

# Programación en R para el análisis de datos

## Expresiones Regulares

---

Nicolás Schmidt

mail: `nschmidt@cienciassociales.edu.uy`

GitHub: `@Nicolas-Schmidt`

Departamento de Ciencia Política

Facultad de Ciencias Sociales

1. Qué es una expresión regular?
2. Expresiones regulares
3. Carguemos datos
4. Funciones de base
5. Funciones de stringr

## Qué es una expresión regular?

---



# expresión regular?

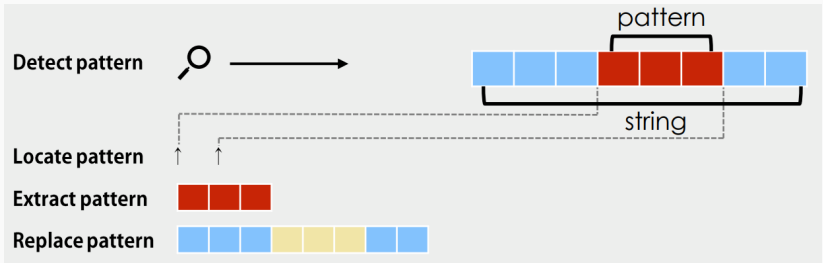
Una expresión regular es un patrón que describe un conjunto de cadenas de caracteres. Es decir, es una secuencia de caracteres que conforma un patrón de búsqueda.

Este patrón de búsqueda a diferencia del reconocimiento de patrones, por ejemplo, la coincidencia buscada es comúnmente exacta.

## Cuándo usarla?

Obviamente, es necesario utilizar una *regular expression* cuando se quiere buscar dentro de datos de tipo `character`. De lo contrario se utilizaría lógica matemática. Esto no quiere decir que no se pueden incorporar

números en la conformación de la expresión regular. Simplemente se deben tener caracteres (numéricos o alfabéticos) almacenado en modo texto.



# Expresiones regulares

---



# Cuantificadores

*	Coincide al menos 0 veces.
+	Coincide al menos 1 vez.
?	Coincide a lo sumo 1 vez.
{n}	Coincide 'n' veces exactamente.
{n,}	Coincide al menos 'n' veces.
{n,m}	Coincide entre 'n' y 'm' veces.

# Ubicación dentro del patrón

^	Coincidencia con el inicio de la cadena.
\$	Coincidencia con el fin de la cadena.
\b	Coincide con la cadena vacía en el borde de una palabra.
\B	Coincide con la cadena vacía siempre que no esté en el borde de una palabra.

# Operadores

.	Coincidencia con cualquier caracter.
[...]	Lista de caracteres.
[^ ...]	Negación de lista de caracteres.
	El operador 'o'.
(...)	Agrupación de expresiones regulares.

# Escapes

Hay caracteres especiales que en R tienen un uso funcional. Por lo que se debe generar escape para poder incorporarlos como elementos de una expresión regular.

\$ \* + . ? [ ] ^ { } \| ()

Para escapar hay que usar doble barra, por ejemplo: `\\$`

<code>[digit:]</code>	Números del 0 al 9: <code>[0-9]</code> .
<code>\D</code>	No dígitos: <code>[^ 0-9]</code> .
<code>[lower:]</code>	letras minúsculas: <code>[a-z]</code> .
<code>[upper:]</code>	letras mayúsculas: <code>[A-Z]</code> .
<code>[alnum:]</code>	caracteres alfabéticos: <code>[[:lower:][:upper:]]</code> o <code>[A-z]</code> .
<code>\w</code>	caracteres de palabras: <code>[[:alnum:]_]</code> o <code>[A-z0-9_]</code> .
<code>\W</code>	no palabra: <code>[^ A-z0-9_]</code> .
<code>[blank:]</code>	espacio y tabulación.
<code>[space:]</code>	tabulación, nueva línea, tabulación vertical...
<code>\s</code>	espacio.
<code>\S</code>	no espacio.
<code>[punct:]</code>	caracteres de puntuación

Importante:

- Las clases '[...]' deben usarse entre corchetes: '[[...]]'
- El '\' es un caracter especial. Por lo de que se debe usar doble: '\\D

# Carguemos datos

---

# Cargando datos

```
## Warning: package 'tidyverse' was built under R version 3.6.3
```

```
library(gapminder)
paises <- gapminder$country %>% as.character() %>% unique()
length(paises)
```

```
## [1] 142
```

```
paises[1:40]
```

```
## [1] "Afghanistan"      "Albania"
## [3] "Algeria"          "Angola"
## [5] "Argentina"        "Australia"
## [7] "Austria"          "Bahrain"
## [9] "Bangladesh"       "Belgium"
## [11] "Benin"            "Bolivia"
## [13] "Bosnia and Herzegovina" "Botswana"
## [15] "Brazil"           "Bulgaria"
## [17] "Burkina Faso"     "Burundi"
## [19] "Cambodia"         "Cameroon"
## [21] "Canada"           "Central African Republic"
## [23] "Chad"             "Chile"
## [25] "China"            "Colombia"
## [27] "Comoros"          "Congo, Dem. Rep."
## [29] "Congo, Rep."      "Costa Rica"
## [31] "Cote d'Ivoire"    "Croatia"
## [33] "Cuba"             "Czech Republic"
## [35] "Denmark"          "Djibouti"
## [37] "Dominican Republic" "Ecuador"
## [39] "Egypt"            "El Salvador"
```



# Funciones de base

---

<code>grep()</code>	Devuelve los indices en los que se encuentra el patron
<code>grepl()</code>	Devuelve un vector logico con TRUE donde hay coincidencia
<code>regexpr()</code>	Devuelve los indices de la cadena en donde comineza la
<code>gregexpr()</code>	coincidencia y la duracion de la misma
<code>sub()</code>	Busca y reemplaza en donde hay
<code>gsub()</code>	coincidencia

# Ejemplos

```
vec <- grep("ina", paises)  
length(vec)
```

```
## [1] 5
```

```
vec
```

```
## [1] 5 13 17 25 56
```

# Ejemplos

```
a <- grep("A", paises)
length(a)
```

```
## [1] 10
```

```
aA <- grep("A", paises, ignore.case = TRUE, value = TRUE)
aA[1:10]
```

```
## [1] "Afghanistan" "Albania"      "Algeria"      "Angola"       "Argentina"
## [6] "Australia"   "Austria"     "Bahrain"     "Bangladesh"  "Bolivia"
```

```
aA <- grep("[Aa]", paises)
aA[1:10]
```

```
## [1] 1 2 3 4 5 6 7 8 9 12
```

# Ejemplos

```
grep(países, pattern = "^U", value = TRUE)
```

```
## [1] "Uganda"          "United Kingdom" "United States"  "Uruguay"
```

```
grep(países, pattern = "c$", value = TRUE)
```

```
## [1] "Central African Republic" "Czech Republic"  
## [3] "Dominican Republic"      "Slovak Republic"
```

```
grep(países, pattern = "\\s", value = TRUE)
```

```
## [1] "Bosnia and Herzegovina" "Burkina Faso"  
## [3] "Central African Republic" "Congo, Dem. Rep."  
## [5] "Congo, Rep."           "Costa Rica"  
## [7] "Cote d'Ivoire"         "Czech Republic"  
## [9] "Dominican Republic"    "El Salvador"  
## [11] "Equatorial Guinea"     "Hong Kong, China"  
## [13] "Korea, Dem. Rep."      "Korea, Rep."  
## [15] "New Zealand"           "Puerto Rico"  
## [17] "Sao Tome and Principe" "Saudi Arabia"  
## [19] "Sierra Leone"         "Slovak Republic"  
## [21] "South Africa"          "Sri Lanka"  
## [23] "Trinidad and Tobago"   "United Kingdom"  
## [25] "United States"         "West Bank and Gaza"  
## [27] "Yemen, Rep."
```

