



UNIVERSIDAD
NACIONAL DEL OESTE

Explotación de Datos

ACTIVIDAD N^o 1

Correlación de Tasa de Movilidad con Casos

PROFESORES:

*Dejean, Gustavo
Españadero, Juan
Mendoza, Dante*

INTEGRANTES GRUPO B:

*Benitez, Nicolas
Garcia Ravlic, Ignacio Agustin
Rechimon, Pablo Hernan
Rodriguez, Miguel Angel*

FECHA DE ENTREGA:

29 de Agosto de 2020

Índice

1	Resumen	2
2	Introducción	3
2.1	Problemática	3
2.2	Datos a utilizar	3
2.3	Objetivo	3
3	Desarrollo	4
3.1	Análisis de los datos	4
3.2	Preparación de los datos	4
3.3	Visualización de datos	5
3.4	Análisis de las visualizaciones	5
4	Conclusión	7
5	Anexo	8
5.1	Código en R	8

Graficos

Fig. 1	Comparación de Time Series de Movilidad 1	5
Fig. 2	Comparación de Time Series de Movilidad 2	5
Fig. 3	Time Series de Movilidad con aumento de Casos	5
Fig. 4	Dispersion entre Movilidad y Casos	6
Fig. 5	Dispersion entre Movilidad y Casos completo	6
Fig. 6	Comparacion de confirmados c/10.000 habitantes	6

1 Resumen

En base al contexto de pandemia decidimos estudiar la correlación entre la tasa de movilidad diaria y los casos diarios de cada 10000 habitantes para afirmar o negar que *"a mayor movilidad aumentan los casos diarios"*, para esto contamos con 3 datasets, los cuales manipulamos a través del lenguaje R, dando como resultado que podamos afirmar que el aumento de la movilidad provoque, en alta probabilidad, un aumento de casos.

Palabras Clave: *Analisis de Datos - Programacion - Estudio de Correlacion - COVID-19 - Movilidad Peatonal - Diagramas de Dispersion - Time Series*

2 Introducción

2.1 Problemática

Dado el contexto de pandemia, decidimos investigar y analizar la correlación entre la tasa de movilidad y los casos diarios de Covid-19 por cada 10000 habitantes en la región sudamericana, específicamente en Chile, Brasil y Argentina.

2.2 Datos a utilizar

Para esto contamos con tres datasets:

- Tasa de Movilidad Mundial en la Pandemia, provisto por Apple
- Confirmados Globales de COVID-19, provisto por Hopkins University
- Habitantes del mundo, provisto por las Organización de las Naciones Unidas

2.3 Objetivo

Confirmar que a mayor tasa de movilidad, mayor cantidad de contagios diarios de Covid-19.

3 Desarrollo

3.1 Análisis de los datos

Analizamos los distintos dataset para verificar la correcta carga de datos y poder empezar a prepararlos. Para esto utilizamos las siguientes funciones:

- `str()`
- `colnames()`
- `View()`
- `head()`
- `tail()`

A partir de esto pudimos observar que algunos se encontraban en formato "horizontal", que tenían columnas que no nos interesaban y que los tipos de datos no eran los apropiados, tales observaciones las tuvimos en cuenta a la hora de preparar los datos.

3.2 Preparación de los datos

Preparamos los datos en base a nuestros requerimientos. Para ello, obtenemos datos estadísticos, modificamos el formato de la tabla de "horizontal" a "vertical", los tipos de datos, eliminamos y renombramos columnas, filtramos los datos creando subconjuntos de interés. Para esto se utilizaron las funciones:

- `summary()`
- `sd()`
- `gather()`
- `setnames()`
- `as.Date()`, `as.factor()`, `as.character()`
- `select()`
- `subset()`
- `sqldf()`

Es importante aclarar que para esto hubo que importar las bibliotecas *RSQLite* - *sqldf* - *lubridate* - *tidyr* - *dplyr*. También tuvimos que calcular los Casos Diarios por cada 10.000 Habitantes, para esto usamos la estructura iterativa `while` para calcular esa diferencia entre día y día.

Una vez terminada la preparación nos encontrábamos con un *data frame* final, el cual contenida la fecha y de los países de interés sus casos nuevos cada día y su tasa de movilidad, por ende nos encontrábamos en posición de hacer un análisis a través de las visualizaciones.

