

### SEMINARI 3.2. *Data Massaging* (Respostes)

### 1. OBJECTIUS

Aquest seminari introdueix al Data Massaging.

### 2. PART 1. Data Massaging

En aquesta primera part del seminari anem a fer alguns canvis en el dataset starwars que ja vam explorar en el seminari 3 de la llibreria tidyverse. Si heu obert R de nou, recordeu carregar la llibreria tidyverse.

Starwars és un dataset on cada fila és una observació i cada columna és una variable

```
> starwars
# A tibble: 87 x 14
                                                                         height mass hair color
                                                                                                                                                    homeworld species films
                                                                                                                                                                                              vehicles starships
                                                                                                                  <chr>
                                                                                                                                                                              <chr [5]> <chr [2]> <chr [2]>
                                                                        <chr>
fair
gold
white, blue
white
light
                                                                                                                19 male
112 none
 1 Luke Skywalker
                                                                                          blue
                                                                                          blue
yellow
red
yellow
                                                                                                                                                                              cchr [5]> cchr [2]> cchr
cchr [6]> cchr [0]> cchr
cchr [7]> cchr [0]> cchr
cchr [4]> cchr [0]> cchr
cchr [5]> cchr [1]> cchr
cchr [3]> cchr [0]> cchr
    C-3P0
                                                                                                                                                                   Droid
 2 C-3PO
3 R2-D2
4 Darth Vader
5 Leia Organa
6 Owen Lars
                                                                                                                                                                   Droid
Human
Human
                                            32 <NA>
136 none
49 brown
120 brown, grey
                                                                                          brown
blue
blue
red
                                                                         light
                                                                                                                                                                   Human
   Beru Whitesun lars
                                               75 brown
32 <NA>
                                                                         light
                                                                                                                                    feminine Tatooine masculine Tatooine
                                                                                                                                                                   Human
                                                                                                                                                                               <chr [3]> <chr [0]> <chr [0]>
<chr [1]> <chr [0]> <chr [0]>
 8 R5-D4
                                                                         white, red
                                                                                          red
                                                                                                                                                                   Droid
                                               84 black light
77 auburn, white fair
9 Biggs Darklighter
10 Obi-Wan Kenobi
                                                                                                                   24 male
57 male
                                                                                                                                    masculine Tatooine
masculine Stewjon
```

Primer ens familiaritzarem amb algunes funcions de la llibreria dyplr com:

- select: Seleccionar columnes
- filter: Filtrar
- arrange: Ordenar un conjunt de dades
- group by: Agrupar per alguna variable
- summarize/summarise: Especificar algunes funciones d'agregats:
  - o n(): comptar;
  - sum(): sumar variables numèriques;
  - o mean(): la mitjana de variables numèriques entre altres
- mutate: modificar, transformar o agregar variables del conjunt de dades
- pipes %>%: combinar operacions

### **EXERCICIS:**

En primer lloc, carreguem la llibreria dyplr

- > library(tidyverse)
- > library(dplyr)

### 1.- Seleccioneu només el nom i gènere del conjunt de dades starwars

Primer cridem les llibreries i desprès fem us de la funció select(), ja que estem seleccionant variables/columnes de la taula.

**UAB** 



9 Biggs Darklighter

# ... with 77 more rows

10 Obi-Wan Kenobi

>

```
select(.data, ...)
                Extract columns as a table. Also select_if().
                select(iris, Sepal.Length, Species)
   Use these helpers with select (),
   e.g. select(iris, starts_with("Sepal"))
  contains(match)
                    num_range(prefix, range) :, e.g. mpg:cyl
  ends_with(match) one_of(...)
                                            -, e.g, -Species
  matches(match)
                    starts with(match)
> select(starwars, name, gender)
# A tibble: 87 x 2
   name
                           gender
   <chr>
                           <chr>>
 1 Luke Skywalker
                           masculine
 2 C-3PO
                           masculine
 3 R2-D2
                           masculine
 4 Darth Vader
                           masculine
 5 Leia Organa
                           feminine
 6 Owen Lars
                           masculine
 7 Beru Whitesun lars feminine
 8 R5-D4
                           masculine
```

### 2.- Seleccioneu els personatges que són humans. Després seleccioneu els que no

En la columna *species*, trobem algunes files que són "Human". Per tant, seleccionem les observacions/files amb la funció *filter* com hem vist a classe

masculine

masculine

# Manipulate Cases EXTRACT CASES Row functions return a subset of rows as a new table. filter(.data, ...) Extract rows that meet logical criteria. filter(iris, Sepal.Length > 7) filter(starwars, species=="Human") #Or also: starwars %>% filter(species=="Human")

