

OPSO79-1-UCSH2021

Corrección prueba 2, Estadística
descriptiva II y transformación
avanzada de data frames

29/10/2021

Transformación avanzada de datos

Un poco más de agrupación, pivotear y combinar data frames

Introducción

La sesión subsiguiente veremos en detalle como elaborar gráficos elegantes en R.

Antes es necesario revisar algunos últimos aspectos sobre transformación de datos.

La clave para elaborar buenos gráficos en R es tener una data frame coherente con el gráfico que queremos

Por ejemplo,

- si queremos graficar N de hogares por región, no nos servirá una base de datos de personas.
- si queremos graficar mediante barras el porcentaje de personas que reciben mas y menos del sueldo mínimo, la variable numérica salario debe ser categorizada
- Si queremos graficar 2 variables, distinguiendo la relación por una tercera, necesitamos tener una base en formato *longer* (hacia abajo), no *wider* (hacia el lado)

Introducción

A continuación veremos herramientas que nos permitirán lidiar con estos y otros problemas:

- Funciones de agrupación (`group_by()`, `summarise()`).
- Funciones para pivotar la data (`pivot_longer()`, `pivot_wider()`).
- Funciones para combinar data (`merge()`, `rbind()` y `cbind()`)

Aplicaremos estas funciones a los datos del paquete `Gapminder`, a datos del Banco Mundial (de donde venían los de Afganistán), entre otros.

Agrupación de datos

profundización función `group_by()`

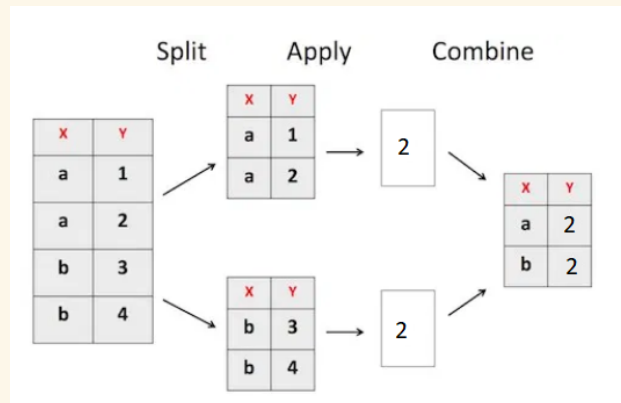
group_by() y summarise()

En conjunto nos permiten resumir información para cada grupo de una variable

Podemos obtener edad promedio por sexo, número de personas en cada región, ingresos por hogar, etcétera.

Estrategia **split-apply-combine**.

Esta estrategia sucede tras bambalinas (no la vemos). Solo observamos el resultados.



group_by() y summarise()

```
data <- readRDS("data/Latinobarometro_2020_Esp_Rds_v1_0.rds")
```

Conteo de frecuencias

```
data %>% group_by(sexo) %>% summarise(n=n())
```

```
## # A tibble: 2 x 2
##   sexo      n
##   <hvn_lbl> <int>
## 1 1      9667
## 2 2     10537
```

Obtención de estadísticos para cada grupo

```
data %>% group_by(sexo) %>% summarise(edad=mean(edad))
```

```
## # A tibble: 2 x 2
##   sexo      edad
##   <hvn_lbl> <dbl>
## 1 1      41.6
```


group_by() y mutate()

Con summarise "perdemos" la data original. Esta es resumida a una más pequeña.

```
data2 <- data %>% group_by(sexo) %>% summarise(edad=mean(edad))  
dim(data2)
```

```
## [1] 2 2
```

```
dim(data)
```

```
## [1] 20204 408
```

Pero en ocasiones queremos una medida de resumen sin perder la data, para poder generar nuevos cálculos.

