# R: A Hitchhikers Guide to Reproducible Research

## - Take a parachute and jump (into the tidyverse)

Brendan Palmer,

Clinical Research Facility - Cork &
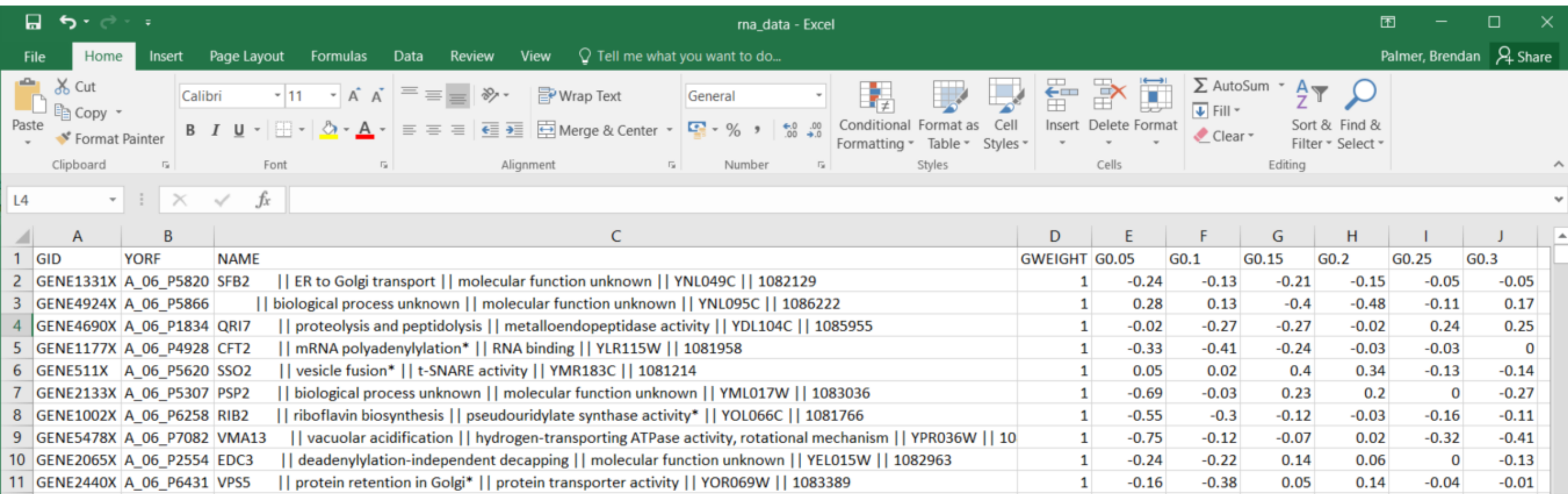
School of Public Health

@B_A_Palmer

CRF-C
HRB Clinical Research Facility Cork

UCC | School of Public Health
University College Cork, Ireland
Coláiste na hOllscoile Corcaigh

# Tidyverse works best with tidy data

- Each variable forms a column
- Each observation forms a row

**Problems with the example RNA data set…**

# Tidyverse works best with tidy data



- Multiple variables are stored in one column
    - e.g. column "NAME" contains values such as;

G0.05 - letter identifies a compound
    - number is the concentration of that compound

# Tidyverse code structure has two main forms

(1)

`new_object` `<-` `function(` `input_data,`

`data_to_b_modified,` `arguments_to_function` `)`

(2)

`new_object` `<-` `input_data` `%>%` ← `magrittr / pipe operator`

`function(` `data_to_b_modified,` `arguments_to_function` `)`

# Line by line

# Line by line

# Line by line

Source on Save → Run → Source

```r
12
13  raw_gene_df <- read_delim("Brauer2008_DataSet1.tds", delim = "\t")
14
15  separated_gene_df <- separate(raw_gene_df, NAME,
16                                c("name", "BP", "MF", "systematic_name",
17                                  "number"),
18                                sep = "\\|\\|")
19
20  mutated_gene_df <- mutate_at(separated_gene_df,
21                               vars(name:systematic_name),
22                               funs(trimws)
23                               )
24
25  selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
26
27  gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
28
29  nearly_there_df <- separate(gathered_gene_df, sample,
30                              c("nutrient", "rate"), sep = 1, convert = TRUE)
31
32  nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
33                      S = "Sulfate", N = "Ammonia", U = "Uracil")
34
35  cleaned_genes_df <- mutate(nearly_there_df,
36                             nutrient = plyr::revalue(nutrient, nutrient_names)
37                             ) %>%
38      filter(!is.na(expression), systematic_name != "")
```