Al primer seminari vam veure que les cadenes de caràcters podien escriure's entre "Human" o 'Human'.





```
starwars %>% filter(species=="Human")
# A tibble: 35 x 14
                                                           skin color eye color birth year sex
   name
   <chr>
                           <int> <dbl> <chr
                                                           <chr>
                                                                                          <dbl> <chr>
                                                                                                         <chr>
                                                                                                                      <chr>
                                                                                                                                  <chr>
                                                                                                                                           st>
                                                                                                                                                       st>
                                                                                                                                                                    st>
 l Luke Skywalker
                             172
                                     77 blond
                                                           fair
                                                                       blue
                                                                                           19
                                                                                                male
                                                                                                         masculine Tatooine
                                                                                                                                 Human
                                                                                                                                           <chr [51> <chr [21> <chr [21>
 2 Darth Vader
3 Leia Organa
                                                                                                                                           <chr [4]> <chr
<chr [5]> <chr</pre>
                                                           white
                                                                        yellow
                                                                                                         masculine Tatooine
                                                                                                 female feminine Alderaan
male masculine Tatooine
                                    120 brown, grey
                                                           light
                                                                                                                                  Human
                                                                                                                                           <chr [3]> <chr
                                                                                                                                                                   <chr
                                                                                                 female feminine
 5 Beru Whitesun lars
                                      75 brown
                                                           light
                                                                        blue
                                                                                                                     Tatooine
                                                                                                                                  Human
                                                                                                                                           <chr [3]> <chr
                                                                                                                                                              [0]> <chr
 6 Biggs Darklighter
7 Obi-Wan Kenobi
8 Anakin Skywalker
9 Wilhuff Tarkin
                                      84 black
                                                                                                 male
                                                                                                         masculine Tatooine
                                                                                                                                           cohr [115 cohr
                                      77 auburn, whit
84 blond
                                                           fair
                                                                                                 male
                                                                                                         masculine Tatooine
                                                                                                                                  Human
                                                                                                                                           <chr [3]> <chr [2]> <chr 
<chr [2]> <chr [0]> <chr
                                                                                                                                                                         [3]>
                              180
                                      NA auburn, grey
                                                          fair
                                                                        blue
                                                                                                 male
                                                                                                         masculine Eriadu
                                                                                                                                  Human
10 Han Solo
                              180
                                      80 brown
                                                           fair
                                                                                                 male
                                                                                                         masculine Corellia
                                                                                                                                 Human
                                                                                                                                           <chr [41> <chr [01> <chr [21>
   .. with 25 more rows
```

Per veure els que no, utilitzem la negació en R!

```
Logical and boolean operators to use with filter()
                                                %in%
      <
                    <=
                                  is.na()
                                                                            xor()
                                               !
                   >=
                                  !is.na()
    See ?base::Logic and ?Comparison for help.
> filter(starwars, !species=="Human")
   starwars %>% filter(!species=="Human") #alternatively
   #or:
   starwars %>% filter(species!="Human")
  #or:
> filter(starwars, species!="Human")
                     height mass hair color skin color
                                                       eye color birth year sex
                                                                                      gender homeworld species
                                                                                                                           vehicles starships
                     __year sex
<dbl> <chr>
112 none
33 none
NA none
200 male
44 male
600 herman
896 male
   <chr>
                                                                                                                           <chr [0]> <chr [0]>
<chr [0]> <chr [0]>
                                                                                      <chr> <chr> <chr> masculine Tatooine Droid
                                                                                                                   t>
                                                       yellow
                                         gold
  C-3F0
R2-D2
R5-D4
Chewbacca
Greedo
Jabba Desilijic Tiure
                                                                                      masculine <NA> Droid
masculine Trandosha Trandoshan
  Bossk
                                                       orange
                                                                                      masculine Mon Cala Mon Calamari
```

- Seleccioneu només el nom i gènere dels personatges amb gènere femení
  - a) Pas a pas: Assigneu primer a una variable temporal (per exemple de nom dadesub) la selecció feta en l'exercici 1. Després seleccioneu del subconjunt dadesub les observacions/files que corresponen al gènere femení i assigneu el subconjunt resultant a una altra variable (per exemple de nom femenins)
  - b) Feu ús de les pipes que hem vist a classe, per tal de combinar les dues operacions fetes en l'apartat a. Assigneu el subconjunt resultant a una variable (per exemple de nom feminisub) i verifiqueu que correspon al mateix subconjunt assignat a femenins en l'apartat a). NOTA: Sempre que puguem utilitzarem pipes com en aquest apartat, és a dir, no crearem variables temporals extres com em fet en l'apartat a).