# Ejemplos

```
grep(países, pattern = "[[:punct:]]", value = TRUE)
```

```
## [1] "Congo, Dem. Rep." "Congo, Rep."      "Cote d'Ivoire"      "Guinea-Bissau"  
## [5] "Hong Kong, China" "Korea, Dem. Rep." "Korea, Rep."       "Yemen, Rep."
```

```
grep(países, pattern = "[[:punct:]]", value = TRUE)
```

```
## [1] "Congo, Dem. Rep." "Congo, Rep."      "Cote d'Ivoire"      "Guinea-Bissau"  
## [5] "Hong Kong, China" "Korea, Dem. Rep." "Korea, Rep."       "Yemen, Rep."
```

# Ejemplos

```
strings <- c("a", "ab", "acb", "accb", "acccb", "accccb")
strings
```

```
## [1] "a"      "ab"     "acb"    "accb"   "acccb"  "accccb"
```

```
grep("ac*b", strings, value = TRUE)
```

```
## [1] "ab"     "acb"    "accb"   "acccb"  "accccb"
```

```
grep("ac+b", strings, value = TRUE)
```

```
## [1] "acb"    "accb"   "acccb"  "accccb"
```

```
grep("ac?b", strings, value = TRUE)
```

```
## [1] "ab"     "acb"
```

# Ejemplos

```
strings <- c("a", "ab", "acb", "accb", "acccb", "accccb")
strings
```

```
## [1] "a"      "ab"     "acb"    "accb"   "acccb"  "accccb"
```

```
grep("ac{2}b", strings, value = TRUE)
```

```
## [1] "accb"
```

```
grep("ac{2,}b", strings, value = TRUE)
```

```
## [1] "accb"   "acccb"  "accccb"
```

```
grep("ac{2,3}b", strings, value = TRUE)
```

```
## [1] "accb"   "acccb"
```



# Ejemplos

```
strings <- c("abcd", "cdab", "cabd", "c abd")  
strings
```

```
## [1] "abcd" "cdab" "cabd" "c abd"
```

```
grep("ab", strings, value = TRUE)
```

```
## [1] "abcd" "cdab" "cabd" "c abd"
```

```
grep("^ab", strings, value = TRUE)
```

```
## [1] "abcd"
```

# Ejemplos

```
strings
```

```
## [1] "abcd" "cdab" "cabd" "c abd"
```

```
grep("ab$", strings, value = TRUE)
```

```
## [1] "cdab"
```

```
grep("\\bab", strings, value = TRUE)
```

```
## [1] "abcd" "c abd"
```

# Ejemplos

```
strings <- c("^ab", "ab", "abc", "abd", "abe", "ab 12")
strings
```

```
## [1] "^ab"    "ab"     "abc"    "abd"    "abe"    "ab 12"
```

```
grep("ab.", strings, value = TRUE)
```

```
## [1] "abc"    "abd"    "abe"    "ab 12"
```

```
grep("ab[c-e]", strings, value = TRUE)
```

```
## [1] "abc" "abd" "abe"
```

```
grep("ab[^c]", strings, value = TRUE)
```

```
## [1] "abd"    "abe"    "ab 12"
```

# Ejemplos

```
strings <- c("^ab", "ab", "abc", "abd", "abe", "ab 12")
strings
```

```
## [1] "^ab" "ab" "abc" "abd" "abe" "ab 12"
```

```
grep("^ab", strings, value = TRUE)
```

```
## [1] "ab" "abc" "abd" "abe" "ab 12"
```

```
grep("\\^ab", strings, value = TRUE)
```

```
## [1] "^ab"
```

```
grep("abc|abd", strings, value = TRUE)
```

```
## [1] "abc" "abd"
```

```
gsub("(ab) 12", "\\1 34", strings)
```

```
## [1] "^ab" "ab" "abc" "abd" "abe" "ab 34"
```

# Funciones de stringr

---

El paquete `stringr` ofrece un conjunto de funciones para trabajar con expresiones regulares que agregan y mejoran en algunos casos funciones existentes en el paquete base. Asimismo, en muchas circunstancias el rendimiento es mejor.

De todas maneras la gran virtud refiere a que las funciones tienen un nombre intuitivo. Todas las funciones comienzan con el prefijo `str_*`