

Diplomado de Análisis Estadístico usando R

Módulo 3: Exploración de datos

Profesor: Víctor Macías E.

En estas notas trabajaremos con:

- Fechas y horas
- Factores
- Strings

Para obtener información adicional a la discutida en esta guía, se recomienda revisar:

<https://www.rstudio.com/resources/cheatsheets/>

1. Fechas y horas

- Para trabajar con fechas y horas usaremos el paquete *lubridate*:

<https://cran.r-project.org/web/packages/lubridate/lubridate.pdf>

```
library(lubridate)
```

- La norma ISO 8601 especifica la notación estándar utilizada para representar fechas y horas:
 - Fechas representadas como *YYYY-MM-DD*
 - Horas representadas como *HH:MM:SS*
- Hay tres tipos de datos que capturan la dimensión temporal de los datos:
 - *date* es la fecha. Por ejemplo, 2020-09-28
 - *time* es una hora dentro de un día. Por ejemplo, 08:48:12
 - *date-time* combina *date* y *time*. Son también conocidos como *POSIXct*. Por ejemplo, 2020-09-28 08:48:12
- Una fecha se puede crear de:
 - Un *string*.
 - Componentes individuales de *date-time*.
 - Otra fecha o fecha-hora.

A continuación se presentan varios ejemplos de cómo crear un *date*, *time* o *date-time*:

```
as.Date("2020-09-27") # forma 1
```

```
## [1] "2020-09-27"
```

```
ymd("2020-09-27") # forma 2
```

```
## [1] "2020-09-27"
```

```
ymd(20200927) # forma 3 - Note que en esta forma no se usó un string
```

```
## [1] "2020-09-27"
```

```
dmy("27-09-2020") # forma 4
```

```
## [1] "2020-09-27"
```

```
as.POSIXct("2020-09-28 20:13:00", tz = "America/El_Salvador") # forma 1
```

```
## [1] "2020-09-28 20:13:00 CST"
```

```
ymd_hms("2020-09-21 14:04:32", tz = "America/El_Salvador") # forma 2
```

```
## [1] "2020-09-21 14:04:32 CST"
```

```
Sys.timezone()
```

```
## [1] "America/Santiago"
```

```
dmy(31012017) # forma 1
```

```
## [1] "2017-01-31"
```

```
dmy("31-Enero-2017") # forma 2
```

```
## [1] "2017-01-31"
```

Para convertir un objeto a *date* o *datetime* usar *as_date()* o *as_datetime()*.

```
as_date("2020-09-29 20:49:34")
```

```
## [1] "2020-09-29"
```

```
as_datetime("2020-09-29")
```

```
## [1] "2020-09-29 UTC"
```

```
as_date(0) # En este caso 0 es un número ¿Qué representa?
```

```
## [1] "1970-01-01"
```

```
as_date(7)
```

```
## [1] "1970-01-08"
```

```
as_date(-2)
```

```
## [1] "1969-12-30"
```

- Objetos que son definidos como fechas son almacenados internamente como el número de días desde el 1 de enero de 1970. Al definir fechas se pueden realizar operaciones aritméticas entre éstas, como por ejemplo:

```
today()
```

```
## [1] "2020-11-24"
```

```
now()
```

```
## [1] "2020-11-24 11:06:45 -03"
```

```
ymd("2020-11-18") + 2
```

```
## [1] "2020-11-20"
```

```
ymd("2020-12-31") - ymd("2020-11-18")
```

```
## Time difference of 43 days
```

```
ymd("2020-11-18") + weeks(5)
```

```
## [1] "2020-12-23"
```

```
today() + 1:5*weeks(1)
```

```
## [1] "2020-12-01" "2020-12-08" "2020-12-15" "2020-12-22" "2020-12-29"
```

```
dmy("31-12-2019") %m+% (1:12*months(1))
```

```
## [1] "2020-01-31" "2020-02-29" "2020-03-31" "2020-04-30" "2020-05-31"
```

```
## [6] "2020-06-30" "2020-07-31" "2020-08-31" "2020-09-30" "2020-10-31"
```

```
## [11] "2020-11-30" "2020-12-31"
```