3.3 Visualización de datos

Para la visualización de los datos, hemos utilizado las siguientes bibliotecas:

- ggplot2 para graficar scatterplots.
- plotly para gráficos dinámicos.

3.4 Análisis de las visualizaciones

- En la Fig. 1 se puede ver la evolución de la Tasa de Movilidad en el Tiempo, precisamente de los países Argentina, Brasil y Chile. Podemos observar que en un primer segmento la tasa se encuentra arriba del 60% y luego tiene una notable caída de la cual aún se está en recuperación, esto producto de las medidas de prevención tomadas por los diferentes países. A partir de esto podemos decir que las medidas fueron más duras en Argentina y más leves en Brasil.

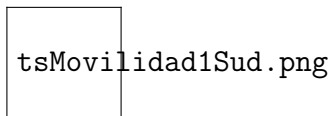


Fig. 1: Comparación de Time Series de Movilidad 1

- En la Fig. 2 se puede observar lo mismo, pero realizado a través de plotly, lo cual permite interactuar.

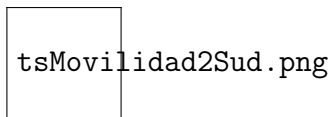


Fig. 2: Comparación de Time Series de Movilidad 2

- En la Fig. 3 podemos observar que hay 3 dimensiones a partir de los datos de Argentina únicamente, tenemos la evolución de la movilidad y los casos en el tiempo. Podemos comprobar que a mayor fecha, mayor movilidad y junto a esto, mayor cantidad de casos.

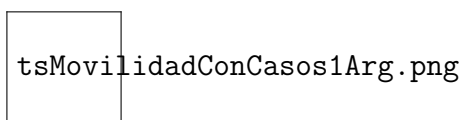


Fig. 3: Time Series de Movilidad con aumento de Casos

- En la Fig. 4 tenemos la dispersión entre la tasa de movilidad y los casos cada 10.000 habitantes, podemos ver que en cada país la correlación es diferente, el caso de Chile se ve un tanto especial, ya que la movilidad se ve casi constante

pero aun así tiene altos niveles de casos diarios, en cambio entre Brasil y Argentina vemos que la "curva" es bastante similar, donde vemos que el aumento de la movilidad va acompañado del aumento de casos.

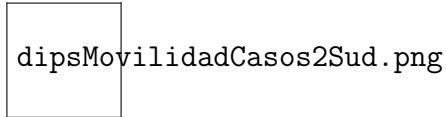


Fig. 4: Dispersion entre Movilidad y Casos

- En la Fig. 5 tenemos una extensión de la figura anterior donde podemos observar unos outliers pertenecientes a los datos de Chile:

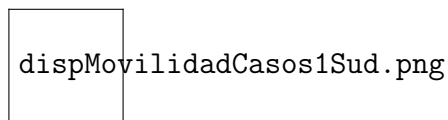


Fig. 5: Dispersion entre Movilidad y Casos completo

- Por último, en la Fig. 6 tenemos la evolución en el tiempo de cantidad de casos confirmados cada 10.000 habitantes desde el 12/03/2020 en Argentina, Chile y Brasil, aquí también presentes los ya nombrados outliers de Chile, pero además podemos ver más picos bruscos tanto en Chile como en Brasil.

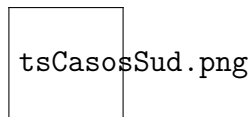


Fig. 6: Comparacion de confirmados c/10.000 habitantes

4 Conclusión

Se puede observar en las figuras, sobre todo en los gráficos de dispersión de los tres países que la relación entre la tasa de movilidad, los casos diarios por 10000 habitantes y su evolución en el tiempo va siguiendo la misma tendencia en los tres países, diagonal hacia arriba, o sea a medida que aumenta la tasa, aumentan los casos. Por lo cual podemos afirmar que el aumento de la movilidad provoque, en alta probabilidad, un aumento de casos.

