

Apéndices del Libro Pobreza y Desigualdad en R

2022-02-19

Contents

1	Bienvenidos	5
	Sobre este Proyecto:	5
	El Libro	5
2	Capítulo 2	7
	Set Inicial	7
2.1	Introducción: ejemplo Brasil	8
2.2	Histogramas	13
2.3	Función de distribución	17
2.4	Pareto	18
2.5	Box-plot	20
2.6	Curva de Lorenz	21
2.7	Curva Generalizada de Lorenz	23
2.8	Curva de Incidencia	24
3	Capítulo 3	29
	Set Inicial	29
3.1	Cociente de cuantiles	30
3.2	Replicar programa ratq51	31
3.3	Replicar programa gcuán	36
3.4	Tamaño de los hogares	39
3.5	La distribución intrahogar	41
3.6	Empleo de ponderadores	43

3.7	Diseño muestraL	46
3.8	Fuentes de ingreso	46
4	Capítulo 4	49
	Set Inicial	49
4.1	Indicador FGT	50
4.2	Computar FGT	51
4.3	Pobreza relativa	53
4.4	Descomposición regional de la pobreza	54
4.5	Pobreza según consumo e ingreso	55
4.6	Pobreza por edad	59
4.7	Significatividad estadística	61
5	Capítulo 5	65
	Set Inicial	65
5.1	Pobreza Multidimensional	66
5.2	Índice Bourguignon y Chakravarty (BC) - Pobreza Multidimensional	69
5.3	Índice Alkire y Foster (AF) - Pobreza Multidimensional	71
5.4	Perfiles de Pobreza	72
5.5	Perfiles de Pobreza Condicionados	74

Chapter 1

Bienvenidos

Sobre este Proyecto:

Nos propusimos transcribir al lenguaje R y Python los apéndices del libro “Pobreza y Desigualdad en América Latina” de Gasparini, Cicowiez y Sosa Escudero que originalmente fueron escritos para Stata y que permitían replicar los datos e información presentados por los autores en el texto. Cada capítulo consta de un apéndice con códigos que permiten llevar a la práctica los conceptos desarrollados. Aquí los traducimos a nuevos lenguajes y los presentamos en un formato amigable para permitir que un público más amplio y de diversas disciplinas pueda aprovecharlos.

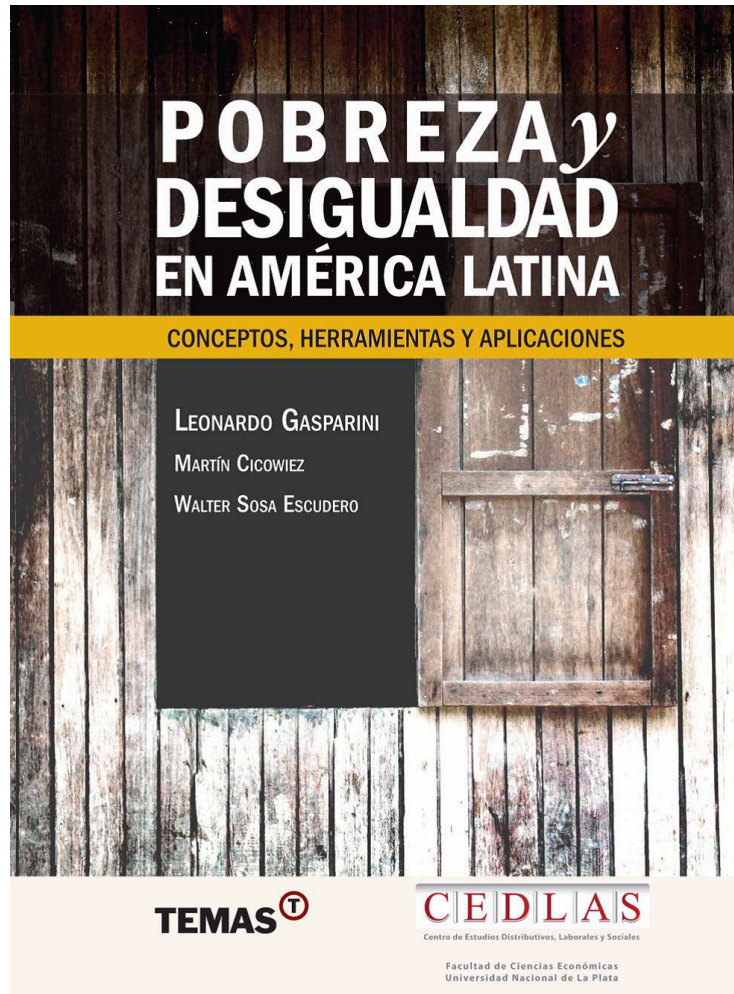
El objetivo de este mini-proyecto no es otro que poner a disposición de un público más amplio estas herramientas útiles y mantener actualizado un material único que permite adentrarse y trabajar sobre la temática de pobreza y desigualdad. Es por eso que este material es de carácter complementario a las explicaciones y detalles conceptuales que se presentan en el libro de texto y los apéndices.

El Libro

Si esta es la primera vez que te encontrás con este libro, antes de empezar con los códigos, dejarnos presentartelo.

El propósito del libro es ayudar al lector interesado en América Latina a que recorra el arduo camino entre los datos y el reporte de resultados rigurosos que puedan contribuir al debate sobre la pobreza y la desigualdad en América Latina. El volumen pone al alcance del lector un conjunto de instrumentos que lo motiven hacia la investigación empírica, y que le permitan producir resultados de la manera más rigurosa posible, para así contribuir a los objetivos últimos de

explicar y cambiar la realidad social de la región. Las discusiones conceptuales y analíticas son ilustradas con ejemplos concretos contruados con datos de los países de la región.



Chapter 2

Capítulo 2

Apéndices del Libro Pobreza y Desigualdad en R

Escrito por: Cristian Bonavida

Last Update: 02/7/2021

Códigos escritos en base a los apéndices del libro “Pobreza y Desigualdad en América Latina” de Gasparini, Cicowiez y Sosa Escudero. El objeto de este material es reproducir la rutina de códigos para STATA presentada en el libro al lenguaje R. Este material es solo de caracter complementario a las explicaciones y detalles conceptuales que se presentan en el libro de texto y los apéndices

Set Inicial

Cargo las librerías, limpio environment, defino el path y atajo para función paste

```
library(dplyr)
library(tidyverse) # Data wrangling
library(tidygraph)
library(readxl)
library(ggplot2)
library(foreign)
library(Hmisc)
```

```
rm(list=ls())      #empiezo limpiando todo

"%+%" <- function(x,y) paste(x,y,sep = "")      # defino un shorcut parar concat de te
data_dir <- "C:/Users/HP/Desktop/CEDLAS - UNLP/Apendices en R/Material libro/encuestas"
```

2.1 Introducción: ejemplo Brasil

###- (pág. 72)

El siguiente código es introductorio y busca guiar al lector en la sintaxis básica a utilizar, mostrando a modo de ejemplo cómo replicar los resultados correspondientes al cuadro 2.1 del texto.

Para ello como al inicio de cada capítulo el primer paso es obtener base de datos, en ese caso la versión procesada de la PNAD (Pesquisa Nacional por Amostra de Domicílios) de Brasil para el año 2007, accesible desde el link de descarga indicado en la sección anterior: datos y set up básico. (Añadir sección y referir allí)

Una vez descargada la encuesta y guardada en el directorio que el lector indicó al definir el objeto `data_dir` se carga la encuesta que se encuentra en formato STATA utilizando, de la librería `foreign`, el comando `read.dta` que permite importar en R bases de datos que están en formatos compatibles con otros programas o lenguajes. Las encuestas de hogares procesadas que se utilizan a lo largo del libro solo contienen observaciones que denominamos coherentes (ver capítulo 3), que en pocas palabras se trata de observaciones válidas que utilizamos en el cálculo de los ingresos familiares. Por ello siempre que se cargue una base luego se debe filtrarla a partir de la columna `cohh`, para quedarse solo con estas observaciones.

Las encuestas de hogares, al igual que cualquier otra base de datos, se organizan en R como objetos tipo tabla que reciben el nombre de `dataframe`, donde las filas representan observaciones o registros y las columnas variables o campos. Con el comando `head` podemos explorar las primeras observaciones de la encuesta (`[,1:9]` indica incluir solo de la columna 1 a 9 para que el output no sea demasiado extenso)

```
#cargo base
bra07 <- read.dta(data_dir %+% "Bra/2007/Bases/bra07_cedlas.dta") %>%
  filter(cohh==1)

head(bra07[,1:9])
```



```
##      edad relacion tarea hstrp pondera pais  ano  encuesta id
## 1    73          1   NA    NA    239  bra 2007 PNAD-2007 14
## 2    52          1   NA    NA    238  bra 2007 PNAD-2007 16
## 3    10          4   NA    NA    238  bra 2007 PNAD-2007 16
## 4    34          1   NA    NA    238  bra 2007 PNAD-2007 34
## 5    11          3   NA    NA    238  bra 2007 PNAD-2007 34
## 6    20          4   NA    NA    238  bra 2007 PNAD-2007 34
```

A su vez, por tratarse de una encuesta, cada observación representa a varios individuos, tantos como indica el factor de expansión o variable de ponderación. En nuestro caso, todas las encuestas que utilizaremos contienen una variable de nombre “*pondera*” que almacena el factor de expansión. Para más detalles sobre el uso de ponderadores, consultar más adelante la sección 3.6 del capítulo 3.

Para expandir las observaciones por su peso muestral debemos emplear alguna de las librerías disponibles creadas para tal fin, ya que “por default” los comandos base de R no estiman estadísticos ponderados, como es el caso por ejemplo del comando `summary`. Si se comparan los valores que arroja respecto a los valores de los mismos estadísticos pero ponderados, es evidente que existen diferencias (excepto obviamente para el mínimo y el máximo). Existen varias librerías que nos permiten obtener los mismos descriptivos ponderados, aquí se presentan los resultados con el uso de `Hmisc` pero también puede explorarse de forma similar las librerías `TAM` y `srvy`.

```
#summary no arroja estadísticas descriptivas ponderadas
summary(bra07$ipcf)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
##      0.0   157.5   310.0   559.1   590.0 66000.0    1057
```

```
#replicar la información que nos da summary pero ponderandola - usando libreria Hmisc
wtd.mean(bra07$ipcf, w=bra07$pondera)
```

```
## [1] 574.3455
```

```
wtd.quantile(bra07$ipcf, probs =c(0.25, 0.50, 0.75) , w=bra07$pondera)
```

```
##      25%   50%   75%
## 165.0 330.0 621.5
```

```
min(bra07$ipcf, na.rm =T)
```

```
## [1] 0
```

```
max(bra07$ipcf, na.rm=T)
```

```
## [1] 66000
```

Mientras que la media no pondera es de 555.5 al ponderar el valor estimado es 569.6. Lo mismo ocurre para la mediana, el primer y el tercer cuartil de la distribución.

Una opción recomendada para usuarios más familiarizados con la sintaxis es definir una función que nos calcula e imprima los resultados para estos estadísticos ponderados. Las funciones son una especie de comando personalizado que nos permiten customizar los cálculos y el output a nuestra necesidad o gusto. (Ver capítulo siguiente sobre uso de funciones)

```
my_weighted_stats <- function(argumento1, argumento2){
  print(paste("Mean: ", wtd.mean(argumento1, w=argumento2), sep=""))
  print(paste("Sd: ", sqrt(wtd.var(argumento1, w=argumento2)), sep=""))
  print(wtd.quantile(bra07$ipcf, probs =c(0.25, 0.50, 0.75) , w=bra07$pondera))
  print(paste("Min: ", min(argumento1, na.rm=TRUE), sep=""))
  print(paste("Max: ", max(argumento1, na.rm=TRUE), sep=""))
}
```

Una vez definida la función debemos llamarla por el nombre indicando sus argumentos. En este caso el primero hace referencia a la variable y el segundo al ponderador. La función replica exactamente la información que arrojaría el comando `summ` en stata ponderando las estimaciones.

```
my_weighted_stats(bra07$ipcf, bra07$pondera)
```

```
## [1] "Mean: 574.345522674979"
## [1] "Sd: 970.243045338968"
##   25%   50%   75%
## 165.0 330.0 621.5
## [1] "Min: 0"
## [1] "Max: 66000"
```

Una de las principales ventajas de las funciones es que nos permiten replicar los resultados para distintas encuestas o también para distintos subconjuntos de los datos. Por ejemplo, en el caso de la PNAD 2007 de Brasil, la variable *region* puede tomar los valores 1, 2, 3, 4 o 5 dependiendo de si la observación corresponde a la región Norte, Nordeste, Sudeste, Sur o Centro-Oeste, respectivamente. Así, las líneas siguientes pueden utilizarse para computar los estadísticos para dichas regiones, replicando aquí las columnas “Norte” y “Nordeste” del cuadro 2.1

```
my_weighted_stats(bra07$ipcf[bra07$region==1], bra07$pondera[bra07$region==1])
```

```
## [1] "Mean: 391.006517522251"
## [1] "Sd: 753.58106580221"
##   25%   50%   75%
## 165.0 330.0 621.5
## [1] "Min: 0"
## [1] "Max: 49591.71875"
```

```
my_weighted_stats(bra07$ipcf[bra07$region==2], bra07$pondera[bra07$region==2])
```

```
## [1] "Mean: 344.654903240737"
## [1] "Sd: 655.945582968748"
##   25%   50%   75%
## 165.0 330.0 621.5
## [1] "Min: 0"
## [1] "Max: 30120"
```

Ahora el argumento no son todas las filas de las columnas *ipcf* y *pondera* sino solo aquellas que cumplen con la condición de pertenecer a la región 1 y 2 en cada caso, lo que se instrumenta con el subscript `[bra07$region==x]`.

Siguiendo con el cuadro 2.1 del texto, si quisiéramos replicar el coeficiente de variación, ahora que ya sabemos cómo obtener estadísticas ponderadas, se vuelve muy sencillo. Solo basta con definir un objeto que *cv* que contenga el cociente entre la media y el desvío estándar.

```
cv=sqrt(wtd.var(bra07$ipcf, w=bra07$pondera)) / wtd.mean(bra07$ipcf, w=bra07$pondera)
cv
```

```
## [1] 1.689302
```

Y si quisiéramos estimarlo para una región en particular solo deberíamos aplicar el subscript de la misma forma que lo hicimos antes en la función.

```
cv= sqrt(wtd.var(bra07$ipcf[bra07$region==1], w=bra07$pondera[bra07$region==1])) /
      wtd.mean(bra07$ipcf[bra07$region==1], w=bra07$pondera[bra07$region==1])
cv
```

```
## [1] 1.927285
```

La primera columna del cuadro 2.1 muestra además la población de referencia o número de observaciones expandidas. Para calcularlo basta con sumar la

columna pondera, es decir sumar a cada persona aumentándola por su factor de expansión. De esta forma vemos que la población de referencia de la PNAD 2007 es 187 millones de personas aproximadamente.

```
#número de observaciones expandidas
sum(bra07$pondera)
```

```
## [1] 187441902
```

El ejemplo del texto finaliza con el cómputo de la pobreza en Brasil para el año 2007, utilizando una línea de pobreza de 130 reales mensuales. Una forma sencilla de computar la proporción de individuos con ingresos mensuales menores a 130 reales es mediante el bloque de código siguiente, donde primero se genera una variable binaria “*pobre*” a partir del comando `ifelse`, que asignará el valor 1 en caso de que el individuo no supere el umbral monetario y 0 en caso contrario, y luego se computa simplemente su media e imprime el resultado redondeado a 2 decimales.

```
#identificar individuos pobres
bra07 <- bra07 %>% mutate(pobre=ifelse(ipcf<129.883, 1, 0))

#total país
pobres_pais = wtd.mean(bra07$pobre, bra07$pondera)
print(paste("shr pobres=", round(pobres_pais, d=2)))
```

```
## [1] "shr pobres= 0.18"
```

El mismo resultado podría obtenerse como el cociente de la suma ponderada de personas pobres sobre el total de población.

```
#otra forma
pobres_pais = sum(bra07$pondera[bra07$pobre==1], na.rm = TRUE)/sum(bra07$pondera)
print(paste("shr pobres=", round(pobres_pais, d=2)))
```

```
## [1] "shr pobres= 0.18"
```

Para calcular la tasa de pobreza para una región en vez de todo el país, basta simplemente con filtrar las observaciones.

```
## Region Norte
pobres_norte = wtd.mean(bra07$pobre[bra07$region==1], bra07$pondera[bra07$region==1])
print(paste("shr pobres=", round(pobres_norte, d=2)))
```

```
## [1] "shr pobres= 0.26"
```

2.2 Histogramas

###- (pág 76)

