

Practica N° 1. Clase 2

Leyendo Datos

Vamos a cargar datos a R y a explorarlos.

1. Entremos a RStudio
2. Instalaremos un paquete con datos y lo cargaremos para usarlo:

```
> devtools::install_github("cienciadedatos/datos")
> install.packages("nycflights13")
```

```
> library(datos)
```

El paquete `datos` y `nycflights13` contiene información de vuelos, aeropuertos, clima entre otras cosas, vamos a utilizarlos para practicar conceptos de manejo de datos:

3. Para ver que contiene un conjunto de datos podemos escribir su nombre en la consola:

```
> vuelos
```

Tibble: es un tipo de data.frame

Cantidad de filas que tiene el tibble

Cantidad de columnas

columnas/variables

Tipo de dato

10 primeras filas/casos

```
# A tibble: 336,776 x 19
  año   mes día horario_salida salida_programa~ atraso_salida horario_llegada
  <int> <int> <int>          <int>          <int>          <dbl>          <int>
1  2013     1   1         517             515             2             830
2  2013     1   1         533             529             4             850
3  2013     1   1         542             540             2             923
4  2013     1   1         544             545            -1            1004
5  2013     1   1         554             600            -6             812
6  2013     1   1         554             558            -4             740
7  2013     1   1         555             600            -5             913
8  2013     1   1         557             600            -3             709
9  2013     1   1         557             600            -3             838
10 2013     1   1         558             600            -2             753
# ... with 336,766 more rows, and 5 more variables: tiempo_vuelo <dbl>, distancia <
```

4. Para poder analizar mejor el set de datos existen otras funciones, probemos con:

```
> view(vuelos)
```

The screenshot shows the RStudio environment with the 'vuelos' data frame loaded. The top toolbar includes a 'Filter' button and a search icon. The data frame is displayed in a table view with the following columns: año, mes, día, horario_salida, salida_programada, atraso_salida, horario_llegada, llegada_programada, and atraso_llegada. The status bar at the bottom indicates 'Showing 1 to 10 of 336,776 entries'.

Annotations in the image include:

- An orange arrow pointing to the column headers with the text 'columnas/variables'.
- A green arrow pointing to the search icon in the top toolbar with the text 'filtro'.
- A blue arrow pointing to the row count '336,776' in the status bar with the text 'Cantidad de filas'.

5. Otra forma de analizar el set de datos es usando la función `glimpse` del paquete `dyplr`

```
> glimpse(vuelos)
```

Diagrama de flujo que muestra la estructura de los datos:

- Cantidad de columnas:** 336,776 (indicado por una flecha verde).
- Cantidad de filas:** 19 (indicado por una flecha azul).
- Tipo de dato:** Se detallan los tipos de datos para cada variable (indicado por una flecha azul).

Variables y sus tipos de datos:

- observations: 19
- variables: 19
- \$ año: <int>
- \$ mes: <int>
- \$ día: <int>
- \$ horario_salida: <int>
- \$ salida_programada: <int>
- \$ atraso_salida: <dbl>
- \$ horario_llegada: <int>
- \$ llegada_programada: <int>
- \$ atraso_llegada: <dbl>
- \$ aerolinea: <chr>
- \$ vuelo: <int>
- \$ codigoCola: <chr>
- \$ origen: <chr>
- \$ destino: <chr>
- \$ tiempo_vuelo: <dbl>
- \$ distancia: <dbl>
- \$ hora: <dbl>
- \$ minuto: <dbl>
- \$ fecha_hora: <dtm>

columnas/variables

- Otra opción para ver que contiene un set de datos es utilizar la función `kable` del paquete `knitr`, instalemos el paquete `knitr`, cárgalo y probemos la función `kable` con el set de datos de aerolíneas (no usamos vuelos porque tiene más de 300.000 casos y tardaría mucho en generar la salida)

```
> library(knitr)
> kable(aerolineas)
```

```
> kable(aerolineas)
```

| codigo_carrier | nombre | columnas/variables |
|----------------|-----------------------------|--------------------|
| 9E | Endeavor Air Inc. | |
| AA | American Airlines Inc. | |
| AS | Alaska Airlines Inc. | |
| B6 | JetBlue Airways | |
| DL | Delta Air Lines Inc. | |
| EV | ExpressJet Airlines Inc. | |
| F9 | Frontier Airlines Inc. | |
| FL | AirTran Airways Corporation | |
| HA | Hawaiian Airlines Inc. | |
| MQ | Envoy Air | |
| OO | Skywest Airlines Inc. | |
| UA | United Air Lines Inc. | |
| US | US Airways Inc. | |
| VX | Virgin America | |
| WN | Southwest Airlines Co. | |
| YV | Mesa Airlines Inc. | |

