended to provide medical or legal advice. Statements of fact, positions taken and opinions expressed are those of the speaker individually and, unless expressly stated to the contrary, do not necessarily reflect the opinion or position of the speaker's employer, Astellas, or any of its subsidiaries and/or related entities.

Background

Background

- Astellas in June 2024
- 5502-CL-0001 study selected for End to End (E2E) R:

A Phase 1 Combined Single and Multiple Ascending Oral Dose Study to Evaluate the Safety, Tolerability, Pharmacokinetics, Food Effect and Biomarker Response of ASP5502 in Healthy Adults and Participants With Primary Sjögren's Syndrome

- E2E R SDTM, ADaM, TLGs & eSub
- Why E2E R? Because there isn't a better time. Pharmaverse packages. Many companies have adopted it already. More and more programmers have R skills.

Renvironment

Development of the infrastructure

Development of an R Environment



R Language and Packages

R Environment

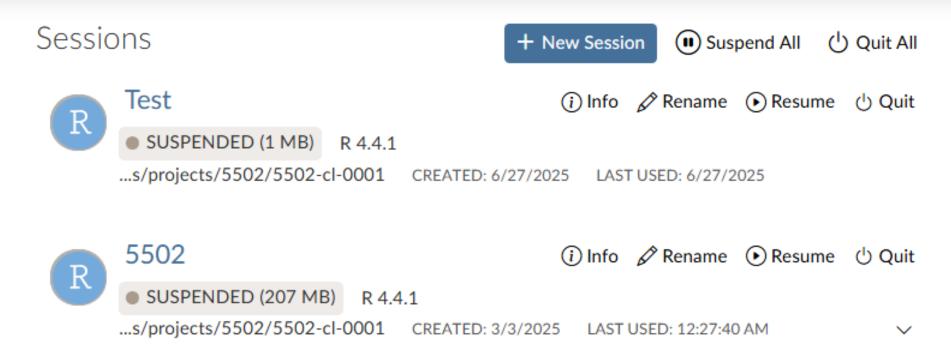
Posit Workbench®

- Vendor supplied tool
- Kubernetes system
- User friendly interface

RStudio® (IDE)

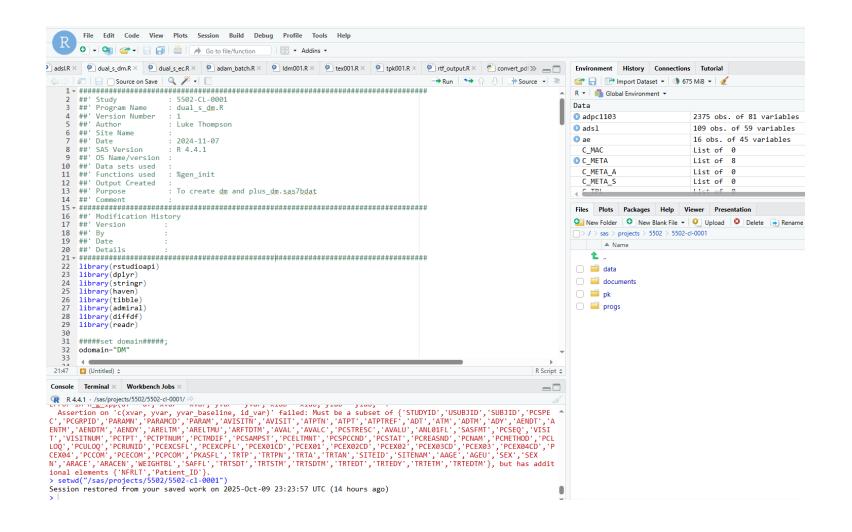
- Vendor supplied tool
- Internal Astellas website
- Connects to the Astellas Statistical Computing Environment (SCE)

posit Workbench



R Environment

Figure: a screenshot of 2 suspended sessions in Posit Workbench.