Siguiendo al apéndice del libro mostramos cómo puede graficarse un histograma de la distribución del ipcf en México para el año 2006 (figuras 2.2 a 2.8 del texto). Al igual que en el ejemplo de Brasil, el lector puede obtener la ENIGH (Encuesta Nacional de Ingresos y Gastos de los Hogares) mexicana de 2006 del sitio web del libro. La misma también cuenta con las variables ipcf y pondera que utilizaremos en lo que resta de este apéndice

```
#carga base
mex06 <- read.dta(data_dir %+% "Mex/2006/bases/mex06_cedlas.dta")

#elimino observaciones incoherentes o con ingreso missing
mex06 <- mex06 %>% filter(cohh==1, !is.na(ipcf))
```

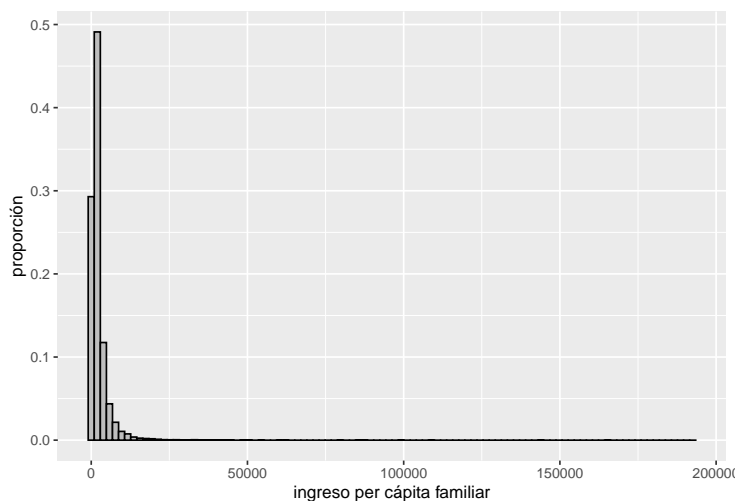
Para graficar optamos por la librería más extendida en uso en R, `ggplot2`, ya que tiene una enorme variedad de gráficos y mantiene un sintaxis común. Para graficar existen básicamente 3 ítems o funciones comunes que todos los gráficos deben tener:

- **`ggplot()`**: Es la función que da inicio al gráfico, indicando que un objeto grafico va a generarse.
 - `data`: como primer argumento debemos pasarle el *dataframe* donde se encuentran las variables a graficar
 - `aes`: es el segundo argumento donde indicamos el o los ejes (x e y) y la diversificación de colores o formas según los grupos en los que se dividen los datos (por ejemplo si fuéramos a graficar líneas de tendencia y quisiéramos obtener una línea de color diferente para cada región)
- **`geom_...()`**: donde los puntos se completan con el tipo de grafico que queremos. Existe una función para cada gráfico. En el caso de un gráfico de líneas será por ejemplo “`geom_line`”
- **customización**: existe múltiples funciones y casi infinitas combinaciones para adaptar el grafico en sus escalas, colores, títulos, dimensiones, fuente de los ejes, leyendas, faceteado, etc.

A continuación, siguiendo estas líneas generales, graficamos un histograma para el ipcf, empleando la función `geom_histogram`. Igual que antes, `[w=pondera]` indica que cada observación de la encuesta debe expandirse según la cantidad de individuos que representa. La opción `bin(100)` especifica que el histograma debe identificar 100 grupos (100 barras). Por último el termino `stat(count) / sum(count)` para estimar valores del eje y en proporciones y no en valores

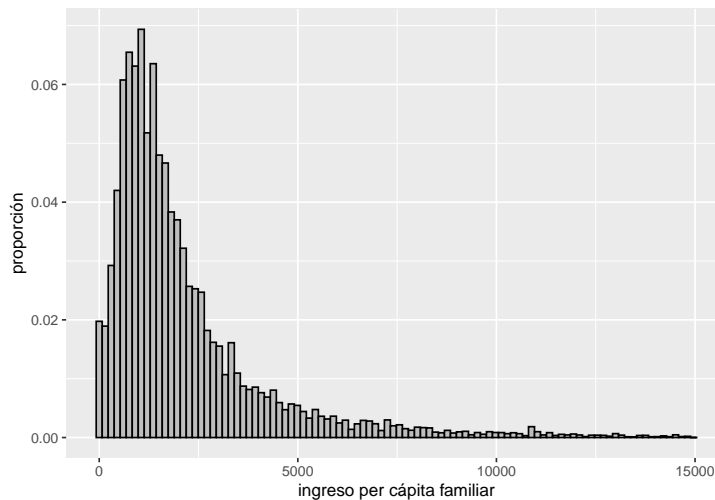
absolutos. Los demás son pequeñas customizaciones de títulos y colores que el lector podrá entender fácilmente.

```
## Figura 2.2 - histograma ipcf
ggplot(data=mex06,
       aes(x=ipcf, weight=pondera)) +
  geom_histogram(bins = 100, aes(y = stat(count) / sum(count)),
                color="black", fill="grey") +
  labs(y="proporción", x="ingreso per cápita familiar")
```



Con las líneas siguientes se filtra la base para individuos con ipcf menor a 1500 para evitar la típica distorsión generada por los valores extremos.

```
## Figura 2.3 - histograma ipcf sin outliers
ggplot(data=mex06 %>% filter(ipcf < 15000),
       aes(x=ipcf, weight=pondera)) +
  geom_histogram(bins = 100, aes(y = stat(count) / sum(count)),
                color="black", fill="grey") +
  labs(y="proporción", x="ingreso per cápita familiar")
```



Las figuras 2.4 y 2.5 del texto pueden replicarse utilizando el bloque de código siguiente. En la primera se grafica el logaritmo de la variable `ipcf`, indicando la función `log` para el eje `x`. En la siguiente se incrementa sucesivamente el número de barras del histograma. Para ello podríamos repetir 4 veces el gráfico cambiando solamente la cantidad de bins. Aquí optamos por presentar un solución más elegante en la que generamos los 4 gráficos a partir de un bucle que itera sobre el número de bins deseado. El lector no familiarizado con los bucles puede optar por la primera opción hasta que en los siguientes capítulos presentemos con más detalle cómo se introducen iteraciones mediante bucles.

```
## Figura 2.4 - histograma logaritmo ipcf
ggplot(mex06,
  aes(x=log(ipcf), weight=pondera)) +
  geom_histogram(bins = 100, aes(y = stat(count) / sum(count)),
    color="black", fill="grey") +
  labs(y="proporción", x="ingreso per cápita familiar")
```

```
## Warning: Removed 1305 rows containing non-finite values (stat_bin).
```

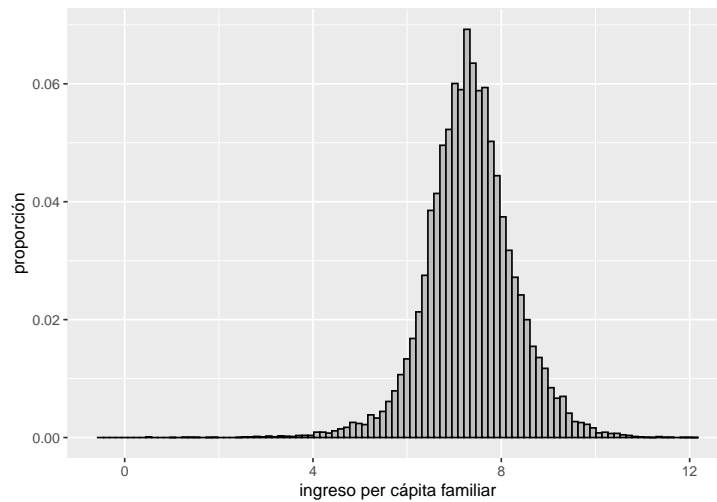


Figura 2.5 - histogramas con diferente cantidad de intervalos

```
n_bins = c(10,50,100,1000)
my_graphs = list()

i=1
for (n in n_bins){

  my_graphs[[i]] <- ggplot(mex06,
                           aes(x=log(ipcf), weight=pondera)) +
    geom_histogram(bins = n, aes(y = stat(count) / sum(count)),
                  color="black", fill="grey") +
    labs(y="proporción", x="ingreso per cápita familiar") +
    ggtitle( n %+" intervalos") +
    theme(plot.title = element_text(size = 12, hjust = 0.5))

  i=i+1
}

g1 <- my_graphs[[1]]
g2 <- my_graphs[[2]]
g3 <- my_graphs[[3]]
g4 <- my_graphs[[4]]

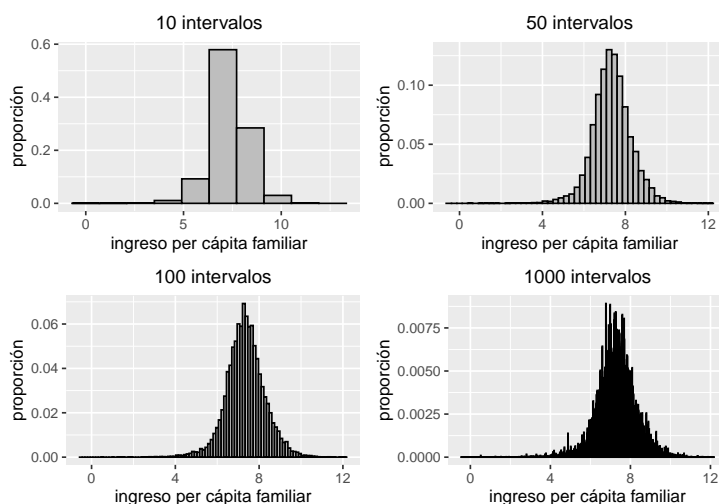
library(gridExtra)
grid.arrange(g1, g2, g3, g4, ncol=2, nrow=2)
```

Warning: Removed 1305 rows containing non-finite values (stat_bin).


```
## Warning: Removed 1305 rows containing non-finite values (stat_bin).
```

```
## Warning: Removed 1305 rows containing non-finite values (stat_bin).
```

```
## Warning: Removed 1305 rows containing non-finite values (stat_bin).
```



2.3 Función de distribución

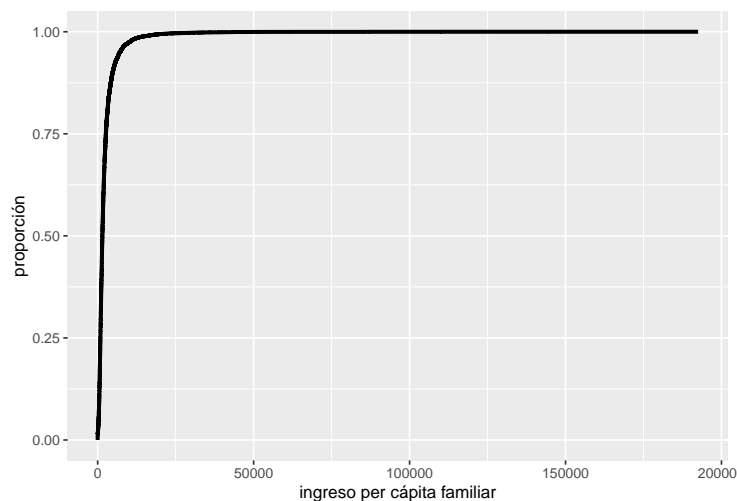
###- (pág. 78)

En este apartado se muestra cómo pueden graficarse las funciones de distribución presentadas en la sección 2.3.2 del cuerpo principal del capítulo. El primer paso para construir una función de distribución es ordenar (de menor a mayor) las observaciones de la encuestas, según la variable de ingreso elegida, *ipcf* en nuestro caso. Esto es lo que hacemos en la primera línea. En la línea siguiente se crea la variable “*shrp*” para almacenar la proporción relativa acumulada de la variable pondera, es decir el porcentaje acumulado de la población que representa cada observación más todas las anteriores. “*shrp*” tendrá entonces un valor de 100 para la última observación. Para obtener la suma acumulada expandida empleamos el comando `cumsum()` y lo dividimos por el total de población expandido que se obtiene con `sum()`. La función de distribución presenta las variables *shrp* e *ipcf* en los ejes vertical y horizontal, respectivamente, como se indica en `aes()`, y se grafica simplemente con una línea, por lo que usamos el comando `geom_line()`, donde solo especificamos su tamaño o grosor.

```
#ordenar según ipcf
mex06 <- mex06 %>% arrange(ipcf)

#población acumulada ordenamiento ipcf
mex06 <- mex06 %>% mutate(shrpop=cumsum(pondera)/sum(pondera))

## Figure 2.9 - función de distribución acumulada
ggplot(mex06, aes(x=ipcf, y=shrpop))+
  geom_line(size=1.2) +
  labs(y="proporción", x="ingreso per cápita familiar")
```



Se deja como ejercicio para el lector elaborar las otras funciones de distribución presentadas en la sección 2.3.2. Por su parte, la curva de Pen (ver figuras 2.12 y 2.13) se construye igual que la función de distribución pero se grafica invirtiendo los ejes.

2.4 Pareto

###- (pág. 78)

En esta sección se muestra cómo replicar la figura 2.14 del texto, que muestra los diagramas de Pareto para las regiones Noroeste y Sur de México. El procedimiento es muy similar al caso de la función de distribución ya que en esencia representan lo mismo, de forma distinta. El gráfico de Pareto muestra para cada valor del ingreso x el porcentaje de la población que recibe ingresos superiores a ese valor x , en una escala doble logarítmica. El cambio de escala genera una suerte de zoom óptico sobre los estratos de mayores ingresos, permitiendo un examen más detallado de esa parte de la distribución.

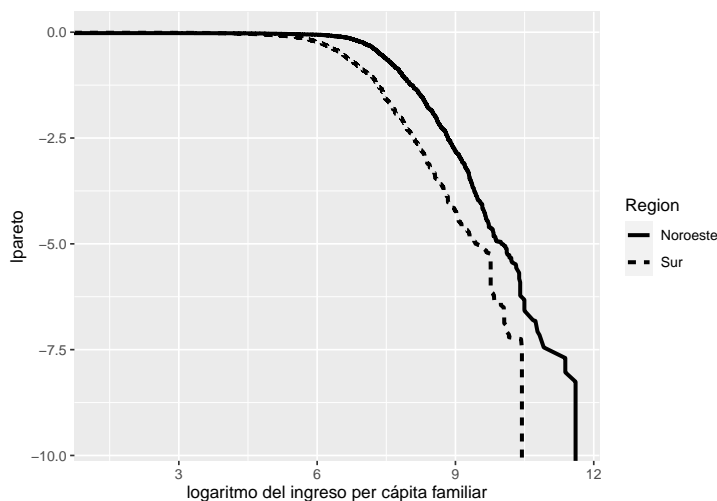
Para graficarlo seguimos los pasos anteriores, pero ahora ordenamos a la población por ingreso dentro de cada región y ya no considerando el total país. Para ello agrupamos las observaciones por región con `group_by()` previo a calcular el share acumulado, de forma tal que este cálculo se haga solo entre individuos de una misma región. En la siguiente línea se genera la variable “*lpareto*” a partir de la variable “*shrpap*”, siguiendo la explicación de la sección 2.3.4 del texto. Finalmente se grafica filtrando la base para las regiones de interés e indicando con la opción *linetype* que queremos un tipo de línea diferente para distinguir por región. Dado que la variable región es numérico empleamos `factor()` para tratarla como una categórica, que identifica a las regiones. Por último, cómo empleamos la opción *linetype* con `scale_linetype()` customizamos las leyendas. Si por ejemplo hubiéramos querido diferenciar las regiones por color, hubiéramos utilizado *color* dentro de `aes()` y luego `scale_color()` para personalizarlo.

```
#población acumulada por región
mex06 <- mex06 %>% group_by(region) %>% mutate(shrpap=cumsum(pondera)/sum(pondera))

mex06 <- mex06 %>% mutate(lpareto=log(1-shrpap))

## Figure 2.14 - Diagrama de Pareto

ggplot(mex06 %>% filter(region==1 | region==6),
       aes(x=log(ipcf), y=lpareto, weight=pondera, linetype=factor(region))) +
  geom_line(size=1.2) +
  scale_linetype(name="Region", labels=c("Noroeste", "Sur")) +
  labs(x="logaritmo del ingreso per cápita familiar")
```

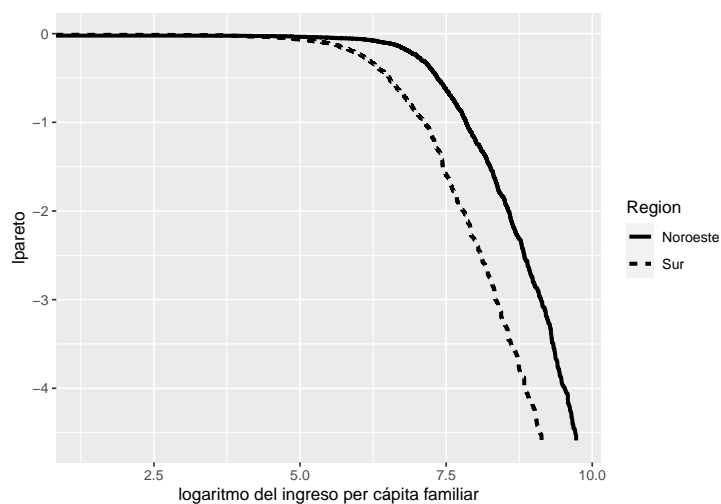


Las líneas siguientes repiten el ejercicio pero dejando de lado al 1% más rico de

las poblaciones regionales

```
cutoff=0.99

ggplot(mex06 %>% filter((region==1 | region==6) & (shrpap<=cutoff)),
  aes(x=log(ipcf), y=lpareto, weight=pondera, linetype=factor(region))) +
  geom_line(size=1.2) +
  scale_linetype(name="Region", labels=c("Noroeste", "Sur")) +
  labs(x="logaritmo del ingreso per cápita familiar")
```



2.5 Box-plot

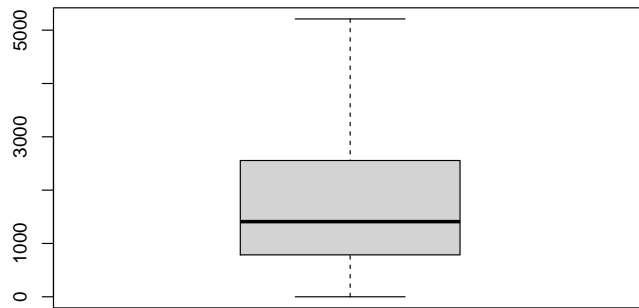
###- (pág. 79)

Aquí se muestra cómo elaborar diagramas de caja o box-plot como los presentados en la sección 2.3.5 del texto. En este caso hacemos una excepción y presentamos una alternativa más directa a `ggplot` ya que el tratamiento de los outliers exige ajuste extras en los ejes. Con el comando `boxplot` base de R podemos obtener un boxplot de forma rápida y sencilla excluyendo valores extremos que distorsionan las escalas.