- Por último, el operador `$` nos permite explorar una sola variable dentro de un marco de datos. Por ejemplo, ejecuta lo siguiente en tu consola:

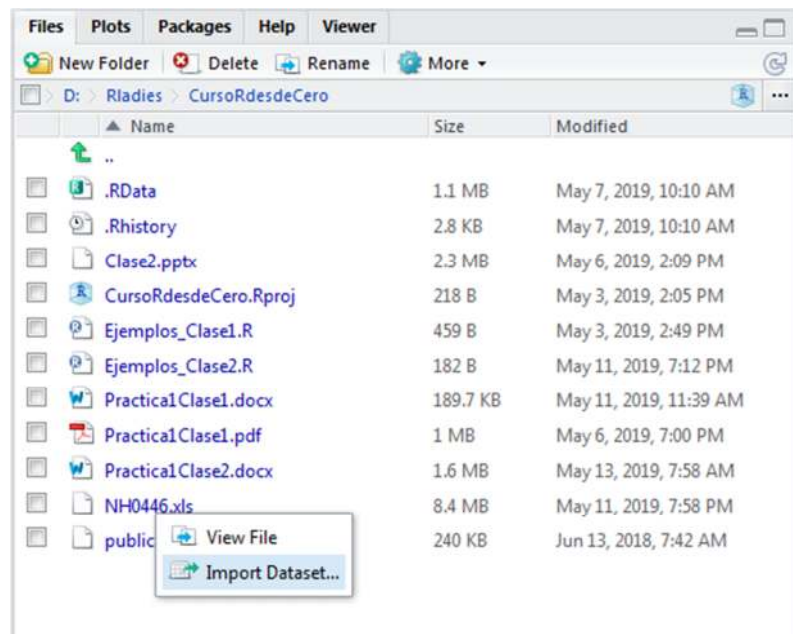
```
> aerolineas
> aerolineas$nombre
```

- Hay una diferencia sutil entre los tipos de variables que podemos tener un set de datos: **variables de identificación y variables de medición/medidas o características**. Por ejemplo, exploremos el marco de datos de los aeropuertos mostrando la salida de `glimpse` (aeropuertos) a continuación:

```
> glimpse (aeropuertos)
```

- ¿Qué columnas sirven para identificar de forma única a cada fila?

10. ¿A qué hace referencia UNA fila en este conjunto de datos de vuelos?
 - a. Datos de una aerolínea
 - b. Datos de un vuelo.
 - c. Datos de un aeropuerto.
 - d. Datos de vuelos múltiples.
11. ¿Cuáles son algunos ejemplos en este conjunto de datos de variables categóricas?
12. ¿Qué las hace diferentes a las variables cuantitativas?
13. ¿Qué ejemplos de variables cuantitativas encontramos en vuelos?
14. Entren a su correo electrónico, a la carpeta compartida de la clase 2 y descarguen los archivos de Excel contenidos en la carpeta.
15. Vamos a cargar una serie de datos externos a R para trabajar con ellos, tenemos tres formas diferentes de hacerlo, la primera en el Panel de Archivos, hacer click sobre el archivo a importar NH0446.xlsx y seleccionar la opción Import DataSet:

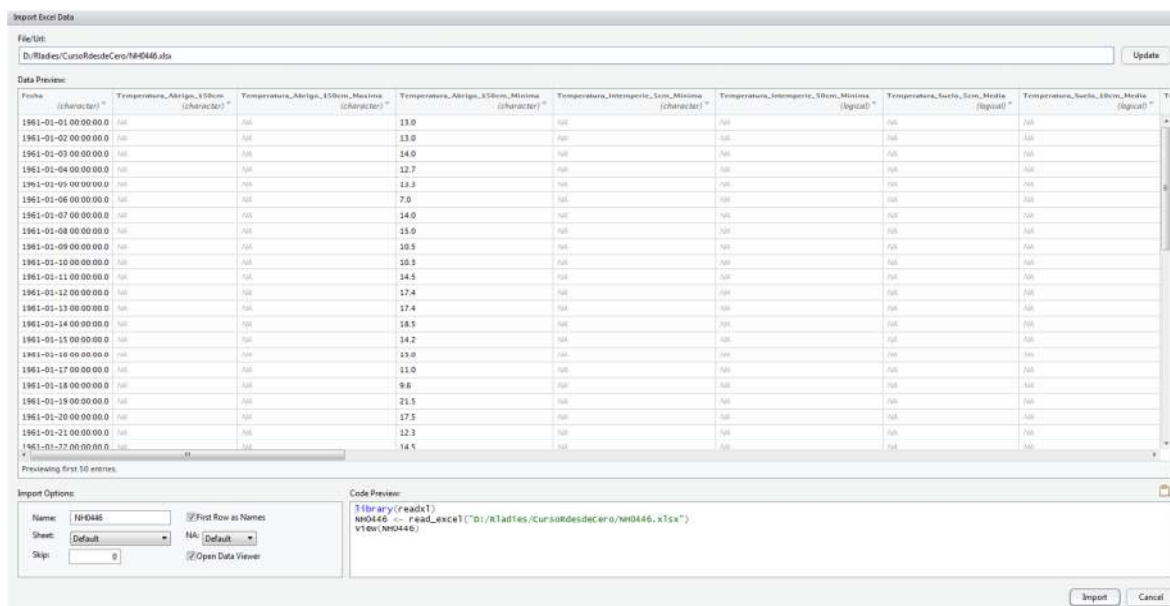




R desde Cero



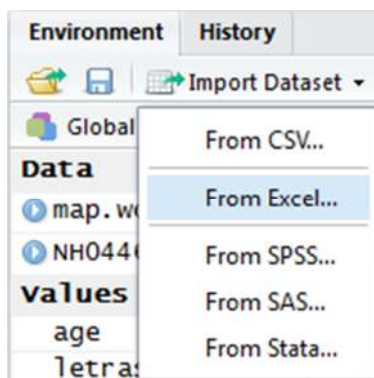
La pantalla muestra una vista previa de los datos a importar, el tipo de dato, el nombre de las columnas, nos muestra el código de R para poder importar ese archivo, vamos a copiar el código y luego presionamos el botón Import:



16. Vamos a crear un nuevo archivo R Script y allí vamos a pegar el código que se generó cuando importamos el archivo NH0446.xlsx:

```
library(readxl)
NH0446 <- read_excel("D:/Rladies/CursoDesdeCero/NH0446.xlsx")
view(NH0446)
```

17. Ahora vamos a importar de otra manera, utilizando el botón Import DataSet del panel de Entorno/Historial, nos presenta varias opciones de tipos de archivos a importar, seleccionamos Excel.





R desde Cero



18. La pantalla para importar es la misma que con la opción anterior, pero en este caso debemos indicar cual es el archivo que tenemos que importar, para eso presionamos en el botón Browse y allí elegimos el archivo publicaciones_propias_2013.xls. Nuevamente nos presenta la misma pantalla. Vamos a copiar el código y presionamos el botón Importar. Pegamos el código debajo del código anterior y guarda el archivo como PracticaClase2.R.

```
library(readxl)
publicaciones_propias_2013 <- read_excel("D:/Rladies/CursoRdesdeCero/publicaciones_propias_2013.xlsx")
view(publicaciones_propias_2013)
```

¿Hay alguna parte que se repita del código anterior?. Acomodar el código para que el paquete `readxl` solo se cargue una vez

¿En qué objetos se guardaron los datos que importamos?

¿Cómo podemos explorar estos dos set de datos?

¿Qué cantidad de columnas y de casos tiene cada set de datos?