La alternativa es **agrupar** sin resumir, sino que **mutando** la data.

```
data2 <- data %>% group_by(sexo) %>% mutate(edad_promedio=mean(edad))  
dim(data2)
```

```
## [1] 20204 409
```

group_by() y mutate()

Veamos un pedazo de la nueva data

```
data2 %>% select(idenpa,sexo,edad,edad_promedio) %>% head()
```

```
## # A tibble: 6 x 4
## # Groups:   sexo [2]
##   idenpa      sexo      edad  edad_promedio
##   <hvn_lbl> <hvn_lbl> <hvn_lbl>      <dbl>
## 1 32         2      63      40.4
## 2 32         1      24      41.6
## 3 32         1      20      41.6
## 4 32         2      54      40.4
## 5 32         1      38      41.6
## 6 32         2      62      40.4
```

Edad promedio aparece en cada observación.

Es el promedio de la edad del grupo (sexo) al que pertenece la observación.

En este caso solo hay valores 41,6 (para los hombres) y 40,4 (para las mujeres)

group_by() y mutate()

¿Cuál es la utilidad?

Sirve para el procesamiento de datos más que para el análisis.

Por ejemplo, identificar casos extraños dentro de un conjunto para luego editarlos.

Países que pertenecen a continentes pobres pero que son **MUY** ricos:

```
library(gapminder)

países_1972 <- gapminder %>%
  filter(year==1972 ) %>%
  group_by(continent) %>%
  mutate(gdpPercap_continente=quantile(gdpPercap,0.90)) %>%
  ungroup()
```

group_by() y mutate()

```
países_1972 %>%  
  filter(continent %in% c("Africa", "Americas") &  
         gdpPercap > gdpPercap_continente) %>%  
  arrange(-gdpPercap) %>% select(-year, -continent)
```

country	lifeExp	pop	gdpPercap	gdpPercap_continente
United States	71.340	209896000	21806.036	10080.371
Libya	52.773	2183877	21011.497	4139.005
Canada	72.880	22284500	18970.571	10080.371
Gabon	48.690	537977	11401.948	4139.005
Venezuela	65.712	11515649	10505.260	10080.371
South Africa	53.696	23935810	7765.963	4139.005
Angola	37.928	5894858	5473.288	4139.005
Reunion	64.274	461633	5047.659	4139.005
Algeria	54.518	14760787	4182.664	4139.005

Pivotear los datos

funciones `pivot_wider()` y `pivot_longer()`

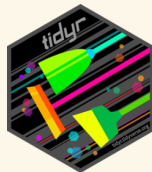
Pivotear los datos

Alargamiento o ensanchamiento de una data frame.

Alargamiento: incremento en el número de filas y decrecimiento del número de columnas

Ensanchamiento: incremento en el número de columnas y decrecimiento del número de filas

Para esto utilizaremos las funciones `pivot_wider()` y `pivot_longer()` del paquete `tidyr`



wide				vs	long		
					ID	ID2	A
					1	a1	
					2	a1	
					3	a1	
ID	a1	a2	a3		1	a2	
					2	a2	
					3	a2	
					1	a3	
					2	a3	
					3	a3	

Función pivot_wider()

Esta función se utiliza para ordenar un dataframe de forma tal de mostrar categorías de una variable como columnas de un dataframe.

Incrementa el número de las columnas y disminuye el número de las filas.

Es útil para la presentación de cuadros de resumen con doble entrada.

sexo	posicion_politica	n
1	centro	3740
1	der	2161
1	izq	2245
1	ninguna	979
2	centro	3582
2	der	2187
2	izq	2327
2	ninguna	1415

Función pivot_wider()

Ahora vemos las categorías de sexo hacia la derecha

posicion_politica	1	2
centro	3740	3582
der	2161	2187
izq	2245	2327
ninguna	979	1415

Pasamos de un formato largo a uno ancho

```
library(tidyr)
```

```
data %>%  
  filter(!is.na(posicion_politica)) %>%  
  group_by(sexo, posicion_politica) %>%  
  summarise(n=n()) %>%  
  pivot_wider(names_from = sexo,  
              values_from = n)
```


Función pivot_wider()

Básicamente dos argumentos:

- *names_from*: categorías que se quiere convertir en columnas
- *values_from*: columna desde la cual extraer los valores

Además, podemos usar el argumento *names_prefix* cuando tenemos números

```
data %>% filter(!is.na(posicion_politica)) %>%  
  group_by(sexo, posicion_politica) %>%  
  summarise(n=n()) %>%  
  pivot_wider(names_from = sexo,  
              values_from = n,  
              names_prefix = "sexo_")
```

```
## # A tibble: 4 x 3  
##   posicion_politica sexo_1 sexo_2  
##   <chr>             <int> <int>  
## 1 centro           3740  3582  
## 2 der              2161  2187  
## 3 izq             2245  2327  
## 4 ninguna          979  1415
```

Función pivot_longer()

Esta función se puede considerar como la opuesta a pivot_wider().