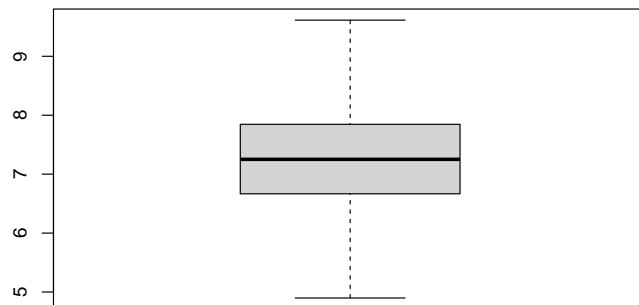
#Una opción mas directa fuera de ggplot

Figure 2.15 - Box Plot

```
boxplot(mex06$ipcf, outline = FALSE)
```



```
## Figure 2.16 - Box Plot en log  
boxplot(log(mex06$ipcf), outline = FALSE)
```



2.6 Curva de Lorenz

###- (pág. 80)

En este apartado se muestra cómo pueden construirse las curvas de Lorenz introducidas en la sección 2.3.6 del capítulo. El primer paso consiste en ordenar a los individuos de menor a mayor según su ingreso, en nuestro caso contenido

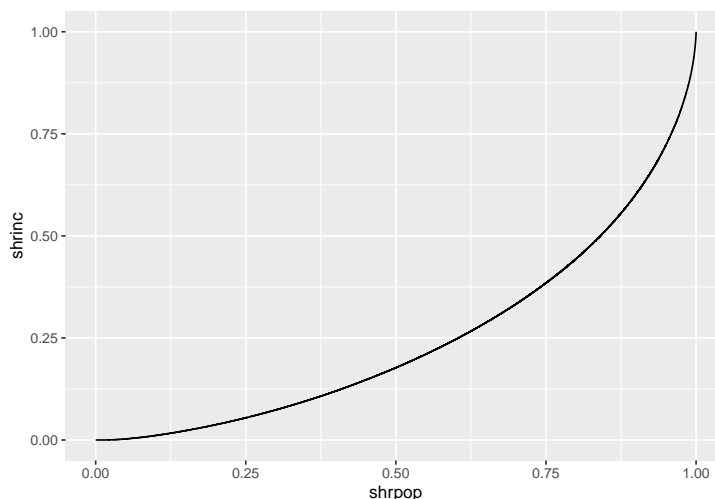
en la variable `ipcf`. Para ellos aplicamos la función `arrange()`. Notar que previo a esto se aplica la función `ungroup()`, lo cual es siempre recomendable cuando en la misma base previamente se había realizado una agrupación, para de alguna forma resetear el agrupamiento. Las líneas siguientes generan la variable “*shrpap*” de la misma forma en la que fue generada más arriba, para construir el eje vertical de la función de distribución: contiene la proporción de la población que se acumula hasta cada observación de la encuesta. La diferencia es que ahora también generaremos una variable “*shrinc*” que contiene la proporción del ingreso que se acumula hasta cada observación de la encuesta y se estima expandiendo el ingreso por el factor de expansión, es decir multiplicando ambas columnas como puede verse en el código.

La curva de Lorenz nos muestra qué porcentaje de la población acumula un dado porcentaje del ingreso total. Para visualizarlo indicamos entonces los ejes respectivos para “*shrpap*” y “*shrinc*” y luego aplicamos un gráfico de línea como `geom_line()`.

```
#ordenar según ipcf
mex06 <- mex06 %>% ungroup() %>% arrange(ipcf)

#población e ingreso acumulado
mex06 <- mex06 %>% mutate(shrpap=cumsum(pondera)/sum(pondera),
                          shrinc=cumsum(ipcf*pondera)/sum(ipcf*pondera))

## Figura 2.18 - curva de Lorenz
ggplot(mex06, aes(x=shrpap, y=shrinc)) +
  geom_line()
```

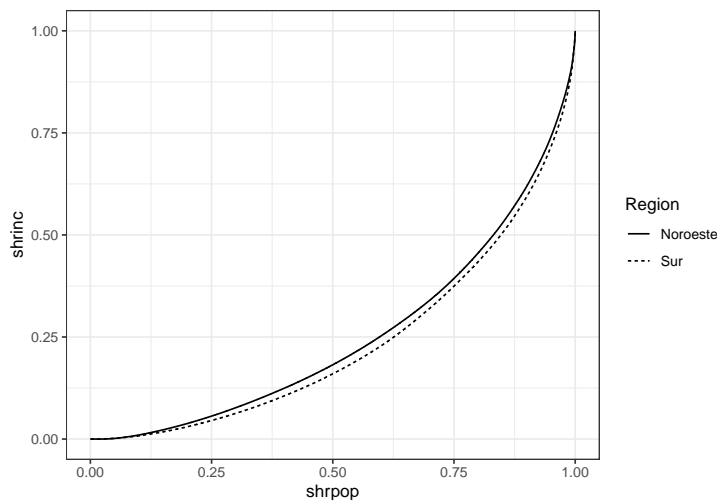


En las líneas siguientes se comparan las curvas de Lorenz para dos regiones de México. El código sigue los mismos pasos que antes pero ahora los cálculos de

población e ingreso acumulado se realizan por región, agregando un `group_by()`. En el gráfico filtramos las regiones en cuestión y las diferenciamos con distintos tipos de líneas con la opción `linetype`, al igual que antes. Con `theme_bw()` modificamos un poco el aspecto visual del gráfico.

```
#ordenar según ipcf + región y acumular por región
mex06 <- mex06 %>% arrange(region, ipcf) %>%
  group_by(region) %>%
  mutate(shrpop=cumsum(pondera)/sum(pondera),
         shrinc=cumsum(ipcf*pondera)/sum(ipcf*pondera))

## Figura 2.19 - curva de lorenz por regiones
ggplot(mex06 %>% filter(region==1 | region==6),
       aes(x=shrpap, y=shrinc, linetype=factor(region))) +
  geom_line() +
  scale_linetype(name="Region", labels=c("Noroeste", "Sur")) +
  theme_bw()
```



2.7 Curva Generalizada de Lorenz

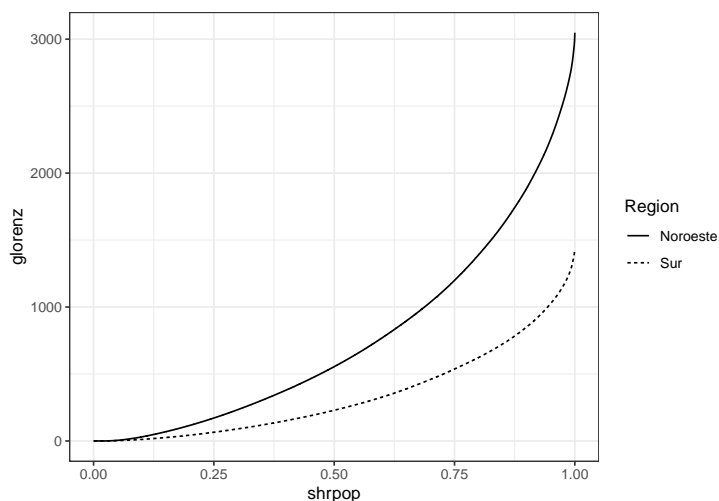
###- (pág. 81)

La curva generalizada de Lorenz se construye a partir de la curva de Lorenz pero multiplicando su eje vertical por el ingreso promedio (ver sección 2.3.6) en el cuerpo del capítulo. Las líneas no se modifican respecto de las utilizadas para estimar la curva de Lorenz, salvo en que ahora la variable *shrpap* la generamos ahora la base de la variable *pondera*. Para entender el álgebra detrás de esta

forma de calcular la curva generalizada de Lorenz, el lector puede remitirse a los apéndices del libro (pagina 81-82) donde queda claramente explicitado.

```
#ordenar según ipcf + región y acumular por región
mex06 <- mex06 %>% arrange(region, ipcf) %>%
  group_by(region) %>%
  mutate(shrpop=cumsum(pondera)/sum(pondera),
         glorenz=cumsum(ipcf*pondera)/sum(pondera))

## Figura 2.20 - curva de Lorenz generalizada por regiones
ggplot(mex06 %>% filter(region==1 | region==6),
       aes(x=shrpop, y=glorenz, linetype=factor(region))) +
  geom_line() +
  scale_linetype(name="Region", labels=c("Noroeste", "Sur")) +
  theme_bw()
```



2.8 Curva de Incidencia

###- (pág. 83)

En este apartado se muestra cómo pueden estimarse las curvas de incidencia del crecimiento que aparecen en la figura 2.21 del texto. A modo de ejemplo, se computa la curva de incidencia del crecimiento para Argentina, 1992-2006, utilizando percentiles del ingreso per cápita familiar. Para ello el lector deberá descargar dichas bases en su directorio de trabajo.

La idea básica será generar un bucle que en cada repetición cargue una base, ordene por ipcf, genere el porcentaje de población acumulado, en base a este

identifique a qué percentil de ingreso corresponde y finalmente agrupando por percentiles estime la media del ingreso en cada uno. Dado que tenemos dos bases distintas, el bucle iterará dos veces.

Para implementarlo se propone utilizar una lista (objeto tipo *list*), en la que la EPH de cada año va a ser un elemento distinto. Así el primer elemento de esa lista corresponderá a toda la base de 1992 y el segundo, a toda la base del 2006. Para señalar el primer y segundo elemento nos valemos de un contador *i* que en la primer iteración del bucle valdrá uno, y al que iremos sumando de uno en uno, por lo que en la segunda tomará valor dos. Por eso la expresión `bases_mod[[i]]` será equivalente a decir “elemento 1 de la lista *bases_mod*” en la primera iteración y “elemento 2 de la lista *bases_mod*” en la segunda iteración. Con esto logramos que en vez de ordenar, generar y cambiar variables sobre un dataframe estático, lo hagamos sucesivamente sobre dataframes distintos contenidos en un lista “madre”, que aquí llamamos *bases_mod*.

Nuestro bucle irá iterando sobre *j*, que al igual que *i* cambiará de valor en cada vuelta. En la primera será 92 y en la segunda 06 y nos servirá para indicar con qué base trabajar. Así `read.dta()` tiene como argumento la función `paste()` que une el prefijo “*arg*” al año de la base correspondiente (*j*), generando, junto con el prefijo “*dta*”, un string igual al nombre de la base completa. De esta forma logramos cambiar el nombre en cada iteración para cargar bases de años diferentes.

Dado que estamos trabajando con años distintos debemos prever que necesitaremos ajustar los valores de ingresos por inflación. Por eso en caso de que la base sea la del 92, se multiplica el *ipcf* por el factor de ajuste (1.7865). Esto se debe hacer solo en el caso de la base de 1992 y la manera de instrumentarlo es diciéndole a R que ese cambio debe realizarse solo cuando *j* sea igual a 92 `-if (j=="92")-`.

Una vez cargada la base, logrado el ajuste por inflación, aplicamos los mismos pasos que antes para ordenar por *ipcf* a la población y calcular el porcentaje acumulado. Con ese porcentaje acumulado generaremos los percentiles anidando un bucle a nuestro bucle inicial, que por cada base hará 100 repeticiones. En cada una de ellas asignará los percentiles en base a *shrpap* ayudándose de la variable *z* multiplicada por 0.01. Esto porque si queremos generar percentiles ($n=100$), necesitamos 100 cuantiles por lo que cada cuantil se asigna de a intervalos de población acumulada iguales a 0.01 (1/100). Así por ejemplo, cuando *z* sea igual a 20, $(z-1)*0.01$ será igual a 0.19 y $z*0.01$ a 0.20, por lo que caerán en este cuantil veinte, todos aquellos individuos que ordenados por ingreso, estén entre el 19 y el 20 por ciento de población acumulada. Dado que empleamos el comando `ifelse()`, para las observaciones restantes que caen fuera de este intervalo, se mantendrá el valor que ya traía la variable “*percentil*”. Es por esto que antes del bucle, la generamos como vacía, para luego ir rellenándola sucesivamente. *z* continuará incrementándose de a uno, hasta llegar a 100, donde asigne el último percentil a aquellas personas que quedaron ultimas en nuestra ordenación, ya que poseían los ingresos mas alto. Al considerar a este 1% de

la población faltante, se habrá acumulado al 100%. Finalmente el último paso, consiste en reducir todo el dataframe a 100 observaciones correspondientes a los 100 cuantiles creados y calcular la media del ingreso de las personas que pertenezcan a cada percentil.

```
bases_mod <- list() #lista vacia en la que se guardaran los cambios

i <- 1 #contador para iterar sobre los elementos de la lista

#j tomará valor 92 y 06
for (j in c("92","06")){

  #carga base
  bases_mod[[i]] <- read.dta(data_dir %+% paste("Arg/bases/arg",j,"_cedlas.dta", sep = " "),
                             filter(cohh==1 & !is.na(ipcf))

  #ajuste por inflación para el año 92 solamente
  if (j=="92"){

    bases_mod[[i]] <- bases_mod[[i]] %>% mutate(ipcf=ipcf*2.099) #ver comentarios ;

  }

  #ordenar según ipcf
  bases_mod[[i]] <- bases_mod[[i]] %>% arrange(ipcf)

  #computar porcentaje de población
  bases_mod[[i]] <- bases_mod[[i]] %>% mutate(shrpop=cumsum(pondera)/sum(pondera))

  #identificar percentil de ipcf
  bases_mod[[i]] <- bases_mod[[i]] %>% mutate(percentile=0) #esto es equivalente a ge

  for (z in 1:100){

    bases_mod[[i]] <- bases_mod[[i]] %>% mutate(percentile=ifelse((shrpop>(z-1)*0.01,
                                                                    percentile,0))

  }

  bases_mod[[i]] <- bases_mod[[i]] %>% group_by(percentile) %>% summarise(ipcf=mean(ipcf))

  i=i+1
}
```

Cuando el bucle finaliza, podemos “recuperar” los elementos almacenados en nuestra lista *bases_mod*. El primero corresponde a los ingresos promedios de

cada percentil construidos con la base de 1992 y el segundo en base a los datos de 2006.

```
#recupero las bases de cada año como respectivos elementos de la lista
arg92 <- bases_mod[[1]]
arg06 <- bases_mod[[2]]
```

```
head(arg92)
```

Para finalizar la estimación, debemos calcular el cambio porcentual en el ingreso promedio de cada percentil entre estos años. Para ello necesitaremos juntar ambos dataframes en un solo, a partir de la variable percentil, que se encuentra en ambos. Con el comando `merge()` unificamos las bases y le asignamos un sufijo que nos permita identificar cada columna. Hecho esto ya podemos calcular el cambio porcentual en el ipcf promedio y graficarlo para cada percentil.

```
#junto bases, asigno nombres y calculo el cambio porcentual
change <- merge(arg92, arg06, by="percentile", suffix=c("_92", "_06")) %>%
  mutate(change=( (ipcf_06/ipcf_92) -1)*100)
```

```
## Figura 2.21 - Curva de incidencia ipcf
```

```
ggplot(change %>% filter(percentile>2), #elimino los primeros dos porque son outliers debido a
  aes(x=percentile, y=change))+
  geom_line(size=1.2, color="darkcyan") +
  geom_hline(yintercept = 0, linetype="dashed") +
  theme_bw()
```

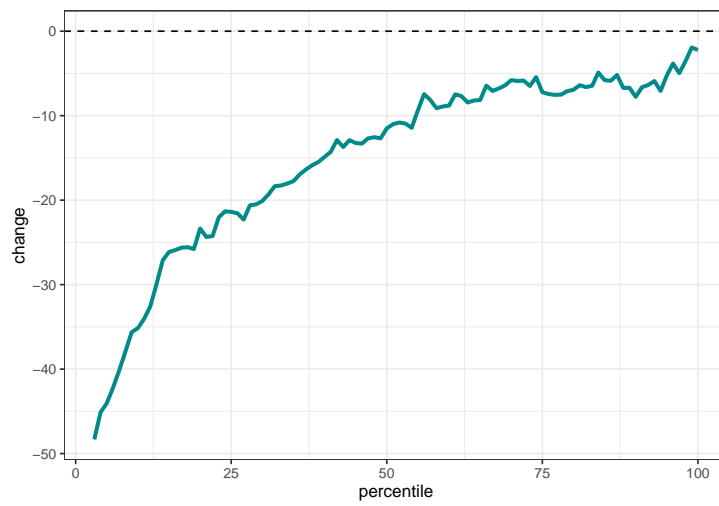


Figure 2.1: Esta figura presenta resultados en magnitudes diferentes a los reportados originalmente en el capítulo del libro.

