## Más funciones de dplyr Sesión 2

Natalie Julian - www.nataliejulian.com

Estadística UC y Data Scientist en Zippedi Inc.

#### Datos mtcars

Trabajaremos en esta sesión complementaria con los datos mtcars. Para acceder a estos basta con utilizar:

data(mtcars)

Y se guardará un objeto de tipo dataframe llamado mtcars:

class(mtcars)
[1] "data.frame"

1 | 15

## count y n(): Contando casos

#### Contar casos

La variable vs de los datos mtcars toma los valores 1 (si el motor del auto es recto) y 0 (si el motor del auto es en forma de V). ¿Cómo podríamos contar cuántos casos por tipo de motor hay?

Usualmente podemos utilizar la función table() que entrega una tabla de frecuencias de la variable vs:

```
table(mtcars$vs)
```

0 1

#### Contar casos

#### Con dplyr podemos utilizar las funciones n() y count():

```
mtcars%>%
  group_by(vs)%>%
  count()
# A tibble: 2 x 2
    ٧S
        n
  <dbl> <int>
        14
mtcars%>%
  group_by(vs)%>%
  summarise(n=n())
# A tibble: 2 x 2
    ٧S
  <dbl> <int>
         18
          14
```

3 | 15



## Group by

Cuando agrupamos, las estadísticas se calculan en base a esta agrupación, por ejemplo:

```
mtcars%>%
           group by(vs)%>%
           mutate(prom=mean(disp)) #Calcula promedio de disp para vs=0 y vs=1
 # A tibble: 32 x 12
 # Groups:
                                                                        vs [2]
                                                                 cyl disp
                                                                                                                                               hp
                                                                                                                                                                        drat
                                                                                                                                                                                                                                                                                                                                                                                                 carb
                              mpg
                                                                                                                                                                                                                       wt
                                                                                                                                                                                                                                               gsec
                                                                                                                                                                                                                                                                                                 VS
                                                                                                                                                                                                                                                                                                                                     am gear
                  <dbl> 
                     21
                                                                                               160
                                                                                                                                           110
                                                                                                                                                                   3.9
                                                                                                                                                                                                           2.62 16.5
                                                                                                                                                                                                                                                                                                                                                                                                                                     307.
                                                                              6
                     21
                                                                                              160
                                                                                                                                          110
                                                                                                                                                                  3.9
                                                                                                                                                                                                            2.88 17.0
                                                                                                                                                                                                                                                                                                                                                                                                                   4 307.
                   22.8
                                                                             4 108
                                                                                                                                               93 3.85
                                                                                                                                                                                                          2.32 18.6
                                                                                                                                                                                                                                                                                                                                                                                                                                    132.
      4 21.4
                                                                              6 258
                                                                                                                                          110 3.08 3.22 19.4
                                                                                                                                                                                                                                                                                                                                                                                                                   1 132.
```

Pero, ¿qué pasa si quisiera ahora añadir una columna con el promedio de wt de todos los registros?

## Luego de group by, R piensa en ese nivel de agregación

Si sólo añadimos mutate(promwt=mean(wt)) no obtendremos lo que queremos, pues se está respetando la agrupación anterior:

```
mtcars%>%
 group_by(vs)%>%
 mutate(prom=mean(disp)) %>% #Calcula promedio de disp para vs=0 v vs=1
 mutate(promwt=mean(wt)) #Calcula promedio de wt para vs=0 y vs=1
# A tibble: 32 x 13
# Groups:
         vs [2]
         cyl disp
                       drat
    mpg
                   hp
                             wt
                                 gsec
                                        VS
                                             am gear
                                                     carb prom promwt
  <dh1>
  2.1
             160
                   110
                      3.9
                            2.62
                                16.5
                                                          307.
                                                                3.69
          6
  21
          6 160
                   110
                      3.9
                            2.88 17.0
                                                       4 307.
                                                                3.69
  22.8
          4 108
                   93 3.85
                            2.32 18.6
                                                          132.
                                                                2.61
4 21.4
          6 258
                   110
                      3.08 3.22 19.4
                                                          132.
                                                                2.61
```

