

Spanish TDT movies report

Guillermo Diaz

2022-09-19

Índice

Basic data	1
Real Stuff with TDT data	5
Create tables the handy way	6
Para escribir una cita en bloque	6

Link

This line of code is quite magical. Change any `opts_chunk` here and it will change the entire document. You can set any option you want: `echo`, `include`, `warning` or `messages`.

```
knitr::opts_chunk$set(  
  echo = TRUE,  
  message = FALSE,  
  warning = FALSE  
)
```

Basic data

The data set contains 127 instances and 9 variables. From now on, we won't be using the variable `Description`, which included a brief synopsis of each film. Let's see a few examples of the data.

```
library(knitr)  
df_no_desc <- df %>%  
  select(1:8)  
  
df_no_desc %>%  
  head() %>%  
  knitr::kable()
```

date_time	channel	sp_title	original_title	year	genre	country	length
2022-09-13 00:19:00	Paramount Network	Vanilla Sky	Vanilla Sky	2001	Drama	NA	NA

date_time	channel	sp_title	original_title	year	genre	country	length
2022-09-13 00:42:00	Neox	Ruslan: la venganza del asesino	Driven to Kill	2009	Acción	NA	NA
2022-09-13 01:10:00	laSexta	La mujer del pastor	The Pastor's Wife	2011	Drama	NA	NA
2022-09-13 13:10:00	La 2	El sonido de un tambor	Cimarron: The Sound of a Drum	1968	Western	NA	NA
2022-09-13 16:05:00	TRECE	Comando secreto	The Secret War of Harry Frigg	1968	Comedia	NA	NA
2022-09-13 16:20:00	La 1	La cuchara de Elli	Tessa Hennig - Elli gibt den Löffel ab	2012	Drama	NA	NA

The variables names are formatted to work with them in R, not to be shown in a document. We can clean them.

```
df_clean <- df_no_desc %>%
  dplyr::rename(
    Date = date_time,
    Channel = channel,
    "Spanish title" = sp_title,
    "Original title" = original_title,
    Year = year,
    Genre = genre,
    Country = country,
    Length = length
  )

df_clean %>%
  head() %>%
  knitr::kable()
```

Date	Channel	Spanish title	Original title	Year	Genre	Country	Length
2022-09-13 00:19:00	Paramount Network	Vanilla Sky	Vanilla Sky	2001	Drama	NA	NA
2022-09-13 00:42:00	Neox	Ruslan: la venganza del asesino	Driven to Kill	2009	Acción	NA	NA
2022-09-13 01:10:00	laSexta	La mujer del pastor	The Pastor's Wife	2011	Drama	NA	NA
2022-09-13 13:10:00	La 2	El sonido de un tambor	Cimarron: The Sound of a Drum	1968	Western	NA	NA
2022-09-13 16:05:00	TRECE	Comando secreto	The Secret War of Harry Frigg	1968	Comedia	NA	NA
2022-09-13 16:20:00	La 1	La cuchara de Elli	Tessa Hennig - Elli gibt den Löffel ab	2012	Drama	NA	NA

This is so much better. I want to see some examples with **Country** and **Length** data.

```
df_clean %>%
  drop_na() %>%
  head() %>%
  knitr::kable()
```

Date	Channel	Spanish title	Original title	Year	Genre	Country	Length
2022-09-18 00:26:00	Neox	Tenemos que hablar	Tenemos que hablar	2016	Comedia	España	91 min
2022-09-18 00:35:00	Antena 3	Suplantación de identidad	The Cheating Pact	2013	Suspense	Estados Unidos	85 min
2022-09-18 00:53:00	Cuatro	Colonia V	The Colony	2013	Ciencia ficción	Canadá	95 min
2022-09-18 01:15:00	La 1	Amor, ladrón, diamantes	Liebe, Diebe, Diamanten	2015	Drama	Alemania	90 min
2022-09-18 01:30:00	TRECE	Sol naciente	Rising Sun	1993	Suspense	Estados Unidos	129 min
2022-09-18 01:45:00	Paramount Network	Shame	Shame	2011	Drama	Reino Unido	101 min

I don't like to see the unit in the **Length** column. I noticed that film genres are in spanish. Let's clean both columns.

¿How many unique values there are in each variable?

```
library(purrr)
df %>% map_dbl(
  n_distinct
)
```

```
##      date_time      channel  sp_title original_title      year
##          115           10        123          123          51
##      genre      country    length  description
##          16           7         33          123
```

So, if there are 127 instances, why do we only have 123 movies? I guess some of them were broadcasted more than once. These are the ones:

```
data %>%
  group_by(sp_title) %>%
  summarise(emisiones = n()) %>%
  filter(emisiones > 1)
```

```
## # A tibble: 4 x 2
##   sp_title      emisiones
##   <chr>          <int>
## 1 ¡Viven!          2
## 2 Cómo entrenar a tu dragón 2
## 3 Diario de Greg 3: Días de perros
## 4 Indiana Jones y la última cruzada 2
```