27:1   # Section 1: Data import, tidying and transformation ↕                          R Script ↓

**Console**   **Terminal** ×

~/R_Users_Workshop/8_weeks_Oct-Dec_17/Workshop_1/workshop_1_project/

```
Parsed with column specification:
cols(
  .default = col_double(),
  GID = col_character(),
  YORF = col_character(),
  NAME = col_character(),
  GWEIGHT = col_integer()
)
See spec(...) for full column specifications.
> separated_gene_df <- separate(raw_gene_df, NAME,
+                               c("name", "BP", "MF", "systematic_name",
+                                 "number"),
+                               sep = "\\|\\|")
> mutated_gene_df <- mutate_at(separated_gene_df,
+                              vars(name:systematic_name),
+                              funs(trimws)
+                              )
> selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
>
```

**Environment**   **History**   **Connections**

Import Dataset                                        Grid

Global Environment

| Name | Type | Length | Size | Value |
|---|---|---|---|---|
| mutated_gene_df | tbl_df | 44 | 3.5 MB | 5537 obs. of 44 variables |
| raw_gene_df | tbl_df | 40 | 3.3 MB | 5537 obs. of 40 variables |
| selected_gene_df | tbl_df | 40 | 2.4 MB | 5537 obs. of 40 variables |
| separated_gene... | tbl_df | 44 | 3.6 MB | 5537 obs. of 44 variables |

**Files**   **Plots**   **Packages**   **Help**   **Viewer**

New Folder   Delete   Rename   More

Home › R_Users_Workshop › 8_weeks_Oct-Dec_17 › Workshop_1 › workshop_1_project

| Name | Size | Modified |
|---|---|---|
| .. | | |
| .RData | 2.5 KB | Oct 2, 2017, 1:49 PM |
| .Rhistory | 20.3 KB | Dec 6, 2017, 3:43 PM |
| Brauer2008_DataSet1.csv | 1.6 MB | Sep 27, 2017, 11:32 PM |
| Brauer2008_DataSet1.tds | 1.6 MB | Sep 28, 2017, 10:22 AM |
| house_completions.csv | 4 KB | Sep 28, 2017, 1:35 PM |
| irish_population.csv | 315 B | Aug 28, 2017, 4:21 PM |
| raw_house_completions.csv | 16.2 KB | Aug 25, 2017, 3:45 PM |
| workshop_1.Rproj | 217 B | Oct 18, 2018, 12:18 PM |
| ws1_script1_stepwise_Bauer_dataset_analysis.R | 6.1 KB | Dec 5, 2017, 12:19 PM |
| ws1_script2_Bauer_dataset_analysis.R | 2 KB | Dec 6, 2017, 2:33 PM |
| ws1_script3_house_completions.R | 2.4 KB | Oct 2, 2017, 3:53 PM |

# Line by line

```
12
13  raw_gene_df <- read_delim("Brauer2008_DataSet1.tds", delim = "\t")
14
15  separated_gene_df <- separate(raw_gene_df, NAME,
16                                c("name", "BP", "MF", "systematic_name",
17                                  "number"),
18                                sep = "\\|\\|")
19
20  mutated_gene_df <- mutate_at(separated_gene_df,
21                               vars(name:systematic_name),
22                               funs(trimws)
23                               )
24
25  selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
26
27  gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
28
29  nearly_there_df <- separate(gathered_gene_df, sample,
30                              c("nutrient", "rate"), sep = 1, convert = TRUE)
31
32  nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
33                      S = "Sulfate", N = "Ammonia", U = "Uracil")
34
35  cleaned_genes_df <- mutate(nearly_there_df,
36                             nutrient = plyr::revalue(nutrient, nutrient_names)
37                             ) %>%
38    filter(!is.na(expression), systematic_name != "")
```

