



# Introducción a R

## Reshape Data - change the layout of values in a table

Use **gather()** and **spread()** to reorganize the values of a table into a new layout. Each uses the idea of a key column: value column pair.

**gather()**(data, key, value, ..., na.rm = FALSE,  
convert = FALSE, factor\_key = FALSE)

Gather moves column names into a key column, gathering the column values into a single value column.

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K

→

country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

key value

*gather(table4a, `1999`, `2000`,  
key = "year", value = "cases")*

**spread()**(data, key, value, fill = NA, convert = FALSE,  
drop = TRUE, sep = NULL)

Spread moves the unique values of a key column into the column names, spreading the values of a value column across the new columns that result.

table2

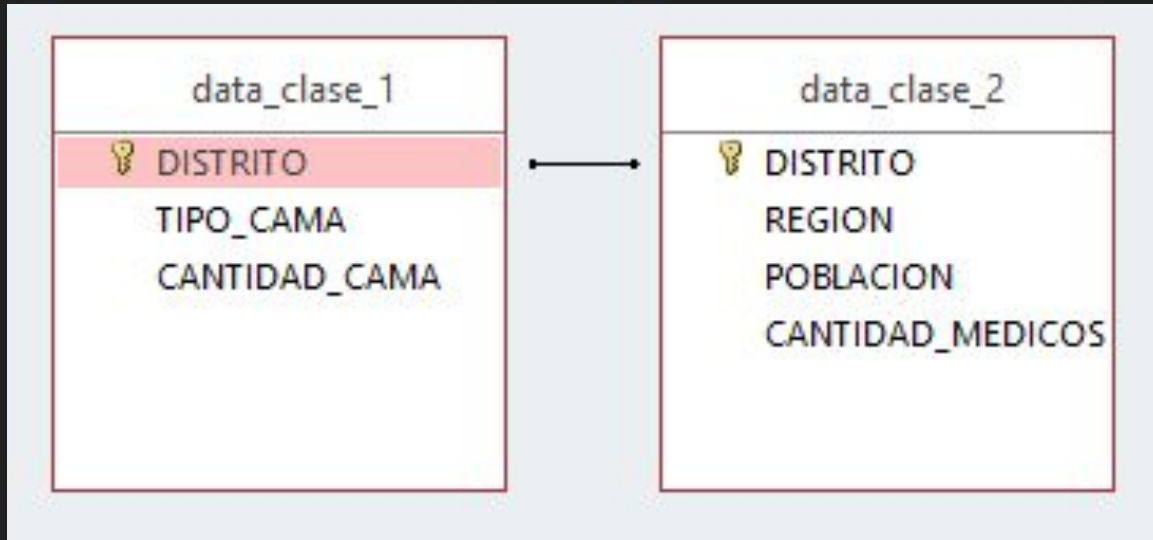
country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T

key value

*spread(table2, type, count)*

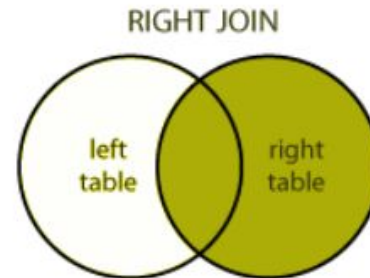
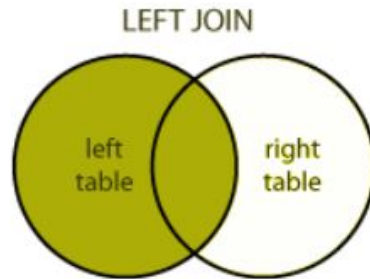
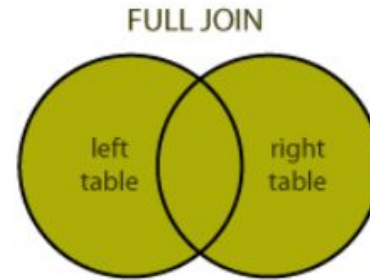
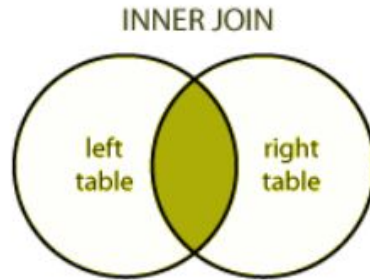
## PRIMARY KEY

*Permite identificar en un dataset a una columna que identifica de forma única a cada fila de una tabla*