Esta función incrementa el número de filas y disminuye el número de columnas.

Los dataframes obtenidos por esta función son más fáciles de manipular, pero no de visualizar

```
df1 <- data.frame(region = c(1, 2),  
                  hombres = c(100, 200),  
                  mujeres = c(50, 300))
```

```
df1
```

```
##   region hombres mujeres  
## 1      1     100      50  
## 2      2     200     300
```

Función pivot_longer()

```
df1 %>%  
  pivot_longer(cols = c(hombres,mujeres) )
```

```
## # A tibble: 4 x 3  
##   region name      value  
##   <dbl> <chr>    <dbl>  
## 1      1 hombres    100  
## 2      1 mujeres     50  
## 3      2 hombres    200  
## 4      2 mujeres    300
```

El argumento principal es cols:

- *cols*: columnas a las que se le aplicará la operación (que se convertirán en categorías de una nueva variable)

Función pivot_longer()

Además, se pueden especificar los nombres de las columnas "name" y "value"

- *names_to*: indica el nombre de la variable que será creada para "guardar" los nombres de las categorías.
- *values_to*: indica el nombre de la variable que será creada para "guardar" los valores asociados a las categorías.

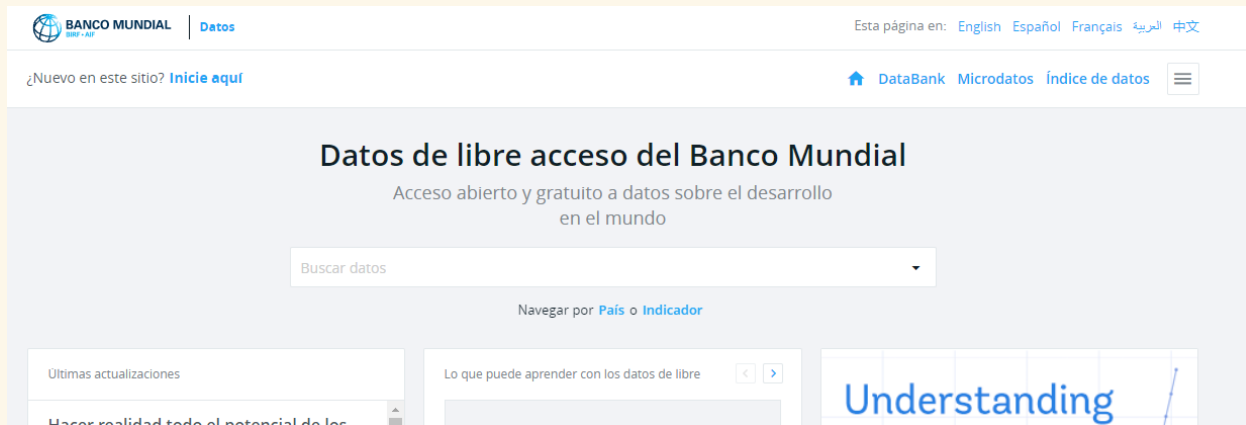
```
df1 %>%  
  pivot_longer(cols = c(hombres,mujeres) ,  
               names_to = "sexo", values_to = "total_sexo")
```

```
## # A tibble: 4 x 3  
##   region sexo    total_sexo  
##   <dbl> <chr>      <dbl>  
## 1     1  hombres      100  
## 2     1  mujeres       50  
## 3     2  hombres      200  
## 4     2  mujeres      300
```

Pivotear los datos

Relevante para visualizar (ggplot2) y para trabajar datos importados

Por ejemplo, descarguemos los datos de Afganistán que usamos clases atrás. Esta vez sin trampa.



