

# MANIPULATING DATA WITH DPLYR

BASED ON R-ECOLOGY LESSON - DATA CARPENTRY

Marco Chiapello, PhD

March 29, 2017



# WHAT IS DPLYR?

- **dplyr** is a package for making data manipulation easier
  - ▶ Packages in R are basically sets of additional functions that let you do more stuff

# WHAT IS DPLYR?

- **dplyr** is a package for making data manipulation easier
  - ▶ Packages in R are basically sets of additional functions that let you do more stuff
- dplyr provides **easy tools** for the most common **data manipulation** tasks

# WHAT IS DPLYR?

- **dplyr** is a package for making data manipulation easier
  - ▶ Packages in R are basically sets of additional functions that let you do more stuff
- dplyr provides **easy tools** for the most common **data manipulation** tasks
- dplyr addresses this by porting much of the computation to C++

# WHAT IS DPLYR?

- **dplyr** is a package for making data manipulation easier
  - ▶ Packages in R are basically sets of additional functions that let you do more stuff
- dplyr provides **easy tools** for the most common **data manipulation** tasks
- dplyr addresses this by porting much of the computation to C++
- An additional feature is the ability to work directly with data stored in an **external database**

# WHAT IS DPLYR?

**Before start to dig into dplyr functions we learn how to import data into R**

- To download the data, run the following:

```
download.file("https://ndownloader.figshare.com/files/2292169",  
              "portal_data_joined.csv")
```

- You are now ready to load the data:

```
surveys <- read.csv('portal_data_joined.csv')
```

# WHAT IS DPLYR?

**Before start to dig into dplyr functions we learn how to import data into R**

- To download the data, run the following:

```
download.file("https://ndownloader.figshare.com/files/2292169",  
             "portal_data_joined.csv")
```

- You are now ready to load the data:

```
surveys <- read.csv('portal_data_joined.csv')
```

- Converts data to tbl class. tbl's are easier to examine than data frames. R displays only the data that fits onscreen

```
surveys <- tbl_df(surveys)
```



# WHAT IS DPLYR?

We're going to learn some of the *most common dplyr functions*:

- **select**
- **filter**
- **arrange**
- **mutate**
- **group\_by**
- **summarize**

# SELECT

- To select columns of a data frame, use `select()`
  - ▶ The **first argument** to this function is the **data frame**
  - ▶ The **subsequent arguments** are the **columns to keep**

```
select(surveys, plot_id, species_id, weight)
```

# SELECT

- To select columns of a data frame, use `select()`
  - ▶ The **first argument** to this function is the **data frame**
  - ▶ The **subsequent arguments** are the **columns to keep**

```
select(surveys, plot_id, species_id, weight)
```

- `select` is much more powerful than just select the interest columns
  - ▶ You can remove one column

```
select(surveys, -weight)
```

# SELECT

- To select columns of a data frame, use `select()`
  - ▶ The **first argument** to this function is the **data frame**
  - ▶ The **subsequent arguments** are the **columns to keep**

```
select(surveys, plot_id, species_id, weight)
```

- `select` is much more powerful than just select the interest columns
  - ▶ You can remove one column

```
select(surveys, -weight)
```

\* Select columns whose name contains a character string

```
select(surveys, contains("ec"))
```

# SELECT

- Select columns whose name starts with a character string

```
select(surveys, starts_with("s"))
```

# SELECT

- Select columns whose name starts with a character string

```
select(surveys, starts_with("s"))
```

- Select all columns between Sepal.Length and Petal.Width (inclusive).

```
select(surveys, plot_id:weight)
```

# SELECT

- Select columns whose name starts with a character string

```
select(surveys, starts_with("s"))
```

- Select all columns between Sepal.Length and Petal.Width (inclusive).

```
select(surveys, plot_id:weight)
```

- Select every column

```
select(surveys, weight, everything())
```