R Environment

Figure: a screenshot of the Rstudio environment

R language and Packages

Importance of Validation

- Required for regulatory compliance, data integrity, quality assurance, standardisation, & traceability
- SAS has built-in validation Due to reputation, pre validated by vendor.
- R requires validation from individual or companies, because it is open source. E.g. package quality may differ (documentation, validation levels) and it is ever-changing.
- Possible to do inhouse, but workload too demanding

R language and Packages

OpenVal® by Atorus Research

- Comprehensive suite of validated packages
- Release every 6 months and packages can be added, updated or removed
- Negotiation possible
- Contained many of the packages we required including some Pharmaverse packages

Utility Macros

Utility Macros

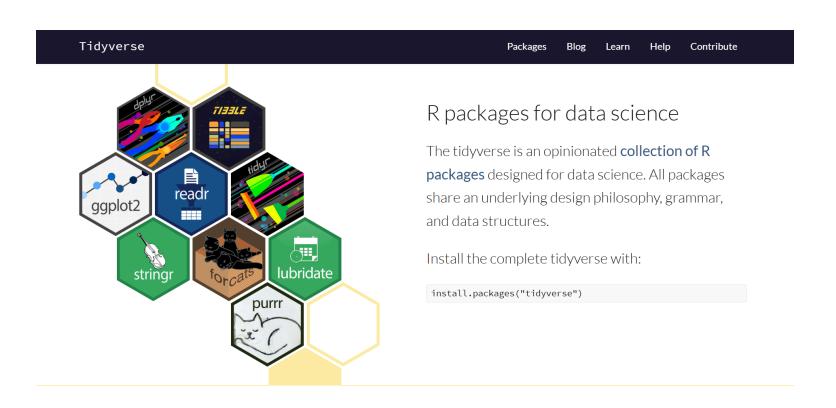
Setup

- Maintain the current workflow we had for SAS
- Utilise code and templates we already had and convert them into SAS.
- We replicated utility SAS Macros into R functions
- Functions to initialise program names, load libraries, read in deliverable files, create formats, ect.
- Tidyverse dplyr, stringr, haven & magnitr

Useful R Resources

Tidyverse websites

Useful tools to begin R programming



Tidyverse websites

Useful tools to begin R programming

Data transformation with dplyr:: cheatsheet

pipes Each variable is in Each observation, or x |> f(y) its own column case, is in its own row becomes f(x, y)

dplyr functions work with pipes and expect tidy data. In tidy data:

Summarize Cases

Apply summary functions to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).



Group Cases

Use group_by(.data, ..., .add = FALSE, .drop = TRUE) to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.



Use rowwise(.data, ...) to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidyr cheat sheet for list-column workflow.



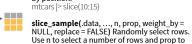
g_mtcars <- mtcars |> group_by(cyl) ungroup(g_mtcars)

Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.





select a fraction of rows. mtcars |> slice_sample(n = 5, replace = TRUE)

slice min(.data, order by, ..., n, prop. with ties = TRUE) and slice max() Select rows with the lowest and highest values. mtcars |> slice_min(mpg, prop = 0.25)

!is.na()

slice_head(.data, ..., n, prop) and slice_tail() Select the first or last rows. mtcars |> slice head(n = 5)

Logical and boolean operators to use with filter() is.na()

See ?base::Logic and ?Comparison for help.

ARRANGE CASES



ADD CASES



Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.



relocate(.data, ..., .before = NULL, .after = NULL) Move columns to new position.

mtcars |> relocate(mpg, cyl, .after = last_col())

Use these helpers with select() and across()

e.g. mtcars |> select(mpg:cyl)

contains(match) num_range(prefix, range) ;, e.g., mpg:cyl ends_with(match) all_of(x)/any_of(x, ..., vars) !, e.g., !gear starts_with(match) matches(match)

MANIPULATE MULTIPLE VARIABLES AT ONCE

df < -tibble(x 1 = c(1, 2), x 2 = c(3, 4), y = c(4, 5))



c across(.cols) Compute across columns in row-wise data.