5 Anexo

5.1 Código en R

```

1 # Correlacion entre Casos Diarios Relativos y Movilidad urbana en
  Argentina vs Chile vs Brasil
2
3 # creado: 2020-08-23    v.    2020-8-28
4 # ultima modificacion: se hicieron los resúmenes y se terminaron
  las vistas
5
6 # Autor: GAD, Ignacio Garcia, Pablo Rechimon, Miguel Rodriguez,
  Nicolas Benitez
7
8 # archivo MOVILIDAD:  https://covid19.apple.com/mobility
9 # archivo HABITANTES: https://drive.google.com/file/d/1
  wi9LrbbJXqwmNhnTSvJT1pcSEE_z-nd5/view?usp=sharing
10
11 #####
12
13
14 ##### IMPORTAR BIBLIOTECAS
15 check_packages <- function(packages) {
16   if (all(packages %in% rownames(installed.packages())) {
17     TRUE
18   } else{
19     cat(
20       "Instalar los siguientes packages antes de ejecutar el
  presente script\n",
21       packages[!(packages %in% rownames(installed.packages()))],
22       "\n"
23     )
24   }
25 }
26 packages_needed <- c("ggplot2", "plotly", "sqldf", "lubridate",
27   "tidyr", "data.table", "readr", "dplyr", "
  RColorBrewer" )
28 check_packages(packages_needed)
29
30 library(gsubfn)
31 library(proto)
32 library(RSQLite)
33 library(sqldf)
34 library(data.table)
35 library(lubridate)
36 library(tidyr)
37 library(readr)
38 library(ggplot2)
39 library(plotly)
40 library(dplyr)
41 library(RColorBrewer)
42
43 options(scipen = 6) #para evitar notacion cientifica
44

```

```

45
46
47 ##### LEER DATOS DEL COVID
48 # URL con datos del COVID-19
49 URL      <- "https://raw.githubusercontent.com/CSSEGISandData/
      COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/"
50
51 # time_series_covid19_confirmed_global.csv     este es el archivo
      que a leer
52 url_archivo <- paste(URL,"time_series_covid19_confirmed_global.csv
      ", sep = "")
53
54 COVID_19_h  <- read.csv(url_archivo, sep = ",", header = T)
55
56 #Analizamos lo leído y corroboramos que se leyó bien
57 str(COVID_19_h)
58
59
60 ##### PREPARAR SUBSET DE DATOS PARA ARGENTINA,
      BRASIL y CHILE
61 colnames(COVID_19_h)
62 #Eliminamos columnas que no nos interesan
63 COVID_19_h <- select(COVID_19_h, -c(Lat, Long, Province.State))
64
65 #Cambiamos el nombre de una columna a algo mas entendible
66 setnames(COVID_19_h, "Country.Region", "pais")
67
68 #Pasamos de formato horizontal a vertical
69 COVID_19 <- COVID_19_h %>% gather(fecha, acumCasos, 2:ncol(COVID_19
      _h))
70
71
72 str(COVID_19)
73 #Cambiamos el formato de la fecha
74 COVID_19$fecha <- as.Date(COVID_19$fecha, format = "%m.%d.%y")
75
76
77 #Quedarse solo con datos de ARG
78 ARG_target = "Argentina"
79 C19_ARG <- subset(COVID_19, pais == ARG_target)
80
81 #Quedarse solo con datos de BR
82 BR_target = "Brazil"
83 C19_BR <- subset(COVID_19, pais == BR_target)
84
85 #Quedarse solo con datos de CHL
86 CHL_target = "Chile"
87 C19_CHL <- subset(COVID_19, pais == CHL_target)
88
89 C19_Sudamerica <- sqldf("
90     SELECT arg.fecha fecha, arg.acumCasos acumCasosARG,
91           br.acumCasos acumCasosBR, chl.acumCasos acumCasosCHL
92     FROM C19_ARG arg, C19_BR br, C19_CHL chl
93     WHERE arg.fecha==br.fecha AND arg.fecha==chl.fecha")
94
95 #Analizamos la salida del join entre los casos de cada pais