Chapter 3

Capítulo 3

Apéndices del Libro Pobreza y Desigualdad en R

Escrito por: Cristian Bonavida

Last Update: 02/7/2021

Códigos escritos en base a los apéndices del libro “Pobreza y Desigualdad en América Latina” de Gasparini, Cicowiez y Sosa Escudero. El objeto de este material es reproducir la rutina de códigos para STATA presentada en el libro al lenguaje R. Este material es solo de caracter complementario a las explicaciones y detalles conceptuales que se presentan en el libro de texto y los apéndices

Set Inicial

Cargo las librerías, limpio environment, defino el path y atajo para función paste

```
library(dplyr)
library(tidyverse) # Data wrangling
library(tidygraph)
library(readxl)
library(ggplot2)
library(foreign)
library(TAM)
```

```
rm(list=ls())      #empiezo limpiando todo

"%+%" <- function(x,y) paste(x,y,sep = "")      # defino un shorcut para concat de te
data_dir <- "C:/Users/HP/Desktop/CEDLAS - UNLP/Apendices en R/Material libro/encuestas
```

3.1 Cociente de cuantiles

###- (pág. 151-152)

El siguiente bloque de código puede utilizarse para computar el cociente de quintiles extremos presentado en el cuadro 3.2 del texto del libro. Para iniciar importo la base desde el formato STATA utilizando de la librería `foreign` el comando `read.dta`. y hago una limpieza de la base. Luego ordeno las observaciones según la variable “`ipcf`”, filtro ingresos nulos y por último computo el porcentaje acumulado por población.

```
#carga base
mex06 <- read.dta(data_dir %+% "Mex/2006/bases/mex06_cedlas.dta")

#elimino observaciones incoherentes o con ingreso missing
mex06 <- mex06 %>% filter(cohh==1, !is.na(ipcf))

#ordenar por ipcf, filtrar ingresos nulos y computar % de población
df <- mex06 %>% arrange(ipcf) %>% filter(ipcf>0) %>%
  mutate(shrpop=cumsum(pondera)/sum(pondera))
```

A partir de allí genero la variable quintil en el dataframe, que vale 1 para el 20% más pobre de la población, 2 para el 20% siguiente, y así sucesivamente. El comando `ifelse` permite evaluar el porcentaje acumulado de población en cada observación y le otorga el valor del quintil correspondiente en caso afirmativo, o deja la variable inalterada en caso negativo.

```
df$quintil = 0
df$quintil = ifelse(df$shrpob<=0.2, 1, df$quintil)
df$quintil = ifelse(df$shrpob> 0.2 & df$shrpob<= 0.4, 2, df$quintil)
df$quintil = ifelse(df$shrpob> 0.4 & df$shrpob<= 0.6, 3, df$quintil)
df$quintil = ifelse(df$shrpob> 0.6 & df$shrpob<= 0.8, 4, df$quintil)
df$quintil = ifelse(df$shrpob> 0.8 & df$shrpob<= 1, 5, df$quintil)
```

Para terminar computamos el ingreso promedio ponderado para las observaciones del quintil 1 y 5. Aquí la expresión `[df$quintil==1]` actúa como filtro.

De esta manera el comando le pide a R considerar los valores de la variable `ipcf` y pondera que corresponden a observaciones que cumplen la condición, en este caso de pertenecer al quintil 1. Por último calculamos el ratio de estos dos valores e imprimimos el resultado. Para mayor claridad siempre se recomienda redondear los valores a imprimir, utilizando el comando `round`.

```
ipcf_q1 = weighted.mean(df$ipcf[df$quintil==1], df$pondera[df$quintil==1], na.rm = TRUE)
ipcf_q5 = weighted.mean(df$ipcf[df$quintil==5], df$pondera[df$quintil==5], na.rm = TRUE)
ratq51 = ipcf_q5/ipcf_q1

paste("ratq51 = ", ratq51)
```

```
## [1] "ratq51 = 13.4341399453849"
```

```
paste("ratq51 = ", round(ratq51, digits = 2))
```

```
## [1] "ratq51 = 13.43"
```

3.2 Replicar programa ratq51

###- (pág. 153-154)

En las siguientes líneas de código proponemos escribir una función que permite computar el cociente de quintiles extremos muy fácilmente; simplemente se lo invoca indicando la base de datos, la variable de la que se quieren computar quintiles, y opcionalmente la variable de ponderación y la condición `if`. Una función en R es similar a lo que en STATA llamaríamos un programa, es decir un conjunto de instrucciones que se guardan en la memoria y que pueden ser fácilmente replicables en distintas bases de datos y para distintas variables.

La función puede resultar un poco extensa por lo que aquí enumeraremos las líneas para mayor claridad. Para iniciar la función, en la primera línea debemos especificar sus argumentos, que en este caso son: un dataframe (objeto “`df`”), una variable (objeto “`varname`”) y opcionalmente una variable de ponderación y una condición (objetos “`var_pondera`” y “`condicion`”). Para llamar a esta función entonces siempre tendremos que, como mínimo, especificar los dos objetos iniciales.

La siguiente línea se asegura que la base de datos “`df`” que hemos especificado siempre sea almacenado como un objeto del tipo *data frame* en un objeto llamado “`aux`” que creamos en el entorno de la función. La noción de entorno es importante para comprender el resto del código. A los fines de nuestro código

diremos que un entorno es un espacio determinando en el que se almacenan distintas clases de objetos que R acepta y, en el cual estos tienen un significado, un valor o simplemente existen como objetos. En R pueden “coexistir” entornos distintos, a los cuales se le asigna un nivel o prioridad¹. Volviendo a nuestro código, cuando a R le especificamos, por ejemplo, el objeto “varname”, R buscará si este está definido en el *entorno o environment* general del programa como prioridad, y no en el entorno o espacio de la función, que más acotado que el anterior. Por lo tanto aquí necesitamos evitar este comportamiento y lograr que R no busque evaluar al objeto por fuera del ámbito de la función. El comando **substitute** le indica a R que considere al objeto con su nombre literal. Así cuando “varname” sea definida, por ejemplo, como la variable “ipcf”, R no buscará a qué es equivalente el objeto “ipcf” en su memoria sino que tomará el nombre dado. Luego la función **eval** le indica a R que ese nombre debe evaluarlo, es decir buscarlo o identificarlo, dentro del dataframe “aux” que hemos definido en nuestra función. De esta forma nos aseguramos de que R considera correctamente las variables dentro de la base de datos que estamos especificando, y por lo tanto que nuestros argumentos sean válidos dentro del ámbito de la función. Es decir que no entren en conflicto con otros objetos similares o definidos previamente en el environment mas amplio del programa.

Utilizando estas dos funciones, en la línea 9 almacenamos todos los valores de ipcf a partir de los cuales luego ordenamos la base en la siguiente línea. Como la condición es un argumento opcional debemos evaluar si esta fue indicada, en términos de R si no es un objeto nulo (línea 15). Si no lo es, luego en la línea 17 verificamos que esté correctamente especificada, escribiendo una sentencia que será VERDADERA cuando al evaluar de forma silenciosa la condición, de como resultado un error (“try-error”). En ese caso entonces se imprime un mensaje que ayuda a identificar donde se encuentra el problema. En caso de que el resultado no sea del tipo error, filtramos el objeto “aux” de acuerdo a la especificación dada. Para ellos nos valemos de la función **parse** que convierte la condición que era del tipo “string” al tipo “expression”, lo que la vuelve compatible de ser evaluada como sentencia, nuevamente con la función **eval** dentro de **filter**.

De forma similar procedemos con el ponderador en la línea 28. Si no fue especificado (objeto nulo), entonces vale 1 para todas las observaciones, a partir de crear una variable que llamamos “w” que pasará a ser nuestra variable ponderador, pero que será inocua a los fines del cálculo. En el objeto “var_pondera_store” almacenamos esta secuencia de valores 1. En cambio, si el ponderador fue especificado entonces ese objeto almacena los valores correspondientes a la vari-

¹Para los usuarios que no esten tan familiarizados con las nociones de *Environment* y *Non-Standard Evaluation* se recomienda revisar las siguientes referencias para un tratamiento mas detallado.

*<http://adv-r.had.co.nz/Environments.html#environments>

*<http://adv-r.had.co.nz/Computing-on-the-language.html>

*<https://advanced-r-solutions-ed1.netlify.app/non-standard-evaluation.html#non-standard-evaluation-in-subse>

able señalada, sin entrar en conflicto con un posible objeto del mismo nombre definido fuera del espacio de la función.

En la línea 43 actualizamos el objeto que almacena los valores de ipcf (en caso de que se hayan filtrado valores de acuerdo a la condición), y en la siguiente ordenamos el dataframe “aux”, para calcular como antes el porcentaje acumulado de población. Por último definimos inicialmente a la variable “quintil” como valores 0. En la línea 50 replicamos la asignación del quintil que se realizó antes con el comando `ifelse` pero ahora de forma iterativa, cambiando en cada iteración los valores del quintil y de población acumulada.

Por último al igual que antes calculamos el ratio entre quintiles extremos filtrando las observaciones del primer y quinto quintil. La única diferencia es que aquí nuevamente debemos asegurarnos de que los valores promedios y los ponderadores se evalúen para estas observaciones y sólo en el contexto o espacio de la función. Para finalizar imprimimos los resultados.

```

1 ratq51 <- function(df, varname, var_pondera=NULL, condicion=NULL) {
2
3
4   aux <- as.data.frame(df)
5
6
7   #substitute() toma literal el nombre sin asignarle ningun valor posible y
8   #eval() evalua ese nombre en el contexto del data frame que es el 1er argumento definido
9   varname_store <- eval(substitute(varname), aux)
10  aux <- aux %>% arrange(varname_store)
11
12
13
14  #la condición es un argumento opcional
15  if(!is.null(substitute(condicion))) {
16
17    if(is(try(eval(substitute(condicion))), silent = TRUE ), "try-error"))
18
19    stop("ERROR: la condicion debe especificarse como character (entre comillas)")
20
21    aux <- aux %>% filter(eval(parse(text=condicion))) #convierte la condicion de tipo string a
22                                                    #expression y luego la evalua dentro del contexto de la funcion
23  }
24
25
26
27  #set pondera= 1 si no está especificado
28  if(is.null(substitute(var_pondera))) {
29
30    aux <- aux %>% mutate(w=1)

```

```

31     var_pondera=substitute(w)
32     var_pondera_store <- eval(substitute(var_pondera), aux)
33
34   } else {
35
36     var_pondera_store <- eval(substitute(var_pondera), aux)
37
38   }
39
40
41
42   #ordeno y genero quintiles
43   varname_store <- eval(substitute(varname), aux)
44   aux <- aux %>% arrange(varname_store) %>%
45     mutate(shrpop=cumsum(var_pondera_store)/sum(var_pondera_store)) %>%
46     mutate(quintil=0)
47
48   for (i in 1:5) {
49
50     aux <- aux %>% mutate(quintil=ifelse((shrpop>(i-1)*0.2 & shrpop<=i*0.2), i, quintil))
51
52   }
53
54
55
56   #ingreso promedio quintil 1
57   aux_1 <- aux %>% filter(quintil==1)
58   media_q1 = weighted.mean(eval(substitute(varname),aux_1), eval(substitute(var_pondera),aux_1))
59   media_q1 = round(media_q1, digits = 2)
60   #ingreso promedio quintil 5
61   aux_5 <- aux %>% filter(quintil==5)
62   media_q5 = weighted.mean(eval(substitute(varname),aux_5), eval(substitute(var_pondera),aux_5))
63   media_q5 = round(media_q5, digits = 2)
64   #ratio
65   ratq_51 = media_q5/media_q1
66
67
68   print(paste("media quintil 5", substitute(varname), "=", media_q5))
69   print(paste("media quintil 1", substitute(varname), "=", media_q1))
70   print(paste("ratio = ", round(ratq_51, digits = 2)))
71
72
73 }

```

Al correr el código, R almacenará en la memoria esta función, a la cual po-

dremos llamar simplemente indicando los argumentos relevantes. En las líneas siguientes se detallan algunos ejemplos de cómo replicar el programa bajo distintas especificaciones.

```
#No especifica condicion
```

```
ratq51(mex06, ipcf, pondera)
```

```
## [1] "media quintil 5 ipcf = 6655.49"
```

```
## [1] "media quintil 1 ipcf = 448.38"
```

```
## [1] "ratio = 14.84"
```

```
#Especifica condicion
```

```
ratq51(mex06, ipcf, pondera, "ipcf>0")
```

```
## [1] "media quintil 5 ipcf = 6707.64"
```

```
## [1] "media quintil 1 ipcf = 499.3"
```

```
## [1] "ratio = 13.43"
```

```
#Especifica condición doble
```

```
ratq51(mex06, ipcf, pondera, "region==4 | region==2")
```

```
## [1] "media quintil 5 ipcf = 6420.24"
```

```
## [1] "media quintil 1 ipcf = 499.01"
```

```
## [1] "ratio = 12.87"
```

```
#Especifica incorrectamente la condición
```

```
ratq51(df=mex06, ipcf, pondera, region==4 & urbano==1)
```

```
## Error in ratq51(df = mex06, ipcf, pondera, region == 4 & urbano == 1): ERROR: la condicion de
```

```
#No especifica pondera
```

```
ratq51(df=mex06, varname=ipcf, condicion="region==4 & urbano==1")
```

```
## [1] "media quintil 5 ipcf = 6384.76"
```

```
## [1] "media quintil 1 ipcf = 594.14"
```

```
## [1] "ratio = 10.75"
```

```
#Identifica los argumentos por explicitamente por nombre
```

```
ratq51(df=mex06, varname=ipcf, var_pondera=pondera, condicion="region==4 & urbano==1")
```

```
## [1] "media quintil 5 ipcf = 6640.93"
```

```
## [1] "media quintil 1 ipcf = 586.57"
```

```
## [1] "ratio = 11.32"
```

3.3 Replicar programa gcuan

###- (pág. 154-155)

El bloque de código a continuación permite identificar cuantiles de cualquier variable. En términos del programa “ratq51”, nos permite generar variables similares a quintil pero que pueden identificar quintiles, deciles, percentiles, etc.

Por esta razón, esta función tendrá más argumentos, aquí además de los anteriores debemos detallar la cantidad de cuantiles a generar (objeto “num”) y la variable que los almacena (“newvar”). Luego el código y la secuencia son idénticos a la de la función anterior, salvo porque aquí definimos de una manera alternativa y más directa a los valores del poderador, pero en esencia replica lo visto antes hasta la línea 43. A partir de aquí el objeto “num” indica cuantos cuantiles deben generarse, hace iterar al bucle “num” cantidad de veces, y define los intervalos de población acumulada de forma equivalente. Por ejemplo, si queremos generar deciles (num=10), necesitamos 10 cuantiles y cada cuantil se asigna de a intervalos de población acumulada iguales a 0.10 (1/10).

Posteriormente en la línea 54 generamos el reporte que imprimiremos como resultado. Este calcula la media ponderada, el desvío standard ponderado y la cantidad de observaciones en base la variable especificada en “varname” y “var_pondera”. Los objetos “my_var” y “my_var2” son variables auxiliares que generamos en el data frame solo con el objeto de facilitarnos el cálculo directo. Para terminar cambiamos el nombre de la variable que hasta ahora generamos como “quintil” por el nombre indicado en el argumento. El comando `names(aux)` trae todos los nombres de este dataframe y con la expresión `[names(aux) == "quintil"]` elegimos de esos nombres, solo el que coincide con la palabra “quintil”. A esta columna específica le asignamos el nombre dado como argumento, tomando la expresión literal con `substitute` pasada al formato string mediante la función `paste`. En resumen la línea 63 sería el equivalente de escribir en la consola `aux$quintil <- "nombre_asignado"`, con la dificultad de que debemos hacerlo dentro del espacio de la función, respetando los argumentos dados. La línea 64 por su parte elimina del output final la variable `shrp`.