### Efecto de ungroup

Si utilizamos ungroup ya no tendremos este problema:

```
mtcars%>%
         group_by(vs)%>%
        mutate(prom=mean(disp)) %>% #Calcula promedio de disp para vs=0 y vs=1
        ungroup()%>%
        mutate(promwt=mean(wt)) #Calcula promedio de wt para todas las observaciones
# A tibble: 32 x 13
                                               cyl disp
                                                                                                       hp drat
                                                                                                                                                                                                                                                                                      carb
                     mpg
                                                                                                                                                            wt qsec
                                                                                                                                                                                                                VS
                                                                                                                                                                                                                                           am gear
                                                                                                                                                                                                                                                                                                               prom promwt
             <dbl> 
              2.1
                                                                  160
                                                                                                    110 3.9
                                                                                                                                                   2.62 16.5
                                                                                                                                                                                                                                                                                                                 307.
                                                                                                                                                                                                                                                                                                                                               3.22
                                                                                                                                                   2.88 17.0
                                                                                                                                                                                                                                                                                                                                         3.22
              21
                                                                160
                                                                                           110 3.9
                                                                                                                                                                                                                                                                                                   4 307.
            22.8
                                       4 108
                                                                                                       93 3.85
                                                                                                                                                  2.32 18.6
                                                                                                                                                                                                                                                                                                   1 132.
                                                                                                                                                                                                                                                                                                                                               3.22
            21.4
                                                                                                                                                                                                                                                                                                                                        3.22
                                                       6 258
                                                                                          110 3.08 3.22 19.4
                                                                                                                                                                                                                                                                                                                132.
     5 18.7
                                                         8 360
                                                                                                   175 3.15 3.44 17.0
                                                                                                                                                                                                                                                                                                                 307.
                                                                                                                                                                                                                                                                                                                                               3.22
```

¿Qué es lo que hacemos con ungroup? De cierta forma, se divide la tubería en dos caminos o subtuberías, la primera rama respeta la agrupación y obtiene resultados en base a la agrupación, y la segunda rama considera los datos no agrupados.

sample\_n y sample\_frac: Muestrear

#### sample\_n

Supongamos que nos interesa muestrear 10 observaciones de los datos, esto es bastante sencillo:

```
mtcars%>%
sample_n(10)
```

También podemos muestrear 10 observaciones luego de realizar una agrupación:

```
mtcars%>%
  group_by(vs)%>%
  sample_n(10)
```

Incluso, podemos asignar prioridad (o pesos) a cada observación dependiendo de una variable en particular:

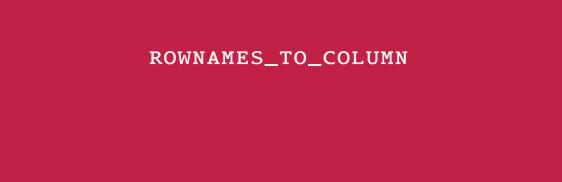
```
mtcars%>%
  group_by(vs)%>%
  sample_n(10, weight = 1/gear)
```

## sample\_frac

Pero...siendo bastante realistas, no siempre tenemos la misma cantidad de observaciones por grupo, por lo cual, dejar un número fijo para muestrear quizás no sea siempre la mejor opción.

Supongamos que en realidad, lo que necesitamos es muestrear el 30% de las observaciones por grupo (es decir, la cantidad de observaciones muestreadas por grupo será proporcional a la cantidad de observaciones totales en cada grupo). Esto lo podemos realizar fácilmente, indicando la fracción que queremos por grupo:

```
mtcars%>%
           group_by(vs)%>%
           sample_frac(0.3, weight = 1/gear)
# A tibble: 9 x 11
# Groups:
                                                                     vs [2]
                                                       cyl disp
                                                                                                                                     hp drat
                                                                                                                                                                                                                                                                                                                                                                                   carb
                                                                                                                                                                                                               wt asec
                                                                                                                                                                                                                                                                                                                                                 gear
           <dbl> 
1 10.4
                                                                       8 460
                                                                                                                                 215
                                                                                                                                                                                                    5 42 17 8
2 18 7
                                                                        8 360
                                                                                                                                                         3.15 3.44 17.0
           19.7
                                                                       6 145
                                                                                                                                                           3.62 2.77
               21
                                                                       6 160
                                                                                                                                                                                                  2.62 16.5
                                                                                                                                                           3.9
               21
                                                                       6 160
                                                                                                                                 110 3 9
                                                                                                                                                                                                                                  17 0
               22.8
                                                                       4 141.
                                                                                                                                       95 3.92 3.15 22.9
               21.5
                                                                     4 120.
               27.3
                                                                        4 79
           32.4
                                                                       4 78.7
                                                                                                                                                                                                                                       19.5
```