```

```

96 str(C19_Sudamerica)
97 C19_Sudamerica
98
99
100 ##### LEER DATOS DE HABITANTES DEL MUNDO
101 data_habitantes <- read.csv("C:/work/Explotacion/Actividad1/
    habitantesMundo.csv",
102                             sep=";", stringsAsFactors=TRUE)
103
104 #Analizamos lo leído y corroboramos que se leyó bien
105 str(data_habitantes)
106
107
108
109 ##### PREPARAR DATOS DE HABITANTES
110 #Quedarse solo con datos de ARG, BR y CHL
111 SUD_habitantes <- subset(data_habitantes, pais == ARG_target |
112                           pais == BR_target |
113                           pais == CHL_target)
114
115 #Eliminar columnas extras
116 colnames(SUD_habitantes)
117 SUD_habitantes <- select(SUD_habitantes, -c(Variant, Index, Notes,
    codigo, Type))
118
119 #Analizamos la salida
120 str(SUD_habitantes)
121 SUD_habitantes
122
123
124 ##### CALCULAR CASOS DIARIOS EN RELACION A LOS
    HABITANTES
125 #Agregar columnas para los casos
126 C19_Sudamerica$casosARG <- 0
127 C19_Sudamerica$casosBR <- 0
128 C19_Sudamerica$casosCHL <- 0
129
130 C19_Sudamerica
131
132 #Contador y tope
133 i <- 1
134 max <- nrow(C19_Sudamerica)
135
136 print(max)#Visualizar las filas
137
138 #Recorrer y calcular casos diarios de cada 10.000 habitantes
139 while (i<=max) {
140     if(C19_Sudamerica$acumCasosARG[i] > 0){
141         C19_Sudamerica$casosARG[i] <- (C19_Sudamerica$acumCasosARG[i] -
            C19_Sudamerica$acumCasosARG[i-1]) * 10000 / SUD_habitantes$
            cantidad[1]
142     }
143     if(C19_Sudamerica$acumCasosBR[i] > 0){
144         C19_Sudamerica$casosBR[i] <- (C19_Sudamerica$acumCasosBR[i] -
            C19_Sudamerica$acumCasosBR[i-1]) * 10000 / SUD_habitantes$
            cantidad[2]

```

```

145 }
146 if(C19_Sudamerica$acumCasosCHL[i] > 0){
147   C19_Sudamerica$casosCHL[i] <- (C19_Sudamerica$acumCasosCHL[i] -
148     C19_Sudamerica$acumCasosCHL[i-1]) * 10000 / SUD_habitantes$
149     cantidad[3]
150 }
151 i <- i+1
152 }
153
154 C19_Sudamerica
155
156 ##### LEER DATOS DE MOVILIDAD
157 data_movilidad_h <- read_csv("apblemobilitytrends-2020-08-26.csv",
158   locale = locale(grouping_mark = ""))
159
160 #Analizamos lo leído y corroboramos que se leyó bien
161 str(data_movilidad_h)
162
163 ##### PREPARAR DATOS
164 colnames(data_movilidad_h)
165 setnames(data_movilidad_h, "region", "pais")
166 setnames(data_movilidad_h, "transportation_type", "transporte")
167
168 data_movilidad_h <- select(data_movilidad_h, -c(geo_type,
169   alternative_name, country, 'sub-region'))
170
171 data_movilidad_h$pais <- as.factor(data_movilidad_h$pais)
172 data_movilidad_h$transporte <- as.factor(data_movilidad_h$
173   transporte)
174
175 #Pasar a formato vertical
176 movilidad <- data_movilidad_h %>% gather(fecha, tasa, 3:ncol(
177   data_movilidad_h))
178
179 #Formatear fecha
180 movilidad$fecha <- as.Date(as.character(movilidad$fecha))
181
182 #Corroboramos que los cambios se efectuaron bien
183 str(movilidad)
184
185 #Quedarse solo con datos de ARG, BR y CHL
186 MOV_ARG <- subset(movilidad, transporte == "walking" & pais == ARG_
187   target)
188 MOV_BR <- subset(movilidad, transporte == "walking" & pais == BR_
189   target)
190 MOV_CHL <- subset(movilidad, transporte == "walking" & pais == CHL_
191   target)
192
193 #Juntar los subsets de movilidad
194 MOV_Sudamerica <- sqldf("
195   SELECT arg.fecha fecha, arg.tasa tasaARG, br.tasa tasaBR, chl.
196   tasa tasaCHL