29:1    Section 1: Data import, tidying and transformation    R Script

**Console**  **Terminal**

~/R_Users_Workshop/8_weeks_Oct-Dec_17/Workshop_1/workshop_1_project/

```
cors(
  .default = col_double(),
  GID = col_character(),
  YORF = col_character(),
  NAME = col_character(),
  GWEIGHT = col_integer()
)
See spec(...) for full column specifications.
> separated_gene_df <- separate(raw_gene_df, NAME,
+                               c("name", "BP", "MF", "systematic_name",
+                                 "number"),
+                               sep = "\\|\\|")
> mutated_gene_df <- mutate_at(separated_gene_df,
+                              vars(name:systematic_name),
+                              funs(trimws)
+                              )
> selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
> gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
>
```

**Environment**  **History**  **Connections**

Global Environment

| Name | Type | Length | Size | Value |
|---|---|---|---|---|
| gathered_gene_df | tbl_df | 6 | 9.8 MB | 199332 obs. of 6 variables |
| mutated_gene_df | tbl_df | 44 | 3.5 MB | 5537 obs. of 44 variables |
| raw_gene_df | tbl_df | 40 | 3.3 MB | 5537 obs. of 40 variables |
| selected_gene_df | tbl_df | 40 | 2.4 MB | 5537 obs. of 40 variables |
| separated_gene… | tbl_df | 44 | 3.6 MB | 5537 obs. of 44 variables |

**Files**  **Plots**  **Packages**  **Help**  **Viewer**

New Folder    Delete    Rename    More

Home > R_Users_Workshop > 8_weeks_Oct-Dec_17 > Workshop_1 > workshop_1_project

| Name | Size | Modified |
|---|---|---|
| .. | | |
| .RData | 2.5 KB | Oct 2, 2017, 1:49 PM |
| .Rhistory | 20.3 KB | Dec 6, 2017, 3:43 PM |
| Brauer2008_DataSet1.csv | 1.6 MB | Sep 27, 2017, 11:32 PM |
| Brauer2008_DataSet1.tds | 1.6 MB | Sep 28, 2017, 10:22 AM |
| house_completions.csv | 4 KB | Sep 28, 2017, 1:35 PM |
| irish_population.csv | 315 B | Aug 28, 2017, 4:21 PM |
| raw_house_completions.csv | 16.2 KB | Aug 25, 2017, 3:45 PM |
| workshop_1.Rproj | 217 B | Oct 18, 2018, 12:18 PM |
| ws1_script1_stepwise_Bauer_dataset_analysis.R | 6.1 KB | Dec 5, 2017, 12:19 PM |
| ws1_script2_Bauer_dataset_analysis.R | 2 KB | Dec 6, 2017, 2:33 PM |
| ws1_script3_house_completions.R | 2.4 KB | Oct 2, 2017, 3:53 PM |

ws1_script1_stepwise_Bauer_dataset_an... *

# Line by line

Source on Save | Run | Source

```r
12
13  raw_gene_df <- read_delim("Brauer2008_DataSet1.tds", delim = "\t")
14
15  separated_gene_df <- separate(raw_gene_df, NAME,
16                                c("name", "BP", "MF", "systematic_name",
17                                  "number"),
18                                sep = "\\|\\|")
19
20  mutated_gene_df <- mutate_at(separated_gene_df,
21                               vars(name:systematic_name),
22                               funs(trimws)
23                               )
24
25  selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
26
27  gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
28
29  nearly_there_df <- separate(gathered_gene_df, sample,
30                              c("nutrient", "rate"), sep = 1, convert = TRUE)
31
32  nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
33                      S = "Sulfate", N = "Ammonia", U = "Uracil")
34
35  cleaned_genes_df <- mutate(nearly_there_df,
36                             nutrient = plyr::revalue(nutrient, nutrient_names)
37                             ) %>%
38     filter(!is.na(expression), systematic_name != "")
```

32:1  # Section 1: Data import, tidying and transformation          R Script

## Environment | History | Connections

Import Dataset | Global Environment

| Name | Type | Length | Size | Value |
|---|---|---|---|---|
| gathered_gene_df | tbl_df | 6 | 9.8 MB | 199332 obs. of 6 variables |
| mutated_gene_df | tbl_df | 44 | 3.5 MB | 5537 obs. of 44 variables |
| nearly_there_df | tbl_df | 7 | 11.3 MB | 199332 obs. of 7 variables |
| raw_gene_df | tbl_df | 40 | 3.3 MB | 5537 obs. of 40 variables |
| selected_gene_df | tbl_df | 40 | 2.4 MB | 5537 obs. of 40 variables |
| separated_gene... | tbl_df | 44 | 3.6 MB | 5537 obs. of 44 variables |

