

# Create your own recipes steps for omics data: the `{scimo}` package

useR! 2024  
Antoine Bichat  
July 9th  
Salzburg



**SERVIER**   
*moved by you*

- **Intro**
- Preprocessing
- Omics data
- Create your first step
- Dependencies without dependencies
- Outro

# Presentation

Data scientist @ Servier

- Exploratory analysis
- Oncology, pediatric cancers, targeted therapies
- R, packages, shiny apps



Workshop at **Agrostat**, Bragança, Portugal, in September 24

*Creation of an end-to-end machine learning pipeline with {tidymodels}*

# Notes

## Open source package

`scimo` is a package developed on my own free time, and not for my employer.

## Code chunk

There will be a lot of code in this presentation, but it is not necessary to look at it in detail on the first read-through.

# The Palmer Archipelago penguins

```
1 library(tidyverse)
2 library(palmerpenguins)
3
4 penguins

# A tibble: 344 × 8
  species island   bill_length_mm bill_depth_mm flipper_length_mm body_mass_g sex   year
  <fct>   <fct>        <dbl>          <dbl>            <int>        <int> <fct> <int>
1 Adelie  Torgersen     39.1           18.7            181       3750 male  2007
2 Adelie  Torgersen     39.5           17.4            186       3800 female 2007
3 Adelie  Torgersen     40.3           18              195       3250 female 2007
4 Adelie  Torgersen      NA             NA              NA         NA <NA> 2007
5 Adelie  Torgersen     36.7           19.3            193       3450 female 2007
6 Adelie  Torgersen     39.3           20.6            190       3650 male  2007
7 Adelie  Torgersen     38.9           17.8            181       3625 female 2007
8 Adelie  Torgersen     39.2           19.6            195       4675 male  2007
9 Adelie  Torgersen     34.1           18.1            193       3475 <NA> 2007
10 Adelie Torgersen      42             20.2            190       4250 <NA> 2007
# i 334 more rows
```



- Intro
- **Preprocessing**
- Omics data
- Create your first step
- Dependencies without dependencies
- Outro

# tidymodels

The tidymodels framework is a collection of packages for modeling and machine learning using tidyverse principles.



# Your first recipe

A recipe is an object that defines a series of roles and steps for data processing.

```
1 library(recipes)
2
3 penguins %>%
4   recipe(flipper_length_mm ~ .)
```

— Recipe —————

— Inputs

Number of variables by role

outcome: 1

predictor: 7

# Steps

```
1 penguins %>%
2   recipe(flipper_length_mm ~ .) %>%
3   step_impute_mean(all_numeric_predictors(), -year) %>%
4   step_normalize(all_numeric_predictors(), -year) %>%
5   step_pca(all_numeric_predictors(), -year, num_comp = 2)
```

— Recipe

---

— Inputs

Number of variables by role

outcome: 1  
predictor: 7

— Operations

- Mean imputation for: all\_numeric\_predictors() and -year
- Centering and scaling for: all\_numeric\_predictors() and -year
- PCA extraction with: all\_numeric\_predictors() and -year

# Estimation

```
1 penguins %>%
2   recipe(flipper_length_mm ~ .) %>%
3   step_impute_mean(all_numeric_predictors(), -year) %>%
4   step_normalize(all_numeric_predictors(), -year) %>%
5   step_pca(all_numeric_predictors(), -year, num_comp = 2) %>%
6   prep()
```

— Recipe

---

— Inputs

Number of variables by role

outcome: 1  
predictor: 7

— Training information

Training data contained 344 data points and 11 incomplete rows.