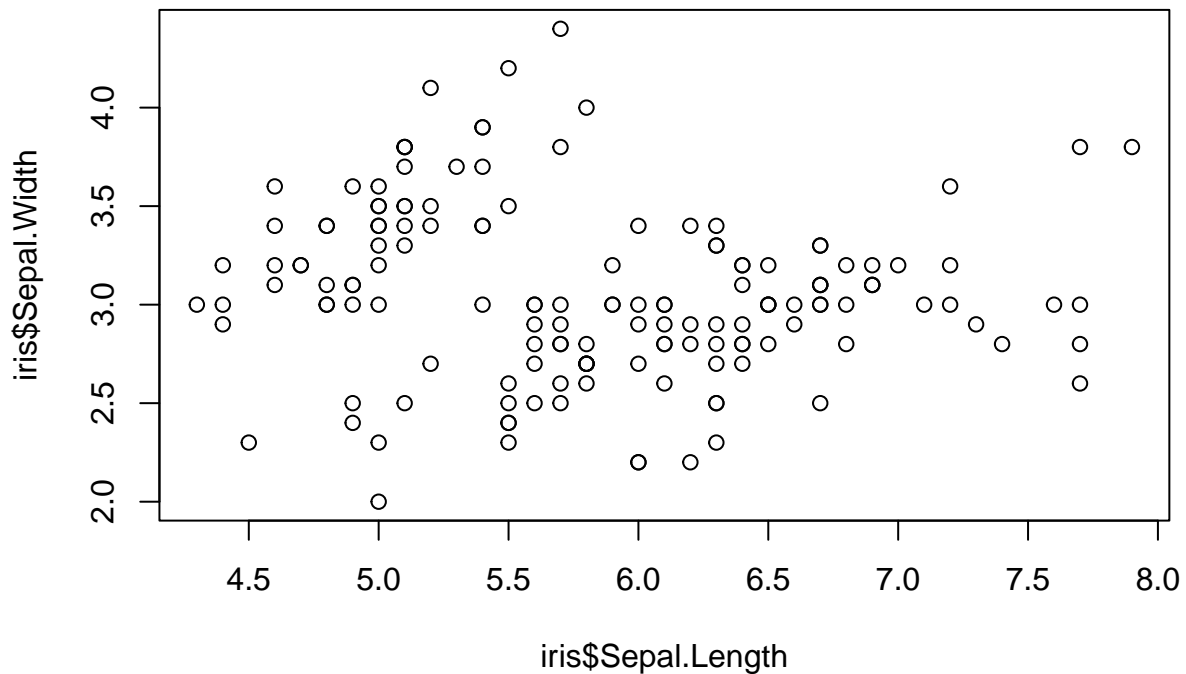
Con `results = 'hold'` lo que hacemos es sacar todos los resultados del tirón, sin partir el cuadro de código cada vez que hay un resultado que mostrar¹.

¹Me marco un pie de página absolutamente escandaloso.

```
data(iris)
a <- 1+1
print(a)
head(iris)
```

```
## [1] 2
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2   setosa
## 2         4.9         3.0         1.4         0.2   setosa
## 3         4.7         3.2         1.3         0.2   setosa
## 4         4.6         3.1         1.5         0.2   setosa
## 5         5.0         3.6         1.4         0.2   setosa
## 6         5.4         3.9         1.7         0.4   setosa
```

Note that the **echo = FALSE** parameter was added to the code chunk to prevent printing of the R code that generated the plot.



```
plot(iris$Sepal.Length, iris$Sepal.Width)
```

```
data("HairEyeColor")
chisq.test(HairEyeColor[, , 2])
```

```
##
## Pearson's Chi-squared test
```

```
##
## data: HairEyeColor[, , 2]
## X-squared = 106.66, df = 9, p-value < 2.2e-16
```

```
library(vcd)
```

¡Cuidado con el nombre de las variables!

```
library(ggplot2)
data(iris)

ggplot(iris, aes(Sepal.Width, Sepal.Length)) +
  geom_point()
```

Real Stuff with TDT data

```
df <- read.csv("https://raw.githubusercontent.com/GuilleDiaz7/Automatic-Web-Scraping-of-Spanish-TDT-Files/master/data/TDT_data.csv",
               fileEncoding = "UTF-8")
library(tidyr)
df_clean <- df %>%
  drop_na()
library(dplyr)
head(
  df_clean %>% select(
    1:8
  )
)
```

```
##           date_time           channel           sp_title
## 1 2022-09-18 00:26:00           Neox       Tenemos que hablar
## 2 2022-09-18 00:35:00       Antena 3 Suplantación de identidad
## 3 2022-09-18 00:53:00         Cuatro         Colonia V
## 4 2022-09-18 01:15:00         La 1   Amor, ladrón, diamantes
## 5 2022-09-18 01:30:00         TRECE         Sol naciente
## 6 2022-09-18 01:45:00 Paramount Network         Shame
##           original_title year           genre           country length
## 1       Tenemos que hablar 2016           Comedia           España   91 min
## 2       The Cheating Pact 2013 Suspense / Thriller Estados Unidos   85 min
## 3           The Colony 2013   Ciencia ficción           Canadá   95 min
## 4 Liebe, Diebe, Diamanten 2015           Drama           Alemania   90 min
## 5           Rising Sun 1993 Suspense / Thriller Estados Unidos 129 min
## 6           Shame 2011           Drama           Reino Unido 101 min
```

Para crear *listas*:

- Un elemento
- Otro elemento

– Otro elemento más

1. Elemento 1
2. Elemento 2
3. Elemento 3

Create tables the handy way

```
library(palmerpenguins)
penguins %>%
  head() %>%
  knitr::kable()
```

Para escribir una cita en bloque

Cara antigua. Es **code block** en la opción *Format* del Editor Visual.