# JOINS

*Unimos dos o más tablas en una única y nueva tabla*



## Formas de pedirle a R que haga la unión:

`left_join()`



`right_join()`



`inner_join()`



`full_join()`





```
select()
```

```
# para seleccionar columnas
```

df

color	value
blue	1
black	2
blue	3
blue	4
black	5



value
1
2
3
4
5

```
select(df, value)
```

`filter()`

# para filtrar filas

df

color	value
blue	1
black	2
blue	3
blue	4
black	5



color	value
blue	1
blue	4

```
filter(df, value %in% c(1, 4))
```



`mutate()`

\* para agregar o modificar valores de columnas no agrupadas

df

color	value
blue	1
black	2
blue	3
blue	4
black	5

→

color	value	double	quadruple
blue	1	2	4
black	2	4	8
blue	3	6	12
blue	4	8	16
black	5	10	20

```
mutate(df, double = 2 * value,  
       quadruple = 2 * double)
```

Función	Características	Aplicación*
<code>sum()</code>	Devuelve la suma	<code>sum(data\$columna)</code>
<code>mean()</code>	Devuelve el promedio	<code>mean(data\$columna)</code>
<code>median()</code>	Devuelve la mediana	<code>median(data\$columna)</code>
<code>max()</code>	Devuelve el maximo valor	<code>max(data\$columna)</code>
<code>min()</code>	Devuelve minimo valor	<code>min(data\$columna)</code>
<code>count()</code>	Devuelve la cantidad de observaciones (frecuencia) con respecto a una variable determinada	<code>count(data\$columna)</code>
<code>table()</code>	Devuelve la cantidad de observaciones (frecuencia) con respecto a una variable determinada (para groupby)	<code>table(data\$columna)</code>
<code>n()</code>	Devuelve la cantidad de observaciones (frecuencia) para una variable no necesariamente determinada (para groupby)	<code>n(data\$columna)</code>
<code>quantile()</code>	Devuelve el valor correspondiente al cuantil deseado (0.25, 0.1, n) el que nosotros deseemos	<code>quantile(data\$columna, quantildeseado)</code>
<code>sd()</code>	Devuelve el valor correspondiente al desvio estandar	<code>sd(data\$columna)</code>

\* OJO CON ESTO! SI YA LE DIJERON A R ANTES CUAL ES EL DATAFRAME AL QUE HACEN REFERENCIA NO DEBEN VOLVER A ACLARARLO. Por ejemplo:

ei: `mutate(data, Promedio = mean(columna))` 

ei: `mutate(data, Promedio = mean(data$columna))` 

`group_by()`

# para agrupar filas

df

color	value
blue	1
black	2
blue	3
blue	4
black	5

→

color	total
blue	8
black	7

```
by_color <- group_by(df, color)
summarise(by_color, total = sum(value))
```

```
summarise()
```

```
# para agregar información de columnas agrupadas (calcula  
# 1 sólo valor por grupo)
```

df

color	value
blue	1
black	2
blue	3
blue	4
black	5

→

total
15

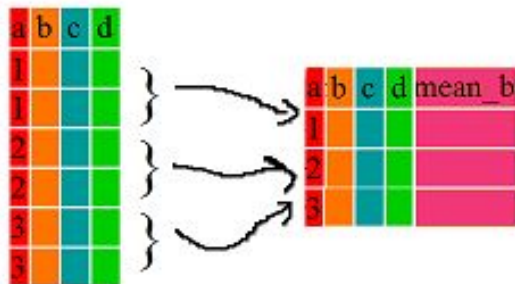
```
summarise(df, total = sum(value))
```

group\_by() %>% mutate()

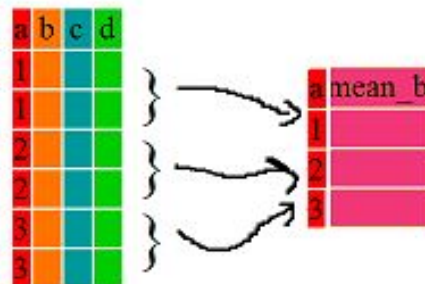
VS.

group\_by() %>% summarize()

```
data %>% group_by(a) %>% mutate(mean_b=mean(b))
```



```
data %>% group_by(a) %>% summarize(mean_b=mean(b))
```



`arrange()`

# ordena las filas en base a una o más columnas

df

color	value
4	1
1	2
5	3
3	4
2	5

→

color	value
5	3
4	1
3	4
2	5
1	2

`arrange(df, desc(color))`