> rowwise() |> mutate(x total = sum(c across(1:2)))

MAKE NEW VARIABLES

Apply vectorized functions to columns. Vectorized functions take vectors as input and return vectors of the same length as output (see back).

vectorized function



mtcars |> mutate(gpm = 1 / mpg) mtcars > mutate(gpm = 1 / mpg, .keep = "none")



rename(.data, ...) Rename columns. Use rename_with() to rename with a function. mtcars | rename(miles_per_gallon = mpg)



CC BY SA Posit Software, PBC • info@posit.co • posit.co • Learn more at dplyr.tidyverse.org • HTML cheatsheets at pos.it/cheatsheets • dplyr 1.1.4 • Updated: 2025-08

xor()

Online courses

Useful R Resources





Online Resources

Useful R Resources





15:2

How would I sort a dataset named data by USUBJID and VISITUM and set a new variable named TEST to 1 when first.usubjid in R as opposed to SAS?

Stella 15:23 Edited



In R, you can sort a dataset and create a new variable based on conditions using the dplyr package, which is part of the tidyverse. Here's how you can do it:

- 1. First, ensure you have the dplyr package installed and loaded.
- 2. Use the $\left[\text{arrange()} \right]$ function to sort the dataset by $\left[\text{usubjid} \right]$ and $\left[\text{visitum} \right]$.
- 3. Use the mutate() function to create a new variable test that is set to 1 for the first occurrence of each usubjid.

Here's an example code snippet:

```
# Load the dplyr package
library(dplyr)

# Assuming your dataset is named 'data'
data <- data %>%
    arrange(usubjid, visitum) %>%  # Sort by usubjid and visitum
    group_by(usubjid) %>%  # Group by usubjid
    mutate(test = ifelse(row_number() == 1, 1, 0)) %>%  # Set test to 1 for the first usubjid
    ungroup()  # Ungroup the data

# View the modified dataset
print(data)
```

In this code:

- arrange() sorts the dataset.
- group_by() groups the data by usubjid .
- mutate() creates the new variable test, setting it to 1 for the first occurrence of each usubjid and 0 otherwise.
- ungroup() is used to remove the grouping after the mutation.

This approach should give you the desired result similar to what you would achieve in SAS. Let me know if you have any further questions!

AI?

Useful R Resources

Pharmaverse

"A connected network of companies and individuals working to promote collaborative development of curated open-source R packages for clinical reporting usage in pharma, in a space where previously we would only ever have worked in silos on our own closed-source and often duplicative solutions."

- There are a wealth of packages available for E2E R
- Highly recommend checking their website: <u>End-to-End Clinical</u> <u>Reporting Packages</u>



Pharmaverse

SDTM – non-pharmaverse

Pharmaverse package called OAK, opted not to use it.

"An EDC (Electronic Data Capture systems) and Data Standard agnostic solution that enables the pharmaceutical programming community to develop CDISC (Clinical Data Interchange Standards Consortium) SDTM (Study Data Tabulation Model) datasets in R"

- Other packages used instead, to continue the replication of standard templates.
- Tidyverse mainly used

```
dm <- dm raw %>%
   dplyr::rename(COUNTRY_ = COUNTRY, SEX = SEX, AGEU = AGEU) %>%
STUDYID = project,
     DOMAIN = odomain,
     SUBJID = stringr::str trim(Subject),
     USUBJID = ifelse(is.na(SUBJIDP),
                      paste0(STUDYID, "-", SUBJID), paste0(STUDYID, "-",
  substr(SUBJID, 1, 5), SUBJIDP)),
     SITEID = stringr::str trim(as.character(`_SITEID`)),
     SEX = stringr::str trim(substr(SEX_, 1, 1)),
     COUNTRY = COUNTRY STD)
```

Figure: a snippet of programming of SDTM DM

Figure: a snippet of programming of SDTM DM

```
# Deriving dthdtc
ddthdtc1 <-
 # Combine with ae 1 where conditions are met
 dplyr::bind rows(
   S RAW[['ae 1']] %>%
     dplyr::mutate(across(where(is.character), ~ na_if(., "")))%>%
     dplyr::filter(toupper(AESDTH) == "YES") %>%
     rename ymd(prefin = "AEENDT", prefout = "DTHDT")
 ) %>%
 # Combine with ds 3 where conditions are met
 dplyr::bind rows(
   S RAW[['ds 3']] %>%
     dplyr::mutate(across(where(is.character), ~ na if(., "")))%>%
     dplyr::filter(toupper(DSDECOD) == "DEATH") %>%
     rename ymd(prefin = "DSSTDT", prefout = "DTHDT")
 ) %>%
 # Combine with ds 4 where conditions are met
 dplyr::bind rows(
   S RAW[['ds 4']] %>%
     dplyr::mutate(across(where(is.character), ~ na if(., "")))%>%
     dplyr::filter(toupper(DSDECOD) == "DEATH") %>%
```