— Operations

- Mean imputation for: bill\_length\_mm, bill\_depth\_mm, body\_mass\_g | Trained
- Centering and scaling for: bill\_length\_mm, bill\_depth\_mm, body\_mass\_g | Trained
- PCA extraction with: bill\_length\_mm, bill\_depth\_mm, body\_mass\_g | Trained

# Application

```
1 penguins %>%
2   recipe(flipper_length_mm ~ .) %>%
3   step_impute_mean(all_numeric_predictors(), -year) %>%
4   step_normalize(all_numeric_predictors(), -year) %>%
5   step_pca(all_numeric_predictors(), -year, num_comp = 2) %>%
6   prep() %>%
7   bake(new_data = NULL)

# A tibble: 344 × 7
  species island    sex    year flipper_length_mm      PC1      PC2
  <fct>   <fct>   <fct>   <int>           <dbl>     <dbl>
1 Adelie   Torgersen male    2007        181 -1.27    -0.0476
2 Adelie   Torgersen female  2007        186 -0.854    0.429
3 Adelie   Torgersen female  2007        195 -1.37    0.157
4 Adelie   Torgersen <NA>    2007        NA  0.000199 -0.0000262
5 Adelie   Torgersen female  2007        193 -1.92    0.00581
6 Adelie   Torgersen male    2007        190 -1.81    -0.827
7 Adelie   Torgersen female  2007        181 -1.16    0.352
8 Adelie   Torgersen male    2007        195 -0.731   -0.522
9 Adelie   Torgersen <NA>    2007        193 -1.86    0.774
10 Adelie  Torgersen <NA>   2007        190 -0.936   -1.03
# i 334 more rows
```

# Step info

```
1 penguins %>%
2   recipe(flipper_length_mm ~ .) %>%
3   step_impute_mean(all_numeric_predictors(), -year) %>%
4   step_normalize(all_numeric_predictors(), -year) %>%
5   step_pca(all_numeric_predictors(), -year, num_comp = 2) %>%
6   prep() %>%
7   tidy(2)
```

```
# A tibble: 6 × 4
  terms      statistic    value id
  <chr>      <chr>        <dbl> <chr>
1 bill_length_mm mean       43.9  normalize_hFPQb
2 bill_depth_mm  mean      17.2  normalize_hFPQb
3 body_mass_g   mean     4202.  normalize_hFPQb
4 bill_length_mm sd        5.44  normalize_hFPQb
5 bill_depth_mm  sd       1.97  normalize_hFPQb
6 body_mass_g   sd       800.  normalize_hFPQb
```

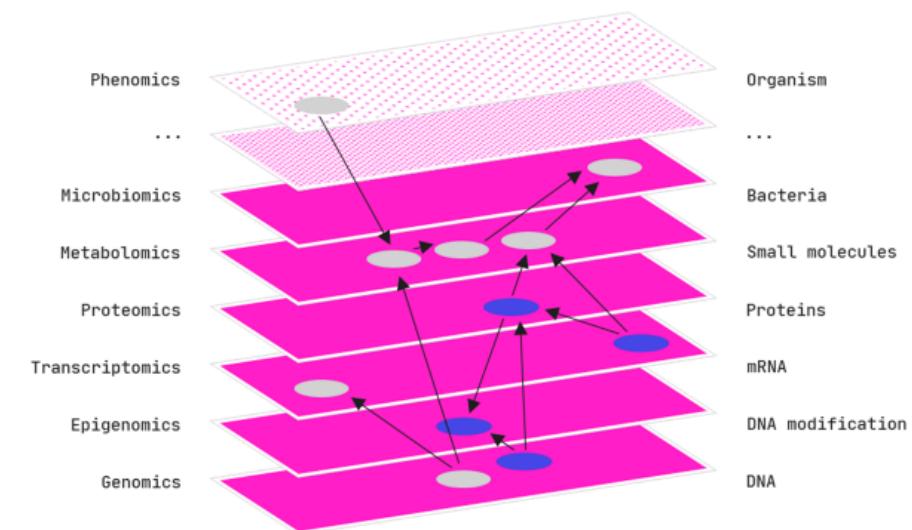
- Intro
- Preprocessing
- **Omics data**
- Create your first step
- Dependencies without dependencies
- Outro

# Omics

Omics data refers to **large datasets** generated through high-throughput sequencing, used to study and understand **complex biological systems**.

- **Genomics:** gene mutations, fusions...
- **Transcriptomics:** gene expression.
- **Proteomics:** protein abundance.
- **Metagenomics:** microorganism abundance.
- ...