#### Filas con nombres

Seguramente ya lo notaste, los datos mtcars corresponden a autos y el modelo de cada auto está como nombre de fila (no explícitamente como columna). Podría ser útil tener esta información como variable. ¿Cómo lograrlo? Con la función de tidyverse rownames\_to\_column:

```
(mtcars<-mtcars%>%
 rownames_to_column())
                      mpg cyl disp hp drat wt qsec vs am gear
              rowname
            Mazda RX4 21.0
                            6 160.0 110 3.90 2.620 16.46
                            6 160.0 110 3.90 2.875 17.02 0 1
        Mazda RX4 Wag 21.0
           Datsun 710 22.8
                                    93 3.85 2.320 18.61 1 1
       Hornet 4 Drive 21.4
                            6 258.0 110 3.08 3.215 19.44 1
    Hornet Sportabout 18.7
                            8 360.0 175 3.15 3.440 17.02
6
              Valiant 18.1
                            6 225.0 105 2.76 3.460 20.22 1
           Duster 360 14.3
                            8 360.0 245 3.21 3.570 15.84
            Merc 240D 24.4
                            4 146.7 62 3.69 3.190 20.00
                                                           0
             Merc 230 22.8
                            4 140.8 95 3.92 3.150 22.90
                                                                4
```

## ¿De qué me podría servir tener el modelo del auto?

Por ejemplo, quizás nos interesa filtrar por ciertas marcas. Por ejemplo, supongamos nos interesan los autos de la marca Toyota o Mazda. Extraemos los registros respectivos a los autos de la siguiente forma:

6 160.0 110 3.90 2.875 17.02 0 1

4 71.1 65 4.22 1.835 19.90 1 1

4 120.1 97 3.70 2.465 20.01 1

mtcars%>%

Mazda RX4 Wag 21.0 Toyota Corolla 33.9

Toyota Corona 21.5

1

## recode: Recodificar variables

#### recode

Recordemos que la variable vs toma los valores 1 si el motor del auto es recto y 0 si el motoro del auto es en forma de V. Podríamos recodificar directamente la variable de 1 y 0 a Motor recto y Motor forma V respectivamente. Podríamos hacerlo con ifelse (¿Recuerdas que lo vimos en el curso R basics?):

```
mtcars%>%
  mutate(vs=ifelse(vs==1, "Motor Recto", "Motor Forma V"))
```

#### O también con recode:

### case\_when

#### Y también podríamos utilizar case\_when:

```
mtcars%>%
 mutate(vs=case_when(vs==1 ~ "Motor Recto", TRUE~ "Motor Forma V"))
              rowname mpg cyl disp hp drat
                                              wt qsec
                                                                 vs am gear carb
            Mazda RX4 21.0 6 160.0 110 3.90 2.620 16.46 Motor Forma V 1
        Mazda RX4 Wag 21.0 6 160.0 110 3.90 2.875 17.02 Motor Forma V 1
           Datsun 710 22.8 4 108.0 93 3.85 2.320 18.61 Motor Recto 1
       Hornet 4 Drive 21.4 6 258.0 110 3.08 3.215 19.44
                                                        Motor Recto 0
    Hornet Sportabout 18.7 8 360.0 175 3.15 3.440 17.02 Motor Forma V 0
6
             Valiant 18.1
                           6 225.0 105 2.76 3.460 20.22
                                                        Motor Recto 0
           Duster 360 14.3 8 360.0 245 3.21 3.570 15.84 Motor Forma V 0
```

Entregando los mismos resultados! ¿Fácil no?

# cut: Agrupar variables numéricas

#### cut

Supongamos que queremos categorizar la variable hp como sigue:

- Si hp  $\in$  (50, 122] indique *bajo*
- si hp  $\in$  (122, 180] indique *alto*
- si hp > 180 indique *potente*