## Files | Plots | Packages | Help | Viewer

New Folder | Delete | Rename | More

Home > R_Users_Workshop > 8_weeks_Oct-Dec_17 > Workshop_1 > workshop_1_project

| Name | Size | Modified |
|---|---|---|
| .. | | |
| .RData | 2.5 KB | Oct 2, 2017, 1:49 PM |
| .Rhistory | 20.3 KB | Dec 6, 2017, 3:43 PM |
| Brauer2008_DataSet1.csv | 1.6 MB | Sep 27, 2017, 11:32 PM |
| Brauer2008_DataSet1.tds | 1.6 MB | Sep 28, 2017, 10:22 AM |
| house_completions.csv | 4 KB | Sep 28, 2017, 1:35 PM |
| irish_population.csv | 315 B | Aug 28, 2017, 4:21 PM |
| raw_house_completions.csv | 16.2 KB | Aug 25, 2017, 3:45 PM |
| workshop_1.Rproj | 217 B | Oct 18, 2018, 12:18 PM |
| ws1_script1_stepwise_Bauer_dataset_analysis.R | 6.1 KB | Dec 5, 2017, 12:19 PM |
| ws1_script2_Bauer_dataset_analysis.R | 2 KB | Dec 6, 2017, 2:33 PM |
| ws1_script3_house_completions.R | 2.4 KB | Oct 2, 2017, 3:53 PM |

## Console | Terminal

~/R_Users_Workshop/8_weeks_Oct-Dec_17/Workshop_1/workshop_1_project/

```r
  GID = col_character(),
  YORF = col_character(),
  NAME = col_character(),
  GWEIGHT = col_integer()
)
See spec(...) for full column specifications.
> separated_gene_df <- separate(raw_gene_df, NAME,
+                               c("name", "BP", "MF", "systematic_name",
+                                 "number"),
+                               sep = "\\|\\|")
> mutated_gene_df <- mutate_at(separated_gene_df,
+                              vars(name:systematic_name),
+                              funs(trimws)
+                              )
> selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
> gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
> nearly_there_df <- separate(gathered_gene_df, sample,
+                             c("nutrient", "rate"), sep = 1, convert = TRUE)
>
```

# Line by line

Source on Save → Run → Source

```r
13   raw_gene_df <- read_delim("Brauer2008_DataSet1.tds", delim = "\t")
14
15   separated_gene_df <- separate(raw_gene_df, NAME,
16                                 c("name", "BP", "MF", "systematic_name",
17                                   "number"),
18                                 sep = "\\|\\\\|")
19
20   mutated_gene_df <- mutate_at(separated_gene_df,
21                               vars(name:systematic_name),
22                               funs(trimws)
23                               )
24
25   selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
26
27   gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
28
29   nearly_there_df <- separate(gathered_gene_df, sample,
30                               c("nutrient", "rate"), sep = 1, convert = TRUE)
31
32   nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
33                       S = "Sulfate", N = "Ammonia", U = "Uracil")
34
35   cleaned_genes_df <- mutate(nearly_there_df,
36                              nutrient = plyr::revalue(nutrient, nutrient_names)
37                              ) %>%
38     filter(!is.na(expression), systematic_name != "")
```

35:1   Section 1: Data import, tidying and transformation       R Script

## Environment / History / Connections

Import Dataset

Global Environment

| Name | Type | Length | Size | Value |
|---|---|---|---|---|
| gathered_gene_df | tbl_df | 6 | 9.8 MB | 199332 obs. of 6 variables |
| mutated_gene_df | tbl_df | 44 | 3.5 MB | 5537 obs. of 44 variables |
| nearly_there_df | tbl_df | 7 | 11.3 MB | 199332 obs. of 7 variables |
| nutrient_names | character | 6 | 984 B | Named chr [1:6] "Glucose" "... |
| raw_gene_df | tbl_df | 40 | 3.3 MB | 5537 obs. of 40 variables |
| selected_gene_df | tbl_df | 40 | 2.4 MB | 5537 obs. of 40 variables |
| separated_gene... | tbl_df | 44 | 3.6 MB | 5537 obs. of 44 variables |