ADaM – pharmaverse (Admiral)

 "(ADaM In R Asset Library) - Modular framework to generate ADaM via R functions relying on community contributions"

```
new vars prefix = "EXST"
 ) %>%
 derive vars dtm(
   dtc = EXENDTC,
   new vars prefix = "EXEN",
   time imputation = "last"
#> The default value of `ignore seconds flag` will change to "TRUE" in admiral
#> 1.4.0.
adsl <- adsl %>%
 derive vars merged(
   dataset add = ex ext,
   filter add = (EXDOSE > 0
     (EXDOSE == 0 &
       str detect(EXTRT, "PLACEBO"))) & !is.na(EXSTDTM),
   new_vars = exprs(TRTSDTM = EXSTDTM, TRTSTMF = EXSTTMF),
   order = exprs(EXSTDTM, EXSEQ),
   mode = "first",
   by vars = exprs(STUDYID, USUBJID)
 derive vars merged(
   dataset add = ex ext,
   filter add = (EXDOSE > 0
     (EXDOSE == 0 &
       str_detect(EXTRT, "PLACEBO"))) & !is.na(EXENDTM),
   new_vars = exprs(TRTEDTM = EXENDTM, TRTETMF = EXENTMF),
   order = exprs(EXENDTM, EXSEQ),
   mode = "last",
   by_vars = exprs(STUDYID, USUBJID)
```

TLG

- TLGs, especially graphs, are where R excels (ggplot2)
- A multitude of packages are available.
- Pharmaverse lists some of these packages that can be used. Image right is non-extensive
- There may be packages not listed such as specific statistical analysis or specific formatting

Listings



Plots



While ggplot2 is a lower level, non-pharma specific plotting package. It is universally accepted as the package for graphics, so included here and as a non-pharma package.



ggplot2 📒 👣

An implementation of the Grammar of Graphics in R, and the most popular plotting package for static plots in R.

Tables



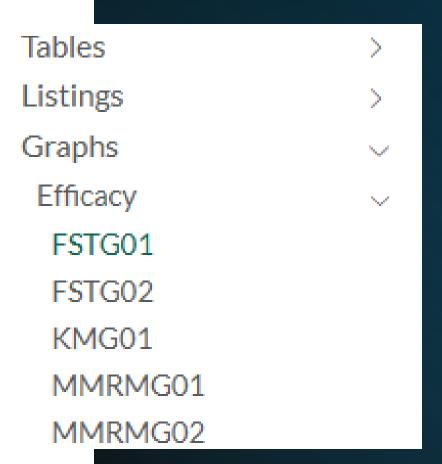


chevron 🗏 😱 CRAN 0.2.12 Stars 16 Contributors 2

Holds TLG template standards to create standard outputs for clinical trials reporting with limited parameterisation.