Esto se puede hacer facilmente con la función cut:

```
mtcars%>%
 mutate(categoriahp = cut(hp.
                           breaks = c(50, 122, 180, Inf),
                           labels = c("Bajo", "Alto", "Potente"),
                           right = TRUE))
              rowname
                       mpg cyl disp hp drat
                                                 wt gsec vs am gear carb categoriahp
            Mazda RX4 21.0
                             6 160.0 110 3.90 2.620 16.46 0
                                                                                 Bajo
         Mazda RX4 Wag 21.0 6 160.0 110 3.90 2.875 17.02 0
                                                                                 Bajo
           Datsun 710 22.8
                             4 108.0 93 3.85 2.320 18.61
                                                                                 Bajo
       Hornet 4 Drive 21.4
                             6 258.0 110 3.08 3.215 19.44
                                                                                 Bajo
    Hornet Sportabout 18.7
                             8 360.0 175 3.15 3.440 17.02
                                                                                 Alto
6
              Valiant 18.1
                             6 225.0 105 2.76 3.460 20.22 1
                                                                                 Baio
           Duster 360 14.3
                             8 360.0 245 3.21 3.570 15.84
                                                                              Potente
            Merc 240D 24.4
                             4 146.7 62 3.69 3.190 20.00
                                                                                 Bajo
```

## ¿Y si no tengo las categorías? ¿Puedo obtenerlas?

Así es! Puedes utilizar la función cut\_width y definir cuál es el ancho del intervalo deseado:

```
(mtcars<-mtcars%>%
 mutate(categoria2hp=cut width(hp. width=60)))
               rowname
                        mpg cyl
                                 disp hp drat
                                                   wt
                                                       gsec vs am gear carb categoria2hp
             Mazda RX4 21.0
                              6 160.0 110 3.90 2.620 16.46 0
                                                                                 (90.150)
2
                              6 160.0 110 3.90 2.875 17.02
         Mazda RX4 Wag 21.0
                                                                                 (90,150]
            Datsun 710 22.8
                              4 108.0 93 3.85 2.320 18.61
                                                                                 (90, 150]
4
        Hornet 4 Drive 21.4
                              6 258.0 110 3.08 3.215 19.44
                                                                                 (90.1507
5
    Hornet Sportabout 18.7
                              8 360.0 175 3.15 3.440 17.02
                                                                                (150,210)
6
               Valiant 18.1
                              6 225.0 105 2.76 3.460 20.22
                                                                                 (90.1507
            Duster 360 14.3
                              8 360.0 245 3.21 3.570 15.84
                                                                           4
                                                                                (210, 270)
8
             Merc 240D 24.4
                                       62 3.69 3.190 20.00
                                                                     4
                                                                                  [30.90]
              Merc 230 22.8
                                      95 3.92 3.150 22.90
                                                                                 (90.1507
10
              Merc 280 19.2
                              6 167.6 123 3.92 3.440 18.30
                                                                                 (90.150]
11
             Merc 280C 17.8
                              6 167.6 123 3.92 3.440 18.90
                                                                           4
                                                                                 (90.150)
```

## Y podemos filtrar por esta nueva categoría:

```
mtcars%>%
  filter(categoria2hp %in% c("[30,90]", "(150,210]"))
                     mpg cyl disp hp drat
                                                wt gsec vs am gear carb categoria2hp
             rowname
                             8 360.0 175 3.15 3.440 17.02
   Hornet Sportabout 18.7
                                                                             (150,210)
            Merc 240D 24.4
                             4 146.7 62 3.69 3.190 20.00
                                                                               [30,90]
          Merc 450SE 16.4
                             8 275.8 180 3.07 4.070 17.40
                                                                             (150.210)
4
          Merc 450SL 17.3 8 275.8 180 3.07 3.730 17.60
                                                                             (150,210)
          Merc 450SLC 15.2
                             8 275.8 180 3.07 3.780 18.00
                                                                             (150, 210]
   Cadillac Fleetwood 10.4
                             8 472.0 205 2.93 5.250 17.98
                                                                             (150, 210]
             Fiat 128 32.4
                                     66 4.08 2.200 19.47
                             4 78.7
                                                                               [30,90]
          Honda Civic 30.4
                             4 75.7 52 4.93 1.615 18.52
                                                                               [30.90]
                                      65 4.22 1.835 19.90
9
      Toyota Corolla 33.9
                                                                               [30.90]
    Pontiac Firebird 19.2
10
                             8 400.0 175 3.08 3.845 17.05
                                                                             (150, 210]
11
            Fiat X1-9 27.3
                             4 79.0 66 4.08 1.935 18.90
                                                                               [30,90]
12
        Ferrari Dino 19.7
                             6 145.0 175 3.62 2.770 15.50
                                                                        6
                                                                             (150,210)
```