## Console / Terminal

~/R_Users_Workshop/8_weeks_Oct-Dec_17/Workshop_1/workshop_1_project/

```r
  NAME = col_character(),
  GWEIGHT = col_integer()
)
See spec(...) for full column specifications.
> separated_gene_df <- separate(raw_gene_df, NAME,
+                                 c("name", "BP", "MF", "systematic_name",
+                                   "number"),
+                                 sep = "\\|\\\\|")
> mutated_gene_df <- mutate_at(separated_gene_df,
+                               vars(name:systematic_name),
+                               funs(trimws)
+                               )
> selected_gene_df <- select(mutated_gene_df, -number, -GID, -YORF, -GWEIGHT)
> gathered_gene_df <- gather(selected_gene_df, sample, expression, G0.05:U0.3)
> nearly_there_df <- separate(gathered_gene_df, sample,
+                               c("nutrient", "rate"), sep = 1, convert = TRUE)
> nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
+                       S = "Sulfate", N = "Ammonia", U = "Uracil")
>
```
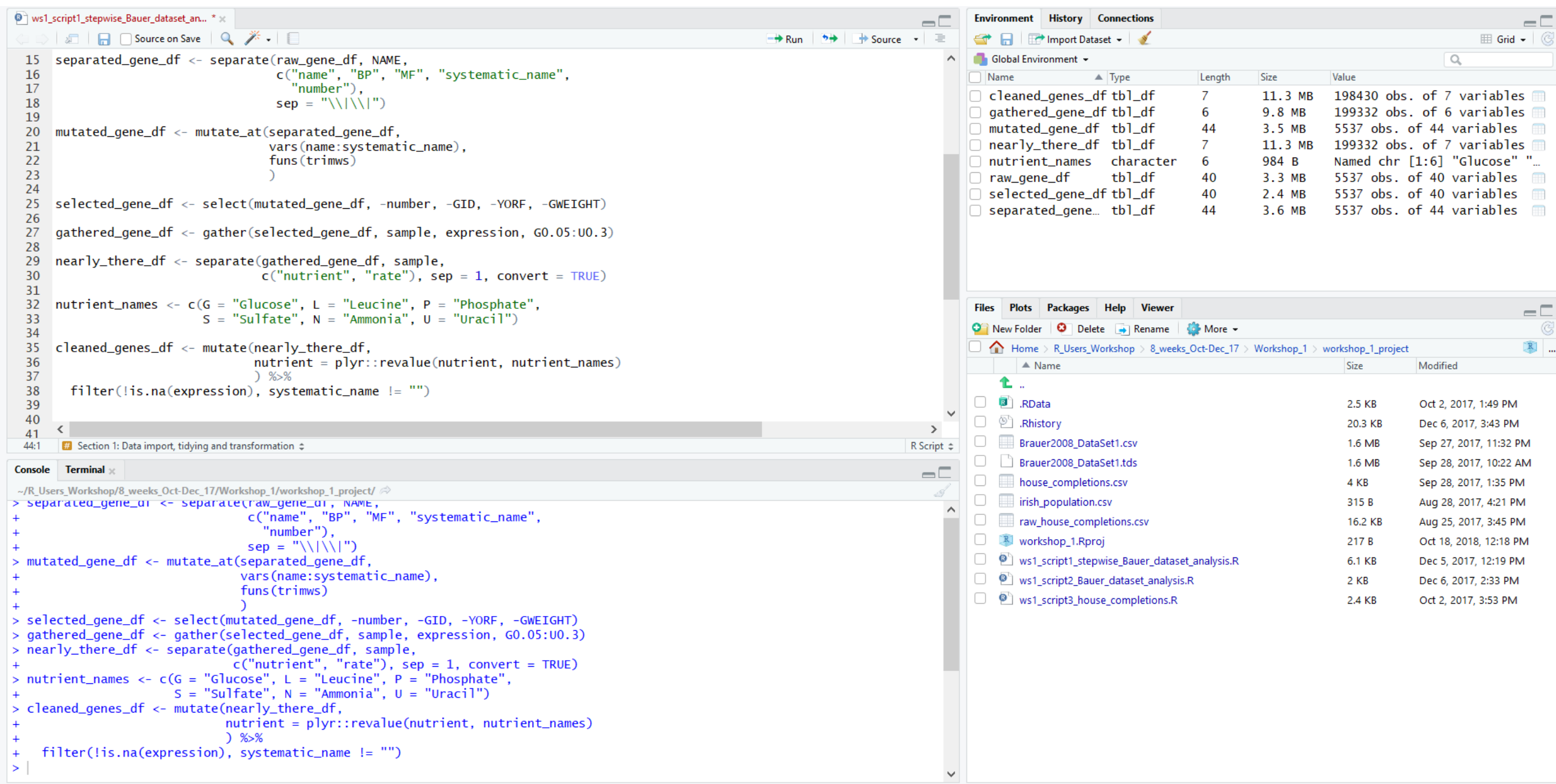
## Files / Plots / Packages / Help / Viewer