- Hay una lista de 594 zonas horarias a nivel mundial, las cuales se pueden obtener `OlsonNames()`. En el caso del continente americano, la lista incluye 167 zonas horarias. Las primeras 10 son:

```
## [1] "America/Adak" "America/Anchorage"
## [3] "America/Anguilla" "America/Antigua"
## [5] "America/Araguaina" "America/Argentina/Buenos_Aires"
## [7] "America/Argentina/Catamarca" "America/Argentina/ComodRivadavia"
## [9] "America/Argentina/Cordoba" "America/Argentina/Jujuy"
## [11] "America/Argentina/La_Rioja" "America/Argentina/Mendoza"
## [13] "America/Argentina/Rio_Gallegos" "America/Argentina/Salta"
## [15] "America/Argentina/San_Juan" "America/Argentina/San_Luis"
## [17] "America/Argentina/Tucuman" "America/Argentina/Ushuaia"
## [19] "America/Aruba" "America/Asuncion"
## [21] "America/Atikokan" "America/Atka"
## [23] "America/Bahia" "America/Bahia_Banderas"
## [25] "America/Barbados" "America/Belem"
## [27] "America/Belize" "America/Blanc-Sablon"
## [29] "America/Boa_Vista" "America/Bogota"
```

Para chequear la zona horaria en que nos encontramos:

```
Sys.timezone()
```

```
## [1] "America/Santiago"
```

2. Factors

- Se usan para representar variables categóricas. También permiten mostrar un vector de strings en un orden que no sea alfabético.
- Para trabajar con factores el paquete ***forcats*** es de gran ayuda:

<https://cran.r-project.org/web/packages/forcats/forcats.pdf>

```
library(forcats)
```

- Dos operaciones comunes realizadas con factores son:
 - Cambiar el orden de los niveles para lo cual puedes usar *fct_reorder()*
 - Cambiar los valores de los niveles usando *fct_recode()*

3. Strings

- Para trabajar con strings usaremos el paquete ***stringr***, cuya descripción se puede encontrar en la siguiente página:

<https://cran.r-project.org/web/packages/stringr/stringr.pdf>

```
library(stringr)
```

- Para crear strings usar comillas:

```
string1 <- "Veinte poemas de amor y una canción desesperada"  
string2 <- 'Veinte poemas de amor y una canción desesperada'
```

```
class(string1)
```

```
## [1] "character"
```

```
class(string2)
```

```
## [1] "character"
```

- Entre las funciones que incluye el paquete `stringr` para manipular `strings` se incluyen:

Símbolo	Descripción
<code>str_length</code>	Retorna el número de caracteres de cada string
<code>str_c</code>	Combinar strings
<code>str_sub</code>	Extrae partes de un string
<code>str_to_upper</code>	Cambia a mayúsculas
<code>str_to_lower</code>	Cambia a minúsculas
<code>str_detect</code>	Entrega el valor=TRUE si el string contiene el patrón buscado y FALSE, si no lo encuentra
<code>str_subset</code>	Retorna los strings que contienen el patrón buscado
<code>str_count</code>	Retorna el número de veces que el patrón buscado ocurre en el string
<code>str_replace</code>	Reemplaza sólo la primera vez que ocurre el patrón
<code>str_replace_all</code>	Reemplaza cada vez que ocurre el patrón
<code>str_trim</code>	Elimina espacio en blanco al comienzo y/o final de un string

- Para manipular strings es muy útil aprender a trabajar con expresiones regulares (*regex*), las que constituyen un modo de describir patrones específicos en strings. *En regex todos son caracteres*. Algunos ejemplos de símbolos usados en las expresiones regulares se presentan en la siguiente tabla:

Símbolo	Descripción
<code>*</code>	cero o más veces
<code>+</code>	uno o más veces
<code>{n}</code>	exactamente <i>n</i> veces
<code>{n,}</code>	al menos <i>n</i> veces
<code>{,m}</code>	como máximo <i>m</i> veces
<code>{n,m}</code>	entre <i>n</i> y <i>m</i> veces
<code>^</code>	Comienzo del string
<code>\$</code>	Término del string
<code> </code>	o
<code>\\d</code> o <code>[0-9]</code>	Cualquier dígito {0,1,2,3,4,5,6,7,8,9}
<code>\\w</code>	Cualquier caracter como letra, número o guión bajo
<code>\\s</code>	espacios en blanco (espacio, tab, nueva línea)
<code>.</code>	Cualquier caracter como número, letra, espacio en blanco
<code>\\.</code>	Símbolo de puntuación (punto seguido, punto aparte)
<code>[a-zA-Z]</code>	Todas las letras
<code>[abc]</code>	Matches letras a, b o c
<code>[^abc]</code>	Cualquier caracter excepto letras a,b,c
<code>^\\d\$</code>	Comienzo string, seguido por un dígito, término del string