Human multi-omics dataset



# Need for specific preprocessing

- Feature **selection** steps:
  - keep genes with the most variability,
  - keep features significantly associated with the outcome.
- Feature **aggregation** steps:
  - compute pathway activity score,
  - sum all abundances belonging to the same clade.
- Feature **normalization** steps:
  - convert absolute counts to proportion.
- Feature **generation** steps:
  - reduce dimension for special data distribution,
  - extract clades from lineages.

# {scimo}

```
1 library(scimo)
```

{scimo} provides extra recipes steps for dealing with omics data, while also being adaptable to other data types.



R-CMD-check passing CRAN 0.0.2

```
step_select_cv()  
step_select_wilcoxon()  
step_aggregate_list()  
step_rownormalize_tss()  
step_taxonomy()
```

...

# Pediatric cancer dataset

Gene expression of 108 CCLE cell lines from 5 different pediatric cancers.

```
1 data("pedcan_expression")
2 pedcan_expression

# A tibble: 108 × 19,197
  cell_line sex event disease   A1BG   A1CF   A2M   A2ML1   A3GALT2   A4GALT   A4GNT   AAAS   AAC
  <chr>     <chr> <chr> <chr>     <dbl>   <dbl>   <dbl>   <dbl>     <dbl>     <dbl>   <dbl>   <dbl>   <dbl>
1 143B      Female Primary Osteosarcoma 3.02 0.0566 2.78 0 0 2.13 0 5.55 2.90
2 A-673     Female Primary Ewing Sarcoma 4.87 0 2.00 3.19 0.0841 4.62 0.189 6.64 4.53
3 BT-12     Female Primary Embryonal Tumor 3.52 0.0286 0.111 0 0 2.32 0.0704 5.78 2.64
4 BT-16     Male Unknown Embryonal Tumor 3.51 0 0.433 0.0144 0 1.54 0.0144 5.48 2.72
5 C396      Male Metastatic Osteosarcoma 4.59 0 0.956 0 0 5.10 0 4.81 2.70
6 CADO-ES1  Female Metastatic Ewing Sarcoma 5.89 0 0.614 0.379 0.0704 6.60 0.151 6.08 3.92
7 CAL-72    Male Primary Osteosarcoma 4.35 0.0426 0.333 0 0 0.614 0 4.53 3.42
8 CBAGPN    Female Primary Ewing Sarcoma 4.87 0.0976 1.33 0.111 0 0.722 0.0704 5.94 3.40
9 CHLA-06   Female Unknown Embryonal Tumor 5.05 0 0.124 0 0 0.848 0.138 4.11 2.22
10 CHLA-10  Female Unknown Ewing Sarcoma 5.05 0.0144 0.949 1.73 0.0704 0.506 0.0704 5.86 4.65
# i 98 more rows
# i 19,184 more variables: AADAC <dbl>,AADACL2 <dbl>,AADACL3 <dbl>,AADACL4 <dbl>,AADAT <dbl>,
# i AAGAB <dbl>,AAK1 <dbl>,AAMDC <dbl>,AAMP <dbl>,AANAT <dbl>,AAR2 <dbl>,AARD <dbl>,AARS1 <dbl>,
# i AARS2 <dbl>,AARSD1 <dbl>,AASDH <dbl>,AASDHPPPT <dbl>,AASS <dbl>,AATF <dbl>,AATK <dbl>,ABAT <dbl>, ...
```

- Intro
- Preprocessing
- Omics data
- **Create your first step**
- Dependencies without dependencies
- Outro

# Coefficient of variation

Omics data are usually large data. For pedcan\_expression,  $p \approx 20000 \gg n \approx 100$ .

We want a step that will keep the top 10% of all variables that have the greater coefficient of variation.

$$CV = \frac{\sigma}{|\mu|}$$

```
1 cv <- function(x, na.rm = TRUE) {  
2   sd(x, na.rm = na.rm) / abs(mean(x, na.rm = na.rm))  
3 }
```