New Folder   Delete   Rename   More

Home > R_Users_Workshop > 8_weeks_Oct-Dec-17 > Workshop_1 > workshop_1_project

| Name | Size | Modified |
|---|---|---|
| .. | | |
| .RData | 2.5 KB | Oct 2, 2017, 1:49 PM |
| .Rhistory | 20.3 KB | Dec 6, 2017, 3:43 PM |
| Brauer2008_DataSet1.csv | 1.6 MB | Sep 27, 2017, 11:32 PM |
| Brauer2008_DataSet1.tds | 1.6 MB | Sep 28, 2017, 10:22 AM |
| house_completions.csv | 4 KB | Sep 28, 2017, 1:35 PM |
| irish_population.csv | 315 B | Aug 28, 2017, 4:21 PM |
| raw_house_completions.csv | 16.2 KB | Aug 25, 2017, 3:45 PM |
| workshop_1.Rproj | 217 B | Oct 18, 2018, 12:18 PM |
| ws1_script1_stepwise_Bauer_dataset_analysis.R | 6.1 KB | Dec 5, 2017, 12:19 PM |
| ws1_script2_Bauer_dataset_analysis.R | 2 KB | Dec 6, 2017, 2:33 PM |
| ws1_script3_house_completions.R | 2.4 KB | Oct 2, 2017, 3:53 PM |

# Line by line

# Nested



```r
nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
                    S = "Sulfate", N = "Ammonia", U = "Uracil")

cleaned_genes_df <-
  filter(
    mutate(
      separate(
        gather(
          select(
            mutate_at(
              separate(
                read_delim("Brauer2008_DataSet1.tds", delim = "\t"),
                NAME,
                c("name", "BP", "MF", "systematic_name",  "number"),
                sep = "\\|\\|"), vars(name:systematic_name),
              funs(trimws)),
            -number, -GID, -YORF, -GWEIGHT),
          sample, expression, G0.05:U0.3),
        sample,
        c("nutrient", "rate"),
        sep = 1, convert = TRUE),
      nutrient = plyr::revalue(nutrient, nutrient_names)),
    !is.na(expression), systematic_name != "")
```

# Nested

```
nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
                    S = "Sulfate", N = "Ammonia", U = "Uracil")

cleaned_genes_df <-
  filter(
    mutate(
      separate(
        gather(
          select(
            mutate_at(
              separate(
                read_delim("Brauer2008_DataSet1.tds", delim = "\t"),
                NAME,
                c("name", "BP", "MF", "systematic_name",  "number"),
                sep = "\\|\\|"), vars(name:systematic_name),
              funs(trimws)),
            -number, -GID, -YORF, -GWEIGHT),
          sample, expression, G0.05:U0.3),
        sample,
        c("nutrient", "rate"),
        sep = 1, convert = TRUE),
      nutrient = plyr::revalue(nutrient, nutrient_names)),
    !is.na(expression), systematic_name != "")
|
```

# Putting the pieces together

# Code structure has two main forms

(1) `new_object` `<-` `function(` `input_data,` `data_to_b_modified,` `arguments_to_function` `)`

(2) `new_object` `<-` `input_data` `function(` `data_to_b_modified,` `arguments_to_function` `)`

