Taller R

Sesion 4

Agosto 2022

Índice

- Limpieza de datos (continuación: variables fecha, numérica y texto)
- topper()
- unique()
- mutate()
- count()

Material de trabajo

Carpeta de drive:

https://drive.google.com/drive/u/0/folders/14QDyRnqnLUjosCBIAzBqXL-qglOMjNJL

Base a descargar ile_completa

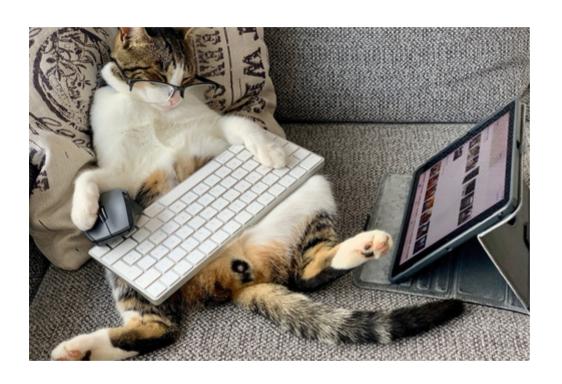
Liga de zoom https://us02web.zoom.us/j/83841692580

ID de reunión: 838 4169 2580

Que hemos visto:

- Estructura de datos
- Paqueterías o librerias
- Importar o subir archivos
- Función unique(), str(), View()
- Simbolo "\$" y pipe %>%
- Renombrar observaciones para limpiar datos
- función count()

Empezamos con la limpieza de datos



Antes de comenzar recuerda:

Llamar librerias:

```
library(readx1)
library(dplyr)
```

Importar o cargar csv:

```
datos<-read.csv("C:/Users/Paola Viridiana/Desktop/ile_completa.csv", encoding = "U")</pre>
```

Recuerda que acá cambia tu ruta o ubicación de documenta y no olvides indicar el encoding UTF-8.

Recordemos:)

unique()

Extrae los valores unicos de una variable o columna.

toupper()

Con esta función se convierte a mayúscula una columna o varias

mutate()

La función mutate() sirve para crear nuevas variables.

case_when

Permite definir una variable, la cual toma un valor particular para cada condición establecida. En caso de no cumplir ninguna de las condiciones establecidas la variable tomara valor NA.

```
# con signo "$" se especifica una variable
datos$mes<- toupper(datos$mes)</pre>
```

Variable year

Exploramos la variable year mostrando sus valores unicos

```
unique(datos$year)
## [1] 2016 2017 2018 2019 NA 2020
```

Parece todo bien... a excepción de los NA o valores perdidos

```
datos %>% count(year)
```

```
## year n
## 1 2016 18086
## 2 2017 17598
## 3 2018 17179
## 4 2019 15251
## 5 2020 6271
## 6 NA 21
```

Donde hay 21 NA's

Variable fingreso

Esta variable hace referencia a la fecha de ingreso, consideremos tener registros del periodo 2016 a 2020, otras fechas se tomaran como valores perdidos o NA.

Acá otra forma de renombrar observaciones: Hay fechas que no tienen formato fecha (ymd) que son números enteros estas observaciones se pasaran a NA.

```
#Remplazo valores no fecha a NA
datos$fingreso[datos$fingreso<=11112] <- NA_character_
datos$fingreso[datos$fingreso>43466] <- NA_character_</pre>
```

Se identifican dos fecha exactas que no entran en el periodo de análisis por lo cual se procede a reclasificarla como NA o valor perdido.

```
#Remplazo fechas fuera de la periodicidad 2016-2020 a NA
datos$fingreso[datos$fingreso=="30/04/2048"] <- NA_character_
datos$fingreso[datos$fingreso<="25/11/2015"] <- NA_character_</pre>
```

Se obseva una gran perdida de información 64,571 NA :c

Variable autorreferida

Esta variable especificar si la paciente es referida de otra unidad o acude directamente. Debe estar clasificada como sí = 1 y no =0.

De todos los valores unicos se identifica los valores que pasan a NA y los pasan a un vector.

```
auto_no<-c(NA, "NE", "N/E")
```

Para remplazar se utiliza el vector nombrado "auto_no"

```
#Con %in% le indicamos que busque dentro de la variable autoref los valores identi
datos$AUTOREF_LIMPIA[datos$autoref %in% auto_no] <- NA_character_</pre>
```

```
# Los valores diff a auto_no que esten en la columna autoref pasan a 1
datos$AUTOREF_LIMPIA[datos$autoref != auto_no] <- 1</pre>
```

```
#Si hay observaciones = a NO se renombran con el valor 0
datos$AUTOREF_LIMPIA[datos$autoref == "NO"] <- 0</pre>
```

Validar nueva variable autorreferida

Valores unicos de la nueva variable, solo se acepta 1,0 y NA

```
unique(datos$AUTOREF_LIMPIA)
## [1] NA "1" "0"
```