# step\_select\_cv( ) in action

```
1 rec_cv <-
2   # recipe(disease ~ ., data = pedcan_expression) %>% # too many variables for a formula
3   recipe(pedcan_expression) %>%                         # <- workaround to avoid stack overflow error
4   update_role(disease, new_role = "output") %>%        #     should be resolved in a future version of recipes
5   update_role(-disease, new_role = "predictor") %>%    #
6   step_select_cv(all_numeric_predictors(), prop_kept = 0.1) %>%
7   prep()
```

```
1 bake(rec_cv, new_data = NULL)
```

```
# A tibble: 108 × 1,923
  cell_line sex event disease A1CF AADAC AADACL3
  <fct>     <fct> <fct> <fct>   <dbl>  <dbl>   <dbl>
1 143B      Female Primary Osteos... 0.0566 0.0704    0
2 A-673      Female Primary Ewing ... 0          0.151    0
3 BT-12      Female Primary Embryo... 0.0286  3.49     0
4 BT-16      Male  Unknown Embryo... 0          0.0286    0
5 C396       Male  Metastatic Osteos... 0          0         0
6 CADO-ES1   Female Metastatic Ewing ... 0          0         0
7 CAL-72     Male  Primary Osteos... 0.0426  0          0
8 CBAGPN     Female Primary Ewing ... 0.0976  0          0
9 CHLA-06    Female Unknown Embryo... 0          0         0
10 CHLA-10   Female Unknown Ewing ... 0.0144  0          0
# i 98 more rows
# i 1,916 more variables: ABCB11 <dbl>, ABCC12 <dbl>,
#   ABRA <dbl>, AC002456.2 <dbl>, AC008397.1 <dbl>,
#   ACOD1 <dbl>, ACSM1 <dbl>, ACSM2B <dbl>, ACSM5 <dbl>, ...
```

```
1 tidy(rec_cv, 1)
```

```
# A tibble: 19,193 × 4
  terms            cv kept id
  <chr>           <dbl> <lgl> <chr>
1 A1BG            0.371 FALSE select_cv_3YfWC
2 A1CF            4.60  TRUE select_cv_3YfWC
3 A2M             1.69 FALSE select_cv_3YfWC
4 A2ML1           2.45 FALSE select_cv_3YfWC
5 A3GALT2         2.37 FALSE select_cv_3YfWC
6 A4GALT          0.979 FALSE select_cv_3YfWC
7 A4GNT           1.53 FALSE select_cv_3YfWC
8 AAAS            0.0934 FALSE select_cv_3YfWC
9 AAC5            0.194 FALSE select_cv_3YfWC
10 AADAC           3.40  TRUE select_cv_3YfWC
# i 19,183 more rows
```

# User interface: step\_select\_cv()

```
1 step_select_cv <- function(recipe, ..., role = NA, trained = FALSE,
2                               n_kept = NULL, prop_kept = NULL,
3                               cutoff = NULL, res = NULL,
4                               skip = FALSE,
5                               id = rand_id("select_cv")) {
6
7   add_step(                                # Add a new step to the existing recipe
8     recipe,
9     step_select_cv_new(                   # Arguments are passed as is
10       terms = enquos(...),
11       role = role,
12       trained = trained,                 # trained = FALSE
13       n_kept = n_kept,
14       prop_kept = prop_kept,
15       cutoff = cutoff,
16       res = res,                      # res = NULL to update later
17       skip = skip,
18       id = id)                      # Random id
19   )
20 }
21 }
```

# Step creation: `step_select_cv_new()`

```
1 step_select_cv_new <- function(terms, role, trained,
2                               n_kept, prop_kept, cutoff,
3                               res, skip, id) {
4
5   step(
6     subclass = "select_cv", # Specify the class to dispatch future methods
7     terms = terms,          # Arguments are passed as is
8     role = role,
9     trained = trained,
10    n_kept = n_kept,
11    prop_kept = prop_kept,
12    cutoff = cutoff,
13    res = res,
14    skip = skip,
15    id = id
16  )
17 }
```