# Piped

ws1_script1_stepwise_Bauer_dataset_an...   ws1_script2_Bauer_dataset_analysis.R*

Source on Save    Run    Source

```
 1  nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
 2                      S = "Sulfate", N = "Ammonia", U = "Uracil"
 3                      )
 4
 5  cleaned_genes_df <- read_delim("Brauer2008_DataSet1.tds", delim = "\t"
 6                      ) %>%
 7
 8    separate(NAME, c("name", "BP", "MF", "systematic_name", "number"), sep = "\\|\\|"
 9            )  %>%
10
11    mutate_at(vars(name:systematic_name), funs(trimws)
12            ) %>%
13
14    select(-number, -GID, -YORF, -GWEIGHT
15            ) %>%
16
17    gather(sample, expression, G0.05:U0.3
18            ) %>%
19
20    separate(sample, c("nutrient", "rate"), sep = 1, convert = TRUE
21            ) %>%
22
23    mutate(nutrient = plyr::revalue(nutrient, nutrient_names)
24            ) %>%
25
26    filter(!is.na(expression), systematic_name != ""
27            )
```

9:18    (Top Level)    R Script

**Console**    Terminal

~/R_Users_Workshop/8_weeks_Oct-Dec_17/Workshop_1/workshop_1_project/

```
+   separate(sample, c("nutrient", "rate"), sep = 1, convert = TRUE
+           ) %>%
+
+   mutate(nutrient = plyr::revalue(nutrient, nutrient_names)
+           ) %>%
+
+   filter(!is.na(expression), systematic_name != ""
+           )
Parsed with column specification:
cols(
  .default = col_double(),
  GID = col_character(),
  YORF = col_character(),
  NAME = col_character(),
  GWEIGHT = col_integer()
)
See spec(...) for full column specifications.
> 
```

**Environment**    History    Connections

Import Dataset

Global Environment

| Name | Type | Length | Size | Value |
|---|---|---|---|---|
| cleaned_genes_df | tbl_df | 7 | 11.3 MB | 198430 obs. of 7 variables |
| nutrient_names | character | 6 | 984 B | Named chr [1:6] "Glucose" "Le... |

**Files**    Plots    Packages    Help    Viewer

New Folder    Delete    Rename    More

Home > R_Users_Workshop > 8_weeks_Oct-Dec_17 > Workshop_1 > workshop_1_project

| Name | Size | Modified |
|---|---|---|
| .. | | |
| .RData | 2.5 KB | Oct 2, 2017, 1:49 PM |
| .Rhistory | 20.3 KB | Dec 6, 2017, 3:43 PM |
| Brauer2008_DataSet1.csv | 1.6 MB | Sep 27, 2017, 11:32 PM |
| Brauer2008_DataSet1.tds | 1.6 MB | Sep 28, 2017, 10:22 AM |
| house_completions.csv | 4 KB | Sep 28, 2017, 1:35 PM |
| irish_population.csv | 315 B | Aug 28, 2017, 4:21 PM |
| raw_house_completions.csv | 16.2 KB | Aug 25, 2017, 3:45 PM |
| workshop_1.Rproj | 217 B | Oct 18, 2018, 12:18 PM |
| ws1_script1_stepwise_Bauer_dataset_analysis.R | 6.1 KB | Dec 5, 2017, 12:19 PM |
| ws1_script2_Bauer_dataset_analysis.R | 2 KB | Dec 6, 2017, 2:33 PM |
| ws1_script3_house_completions.R | 2.4 KB | Oct 2, 2017, 3:53 PM |

# Piped

```
1   nutrient_names <- c(G = "Glucose", L = "Leucine", P = "Phosphate",
2                       S = "Sulfate", N = "Ammonia", U = "Uracil"
3                       )
4
5   cleaned_genes_df <- read_delim("Brauer2008_DataSet1.tds", delim = "\t"
6                                  ) %>%
7
8     separate(NAME, c("name", "BP", "MF", "systematic_name", "number"), sep = "\\|\\|"
9                 )   %>%
10
11    mutate_at(vars(name:systematic_name), funs(trimws)
12               ) %>%
13
14    select(-number, -GID, -YORF, -GWEIGHT
15            ) %>%
16
17    gather(sample, expression, G0.05:U0.3
18            ) %>%
19
20    separate(sample, c("nutrient", "rate"), sep = 1, convert = TRUE
21             ) %>%
22
23    mutate(nutrient = plyr::revalue(nutrient, nutrient_names)
24             ) %>%
25
26    filter(!is.na(expression), systematic_name != ""
27             )
```

# Moral of the story...

You can go from this



To this!!

**Master Builder!!**