Pivotear los datos

¿Cómo vienen los datos?

```
afganistan <- readxl::read_excel("data/API_AFG_DS2_es_excel_v2_3162018.")  
skip = 3)
```

Country Code	2005	2012
AFG	NA	NA
AFG	3.910455	0.5579121
AFG	210.115520	771.2416744
AFG	NA	1.9269141
AFG	NA	0.0000000

¡Las variables vienen como filas!

Pivotear los datos

La data no es un dato ordenado (tidy data)

¿Como graficamos el PIB de Afganistan si no es una variable? Solo podemos tabular años, lo que no tiene sentido:

```
table(afganistan$`1962`)
```

Pivotear los datos

La solución es pivotear los datos. Hacer que las filas pasen a ser variables.

```
# Alargar la data
```

```
afganistan <- afganistan %>% pivot_longer(5:ncol(afganistan)) %>%  
  select(-`Country Name`, `Country Code`, `Indicator Code`)
```

```
# Quitar filas repetidas para evitar errores
```

```
afganistan <- afganistan %>%  
  distinct(`Indicator Name`, value, name)
```

```
# Ensanchar la data
```

```
afganistan <- afganistan %>%  
  tidyr::pivot_wider(names_from = `Indicator Name`,  
                     values_from = value,  
                     values_fn = {sum})
```

```
# Limpiar los nombres
```

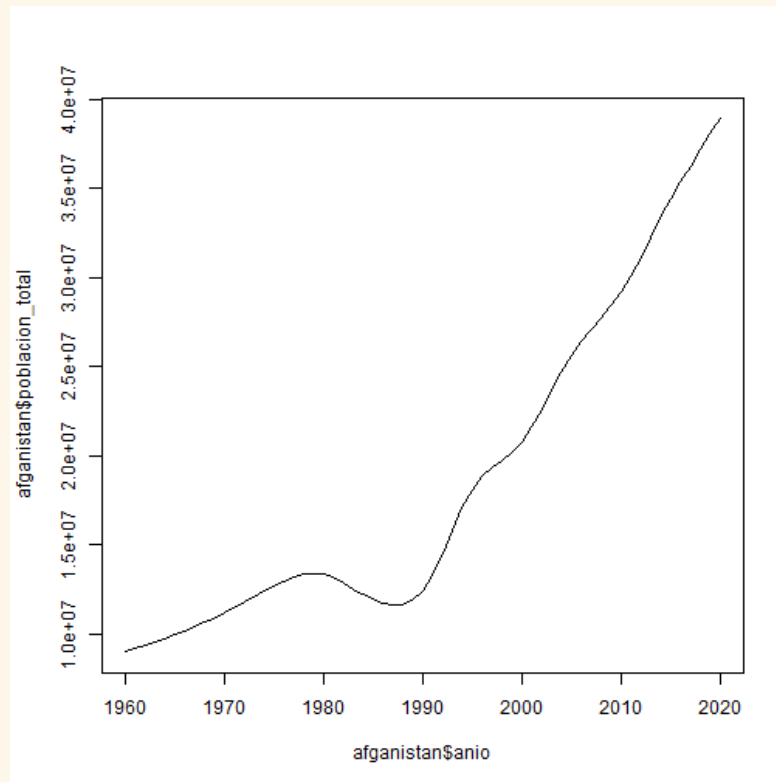
```
afganistan <- afganistan %>%  
  janitor::clean_names() %>% rename(anio=name)
```


Pivotear los datos

anio	ingreso_nacional_bruto_ing_us	poblacion_total
2012	20033093818	31161378
2013	20632806188	32269592
2014	20482514566	33370804
2015	20087077459	34413603
2016	18197299091	35383028
2017	19118263186	36296111
2018	18544615040	37171922
2019	19598008726	38041757
2020	19996141020	38928341

Pivotear los datos

```
plot(afganistan$anio,afganistan$poblacion_total,type = "l")
```



Recursos web utilizados

Xaringan: [Presentation Ninja](#), de Yihui Xie. Para generar esta presentación.

Ilustraciones de Allison Horst

Para reforzar y seguir aprendiendo

Video "[el juego de las estadísticas](#)" (utiliza gapminder)