# Computing: `prep.step_select_cv()`

```
1 prep.step_select_cv <- function(x, training, info = NULL, ...) {  
2  
3   col_names <- recipes_eval_select(x$terms, training, info)      # x is a list containing step info  
4   check_type(training[, col_names], quant = TRUE)                  # Check variables  
5  
6   #####  
7  
8   res_cv <-  
9     training[, col_names] %>%  
10    apply(2, cv) %>%  
11    enframe(name = "terms", value = "cv") %>%  
12    mutate(kept = var_to_keep(.data$cv, x$n_kept, x$prop_kept, x$cutoff, maximize = TRUE))  
13  
14 #####  
15  
16   step_select_cv_new(# Update step in recipe  
17     terms = x$terms,      # Most arguments are passed as is  
18     role = x$role,  
19     trained = TRUE,        # The step is now trained  
20     n_kept = x$n_kept,  
21     prop_kept = x$prop_kept,  
22     cutoff = x$cutoff,  
23     res = res_cv,          # Result to store to use later  
24     skip = x$skip,  
25     id = x$id  
26   )  
27 }
```

# Applying: `bake.step_select_cv()`

```
1  bake.step_select_cv <- function(object, new_data, ...) {  
2  
3    col_names <- object$res$terms          # object is a list containing step info  
4    check_new_data(col_names, object, new_data) # Check variables  
5  
6    #####  
7  
8    col_to_remove <-          # Do things  
9      object$res %>%          # remove unwanted columns  
10     filter(!.data$kept) %>%  
11     pull(.data$terms)  
12  
13    new_data[col_to_remove] <- NULL  
14  
15    #####  
16  
17    new_data # Return updated dataset  
18 }
```

# Getting informations: tidy.step\_select\_cv()

```
1 tidy.step_select_cv <- function(x, ...) {  
2  
3   if (is_trained(x)) {  
4     res <- x$res                      # res contains all necessary information  
5   } else {  
6     term_names <- sel2char(x$terms)  
7     res <-  
8       tibble(  
9         terms = term_names,            # Returns NA when not trained  
10        cv = rlang::na_dbl,  
11        rank = rlang::na_dbl,  
12        kept = rlang::na_lgl  
13      )  
14    }  
15  
16    res$id <- x$id                   # Add the unique random id  
17    res  
18 }
```

# Printing: `print.step_select_cv()`

```
1 print.step_select_cv <- function(x, width = max(20, options()$width - 35), ...) {  
2  
3   title <- "Top CV filtering on "  
4  
5   print_step(  
6     tr_obj = x$res$terms,  
7     untr_obj = x$terms,  
8     trained = x$trained,  
9     title = title,  
10    width = width  
11  )  
12  
13  invisible(x)  
14}
```

# Methods to import

To correctly manage your NAMESPACE

```
1 #' @keywords internal
2 '_PACKAGE'
3
4 ## usethis namespace: start
5 #' @importFrom generics required_pkgs tidy
6 #' @importFrom recipes prep bake
7 #' @importFrom tibble tibble
8 ## usethis namespace: end
9 NULL
```

- Intro
- Preprocessing
- Omics data
- Create your first step
- **Dependencies without dependencies**
- Outro

# Deal with taxonomic lineages

```
1 data("cheese_taxonomy")
2
3 cheese_taxonomy %>%
4   select(asv, lineage)

# A tibble: 74 × 2
  asv      lineage
  <chr>    <chr>
1 asv_01  k_Fungi|p_Aскомикота|c_Dothideomycetes|o_Dothideales|f_Dothioraceae|g_Aureobasidium|s_Aureob...
2 asv_02  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Aspergillus|s_Aspergil...
3 asv_03  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
4 asv_04  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
5 asv_05  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
6 asv_06  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
7 asv_07  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
8 asv_08  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
9 asv_09  k_Fungi|p_Aскомикота|c_Eurotiomycetes|o_Eurotiales|f_Aspergillaceae|g_Penicillium|s_Penicill...
10 asv_10  k_Fungi|p_Aскомикота|c_Saccharomycetes|o_Saccharomycetales|f_Debaryomycetaceae|g_Debaryomyces...
# i 64 more rows
```