Para finalizar especificamos que esta función no solo debe imprimirnos los resultados sino que además debe devolver una base de datos nueva. El comando `return(aux)` le dice a R que el resultado será el objeto “aux” en sí mismo, es decir obtendremos la base original con una nueva columna llamada “newvar” que es la que esta función genera.

```

1 gcuan <- function(df, varname, var_pondera=NULL, condicion=NULL, num, newvar) {
2
3
4   aux <- as.data.frame(df)
5

```

```

6
7 varname_store <- eval(substitute(varname), aux)
8 aux <- aux %>% arrange(varname_store)
9
10
11 #la condición es un argumento opcional
12 if(!is.null(substitute(condicion))) {
13
14     if(is(try(eval(substitute(condicion)), silent = TRUE ), "try-error"))
15
16         stop("ERROR: la condicion debe especificarse como character (entre comillas)")
17
18         aux <- aux %>% filter(eval(parse(text=condicion)))
19
20 }
21
22
23
24 #set pondera= 1 si no está especificado
25 if(is.null(substitute(var_pondera))) {
26
27     var_pondera_store <- c(rep(1, nrow(aux)))
28
29 } else {
30
31     var_pondera_store <- eval(substitute(var_pondera), aux)
32
33 }
34
35
36
37 #ordeno y genero shrpap
38 varname_store <- eval(substitute(varname), aux)
39 aux <- aux %>% arrange(varname_store) %>%
40     mutate(shrpap=cumsum(var_pondera_store)/sum(var_pondera_store)) %>%
41     mutate(quintil=0)
42
43 #genero cuantiles en base a lo indicado
44 shareq = 1/num
45
46 for (i in 1:num) {
47
48     aux <- aux %>% mutate(quintil=ifelse((shrpap>(i-1)*shareq & shrpap<=i*shareq), i, quintil))
49
50 }

```

```

51
52
53   #armo la información que se imprime en la consola
54   show <- aux %>% mutate(my_var=varname_store,
55                          my_var2=var_pondera_store) %>% group_by(quintil) %>%
56                          summarise(mean = weighted_mean(my_var, my_var2),
57                                    std = weighted_sd(my_var, my_var2),
58                                    obs = sum(my_var2))
59
60
61   #renombro variable al nombre indicado como argumento y descarto shrpap
62   names(show)[names(show) == "quintil"] <- paste(substitute(newvar))
63   names(aux)[names(aux) == "quintil"] <- paste(substitute(newvar))
64   aux$shrpap <- NULL
65
66
67   #outputs
68   print.data.frame(show)
69   return(aux)
70
71
72 }
```

A diferencia entonces de la anterior, esta función devuelve no sólo un resultado impreso sino un objeto, por lo que ahora debemos especificar cómo lo nombramos. Si lo llamamos de la misma manera que el argumento que le pasamos a la función simplemente estamos pisando este objeto y añadiéndole una nueva variable. Otra alternativa es crear un nuevo dataframe con otro nombre.

```

#Especifica todo correctamente y crea nuevo dataframe
mex06_bis <- gcuan(df=mex06, varname=ipcf, var_pondera=pondera, condicion="ipcf>0", num
```

```

##   quintil      mean      std      obs
## 1      1  499.2978  197.3755 20372798
## 2      2 1011.4695  136.6778 20383133
## 3      3 1539.5536  172.7480 20385227
## 4      4 2382.6530  345.4059 20380795
## 5      5 6707.6362 6532.2051 20381791
```

```

#No especifica correctamente la condición
mex06_bis <- gcuan(df=mex06, varname=ipcf, var_pondera=pondera, condicion=region==4, num
```

```

## Error in gcuan(df = mex06, varname = ipcf, var_pondera = pondera, condicion = region
```

```
#No especifica condición
```

```
mex06_bis <- gcuan(df=mex06, varname=ipcf, var_pondera=pondera, num=5, newvar=quintil)
```

```
##   quintil      mean      std      obs
## 1        1  448.3757  222.9768 20677642
## 2        2  986.9415  137.2255 20677386
## 3        3 1518.1792  172.0445 20676646
## 4        4 2357.4796  343.4104 20678111
## 5        5 6655.4937 6499.4356 20679061
```

```
#No especifica pondera
```

```
mex06_bis <- gcuan(df=mex06, varname=ipcf, num=5, newvar=quintil)
```

```
##   quintil      mean      std      obs
## 1        1  363.0791  202.9644 16400
## 2        2  900.1146  135.1091 16401
## 3        3 1424.6686  177.6839 16401
## 4        4 2273.6985  346.6132 16401
## 5        5 6401.7186 6319.9421 16401
```

```
#Especifica todo correctamente y reemplaza el dataframe dado
```

```
mex06 <- gcuan(df=mex06, varname=ipcf, var_pondera=pondera, condicion="ipcf>0", num=5, newvar=quintil)
```

```
##   quintil      mean      std      obs
## 1        1  499.2978  197.3755 20372798
## 2        2 1011.4695  136.6778 20383133
## 3        3 1539.5536  172.7480 20385227
## 4        4 2382.6530  345.4059 20380795
## 5        5 6707.6362 6532.2051 20381791
```

3.4 Tamaño de los hogares

###- (pág. 156)

El código siguiente puede utilizarse para computar las estadísticas sobre proporción de hogares unipersonales y multipersonales presentadas en el cuadro 3.4 del texto. Como ya es habitual, en la primera línea indicamos con qué base iremos a trabajar. Luego ordenamos los hogares de forma creciente en base a su identificador “id” y en forma decreciente respecto la variable “jefe”. Dado que esta vale 1 solo para el jefe de hogar y 0 para todo el resto, la primera

observación por “id” corresponderá siempre a la cabeza del hogar. La función `deduplicated` genera una variable que será *FALSE* para la primera observación dentro de cada id, indicando que esta no está duplicada por ser la primera, pero será *TRUE* para todas las que siguen ya que se identificó antes una observación con el mismo id. De esta manera hacemos un “tag” del jefe de hogar y mantenemos todos los demás miembros. En el caso de que nos interesa mantener en nuestra base el resto de las observaciones, estos pasos se simplifican a un simple filtrado por jefe de hogar.

Luego generamos la variable tamaño solo para el jefe de hogar en base la cantidad de miembros que viven con él, utilizando la función `case_when`. Así, por ejemplo, cuando la observación corresponda al jefe y la cantidad de miembros sea igual a 3, en ese caso tamaño valdrá 3. Por último armamos nuestra tabla resultado filtrando a los jefes de hogar, agrupando por tamaño y sumando la cantidad de hogares (porque tenemos una sola observación por hogar) y la frecuencia relativa en cada caso.

```
#indico con qué base de hogares voy a trabajar
df <- mex06

#ordeno por id y jefe (decreciente) e identifico al jefe de hogar
df <- df %>% arrange(id, desc(jefe)) %>%
  mutate(hh=duplicated(id)) #alterantiva para egen hh=tag() en stata

#genero tamaño sólo para los jefes de hogar
df <- df %>% mutate(tamaño=case_when(miembros==1 & hh==FALSE ~ 1,
                                     miembros==2 & hh==FALSE ~ 2,
                                     miembros==3 & hh==FALSE ~ 3,
                                     miembros==4 & hh==FALSE ~ 4,
                                     miembros==5 & hh==FALSE ~ 5,
                                     miembros>=6 & hh==FALSE ~ 6))

#tabla con resultados
table <- df %>% filter(hh==FALSE) %>% group_by(tamaño) %>% summarise(N = sum(pondera))
mutate(freq = 100*N/sum(N))

table

## # A tibble: 6 x 3
##   tamaño      N freq
##   <dbl>   <int> <dbl>
## 1     1 2070859  8.15
## 2     2 3717891 14.6
## 3     3 4815417 18.9
## 4     4 5875572 23.1
## 5     5 4422391 17.4
## 6     6 4516967 17.8
```


El objeto `table` almacena los resultados que se visualizan en la consola. A veces al imprimir los resultados estos no se visualizan con un formato muy amigable. A modo de extensión presentamos dos alternativas para refinar este aspecto. El paquete `formattable` nos permite, entre otras cosas, especificar el tipo de datos y el formato para distintas variables al generar el dataframe. Por su parte `print.data.frame` nos da una visualización más limpia. Incorporando este aspecto el código para nuestra tabla sería:

```
#install.packages("formattable")
library(formattable)

table <- df %>% filter(hh==FALSE) %>% group_by(tamaño) %>% summarise(N = accounting(sum(pondera),
mutate(freq = percent(N/sum(N)))
print.data.frame(table)
```

```
## tamaño      N    freq
## 1      1 2,070,859 8.15%
## 2      2 3,717,891 14.63%
## 3      3 4,815,417 18.94%
## 4      4 5,875,572 23.11%
## 5      5 4,422,391 17.40%
## 6      6 4,516,967 17.77%
```

Además de agregarle esta customización, para los usuarios que utilizan “Rmark-down” habitualmente, es posible mostrar los resultados con una visualización amigable de forma sencilla.

```
rmarkdown::paged_table(table)
```

3.5 La distribución intrahogar

###- (pág. 157)

El fragmento de código siguiente puede utilizarse para generar resultados similares a los presentados en el cuadro 3.7 del texto, que muestra cómo se modifica la desigualdad calculada a través del cociente de deciles extremos cuando cambia la distribución del ingreso hacia el interior del hogar. Cabe recordar que la distribución del ingreso intrahogar se modifica mediante un impuesto proporcional al ingreso per cápita familiar combinado con un subsidio que solo recibe el jefe de hogar. En la implementación, utilizamos quintiles en lugar de deciles ingreso.

En la primer línea cargamos la base, en este caso trabajaremos con la base de Venezuela del año 2006, limpiamos y ordenamos los hogares. Mas abajo creamos el objeto “ty” que toma el valor de la tasa del impuesto aplicada sobre el ipcf, el cual calculamos en línea siguiente generando en el dataframe una nueva columna “impuesto”. Posteriormente calculamos el valor del subsidio a otorgar que surge de sumar el valor de los impuestos para los integrantes de un mismo hogar. En la línea siguiente modificamos la variable subsidio, haciendo que valga 0 para todo miembro distinto al jefe, de esta forma redistribuimos ingreso al interior del hogar: a todos se le hemos quitado una porción “ty” que ahora la recibe solamente el jefe de hogar.

En la anteúltima línea creamos una variable para el nuevo valor de ingreso per cápita familiar, restando el impuesto y sumando el subsidio. Por último hacemos uso de nuestra función del ratio de quintiles para computar el cociente del ingreso promedio de los quintiles 5 y 1 como indicador de desigualdad, a partir del ingreso modificado.

```
ven06 <- read.dta(data_dir %>% "Ven/2006/bases/ven06_cedlas.dta")

#elimino observaciones incoherentes o con ingreso missing y ordeno
df <- ven06 %>% filter(cohh==1, !is.na(ipcf)) %>% arrange(id)

#tasa del impuesto
ty=0.1

#impuesto al ipcf
df$impuesto = df$ipcf*ty      #alternativa a escribir: df <- df %>% mutate(impuesto = ip

#recaudación impuesto total por hogar
df <- df %>% group_by(id) %>% mutate(subsidio = sum(impuesto))

#subsidio solo lo recibe el jefe de hogar
df$subsidio <- ifelse(df$jefe!=1, 0, df$subsidio)

#nuevo ipcf
df$ipcf_star = df$ipcf - df$impuesto + df$subsidio

ratq51(df, ipcf_star, pondera, "ipcf>0")

## [1] "media quintil 5 ipcf_star = 842520.38"
## [1] "media quintil 1 ipcf_star = 79335.26"
## [1] "ratio = 10.62"
```

A modo de comentario, un usuario mas familiarizado con el lenguaje podrá notar que las líneas que computan y asignan el subsidio pueden resumirse en una sola, escribiendo.

```
df <- df %>% group_by(id) %>% mutate(subsidio=ifelse(jefe!=1, 0, sum(impuesto)))
```

3.6 Empleo de ponderadores

###- (pág. 157)

El bloque de código que sigue puede utilizarse para construir un cuadro como el 3.9 del texto, que muestra la relación entre el ingreso per cápita familiar y el valor de la variable de ponderación.

La línea 7 genera la variable shrobs, que contiene el porcentaje acumulado de observaciones, a diferencia de shropop que contiene el acumulado de población usando el ponderador. Aquí la expresión `1:n()` enumera secuencialmente las observaciones en orden de aparición mientras que `n()` computa el total.

```
#indico con que base de hogares voy a trabajar
df <- mex06

#ordenar por ipcf y computo porcentaje acumulado de observaciones
df <- df %>% arrange(ipcf) %>% mutate(shrobs = 1:n()/n())

#identificar quintiles de ipcf
df$quintil = 0
df$quintil = ifelse(df$shrobs<= 0.2, 1, df$quintil)
df$quintil = ifelse(df$shrobs>0.2 & df$shrobs<= 0.4, 2, df$quintil)
df$quintil = ifelse(df$shrobs>0.4 & df$shrobs<= 0.6, 3, df$quintil)
df$quintil = ifelse(df$shrobs>0.6 & df$shrobs<= 0.8, 4, df$quintil)
df$quintil = ifelse(df$shrobs>0.8 & df$shrobs<= 1, 5, df$quintil)

show <- df %>% group_by(quintil) %>% summarise(mean = accounting(mean(ipcf), digits = 0),
                                                means_w = accounting(mean(pondera), digits = 0))
print.data.frame(show)
```

```
##   quintil   mean means_w
## 1      1    413    1,038
## 2      2    926    1,269
## 3      3  1,449    1,345
## 4      4  2,302    1,344
## 5      5  6,457    1,319
```

El bloque de código siguiente calcula las tasas de pobreza con y sin ponderadores para el total del país y para cada una de las regiones de México en 2006

correspondientes al cuadro 3.10 del texto. En la primer línea el objeto “lp” almacena el valor de la línea de pobreza, en base a la cual se genera la variable binaria “pobre”, vale 1 para los individuos pobres (es decir, $ipcf < lp$) y 0 para el resto. La suma de esta variable arroja el número total de personas pobres en la muestra, al multiplicarla por el ponderador obtenemos el número de personas pobres en la población y dividiendo por la población total del país nos devuelve la tasa de incidencia de la pobreza. En este caso el valor es de 13.57%

```
###línea de pobreza us$2.5 Mexico 2006
lp= 633.90918

#identificar individuos pobres
df$pobre = ifelse(df$ipcf<lp, 1, 0)

#total país
sum(df$pobre*df$pondera)*100/sum(df$pondera)

## [1] 13.57454
```

Las siguientes líneas calculan el mismo valor de una manera alternativa, almacenándolo en un dataframe, a partir de contabilizar la frecuencia absoluta y relativa de personas pobres y no pobres.

```
pobreza <- df %>% group_by(pobre) %>% summarise(n = accounting(sum(pondera), digits = 0),
                                                  mutate(tasa_pobreza = percent(n/sum(n))))
print.data.frame(pobreza)

##   pobre      n tasa_pobreza
## 1     0 88,070,783      86.43%
## 2     1 13,832,961      13.57%
```

Para replicar el mismo cálculo pero para las regiones del país, generamos un bucle que itera 8 veces, una vez por cada región. En cada vuelta del bucle se estima, de la misma forma que arriba, el valor de pobreza tanto ponderado como no ponderado. Para este segundo caso el único paso adicional consiste en fijar la variable “pondera” como igual a 1. Como vimos en la salida anterior de la tabla, en la 2da fila de la 3er columna se encuentra la tasa de incidencia de pobreza, con la expresión `pobreza_pondera[2,3]` llamamos a este valor para guardarlo en el objeto “share_p”, luego de redondearlo a dos dígitos con el comando `round`. Lo mismo hacemos para las estimaciones sin ponderar. Por último, imprimimos los resultados acompañados por una leyenda indicativa.

```

#por region
for (i in 1:8){

  pobreza_pondera <- df %>% filter(region==i) %>% group_by(pobre) %>% summarise(n=sum(pondera)) %>%
    mutate(tasa_pobreza=n/sum(n))

  pobreza_s_pondera <- df %>% filter(region==i) %>% mutate(pondera=1) %>% group_by(pobre) %>% sum
    mutate(tasa_pobreza=n/sum(n))

  #recupero el valor llamando a la fila y columna y redondeo
  share_p = round(pobreza_pondera[2,3]*100, digits = 2)
  share_sp = round(pobreza_s_pondera[2,3]*100, digits = 2)

  print(paste("H_ponderado =", share_p, "/ H_sin_ponderar =", share_sp))

}

## [1] "H_ponderado = 7.73 / H_sin_ponderar = 8.21"
## [1] "H_ponderado = 12.23 / H_sin_ponderar = 17.43"
## [1] "H_ponderado = 4.73 / H_sin_ponderar = 10.03"
## [1] "H_ponderado = 11.94 / H_sin_ponderar = 11.18"
## [1] "H_ponderado = 9.55 / H_sin_ponderar = 13.4"
## [1] "H_ponderado = 35.13 / H_sin_ponderar = 44.68"
## [1] "H_ponderado = 20.64 / H_sin_ponderar = 21.14"
## [1] "H_ponderado = 14.51 / H_sin_ponderar = 18.91"