```

```

191 FROM MOV_ARG arg, MOV_BR br, MOV_CHL chl
192 WHERE arg.fecha == br.fecha AND arg.fecha == chl.fecha")
193
194 #Analizamos que el join no afecto y salio bien
195 str(MOV_Sudamerica)
196 MOV_Sudamerica
197
198
199
200 ##### JOIN ENTRE MOVILIDAD Y CASOS
201 SUD_MovC19 <- sqldf("
202     SELECT mov.fecha fecha, c19.casosARG, mov.tasaARG, c19.casosBR,
203           mov.tasaBR, c19.casosCHL, mov.tasaCHL
204     FROM C19_Sudamerica c19, MOV_Sudamerica mov
205     WHERE mov.fecha == c19.fecha")
206
207 #Analizamos que el join no afecto y salio bien
208 str(SUD_MovC19)
209 SUD_MovC19
210
211 #Calculamos y analizamos medidas estadísticas
212 desviaciones <- c( #Desviaciones típicas
213 sd(SUD_MovC19$casosARG, na.rm = T),
214 sd(SUD_MovC19$tasaARG, na.rm = T),
215 sd(SUD_MovC19$casosBR, na.rm = T),
216 sd(SUD_MovC19$tasaBR, na.rm = T),
217 sd(SUD_MovC19$casosCHL, na.rm = T),
218 sd(SUD_MovC19$tasaCHL, na.rm = T))
219
220 desviaciones
221
222
223 resumen <- summary(SUD_MovC19) #Minimos, Maximos, Quartiles y
224 Promedio
225 resumen
226
227
228
229 ##### GUARDAR DATOS
230 write.csv2(C19_Sudamerica, "datos_casosXDia.csv", row.names =
231 FALSE, fileEncoding = "UTF-8")
232 write.csv2(MOV_Sudamerica, "datos_movilidadXDia.csv", row.names =
233 FALSE, fileEncoding = "UTF-8")
234 write.csv2(SUD_MovC19, "datos_correlacion.csv", row.names = FALSE,
235 fileEncoding = "UTF-8")
236
237 ##### GENERAR FIGURA
238 colores <- c("blue", "lightblue", "yellow", "orange", "red", "
239 darkred")
240 #Time Series de Movilidad en Arg, Br y Chl con qplot

```

```

241 fig1 <- ggplot(SUD_MovC19, aes(x = fecha, y = tasaARG)) +
242   geom_line(aes(y = tasaBR), color="darkgreen") +
243   geom_line(aes(y = tasaCHL), color="darkred") +
244   geom_line(size = 0.6) +
245   ggtitle(paste("COVID-19 - Movilidad", sep = "")) +
246   scale_x_date(date_breaks = "7 day", date_labels = "%d %b") +
247   theme(plot.title = element_text(lineheight = 1, face = 'bold')) +
248   ylab("Tasa de Movilidad") +
249   xlab("") +
250   labs(caption = "Fuente de los datos: apple.com/covid19/mobility")
251   +
252   theme_minimal() +
253   theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 7))
254 fig1
255
256 #Time Series de Movilidad en Arg, Br y Chl con plotly
257 fig2 <- plot_ly(SUD_MovC19, x = ~fecha, y = ~tasaARG,
258   name = 'Argentina', type = 'scatter', mode = 'lines') %>%
259   add_trace(y = ~tasaBR, name = 'Brasil', mode = 'lines') %>%
260   add_trace(y = ~tasaCHL, name = 'Chile', mode = 'lines') %>%
261   layout(title = "Time Series de Movilidad",
262     xaxis = list(title = "Fecha", type = "date",
263       tickmode = "linear", tick0 = min(SUD_MovC19$
264         fecha),
265         tickformat = "%d/%m", dtick = 86400*10000,
266         tickangle = 75),
267     yaxis = list(title = "Tasa de Movilidad", tickangle =
268       -45))
269 fig2
270
271 #Time Series de Movilidad en Arg, con color en relacion a cantidad
272 #de casos con plotly
273 fig3 <- qplot(data = SUD_MovC19, x = fecha, y = tasaARG, colour =
274   casosARG, geom = "line") +
275   geom_line(size = 1) +
276   scale_colour_gradient2(low = "lightblue", mid="orange", high = "
277     red",
278     midpoint = max(SUD_MovC19$casosARG)/2) +
279   ggtitle("Time Series de Movilidad con aumento de Casos") +
280   scale_x_date(date_breaks = "7 day", date_labels = "%d/%m") +
281   theme(plot.title = element_text(lineheight = 1, face = 'bold')) +
282   ylab("Tasa de Movilidad") +
283   xlab("") +
284   labs(caption = "Fuente de los datos: apple.com/covid19/mobility")
285   +
286   theme_minimal() +
287   theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 7))
288 fig3
289
290 #Movilidad con calor de Arg con plotly