Conteo de la variable

```
datos %>% count(AUTOREF_LIMPIA)
```

```
## AUTOREF_LIMPIA n
## 1 0 7563
## 2 1 12947
## 3 <NA> 53896
```

Variable edocivil_descripcion: parte 1

Paso 1: extraer valores unicos.

unique(datos\$edocivil_descripcion)

```
## [1] "soltera"
## [4] "divorciada"
## [7] "casada (o)"
## [10] "Viuda"
## [13] "viuda"
## [16] "divociada"
## [19] "CASADA"
## [22] "soltero(a)"
## [25] "union libre"
## [28] "Casada"
## [31] "UNIÓN LIBRE"
## [34] "DIVORCIADO/A"
## [37] "Sin especificar"
```

```
"unión libre"
"N/E"
"separada"
"SOLTERA"
"S/INFORMACION"
"NO RESPONDE"
"SEPARADA"
"divorciado(a)"
"separado (a)"
"Divorciada"
"SOLTERO/A"
"casada"
```

```
"casado(a)"
"no sabe/ sin respuesta"
"Soltera"
"VIUDA"
"SE DESCONOCE"
"UNION LIBRE"
"DIVORCIADA"
"viudo"
NA
"Separada"
"SEPARADO/A"
"SIN ESPECIFICAR"
```

Variable edocivil_descripcion: parte 2

Paso dos: pasar a mayúsculas.

```
# con signo "$" se especifica una variable
datos$edocivil_descripcion<- toupper(datos$edocivil_descripcion)</pre>
```

vemos valores unicos ahora que estan en mayúsculas

unique(datos\$edocivil_descripcion)

```
"UNTÓN LIBRE"
                                                            "CASADO(A)"
    Γ1] "SOLTERA"
                                  "N/F"
## [4] "DIVORCIADA"
                                                            "NO SABE/ SIN RESPUESTA"
  [7] "CASADA (0)"
                                  "SFPARADA"
                                                            "VTUDA"
                                  "SE DESCONOCE"
                                                            "DIVOCIADA"
   [10] "S/INFORMACION"
   [13] "NO RESPONDE"
                                  "UNION LIBRE"
                                                            "CASADA"
## [16] "SOLTERO(A)"
                                  "DIVORCIADO(A)"
                                                            "VTUDO"
## [19] "SEPARADO (A)"
                                                            "SOLTERO/A"
                                  NA
                                  "DIVORCIADO/A"
                                                            "SIN ESPECIFICAR"
## Γ221 "SEPARADO/A"
```

Variable edocivil_descripcion: parte 3

Paso 3:reclasificar obsevaciones con case_when cuando no se espeficica alguna obsevación por default se pasa a NA o valor perdido.

```
datos<- datos %>%
 mutate(
    EDO_CIVIL_DESCRIPCION_LIMPIA = case_when(
     edocivil_descripcion == "CASADO(A)" ~ "CASADA",
     edocivil_descripcion == "SOLTERA" ~ "SOLTERA",
     edocivil_descripcion == "UNIÓN LIBRE" ~ "UNIÓN LIBRE",
     edocivil_descripcion == "DIVORCIADA" ~ "DIVORCIADA",
     edocivil_descripcion == "CASADA (0)" ~ "CASADA",
     edocivil_descripcion == "UNION LIBRE" ~ "UNIÓN LIBRE",
     edocivil_descripcion == "SOLTERO(A)" ~ "SOLTERA",
     edocivil_descripcion == "DIVORCIADO(A)" ~ "DIVORCIADA",
     edocivil_descripcion == "VIUDO" ~ "VIUDA",
     edocivil_descripcion == "SEPARADO (A)" ~ "SEPARADA",
      edocivil_descripcion == "SOLTERO/A" ~ "SOLTERA",
      edocivil_descripcion == "CASADA" ~ "CASADA",
     edocivil_descripcion == "SOLTERO(A)" ~ "SOLTERA",
     edocivil_descripcion == "SEPARADO/A" ~ "SEPARADA",
     edocivil_descripcion == "DIVORCIADO/A" ~ "DIVORCIADA"
```

Esta crea una nueva variable llamada EDO_CIVIL_DESCRIPCION_LIMPIA se añadé a la base en la última columna.

Variable edad

Se calcula las siguientes métricas:

```
summary(datos$edad)

## Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
## 11.00 21.00 25.00 25.73 30.00 53.00 20
```

El promedio de edad de las mujeres es de 26, el mínimo 11 y el máximo 53 años.

Parece que esta variable esta limpia.

Se procede a caluclar rangos de edad por diversión:

Conteo de variable edad

```
datos %>% count(RANGO_EDAD)
```

```
## RANGO_EDAD n
## 1 18 a 25 AÑOS 37980
## 2 26 a 35 AÑOS 26155
## 3 36 a 45 AÑOS 6504
## 4 MAYOR DE 46 AÑOS 46
## 5 MENOS DE 17 AÑOS 2201
## 6 <NA> 1520
```

Hay 1,520 datos perdidos.