```

Este bloque de códigos imprime los resultados en la consola. Si quisiéramos generar un cuadro u objeto que los almacene (dataframe) para luego por ejemplo exportarlo en Excel u otro formato, solo debemos modificar ligeramente el código de arriba. Antes de correr el bucle generamos un dataframe que contiene todos valores 0, pero con 3 columnas (region, pondera, sin_pondera). Al final del bucle guardamos los datos de la región, de pobreza con ponderador y sin ponderador en la columna 1, 2 y 3 respectivamente, pero en cada iteración lo hacemos en una fila distinta, que corresponde a cada región. Por ejemplo en la tercera iteración el código filtra la región a la cual le corresponde el valor 3, y almacena las estimaciones de línea de pobreza en la tercera fila de la tabla “results”.

```

#creo data frame solo con ceros, de 8 filas y 3 columnas
results <- data.frame(region=c(rep(0,8)),
                      pondera=c(rep(0,8)),
                      sin_pondera=c(rep(0,8)))

#por region
for (i in 1:8){

```

```

pobreza_pondera <- df %>% filter(region==i) %>% group_by(pobre) %>% summarise(n=sum(
mutate(tasa_pobreza=

pobreza_s_pondera <- df %>% filter(region==i) %>% mutate(pondera=1) %>% group_by(pob

#recupero el valor llamando a la fila y columna y redondeo
results[i,1] = paste("region", i)
results[i,2] = round(pobreza_pondera[2,3]*100, digits = 2)
results[i,3] = round(pobreza_s_pondera[2,3]*100, digits = 2)
}

```

De esta manera al finalizar el bucle queda armado un objeto de 3 columnas con 8 filas, una para cada región

3.7 Diseño muestral

###- (pág. 160)

En este material no se cubre con ejemplos este apartado pero para los usuarios interesados en adentrarse en el manejo de encuestas de hogares contemplando el diseño muestral, se recomienda explorar el uso del paquete “SURVEY”, similar al paquete “svy” en STATA.

```

#install.packages("survey")
library(survey)
?survey

```

3.8 Fuentes de ingreso

###- (pág. 161)

El bloque de código a continuación muestra cómo computar la importancia que tiene cada fuente de ingresos identificada en las encuestas de hogares (cuadro 3.13). Dentro de las fuentes de ingreso consideramos: laboral (variable *ila*), jubilaciones (*ijubi*), capital (*icap*), transferencias (*itrans*) y otros (*ionl*).

Para comenzar declaramos la base y la limpiamos. Luego sumamos generamos la variable ingreso total (“*itot*”), como la suma de las columnas para cada ingreso.

Hasta ahora todas las veces que aplicamos el comando `sum` lo hicimos sumando una misma columna, por ejemplo para obtener el total de población, sumamos la variable “pondera”. En este caso queremos sumar por fila distintas columnas, para eso debemos anteponer el comando `rowwise`.

Como resultado final necesitamos mostrar la participación en el ingreso total de cada fuente de ingreso. Para ello debemos considerar cada fuente por separado y dividirla por el ingreso total. La forma de hacerlo nuevamente es utilizando un bucle. Entonces el código a escribir debe lograr en cada iteración tomar todos los valores de cada una de las variables que corresponde a la fuente de ingresos, en otras palabras debemos iterar entre columnas distintas. La forma que proponemos aquí consiste en valernos de las listas. Las listas son un objeto que en cada elemento puede almacenar otro objeto de cualquier tipo. Una lista de n elementos puede contener n dataframes distintos en cada uno de ellos, o n columnas de un dataframe, o n vector. Aquí entonces nos valdremos de esta flexibilidad para almacenar en un mismo objeto (lista) múltiples objetos distintos (columnas)

La cuarta línea de códigos declara a “ingresos” como una lista que como elementos contiene a las variables de ingreso de nuestra base de datos. En la línea siguiente creamos un simple vector con los nombres de cada una de estas. Con esto ya podemos generar nuestro bucle, haciéndolo iterar desde 1 hasta n , siendo n la cantidad de fuentes de ingresos distintas que consideramos. Para evitarnos contar manualmente cuantas son, directamente calculamos n con el comando `length` que nos devuelve la cantidad de elementos dentro de la lista “ingresos”. De esta forma si agregamos o quitamos una variable de ingreso no debemos preocuparnos por fijar este valor cada vez.

En la primer línea del código definimos la variable “y” que, en cada iteración, será un elemento distinto del objeto ingreso, es decir una fuente de ingreso distinta, comenzando por “ila” y terminando en “ionl”. A esta variable la expandimos multiplicándola por el ponderador, de la misma manera que expandimos la variable de ingreso total, que ya calculamos mas arriba. Solo nos resta hacer el cociente entre la suma de los ingresos de todas las personas correspondiente a cada fuente, sobre el total de todos los ingresos. El objeto “value” almacena este cálculo, que luego redondeamos a dos dígitos. Nótese la importancia de la opción `na.rm=TRUE` que evita que un valor missing convierta en missing a toda la suma, es decir que indica ignorar los valores missings y preservar el cálculo sobre el resto de las observaciones no missings. Para terminar, “concepto” almacena el nombre de la fuente correspondiente a cada vuelta del bucle, para imprimirlo en la línea final junto con el valor calculado.

```
arg06 <- read.dta(data_dir %+% "Arg/2006/s2/bases/arg06_cedlas.dta")

df <- arg06 %>% filter(cohh==1, !is.na(ipcf))

#sumo los ingresos por fila (rowwise)
```

```

df <- df %>% rowwise %>% mutate(itot = sum(ila, ijubi, icap, itran, ionl, na.rm = TRUE)

#creo una lista en la que cada elemento es un vector distinto de ingreso
ingresos <- list(df$ila, df$ijubi, df$icap, df$itran, df$ionl)
names <- c("laboral", "jubilación", "capital", "transferencias", "otros")

#itero sobre cada uno de esos vectores
for (i in 1:length(ingresos)){

  y = ingresos[[i]]
  y_expand = y * df$pondera
  itot_expand = df$itot * df$pondera

  value <- sum(y_expand, na.rm=TRUE) / sum(itot_expand, na.rm=TRUE) * 100
  value <- round(value, digits = 2)

  concepto = names[i]
  print(paste("shr %", concepto, "= ", value))

}

## [1] "shr % laboral = 80.4"
## [1] "shr % jubilación = 12.24"
## [1] "shr % capital = 1.5"
## [1] "shr % transferencias = 4.77"
## [1] "shr % otros = 1.09"

```


Chapter 4

Capítulo 4

Apéndices del Libro Pobreza y Desigualdad en R

Escrito por: Cristian Bonavida

Last Update: 28/8/2021

Códigos escritos en base a los apéndices del libro “Pobreza y Desigualdad en América Latina” de Gasparini, Cicowiez y Sosa Escudero. El objeto de este material es reproducir la rutina de códigos para STATA presentada en el libro al lenguaje R. Este material es solo de caracter complementario a las explicaciones y detalles conceptuales que se presentan en el libro de texto y los apéndices

Set Inicial

Cargo las librerías, limpio environment, defino el path y atajo para función paste

```
library(dplyr)
library(tidyverse) # Data wrangling
library(tidygraph)
library(readxl)
library(ggplot2)
library(foreign)
library(TAM)
```

```
rm(list=ls())      #empiezo limpiando todo

"%+%" <- function(x,y) paste(x,y,sep = "")      # defino un shorcut parar concat de te
data_dir <- "C:/Users/HP/Desktop/CEDLAS - UNLP/Apendices en R/Material libro/encuestas
```

4.1 Indicador FGT

###- (PÁG. 249-250)

En este apartado se presenta cómo calcular la familia de indicadores FGT. En primer lugar, se muestra cómo puede computarse el indicador FGT de manera relativamente sencilla. Luego, al igual que con el ratio de quintiles y el cálculo de quintiles, se introduce una función para el indicador FGT. Como ejemplo, computamos la pobreza de 2.5 dólares para Ecuador en 2006, utilizando microdatos que provienen de la Encuesta de Condiciones de Vida (ENCOVI).

Luego de cargar y limpiar las bases fijamos los valores de pobreza así como el parámetro “*alfa*” del indicador fgt de aversión a la desigualdad entre los pobres. Luego se computa, para cada individuo pobre, su brecha de pobreza elevada al valor de alfa, que se asigna a una nueva variable que llamamos “each”. La distinción entre individuo pobre y no pobre se operativiza en el comando `ifelse`: en caso de cumplirse la condición de que el ingreso esté por debajo de la línea de pobreza, se computa la brecha, en caso de que sea falso se otorga un valor de cero. Por último se obtiene el valor de FGT como el promedio ponderado de la brecha en toda la muestra.

```
#carga base y elimino observaciones incoherentes o con ingreso missing
ecu06 <- read.dta(data_dir %+% "Ecu/2006/bases/ecu06ecv_cedlas.dta") %>%
  filter(cohh==1, !is.na(ipcf))

df <- ecu06

#linea de pobreza
lp=39.74
#parametro alfa indicador fgt
alfa=0

#computar fgt
df <- df %>% mutate(each = ifelse(ipcf<lp, (1 - ipcf/lp)^alfa, 0 ))
fgt = weighted.mean(df$each, df$pondera, na.rm = TRUE)*100

print("fgt = " %+% round(fgt, d=2))
```

```
## [1] "fgt = 19.02"
```

Otra opción posible es computar los valores directamente como un vector, sin alojarlo como una nueva columna del data frame. Esta opción, desde el punto de vista de la escritura de los códigos, es más eficiente ya que se evita agregar una nueva columna que solo se emplea para el cálculo del indicador. El detalle a considerar es que para poder calcularlo como medida ponderada la cantidad de elementos del vector “each” debe ser exactamente igual la cantidad de observaciones de la columna pondera del dataframe.

```
#alternativa
each <- ifelse(df$ipcf<lp, (1 - df$ipcf/lp)^alfa, 0 )
fgt = weighted.mean(each, df$pondera, na.rm = TRUE)*100

print("fgt = " %+% round(fgt, d=2))
```

```
## [1] "fgt = 19.02"
```

4.2 Computar FGT

###- (pág 250-251)

Como ya vimos en el capítulo 3, las funciones nos permiten replicar un indicador o un cálculo en cualquier base, sobre cualquier variable, imponiendo una condición específica o ponderando por algún factor de expansión. En este caso además agregamos como argumentos los parámetros asociados al FGT, “*alfa*”, para la aversión a la desigualdad, y “*zeta*” para la línea de pobreza. De esta forma podremos replicar y comparar rápidamente el índice para distintos valores que podemos pasarle a estos argumentos. La estructura de la función replica lo visto en las funciones del capítulo 3 hasta la línea 39 donde se computa el FGT con las mismas líneas que empleamos arriba. Adicionalmente se añade la opción de correr la función sin de forma silenciosa, sin imprimir el resultado.

```
1 FGT <- function(df, varname, var_pondera=NULL, condicion=NULL, alfa, zeta, quiet=FALSE) {
2
3   aux <- as.data.frame(df)
4
5   varname_store <- eval(substitute(varname), aux)
6   aux <- aux %>% arrange(varname_store)
7
8 }
```

```

9      #la condición es un argumento opcional
10     if(!is.null(substitute(condicion))) {
11
12         if(is(try(eval(substitute(condicion)), silent = TRUE ), "try-error"))
13
14             stop("ERROR: la condicion debe especificarse como character (entre comillas)")
15
16             aux <- aux %>% filter(eval(parse(text=condicion)))
17
18     }
19
20
21     #set pondera igual a 1 si no está especificado
22     if(is.null(substitute(var_pondera))) {
23
24         var_pondera_store <- c(rep(1, nrow(aux)))
25
26     } else {
27
28         var_pondera_store <- eval(substitute(var_pondera), aux)
29
30     }
31
32
33     #Cómputo de brecha y valor del indicador
34     varname_store <- eval(substitute(varname), aux)
35     aux <- aux %>% mutate(each = ifelse( varname_store < zeta,
36                                         ( 1 - varname_store/zeta)^alfa,
37                                         0 ))
38
39     fgt = weighted.mean(aux$each, var_pondera_store, na.rm = TRUE)*100
40     fgt = round(fgt, digits = 2)
41
42     #output
43     if(substitute(quiet)==TRUE){
44         a=fgt
45
46     } else {
47
48         print(paste("FGT(alfa=", alfa, ",Z=", zeta, ") = ", fgt, sep=""))
49         a=fgt
50
51     }
52
53 }

```

```
#No especifica condicion
```

```
FGT(df=df, varname=ipcf, var_pondera=pondera, alfa=0, zeta=39.740)
```

```
## [1] "FGT(alfa=0,Z=39.74) = 19.02"
```

```
#Especifica condicion
```

```
FGT(df=df, varname=ipcf, var_pondera=pondera, condicion="urbano==1", alfa=0, zeta=39.740)
```

```
## [1] "FGT(alfa=0,Z=39.74) = 11.82"
```

```
#Especifica incorrectamente la condición
```

```
FGT(df=df, varname=ipcf, var_pondera=pondera, condicion=urbano==1, alfa=0, zeta=39.740)
```

```
## Error in FGT(df = df, varname = ipcf, var_pondera = pondera, condicion = urbano == : ERROR: la
```

```
#Especifica opcion "quiet"
```

```
FGT(df=df, varname=ipcf, var_pondera=pondera, alfa=0, zeta=39.740, quiet=TRUE)
```

4.3 Pobreza relativa

```
####- (pág. 251)
```

La estimación de la pobreza relativa implica, como primer paso, el cálculo de una línea de pobreza relativa. A modo de ejemplo, se computa una línea de pobreza igual al 50% del ingreso mediano de Ecuador. El comando `weightedMedian` de la librería `matrixStats` arroja el valor de la mediana. Luego, el cálculo de la pobreza se realiza empleando la función `FGT`.

```
#línea de pobreza del 50% de la mediana del ingreso
```

```
lp = matrixStats::weightedMedian(df$ipcf, df$pondera) * 0.50
```

```
FGT(df=df, varname=ipcf, var_pondera=pondera, condicion="urbano==1", alfa=0, zeta=lp)
```

```
## [1] "FGT(alfa=0,Z=47.484118603645) = 16.13"
```

4.4 Descomposición regional de la pobreza

###- (pág. 252-253)

El código a continuación realiza una descomposición por regiones de la tasa de incidencia de la pobreza (cuadro 4.7 del libro de texto). Para este caso utilizamos la Encuesta Nacional de Ingresos y Gastos de los Hogares de México para el año 2006 con la línea de pobreza de 2.5 dólares, equivalentes a 608.245 pesos mensuales. Luego de cargar, limpiar e indicar la base de datos con la que se trabajará, se fija la línea de pobreza y se genera, a partir de la función de FGT, la tasa de incidencia para el total del país que se almacena en el objeto “p0”. En la línea siguiente, se toma una única observación de cada una de las 8 regiones diferentes de México y se las ordena en orden creciente para almacenarlas como un vector en el objeto “list_rgn”. Sobre este objeto haremos iterar un bucle, para que en cada vuelta los cálculos indicados se hagan para cada una de estas regiones.