TLG

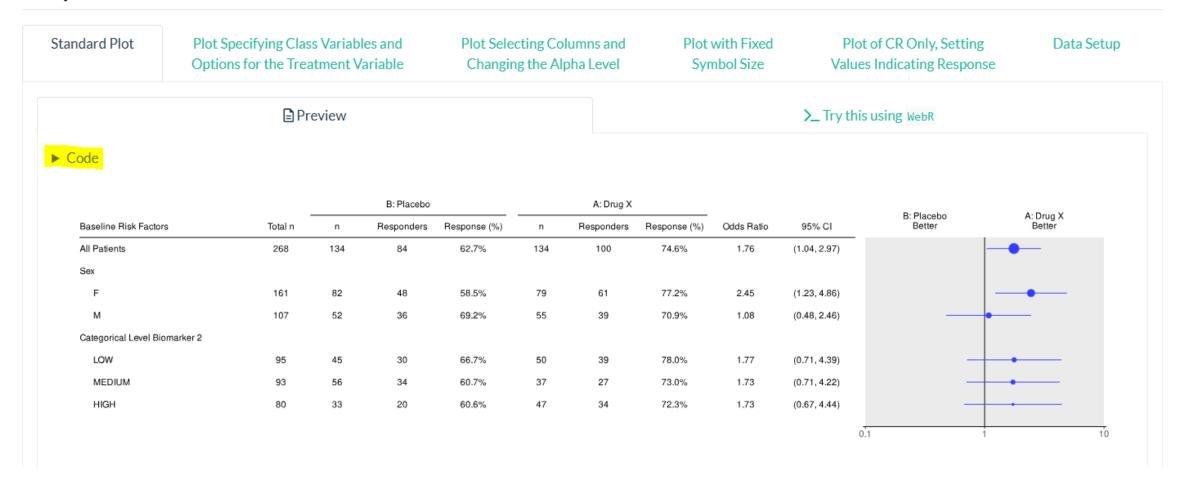
- <u>TLG Catalog</u> provides many template outputs.
- A combination of the TLG catalog, SAS templates, pharmaverse examples & R templates were used as a basis.
- All were modified to be as close to Astellas standards as possible.





Subgroup Analysis of Best Overall Response

Output



```
Preview
                                                                                                                                    >_ Try this using WebR
▼ Code
      df <- extract_rsp_subgroups(</pre>
   2
         variables = list(
           rsp = "is_rsp",
           arm = "ARM",
           subgroups = c("SEX", "BMRKR2"),
           strata = "STRATA2"
   7
   8
         data = anl_rsp_arms_ab,
         conf level = 0.95
   9
  10
 11
      result <- basic_table() %>%
  12
 13
         tabulate_rsp_subgroups(df, vars = c("n_tot", "n", "n_rsp", "prop", "or", "ci"))
  14
  15
      plot <- g_forest(tbl = result)</pre>
  16
      plot
                                                                                               A: Drug X
                                                           B: Placebo
                                                                                                                                                      B: Placebo
                                                                                                                                                                            A: Drug X
     Baseline Risk Factors
                                                                                             Responders
                                                                                                                       Odds Ratio
                                                                                                                                    95% CI
                                                                                                                                                        Better
                                                                                                                                                                             Better
                                      Total n
                                                          Responders
                                                                     Response (%)
                                                                                                         Response (%)
     All Patients
                                                             84
                                                                        62.7%
                                                                                                            74.6%
                                                                                                                                   (1.04, 2.97)
                                       268
                                                 134
                                                                                     134
                                                                                                100
                                                                                                                         1.76
     Sex
                                       161
                                                             48
                                                                        58.5%
                                                                                     79
                                                                                                61
                                                                                                            77.2%
                                                                                                                         2.45
                                                                                                                                   (1.23, 4.86)
       М
                                       107
                                                 52
                                                             36
                                                                        69.2%
                                                                                     55
                                                                                                39
                                                                                                            70.9%
                                                                                                                         1.08
                                                                                                                                   (0.48, 2.46)
```

Image: screenshot of a template output in the TLG catalogue using test data.

66.7%

50

39

78.0%

1.77

(0.71, 4.39)

30

Categorical Level Biomarker 2

95

45

LOW

Conclusion

Conclusion



This is the best time to embark on E2E programming R



Setup of infrastructure can challenging,



Plenty of resources available for personal development Outputs successfully delivered on

this study.



Pharmaverse is in a good condition



We now have the flexibility to use both R & SAS.

Thank You

Any questions?



Acknowledgements

- Astellas 5502-CL-0001 statistics/programming team (Anna Spence & Jing Zhu)
- Astellas Clinical Data and Information Strategy team
- My Manager (Mary Beth Blauwet)
- Open Source Stream leads

Contact Information





EMAIL:

LUKE.THOMPSON@ASTELLAS.COM