```

```

287 fig4 <- plot_ly(data = SUD_MovC19, x = ~fecha, y = ~tasaARG, name
  = "Argentina",
288                 type = 'scatter', mode = 'markers', color = ~SUD_
  MovC19$casosARG "Casos Arg",
289                 size = 1, colors = colores) %>%
290 layout(title = "Time Series de Movilidad con aumento de Casos",
291         xaxis = list(title = "", type = "date",
292                     tickmode = "linear", tick0 = min(SUD_MovC19$
  fecha),
293                     tickformat = "%d/%m", dtick = 86400*10000,
  tickangle = 75,
294                     range=c(min(SUD_MovC19$fecha), max(SUD_MovC19
  $fecha))),
295         yaxis = list(title = "Tasa de Movilidad", tickangle =
  -45))
296 fig4
297
298
299
300 #Filtramos puntos antes de que llegue el primer caso
301 filtered_SUD_MovC19 <- sqldf("SELECT * FROM SUD_MovC19
302                               WHERE casosARG > 0 AND casosBR > 0 AND
  casosCHL > 0")
303
304
305 #Correlacion entre Movilidad y Casos en Arg
306 fig5 <- plot_ly(data = filtered_SUD_MovC19, x = ~casosARG, y = ~
  tasaARG, name = "Argentina",
307                 type = 'scatter', mode = 'markers') %>%
308 layout(title = "Dispercion entre Tasa de Movilidad y Casos",
309         xaxis = list(title = "Casos c/10.000 hab"),
310         yaxis = list(title = "Tasa de Movilidad", tickangle =
  -45))
311 fig5
312
313
314 #Correlacion entre Movilidad y Casos en Arg, Br y Chl
315 fig6 <- plot_ly(data = filtered_SUD_MovC19, x = ~casosARG, y = ~
  tasaARG, name = "Argentina",
316                 type = 'scatter', mode = 'markers') %>%
317 add_trace(x = ~casosBR, y = ~tasaBR, name = 'Brasil', mode = '
  markers') %>%
318 add_trace(x = ~casosCHL, y = ~tasaCHL, name = 'Chile', mode = '
  markers') %>%
319 layout(title = "Dispercion entre Tasa de Movilidad y Casos",
320         xaxis = list(title = "Casos c/10.000 hab"),
321         yaxis = list(title = "Tasa de Movilidad", tickangle =
  -45))
322 fig6
323
324
325 #Correlacion entre Movilidad y Casos en Arg, Br y Chl filtrando
  tasa mayor a 100 y casos mayor a 4
326 fig7 <- plot_ly(data = filtered_SUD_MovC19, x = ~casosARG, y = ~
  tasaARG, name = "Argentina",
327                 type = 'scatter', mode = 'markers') %>%

```

```

328 add_trace(x = ~casosBR, y = ~tasaBR, name = 'Brasil', mode = '
      markers') %>%
329 add_trace(x = ~casosCHL, y = ~tasaCHL, name = 'Chile', mode = '
      markers') %>%
330 layout(title = "Dispercion entre Tasa de Movilidad y Casos",
331         xaxis = list(title = "Casos c/10.000 hab", range= c(0, 4))
      ,
332         yaxis = list (title = "Tasa de Movilidad", tickangle =
      -45, range= c(0, 100)))
333 fig7
334
335
336 #Time Series de casos de Covid-19 en Arg, Br y Chl
337 fig8 <- plot_ly(data = filtered_SUD_MovC19, x = ~fecha, y = ~
      casosARG, name = "Argentina",
338                 type = 'scatter', mode = 'lines') %>%
339 add_trace(y = ~casosBR, name = 'Brasil', mode = 'lines') %>%
340 add_trace(y = ~casosCHL, name = 'Chile', mode = 'lines') %>%
341 layout(title = "Time Series de Confirmados c/10.000 habitantes",
342         xaxis = list(title = "Fecha", type = "date",
343                     tickmode = "linear", tick0 = min(filtered_SUD
      _MovC19$fecha),
344                     tickformat = "%d/%m", dtick = 86400*10000,
      tickangle = 75),
345         yaxis = list (title = "Nro de Casos", tickangle = -45))
346 fig8

```