La primera sentencia dentro de ese bucle, calcula la participación de la región sobre la población y en la siguiente línea se estima su indicador FGT para un valor de $\alpha=0$ y con la línea de pobreza de 2.5 dólares (notar como se instrumenta la condición en la función). Por último la contribución se calcula como el producto entre la participación de cada región en la población total y el cociente entre la tasa de pobreza regional y la tasa de pobreza nacional, se redondea y se indica que se imprima dicho valor.

```
#carga y limpio base
mex06 <- read.dta(data_dir %+"Mex/2006/bases/mex06_cedlas.dta") %>%
  filter(cohh==1, !is.na(ipcf))

df <- mex06

lp=608.245
p0 = FGT(df=df, varname=ipcf, var_pondera=pondera, alfa=0, zeta=lp, quiet=TRUE)

list_rgn = sort(unique(df$region))

for (i in list_rgn){

  #participación region
  shr_rgn = sum(df$pondera[df$region==i]) / sum(df$pondera)
  #fgt region
  p_r <- FGT(df=df, varname=ipcf, var_pondera=pondera, condicion=paste("region==", i, s

  #contribución
  contribut = round( shr_rgn*(p_r/p0)*100 , digits = 2 )
```

```
print(paste("contribución (%) region", i, "=", contribut))

}
```

```
## [1] "contribución (%) region 1 = 6.09"
## [1] "contribución (%) region 2 = 9.52"
## [1] "contribución (%) region 3 = 3.07"
## [1] "contribución (%) region 4 = 13.99"
## [1] "contribución (%) region 5 = 24.22"
## [1] "contribución (%) region 6 = 26.12"
## [1] "contribución (%) region 7 = 13.07"
## [1] "contribución (%) region 8 = 3.95"
```

4.5 Pobreza según consumo e ingreso

###- (pág. 253-254)

El código que sigue puede utilizarse para replicar los resultados sobre pobreza por consumo e ingreso presentados en el cuadro 4.9 del texto para el caso de Nicaragua en 2005. Luego de fijar la tasa de pobreza relevante, el bucle hace que el objeto *i* tome valores desde 0.5 hasta 1.5 a intervalos de 0.1. Estos valores expanden la línea de pobreza que se utiliza para calcular el FGT sobre consumo (variable “cpcf”) e ingreso (variable “ipcf”)

```
nic05 <- read.dta(data_dir %+% "Nic/2005/bases/nic05_cedlas.dta") %>%
  filter(coh_oficial==1)

df <- nic05

#línea de pobreza oficial
lp0=576.5028

for (i in seq(0.5,1.5, by=0.1)) {

  #línea de pobreza
  lp=lp0*i; lp=round(lp, d=2)

  #consumo
  print(paste(i, "*lp = ", lp, "- Consumo", sep=""))
```

```

FGT(df=df, varname=cpcf, var_pondera=pondera, alfa=0, zeta=lp)

#ingreso
print(paste(i, "*lp = ", lp, "- Ingreso", sep=""))
FGT(df=df, varname=ipcf, var_pondera=pondera, alfa=0, zeta=lp)

}

```

```

## [1] "0.5*lp = 288.25- Consumo"
## [1] "FGT(alfa=0,Z=288.25) = 12.21"
## [1] "0.5*lp = 288.25- Ingreso"
## [1] "FGT(alfa=0,Z=288.25) = 16.59"
## [1] "0.6*lp = 345.9- Consumo"
## [1] "FGT(alfa=0,Z=345.9) = 19.03"
## [1] "0.6*lp = 345.9- Ingreso"
## [1] "FGT(alfa=0,Z=345.9) = 22.91"
## [1] "0.7*lp = 403.55- Consumo"
## [1] "FGT(alfa=0,Z=403.55) = 26.23"
## [1] "0.7*lp = 403.55- Ingreso"
## [1] "FGT(alfa=0,Z=403.55) = 28.51"
## [1] "0.8*lp = 461.2- Consumo"
## [1] "FGT(alfa=0,Z=461.2) = 33.63"
## [1] "0.8*lp = 461.2- Ingreso"
## [1] "FGT(alfa=0,Z=461.2) = 33.77"
## [1] "0.9*lp = 518.85- Consumo"
## [1] "FGT(alfa=0,Z=518.85) = 40.06"
## [1] "0.9*lp = 518.85- Ingreso"
## [1] "FGT(alfa=0,Z=518.85) = 38.76"
## [1] "1*lp = 576.5- Consumo"
## [1] "FGT(alfa=0,Z=576.5) = 46.02"
## [1] "1*lp = 576.5- Ingreso"
## [1] "FGT(alfa=0,Z=576.5) = 43.62"
## [1] "1.1*lp = 634.15- Consumo"
## [1] "FGT(alfa=0,Z=634.15) = 52.18"
## [1] "1.1*lp = 634.15- Ingreso"
## [1] "FGT(alfa=0,Z=634.15) = 47.75"
## [1] "1.2*lp = 691.8- Consumo"
## [1] "FGT(alfa=0,Z=691.8) = 56.67"
## [1] "1.2*lp = 691.8- Ingreso"
## [1] "FGT(alfa=0,Z=691.8) = 51.33"
## [1] "1.3*lp = 749.45- Consumo"
## [1] "FGT(alfa=0,Z=749.45) = 60.71"
## [1] "1.3*lp = 749.45- Ingreso"
## [1] "FGT(alfa=0,Z=749.45) = 54.91"
## [1] "1.4*lp = 807.1- Consumo"

```



```
## [1] "FGT(alfa=0,Z=807.1) = 64.76"
## [1] "1.4*lp = 807.1- Ingreso"
## [1] "FGT(alfa=0,Z=807.1) = 58.29"
## [1] "1.5*lp = 864.75- Consumo"
## [1] "FGT(alfa=0,Z=864.75) = 68.21"
## [1] "1.5*lp = 864.75- Ingreso"
## [1] "FGT(alfa=0,Z=864.75) = 61.55"
```

En este caso el código imprime una larga lista de resultados en la consola. Con un cambio menor en el bucle, con las mismas sentencias, podemos generar directamente la tabla a replicar y almacenarla como un dataframe que luego es exportable fácilmente a otros formatos para su presentación. Para ello generamos inicialmente 3 vectores vacíos para cada una de las columnas que tendrá la tabla. Utilizamos un contador auxiliar que irá incrementándose de a uno en cada iteración del bucle, en las cuales iremos guardando los datos que antes se imprimían, ahora como elementos de estos objetos. Así, por ejemplo en la tercer iteración el cálculo de FGT se guardará como el tercer elemento del vector consumo y el vector ingreso. Al finalizar creamos un dataframe de 3 columnas a partir de estos 3 vectores, y por último agregamos la cuarta columna calculando la diferencia.

```
consumo <- c()
ingreso <- c()
linea_pobreza <- c()

j = 1

for (i in seq(0.5,1.5, by=0.1)) {

  #línea de pobreza
  lp=lp0*i; lp=round(lp, d=2)

  linea_pobreza[j] <- paste(i, "*lp", sep="")

  #consumo
  consumo[j] <- FGT(df=df, varname=cpcf, var_pondera=pondera, alfa=0, zeta=lp, quiet = TRUE)
  #ingreso
  ingreso[j] <- FGT(df=df, varname=ipcf, var_pondera=pondera, alfa=0, zeta=lp, quiet = TRUE)

  j=j+1
}

tabla <- data.frame(linea_pobreza, consumo, ingreso) %>% mutate(diferencia = ingreso - consumo)
```

El código siguiente permite replicar la figura 4.13 del texto, que compara las

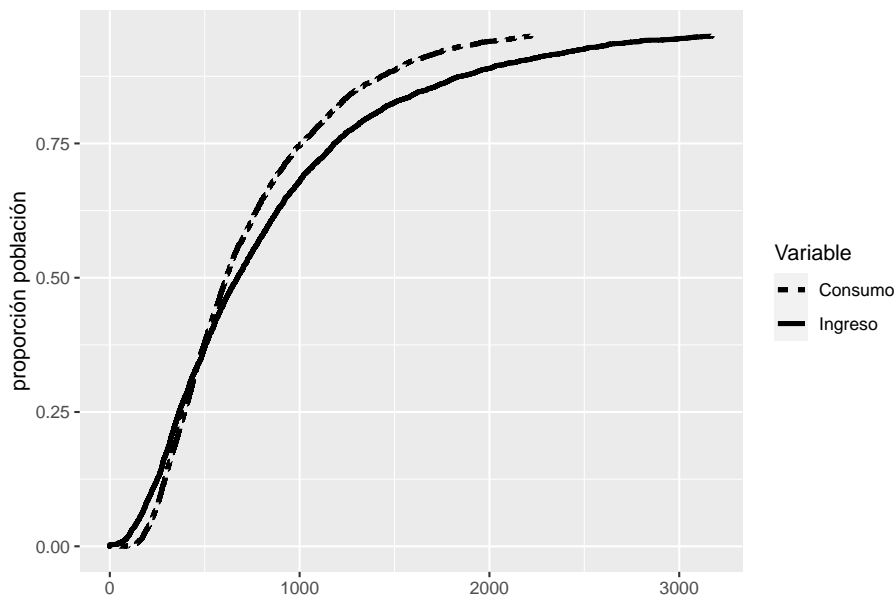
funciones de distribución del ingreso y el consumo per cápita. Para ello se ordena de forma creciente por las variables de ingreso y luego de consumo, calculando en cada caso el share de población. El “*cutoff*” se utiliza para indicar qué porcentaje de las observaciones se mostrará en el gráfico.

##FUNCIÓN DE DISTRIBUCIÓN ACUMULADA

```
#ordenar según ipcf y calcular shrpop
df <- df %>% arrange(ipcf) %>% mutate(shrpop_i = cumsum(pondera)/sum(pondera))
#ordenar según cpcf y calcular shrpop
df <- df %>% arrange(cpcf) %>% mutate(shrpop_c = cumsum(pondera)/sum(pondera))

cutoff=0.95

ggplot(df %>% filter(shrpop_i < cutoff), aes(x=ipcf, y=shrpop_i, linetype="Ingreso"))+
  geom_line(size=1.2) +
  #como la condición es sobre otra variable tengo que volver a indicar dataframe y a
  geom_line(data = df %>% filter(shrpop_c < cutoff), aes(x=cpcf, y=shrpop_c, linetype="Consumo"),
    scale_linetype_manual(name = "Variable", values=c(Ingreso="solid", Consumo="twodash")),
    labs(y="proporción población", x=""))
```



4.6 Pobreza por edad

###- (pág. 254-255)

El bloque de código siguiente muestra cómo puede graficarse la relación entre pobreza y edad (ver figura 4.15 del texto). Para ello luego de cargar, limpiar la base y definir la línea de pobreza, generamos dos objetos “x” e “y” vacíos donde se almacenará los valores del eje x y el eje y. Estos valores se generaran de forma iterativa mediante un bucle que incrementa secuencialmente la edad en 5 años, comenzando en 0 y terminando en 80 años. En cada vuelta estaremos generando el indicador FGT condicionando a las observaciones que caigan dentro de distintos intervalos de edad. Esto equivale a dividir a la población en grupos de edad y para cada uno de ellos calcular el indicador. Al igual que en el código anterior, el contador nos permite almacenar los valores de edad y FGT como elementos de los vectores creados inicialmente.

```
#indico con qué base de hogares voy a trabajar
mex06 <- read.dta(data_dir %+"Mex/2006/bases/mex06_cedlas.dta") %>%
  filter(cohh==1, !is.na(ipcf))
df <- mex06

#línea de pobreza oficial
lp0=608.245

x <- c()
y <- c()

j=1

for (i in seq(0,80,by=5)) {

  print(paste("rango = [", i, ",", i+4, "]", sep=""))
  fgt_edad = FGT(df=df, varname=ipcf, var_pondera=pondera, condicion=paste("edad>=", i, " & edad<=", i+4),
    alfa=0, zeta=lp, quiet = TRUE)

  x[j]=i
  y[j]=fgt_edad

  j=j+1

}
```

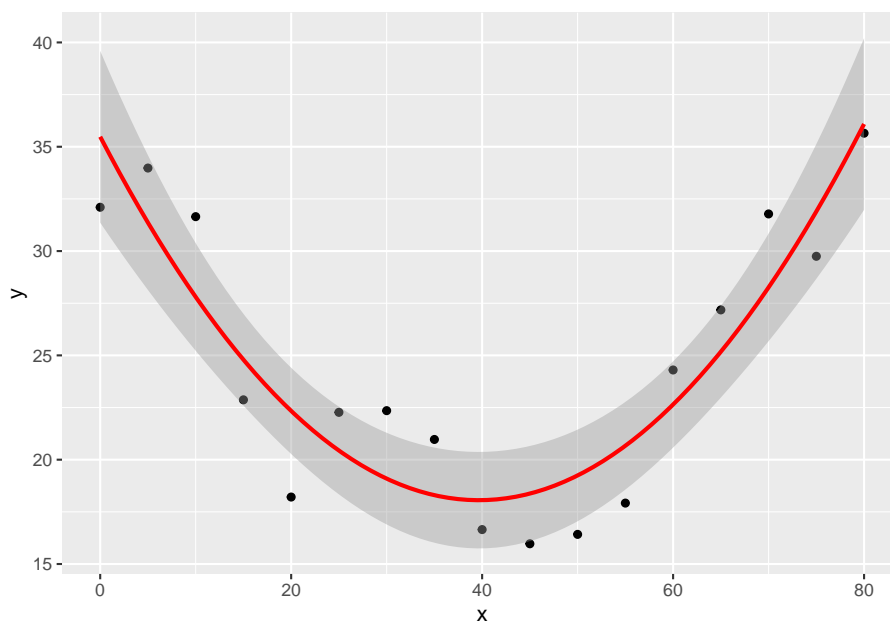
```
## [1] "rango = [0,4]"
## [1] "rango = [5,9]"
## [1] "rango = [10,14]"
## [1] "rango = [15,19]"
```

```
## [1] "rango = [20,24]"
## [1] "rango = [25,29]"
## [1] "rango = [30,34]"
## [1] "rango = [35,39]"
## [1] "rango = [40,44]"
## [1] "rango = [45,49]"
## [1] "rango = [50,54]"
## [1] "rango = [55,59]"
## [1] "rango = [60,64]"
## [1] "rango = [65,69]"
## [1] "rango = [70,74]"
## [1] "rango = [75,79]"
## [1] "rango = [80,84]"
```

Las líneas finales grafican los resultados, superponiendo a las estimaciones de pobreza una línea de regresión polinomial de orden dos. .

```
xst=x^2
aux <- data.frame(x, y, xst)

ggplot(aux, aes(x = x, y = y)) +
  geom_point() +
  geom_smooth(method=lm, formula = y ~ x + I(x^2), colour="red")
```



4.7 Significatividad estadística

###- (pág. 255-256)

En esta última sección recreamos la técnica del *bootstrap* o muestreo para obtener errores estándares e intervalos de confianza para las estimaciones del FGT. La versión más simple del bootstrap requiere (i) tomar una muestra de tamaño N (el tamaño muestral) de la muestra original con reposición, (ii) computar el índice de pobreza deseado y (iii) repetir el procedimiento B veces, con B grande. Esto es lo que haremos mediante un bucle, fijaremos una cantidad de repeticiones, en las que en cada una estaremos tomando un muestreo de la muestra original de igual tamaño. De esta forma la muestra irá cambiando en su composición aleatoriamente y por tanto permitirá generar distintos valores del FGT en cada iteración, a partir de una misma base. El comando en R `sample_n` realiza esta tarea de muestreo, seteando el dataframe, el tamaño de la nueva muestra y la opción con reposición. La expresión `nrow(df)` indica que la nueva muestra tendrá el mismo tamaño que la base original. Cada valor del indicador se almacena en el objeto “store”, sobre el que posteriormente se calcula el desvío, la media y el tamaño. Estos son inputs necesarios para la formula que estima los intervalos de confianza de nuestras estimaciones. En este caso estimamos un intervalo de confianza del 95%.

```
per06 <- read.dta(data_dir %+% "Per/2006/bases/per06_cedlas.dta") %>%
  filter(cohh==1, !is.na(ipcf))

df <- per06

#genero un resamplero del data frame en cada iteración y para ese data frame obtengo el fgt
store <- c()
rep=50

lp = 128.136

for (i in 1:rep) {

  df_sample <- sample_n(df, size=nrow(df), replace=T)

  fgt = FGT(df=df_sample, varname=ipcf, var_pondera=pondera, alfa=0, zeta=lp, quiet = TRUE)
  store[i] = fgt

}
```

```
store
```

```
## [1] 25.06 25.31 25.00 25.24 25.27 25.26 25.32 25.06 25.01 25.16 25.26 25.36
```

```
## [13] 25.28 25.48 25.21 24.98 25.40 25.28 25.38 25.25 25.18 25.27 25.29 25.05
## [25] 25.17 25.18 25.05 25.36 25.08 24.98 25.13 24.91 25.26 25.26 25.13 25.17
## [37] 25.09 25.47 25.24 24.90 25.04 25.13 25.41 25.49 25.02 25.02 25.48 25.31
## [49] 25.26 25.20
```

```
sd=sd(store)
mean=mean(store)
n=length(store)

#con intervalo de confianza del 95%
error <- qt(0.975,df=n-1)*sd/sqrt(n)
left <- mean - error; left
```

```
## [1] 25.15862
```

```
right <- mean + error; right
```

```
## [1] 25.24538
```

A modo de extensión es fácil escribir una función que, tomando los resultados alojados en un vector y el intervalo de confianza deseado, nos devuelve directamente el cálculo del intervalo.

```
#también es posible hacer un función para calcular los intervalos de confianza
ci <- function(vector, intervalo){

  sd=sd(vector)
  mean=mean(vector)
  n= length(vector)

  error <- qt((intervalo+1)/2, df=n-1) * sd/sqrt(n)
  result <- c("lower" = mean - error, "upper" = mean + error)
  return(result)
}
```

```
ci(store, 0.90)
```

```
##      lower      upper
## 25.16581 25.23819
```

```
ci(store, 0.95)
```

```
##      lower      upper
## 25.15862 25.24538
```

```
ci(store, 0.99)
```

```
##      lower      upper  
## 25.14416 25.25984
```


Chapter 5

Capítulo 5

Apéndices del Libro Pobreza y Desigualdad en R

Escrito por: Cristian Bonavida

Last Update: 02/7/2021

Códigos escritos en base a los apéndices del libro “Pobreza y Desigualdad en América Latina” de Gasparini, Cicowiez y Sosa Escudero. El objeto de este material es reproducir la rutina de códigos para STATA presentada en el libro al lenguaje R. Este material es solo de caracter complementario a las explicaciones y detalles conceptuales que se presentan en el libro de texto y los apéndices

Set Inicial

Cargo las librerías, limpio environment, defino el path y atajo para funcion paste

```
library(dplyr)
library(tidyverse) # Data wrangling
library(tidygraph)
library(readxl)
library(ggplot2)
library(foreign)
library(TAM)
library(margins)
```

```
rm(list=ls())      #empiezo limpiando todo

"%+%" <- function(x,y) paste(x,y,sep = "")      # defino un shorcut para concat de te
data_dir <- "C:/Users/HP/Desktop/CEDLAS - UNLP/Apendices en R/Material libro/encuestas"
```

5.1 Pobreza Multidimensional

###- (pág. 334-335)

En este primer apartado se muestra cómo puede replicarse el cuadro 5.1 del libro, sobre tasas de pobreza multidimensional en Nicaragua, Perú y Uruguay. Comenzamos cargando y definiendo la base a utilizar

```
nic05 <- read.dta(data_dir %+% "Nic/2005/bases/nic05_cedlas.dta")
df <- nic05
```

Nos asegurarnos de convertir los missings que puede contener la variable pondera a 0. De esta manera al calcular estimadores ponderados los valores para estas observaciones no tienen ningun peso y evitamos que los missings afecten el nuestros cálculos.

```
df$pondera <- ifelse(is.na(df$pondera), 0, df$pondera)
```

El bloque de código siguiente asigna a todos los miembros del hogar las variables que solo están definidas para el jefe de hogar. Como se trata de variables relacionadas con las características de la vivienda, típicamente se encuentran en las bases de datos de hogares y no de personas.

