

# Data Cleaning with R and Tidyverse

Matt Steele

# Resources

- [Tidyverse Documentation](#)
- [O'Reilly Learning Platform](#)
  - R for Data Science, 2nd Edition
  - R Programming for Statistics and Data Science

# Recap

- Functions and Arguments
- Objects
- Run Code

# Working Directory

The working directory in R is the folder where you are working. Hence, it's the place (the environment) where you have to store your files of your project in order to load them or where your R objects will be saved.

*Session > Set Working Directory > Choose Directory*

```
1 getwd() # show current directory that you are in
2
3
4 setwd("path/to/your/directory") # sets the working directory
```

# Packages

R packages are like toolkits or collections of pre-built functions, data sets, and tools that extend the capabilities of the R programming language.

# Tidyverse

**Tidyverse** is a collection of packages focused on data analysis and data visualizations that share an underlying design philosophy, grammar, and data structures.

```
1 install.packages("tidyverse")
```

```
1 library(tidyverse)
```

# Tidyverse

tibble	lighter and more user-friendly version of data frame
tidyr	create tidy and meaningfully arranged data
readr	better importation of data into R
ggplot	data visualization functions
dplyr	data manipulation tools
lubridate	clean dates and times
purrr	better functional programming
forcats	handle, clean, and manipulate categorical variables
haven	read and write data formats from proprietary sources



# Loading Data with Readr

The `read_csv` function allows you to load data into R in a tibble data frame

```
read_csv("data_set.csv")
```

```
1 sw_df <- read_csv("starwars.csv")
```

# Loading Proprietary Data

- [readXl](#) - this package allows you to read Excel files in a tibble data frame
- [haven](#) - this package allows you to read and export non-proprietary files for SPSS, SAS, and STATA

# The Pipe Operator (|>)

The pipe operator (|> or %>%) allows you to run commands or operation on a single object based on an order of operations

	PC	MAC
Pipe Operator	CTRL + SHIFT + M	CMD + SHIFT + M

# The Pipe Operator (|>)

let's say you want to see the **name**, **height**, **mass**, and **species** of characters who were born on **Tatooine** and then relocate the **mass** variable

```
1 view(sw_df)
2
3 sw_df |> # object we are working on
4   filter(homeworld == "Tatooine") |> # first operation
5   select(name, height, mass, species) |> # second operation
6   relocate(mass, .before = name) # third operation
```

# Order of Operations

With the pipe operator, your code becomes more organized and reads like a step-by-step process:

1. Take the data ...
2. Do this
3. Then do this...
4. And finally, do this.

```
1 # order of operations matter
2
3 sw_df |> # object we are working on
4   select(name, height, mass, species) |> # first operation
5   filter(homeworld == "Tatooine") |> # second operation
6   relocate(mass, .before = name) #third operation
```

# Explore Data - Tidyverse

**view** function: interactively explore the contents of a data frame in a separate viewer window or in the RStudio viewer pane.

```
1 view(sw_df)
```

**glimpse** function: a concise overview of the data, including variable types.

```
1 glimpse(sw_df)
```

# Cleaning Data

<code>filter</code>	retains or filters out observations based on variable criteria
<code>select</code>	retains or filters out variables
<code>arrange</code>	sorts variables
<code>mutate</code>	change variable's observations OR create a new variable and observations using obser
<code>group_by</code>	group observations
<code>summarise</code>	get descriptive statistics about a variable
<code>relocate</code>	change the position of variables in the data frame
<code>rename</code>	change the name of an individual variable
<code>drop_na</code>	remove ALL missing values from a data frame or variable
<code>replace_na</code>	replace missing values with a specified

# Dplyr function: filter

allows you to *select rows* in your data frame that meet specific conditions or criteria in a variable

```
1 sw_df
2
3 # let's filter the data frame so we are seeing characters who have blue eye
4
5 sw_eye <- sw_df |>
6   filter(eye_color == "blue")
7
8 sw_eye
```



# Boolean operators

boolean operators allow you to build criteria in your code

&	AND
	OR
==	EQUAL
!=	NOT EQUAL
<	LESS THAN
>	GREATER THAN
<=	LESS THAN OR EQUAL
>=	GREATER THAN OR EQUAL

# Filter with Boolean

let's filter the data frame for characters who have blue eyes  
and were born after 50 BBY

```
1 sw_eye50 <- sw_df |>
2   filter(eye_color == "blue" & birth_year < 50)
3
4 sw_eye50
```

# Dplyr function: select

allows you to *keep* or *discard* variables

```
1 # keep variables
2
3 sw_select <- sw_df |>
4   select(name, height, mass)
5
6 sw_select
7
8 # remove variables
9
10 sw_not_select <- sw_df |>
11   select(-height, -mass)
12
13 sw_not_select
```