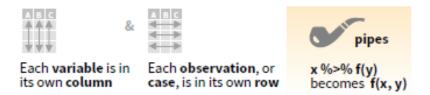
**UAB** 

a) Assignem a la variable temporal dadesub el subconjunt de dades resultant de la selecció feta en l'exercici 1. Desprès seleccionem les observacions/files amb la funció filter com hem vist a l'exercici 2. Fixem-nos que filter selecciona files/observacions que satisfan un cert criteri lògic (en aquest cas seria com pensar if gender='feminine')

```
> dadesub <- select(starwars, name, gender)</pre>
> femenins <- filter(dadesub, gender=='femenine')</pre>
> femenins
> femenins <- filter(dadesub,gender=='feminine')
> femenins
# A tibble: 17 x 2
   name
                      gender
   <chr>
                      <chr>
 l Leia Organa
                    feminine
 2 Beru Whitesun lars feminine
 3 Mon Mothma
 3 Mon room.....
4 Shmi Skywalker feminine feminine
                    feminine
 6 Adi Gallia
                    feminine
 7 Cordé
                     feminine
 8 Luminara Unduli feminine
 9 Barriss Offee
                     feminine
10 Dormé
                     feminine
11 Zam Wesell
                    feminine
12 Taun We
                     feminine
13 Jocasta Nu
                     feminine
                     feminine
14 R4-P17
15 Shaak Ti
                     feminine
16 Rey
                     feminine
17 Padmé Amidala
                    feminine
```

b. Hem vist que podíem seleccionar variables (columnes) i observacions (o files) fent us de *pipes* seguint:

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



Per tant combinem la selecció i el filtratge utilitzant *pip*es de la següent manera:

```
> feminisub <- starwars %>% filter(gender=='feminine')
%>%select(name,gender)
> feminisub
```

En aquest cas, podem intercanviar el filtratge i la selecció obtenint el mateix resultat:

```
> starwars %>%select(name,gender)%>%filter(gender=='feminine')
```





```
> feminisub <- starwars %>% filter(gender=="feminine")%>%select(name,gender)
> feminisub
# A tibble: 17 x 2
   name
                        gender
   <chr>
                         <chr>
 l Leia Organa
                       feminine
 2 Beru Whitesun lars feminine
 3 Mon Mothma feminine
4 Shmi Skywalker feminine
 5 Ayla Secura feminine
6 Adi Gallia feminine
 6 Adi Gallia
 7 Cordé feminine
8 Luminara Unduli feminine
9 Barriss Offee feminine
9 Barriss Offee
10 Dormé
                        feminine
11 Zam Wesell
                        feminine
12 Taun We
                       feminine
                       feminine
feminine
13 Jocasta Nu
14 R4-P17
15 Shaak Ti
                       feminine
16 Rey
17 Padmé Amidala
                        feminine
>
```

4.- Ordeneu els noms dels personatges femenins en ordre descendent. Podríeu fer-ho sense fer us de les variables temporals creades en exercici anterior (és a dir, sense fer ús del subgrup)?

Hem vist:

## ARRANGE CASES arrange(.data, ...) Order rows by values of a column or columns (low to high), use with desc() to order from high to low. arrange(mtcars, mpg) arrange(mtcars, desc(mpg))

### Per tant:

> arrange(feminisub, desc(name)) #partint del subgrup creat

```
> arrange(feminisub,desc(name))
# A tibble: 17 x 2
  name
                     gender
   <chr>
                     <chr>
 <chr>
1 Zam Wesell feminine
2 Taun We
 2 Taun We feminine
3 Shmi Skywalker feminine
                    feminine
 4 Shaak Ti
 5 Rey
                     feminine
 6 R4-P17
                     feminine
                    feminine
feminine
 7 Padmé Amidala
 8 Mon Mothma
 9 Luminara Unduli feminine
10 Leia Organa feminine
12 Dormé
13 Cordé
                      feminine
14 Beru Whitesun lars feminine
15 Barriss Offee feminine
16 Ayla Secura
                      feminine
17 Adi Gallia
                      feminine
>
```

Per tal de no fer us del subgrup feminisub, altre cop hauríem d'utilitzar les pipes:





```
> ordered_fem <- starwars %>%filter(gender=="feminine")
%>%select(name,gender) %>%arrange(desc(name))
> ordered_fem
> ordered_fem <- starwars %>% filter(gender=="feminine")%>%select(name,gender)%>%arrange(desc(name))
 # A tibble: 17 x 2
                      gender
   <chr>
                      <chr>
 1 Zam Wesell
                     feminine
 2 Taun We feminine
3 Shmi Skywalker feminine
 4 Shaak Ti
                     feminine
feminine
 5 Rey
 6 R4-P17
                      feminine
 7 Padmé Amidala feminine
 8 Mon Mothma
                      feminine
 9 Luminara Unduli feminine
10 Leia Organa
                      feminine
                 feminine
11 Jocasta Nu
12 Dormé
                      feminine
13 Cordé
                      feminine
14 Beru Whitesun lars feminine
15 Barriss Offee feminine
16 Ayla Secura feminine
17 Adi Gallia feminine
```

Judit Chamorro Servent Bellaterra, Març 2025