# Deal with taxonomic lineages

```
1 data("cheese_taxonomy")
2
3 cheese_taxonomy %>%
4   select(asv, lineage) %>%
5   mutate(order = yatah::get_clade(lineage, "order"),
6         genus = yatah::get_clade(lineage, "genus"))

# A tibble: 74 × 4
  asv    lineage                                order      genus
  <chr>  <chr>                                 <chr>
1 asv_01 k_Fungi|p_Aскомицота|c_Dothideomycetes|o_Dothideales|f_Dothiorac... Dothideales Aureobasidium
2 asv_02 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Aspergillus
3 asv_03 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
4 asv_04 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
5 asv_05 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
6 asv_06 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
7 asv_07 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
8 asv_08 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
9 asv_09 k_Fungi|p_Aскомицота|c_Eurotiomycetes|o_Eurotiales|f_Aspергillac... Eurotiales Penicillium
10 asv_10 k_Fungi|p_Aскомицота|c_Saccharomycetes|o_Saccharomycetales|f_Deb... Saccharomyceta... Debaryomyces
# i 64 more rows
```



# step\_taxonomy( )

```
1 cheese_taxonomy %>%
2   recipe(~ asv + lineage) %>%
3   step_taxonomy(lineage, rank = c("order", "genus")) %>%
4   prep() %>%
5   bake(new_data = NULL)
```

```
# A tibble: 74 × 3
  asv      lineage_order    lineage_genus
  <fct>    <chr>          <chr>
1 asv_01   Dothideales    Aureobasidium
2 asv_02   Eurotiales     Aspergillus
3 asv_03   Eurotiales     Penicillium
4 asv_04   Eurotiales     Penicillium
5 asv_05   Eurotiales     Penicillium
6 asv_06   Eurotiales     Penicillium
7 asv_07   Eurotiales     Penicillium
8 asv_08   Eurotiales     Penicillium
9 asv_09   Eurotiales     Penicillium
10 asv_10  Saccharomycetales Debaryomyces
# i 64 more rows
```

# The classical way

```
1 #' @importFrom yatah get_clade
2 bake.step_taxonomy <- function(object, new_data, ...) {
3   ...
4   new_col <- paste0(term, "_", rank)
5
6   new_data[[new_col]] <- get_clade(new_data[[term]], rank = rank, same = FALSE)
7   ...
8   return(new_data)
9 }
```

- Need to add {yatah} as a dependency.

# A word about `call2()` and `eval_tidy()`

```
1 head(fruit)  
[1] "apple"      "apricot"     "avocado"     "banana"      "bell pepper" "bilberry"  
1 knitr:::combine_words(head(fruit), and = " or ")  
apple, apricot, avocado, banana, bell pepper, or bilberry
```

```
1 library(rlang)  
2 cl <- call2("combine_words", .ns = "knitr",  
3               words = head(fruit), and = " or ")  
4 cl  
  
knitr:::combine_words(words = c("apple", "apricot", "avocado",  
"banana", "bell pepper", "bilberry"), and = " or ")  
1 eval_tidy(cl)  
apple, apricot, avocado, banana, bell pepper, or bilberry
```

# The dependency-free way

```
1 #' @importFrom rlang eval_tidy call2
2 bake.step_taxonomy <- function(object, new_data, ...) {
3   ...
4   new_col <- paste0(term, "_", rank)
5
6   yatah_call <- call2("get_clade", .ns = "yatah",
7                         lineage = new_data[[term]], rank = rank, same = TRUE)
8   new_data[[new_col]] <- eval_tidy(yatah_call)
9   ...
10  return(new_data)
11 }
```

- {yatah} is no longer needed.
- {rlang} is already a {recipes} dependency.

# required\_pkgs()

```
1 required_pkgs.step_taxonomy <- function(x, ...) {  
2   c("yatah", "scimo")  
3 }
```

- Check if the used package is installed.
- Correctly load the package for parallel processing.
- Used also for other steps in {scimo}, and returns only "scimo".

- Intro
- Preprocessing
- Omics data
- Create your first step
- Dependencies without dependencies
- **Outro**

# Next steps

- Tunable arguments.
- New steps
  - Other tests (`limma`, `DESeq2`...) for feature selection,
  - `PLNmodels` dimension reduction,
  - Batch effect removal,
  - Multi-omics steps (how to define groups with tidy selection?).

# To go further

- Packages
  - `tidymodels`
  - `recipes`
  - `scimo`
- Books
  - Tidy Modeling with R
- Vignette
  - Create your own recipe step function

# Big thanks

**Julie Aubert**

for ideas, discussions and contributions



**Emil Hvitfeldt**

for review and issues

**Sylvain Jonchery**

for the logo