# Dplyr function: mutate

*creates* new variables in your data or *change* existing variables by performing calculations or transformations.

```
1 sw_df
2
3 sw_df <- sw_df |>
4   mutate(bmi = height/mass) |> # run mutate operation
5   relocate(bmi, .after = mass) # relocate variable in data frame
6
7 sw_df
```

# Dplyr function: mutate

NOTE: if you name your variable as an *existing variable*, it will *overwrite* the existing variable. If you give it a *new name*, it will create a *new variable*

```
1 # let's overwrite the old variable
2
3 sw_overwrite <- sw_df |>
4   mutate(height = height/12) # overwrite variable
5
6 sw_overwrite
```

# Dplyr function: arrange

allows you to *sort* variables

```
1 # oldest characters
2
3 sw_df |>
4   arrange(desc(birth_year))
5
6 # characters with the same skin color than the same hair color
7
8 sw_df |>
9   arrange(desc(skin_color), hair_color)
```

# Dplyr function: group\_by & summarise

the **group\_by** function allows you to group common observations in a variable

**summarise** function allows you to get descriptive statistics about the groupings

```
1 sw_group <- sw_df |>
2   group_by(sex) |>
3   summarise(avg_height = mean(height)) # get mean
4
5 sw_group
6
7 write_csv(sw_group, "sw_sex_dv.csv") # export
```

# Recode Data Values

The **as.** function along with **mutate** will allow you to change the data type of a variable. For this example we are going to recode the *character\_id* variable to interpret the data type as a *character* instead of a *double*

```
1 sw_df <- sw_df |>
2   mutate(character_id = as.character(character_id))
3
4 sw_df
```



# Dyplr Function: recode observations

we can rename the values of observations within a variable using the **mutate** function in combination with the **recode** or **recode\_factor** functions

```
1 # let's create a new variable with to give numeric levels instead of value
2
3 sw_df <- sw_df |>
4   mutate(gender_num = recode(gender,
5                               "masculine" = 0,
6                               "feminine" = 1)) |>
7   relocate(gender_num, .after = gender)
8
9
10 view(sw_df)
```

# Dplyr function: rename

allows you rename variables in your data frame

```
1 glimpse(sw_df)
2
3 sw_df <- sw_df |>
4   rename("gender_label" = gender)
5
6 glimpse(sw_df)
```

# Missing Data

missing data in numeric fields can cause an issue when trying to calculate descriptive statistics

```
1 # are there missing NA values
2
3 which(is.na(sw_df$bmi))
4
5 # with missing values we cannot calculate descriptive statistics
6
7 mean(sw_df$bmi)
```

# Tidyr function: drop\_na

removes **all** missing data from data frames or variables

```
1 sw_dropNA <- sw_df |>
2   drop_na()
3
4 mean(sw_dropNA$bmi)
```

we can also just drop NAs from a **variable**

```
1 sw_dropNA_var <- sw_df |>
2   drop_na(bmi)
3
4 mean(sw_dropNA_var$bmi)
```

# Tidyr function: `replace_na`

you can also recode the NA values for observations with *mutate* and *replace\_na*

```
1 # let's replace the NAs the gender_label variable with "unknown"
2
3
4 sw_df <- sw_df |>
5   mutate(gender_label = replace_na(gender_label, "unknown")) |>
6   mutate(gender_label = as_factor(gender_label))
7
8 levels(sw_df$gender_label)
```

# Readr function: write\_csv

the **write\_csv** function allows us to export data frames to a csv file once we are done cleaning it up or when we have done some analysis that we want to export

```
1 # now that we have this date frame cleaned let's save it
2
3 # let's export the file
4
5 write_csv(sw_df, "starwars_clean.csv")
```

# Psych Package

**Psych Package** - built-in functions for factor analysis, reliability analysis, descriptive statistics and data visualization.

```
1 install.packages("psych")  
2 library(psych)
```

```
1 sw_ds <- describe(sw_df)  
2  
3 write_csv(sw_ds, "starwars_ds.csv")
```

# Summarytools Package

**SummaryTools Package** - simplifies data exploration and descriptive statistics generation for data frames and vectors.

```
1 install.packages("summarytools")  
2 library(summarytools)
```

```
1 descr(sw_df)  
2  
3 freq(sw_df$sex)  
4  
5 ctable(sw_df$sex, sw_df$gender)
```



# DataExplorer

**DataExplorer package** - automates and streamlines the process of exploring and visualizing datasets.

```
1 install.packages("DataExplorer")  
2 library(DataExplorer)
```

```
1 create_report(sw_df)
```