En STATA la solución se propone con un bucle. En este caso aplicamos la misma lógica pero nos valemos de la función `across` que nos permite realizar un mismo cálculo para un conjunto de columnas especificadas. Previamente debemos ordenar y agrupar las observaciones por id.

```
df <- df %>% arrange(id) %>% group_by(id) %>%
  mutate(
    across(
      .cols = c(habita, matpreca, agua, banio), #sobre qué columnas aplicar operac
      .fns = mean,      #que operación/función queremos realizar
      na.rm = TRUE,
      .names = "{col}" #como deben llamarse las nuevas variables
    )
  )
```

La expresión `names = "{col}"` indica que el nombre de las nuevas columnas sea el nombre de las variables originales, por lo que las estamos sobrescribiendo.

En las siguiente bloque de código calculamos los indicadores de pobreza para cada una de las dimensiones relevantes, creando una nueva variable en el data frame para cada caso. Empleamos el comando `ifelse` cuando el cálculo es directo y `mutate` cuando se requiere de una variable auxiliar previa.

```
# (1) ipcf < 2.5 USD
df$indic1 <- ifelse(df$ipcf < 564.12, 1 ,0)

# (2) mas de 3 miembros por cuarto
df <- df %>% mutate(
  rat_miembros_cuartos = miembros/habita,
  aux = case_when(
    (rat_miembros_cuartos>3 & !is.na(rat_miembros_cuartos)) ~ 1,
    (rat_miembros_cuartos<=3 & !is.na(rat_miembros_cuartos)) ~ 0)) %>%
  group_by(id) %>% mutate(indic2=max(aux)) %>% select(-aux)

# (3) vivienda construida con material precario
df$indic3 <- ifelse(df$matpreca==1, 1 ,0)

# (4) vivienda sin acceso a agua potable
df$indic4 <- ifelse(df$agua==0, 1 ,0)

# (5) vivienda sin acceso a baño sanitario
df$indic5 <- ifelse(df$banio==0, 1 ,0)

# (6) educación promedio menor a 7 años solo para el jefe y conyuge
df <- df %>% mutate(aedu_avg = ifelse(jefe==1 | conyuge==1, mean(aedu, na.rm=TRUE), NA),
  aux = ifelse(aedu_avg<7 & jefe==1, 1 ,0)) %>% group_by(id) %>%
  mutate(indic6 = max(aux)) %>% select(-aux)
```

Seguidamente obtenemos el porcentaje de personas con privaciones para cada indicador, utilizando la media ponderada.

```
weighted.mean(df$indic1, df$pondera, na.rm=TRUE)*100
```

```
## [1] 42.66527
```

```
weighted.mean(df$indic2, df$pondera, na.rm=TRUE)*100
```

```
## [1] 35.15812
```

```
weighted.mean(df$indic3, df$pondera, na.rm=TRUE)*100
```

```
## [1] 14.43805
```

```
weighted.mean(df$indic4, df$pondera, na.rm=TRUE)*100
```

```
## [1] 37.31393
```

```
weighted.mean(df$indic5, df$pondera, na.rm=TRUE)*100
```

```
## [1] 73.09805
```

```
weighted.mean(df$indic6, df$pondera, na.rm=TRUE)*100
```

```
## [1] 73.3585
```

La variable “npriv” contiene el número de privaciones de cada individuo. Para crearla utilizamos una operación a nivel de fila con el comando `rowSums` que, combinado con `across`, nos permite sumar todas las columnas especificadas. En este caso especificamos todas las columnas que comienzan con el patron “indic”. Notese que a diferencia del uso anterior aqui no se realiza una misma operación repetida para cada columna sino que se especifican las columnas que se incluyen como argumento de la operación suma.

```
#contar condiciones:por fila sumo todas las columnas que comienzan con indic
df <- df %>% mutate(npriv = rowSums(across(starts_with("indic"))))
```

Una forma alternativa más intuitiva sería especificar manualmente las columnas a sumar, pero se vuelve poco eficiente en el caso de que estas sean numerosas, por lo que la posibilidad de identificar columnas por patrones se vuelve particularmente atractiva.

```
#manera alternativa
df$npriv = df$indic1 + df$indic2 + df$indic3 + df$indic4 + df$indic5 + df$indic6
```

A partir de la variable “npriv” se generan las variables “pobre1” a “pobre6” que valen 1 de acuerdo con la cantidad de privaciones que sufre cada individuo. Por ejemplo, la variable “pobre4” vale 1 para los individuos que tienen 4 o más privaciones, y 0 en caso contrario.

En cada iteración se concatena el prefijo *pobre* con el contador *i*, dándole nombre a cada nueva columna del data frame. Luego se calcula el porcentaje como la media ponderada de esta columna, se redondea e imprime el resultado.

```
#condición de pobreza segun cantidad de privaciones
for (i in 1:6){

  df[paste("pobre",i,sep="")] <- ifelse(df$npriv>=i, 1, 0)

  p = weighted.mean(df[paste("pobre",i,sep="")], df["pondera"], na.rm=TRUE)*100
  print(paste(i, " privaciones = ", round(p, d=2), "%", sep = ""))

}
```

```
## [1] "1 privaciones = 86.75%"
## [1] "2 privaciones = 73.02%"
## [1] "3 privaciones = 56.6%"
## [1] "4 privaciones = 37.67%"
## [1] "5 privaciones = 19.2%"
## [1] "6 privaciones = 4.51%"
```

5.2 Índice Bourguignon y Chakravarty (BC) - Pobreza Multidimensional

###- (pág. 335-336)

El código a continuación permite reproducir el cuadro 5.2 del texto sobre pobreza multidimensional computada con el índice de Bourguignon y Chakravarty (BC). El cómputo de dicho índice se realiza empleando solo las observaciones que tienen información para las tres dimensiones consideradas en el texto; por lo que se eliminan las observaciones con missing en al menos una de esas dimensiones. Se eligen los valores para los parametros relevantes y se fija el número de dimensiones a considerar

```
df <- df %>% filter(!is.na(ipcf), !is.na(aedu_avg), !is.na(rat_miembros_cuartos))

theta=1
alpha=2

dim_t=3    #total dimensiones
```

Posteriormente almacenamos los valores de las observacion en una lista, donde cada elemento contiene todos los valores de cada una de las 3 variables a considerar. En vectores separados almacenamos los umbrales y los pesos

```
dimension <- list( df$ipcf,           # (1) ipcf
                  1/df$rat_miembros_cuartos, # (2) ratio de miembros por cuarto
                  df$aedu_avg           # (3) educación promedio de jefe y conyug
                  )

umbral <- c(564.119195, 1/3, 7)      #valores para los umbrales de cada dimensiones
wt <- c(1, 1, 1)                   #wt correspondiente
```

El objeto `brechas` se define como vacío y cada uno de sus elementos se genera en las iteraciones sucesivas del bucle al comparar cada valor de la variable contra los umbrales fijados. El objeto “`suma_brechas`” se crea como un vector único con valores 0, y luego se reemplaza iterativamente para computar la suma de brechas. El bucle itera n veces en total, siendo n la cantidad de dimensiones relevadas. En cada iteración replica la fórmula de BC para cada dimensión.

```
#defino la lista brechas como vacia para generar cada uno de sus elementos en el bucle
brecha <- list()

suma_brechas <- c(rep(0, nrow(df)))

for (i in 1:dim_t) {

  #generar brechas a partir de valores de las dimensiones vs umbrales
  brecha[[i]] <- ifelse(dimension[[i]]<umbral[i], 1-dimension[[i]]/umbral[i], 0)

  #construir brechas ponderadas
  brecha[[i]] <- wt[i]/dim_t * (brecha[[i]]^theta)

  #computar suma de las brechas. Suma_brechas será = 0 solo si todas las brechas son 0
  suma_brechas = suma_brechas + brecha[[i]]
}
```

Finalmente se calcula, para cada individuo, la suma de las brechas ponderadas elevadas a la potencia θ , siempre que la suma de las brechas sea distinta de cero. Por último, se computa el índice BC como el cociente entre la suma ponderada de las brechas individuales almacenadas en la variable `suma_brechas` y la población de referencia.

```
suma_brechas = ifelse(suma_brechas!=0, suma_brechas^(alpha/theta), suma_brechas)

BC = round(sum(suma_brechas*df$pondera)/sum(df$pondera), d=3)

print(paste("BC =", BC))
```

5.3. INDICE ALKIRE Y FOSTER (AF) - POBREZA MULTIDIMENSIONAL71

```
## [1] "BC = 0.08"
```

5.3 Índice Alkire y Foster (AF) - Pobreza Multidimensional

###- (pág. 337-338)

A continuación se replica la formula de Alkire y Foster que permite replicar el cuadro 5.3 del texto. Las primeras líneas son idénticas al caso anterior, cambiando los parámetros de interés y agregando la lista “pobre” como objeto vacío.

```
k=2
alpha=2

dim_t=3

dim <- list( df$ipcf,
             1/df$rat_miembros_cuartos,
             df$aedu_avg
)

umbral <- c(564.119195, 1/3, 7)
wt <- c(1, 1, 1)

#defino la lista "brechas" Y pobre como vacia para generar cada uno de sus elementos en el bucle
brecha <- list()
pobre <- list()
```

Nuevamente el bucle itera sobre las n dimensiones fijada generando los valores de brecha para cada observación y ahora también completando el objeto binario “pobre” según el valor que toma la brecha. Al finalizar el bucle se construye el objeto “npriv” que contiene el número de dimensiones en que cada individuo fue identificado como pobre (vale cero para los individuos no pobres). El objeto “pobre_k” vale 1 para los individuos que son pobres en, al menos, k dimensiones.

```
for (i in 1:dim_t) {

  #generar brechas a partir de valores de las dimensiones vs umbrales
  brecha[[i]] <- ifelse(dim[[i]]<umbral[i], (1-dim[[i]]/umbral[i])^alpha, 0)
```

```

#identificar si es pobre en dimensión i
pobre[[i]] <- ifelse(brecha[[i]]!=0, 1, 0)

}

#identificar pobres en al menos k dimensiones
npriv = pobre[[1]] + pobre[[2]] + pobre[[3]]
pobre_k = ifelse(npriv>=k, 1, 0)

```

La línea siguiente genera la variable “suma_brechas” que, como antes, se emplea luego para almacenar la suma de las brechas en cada una de las dimensiones consideradas. El bucle constuye para cada dimensión los objetos necesarios para replicar la formula de AK

```

for (i in 1:dim_t) {

  #brechas positivas solo si el número de privaciones mayor a k
  brecha[[i]] <- ifelse(pobre_k!=1, 0, brecha[[i]])

  #construir brechas ponderadas
  brecha[[i]] <- wt[i] * brecha[[i]]

  #computar suma de las brechas. Suma_brechas será = 0 solo si todas las brechas son
  suma_brechas = suma_brechas + brecha[[i]]

}

```

Por último se computa y se redondean los valores del índice de AF

```

AK = round(sum(suma_brechas*df$pondera)/(dim_t*sum(df$pondera)), d=5)

print(paste("AK =", AK))

```

```
## [1] "AK = 0.12319"
```

5.4 Perfiles de Pobreza

###- (pág. 338-339)

El bloque de código a continuación puede emplearse para computar el perfil de pobreza monetaria para vivienda y servicios que se muestra en el cuadro 5.8. El código del ejemplo se aplica a la EPH (Encuesta Permanente de Hogares) de Paraguay para el año 2007. Luego de cargar y definir la base, la tercer línea de código genera la variable “hh” que vale 1 para una única observación de cada hogar, a partir de indentificar las observaciones duplicadas. El comando `duplicated` asigna valor `FALSE` a la primera observación del por hogar, y `TRUE` a todo el resto.

```
pry07 <- read.dta(data_dir %+% "Par/2007/bases/par07_cedlas.dta")
df <- pry07

df$hh <- ifelse(duplicated(df$id)==FALSE, 1,0)
```

Luego generamos en el data frame la variable indicativa de pobreza monetaria y a partir de ella computamos para el grupo de pobres y no pobres, el promedio de las variables “habita”, “matpreca”, “agua”, “banio” y “elect”, indicado la proporción de personas que cuenta con estos servicios.

```
df$pobre <- ifelse(df$ipcf<205970.366, 1, 0)

df %>% filter(hh==1 & !is.na(pobre)) %>% group_by(pobre) %>%
  summarise( mean_habita = weighted.mean(habita, pondera, na.rm = TRUE),
             mean_matpreca = weighted.mean(matpreca,pondera, na.rm = TRUE),
             mean_banio = weighted.mean(banio, pondera, na.rm = TRUE),
             mean_agua = weighted.mean(agua, pondera, na.rm = TRUE),
             mean_elect = weighted.mean(elect, pondera, na.rm = TRUE))
```

```
## # A tibble: 2 x 6
##   pobre mean_habita mean_matpreca mean_banio mean_agua mean_elect
##   <dbl>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
## 1     0         3.36         0.0108         0.761         0.975         0.978
## 2     1         2.53         0.0414         0.364         0.905         0.898
```

Con el comando `ttest` buscamos evaluar la significatividad estadística de estas diferencias de medias entre pobres y no pobres para las variables incluidas. Para esto elegimos el nivel de confianza y definimos un bucle que itera sobre cada variable respectiva del data frame. Dentro de él construimos un objeto “x” que contiene los valores para estas variables sólo para una única observación por hogar y para el grupo de pobres y otro objeto “y” con los mismos datos para el grupo de no pobres. A partir de estos objetos se evalúan las significatividad de la diferencia de medias entre ambos grupos y se reporta si su p-valor es mayor al nivel de confianza fijado.

```

set_confidence = 95
confidence = 1 - set_confidence/100

dim <- c("habita", "matpreca", "banio", "agua", "elect")
for (i in dim){

  print(i)
  x <- df[df$pobre==1 & df$hh==1, colnames(df)==i]
  y <- df[df$pobre==0 & df$hh==1, colnames(df)==i]

  test = t.test(x,y)
  print(paste("No significant mean diff:", test$p.value > confidence))

}

```

```

## [1] "habita"
## [1] "No significant mean diff: FALSE"
## [1] "matpreca"
## [1] "No significant mean diff: FALSE"
## [1] "banio"
## [1] "No significant mean diff: FALSE"
## [1] "agua"
## [1] "No significant mean diff: FALSE"
## [1] "elect"
## [1] "No significant mean diff: FALSE"

```

5.5 Perfiles de Pobreza Condicionados

###- (pág. 340-341)

El bloque de código siguiente permite replicar el cuadro 5.12, que muestra perfiles condicionados de pobreza. En el ejemplo se emplea la encuesta de México para el año 2006. Luego de cargar la base eliminamos las observaciones incoherentes y al igual que antes generamos la variable indicativa de pobreza monetaria línea de 2.5 dólares. Las líneas siguientes agregan nuevas variables al data frame que suman la cantidad de individuos que pertenecen a distintos grupos etarios dentro de cada hogar, calculan el ratio de miembros por habitaciones y los valores de educación y edad al cuadrado

```

#cargo base
mex06 <- read.dta(data_dir %+"Mex/2006/bases/mex06_cedlas.dta")

```

```
df <- mex06

df <- df %>% filter(cohh==1)
df$pobre <- ifelse(df$ipcf<608.24533, 1, 0)

#número de miembros en cada grupo
df <- df %>% arrange(id) %>% group_by(id) %>%
  mutate( miembros_edad_0015 = sum(ifelse(edad<=15, 1, 0)),
          miembros_edad_1625 = sum(ifelse(edad %in% (16:25), 1, 0)),
          miembros_edad_2640 = sum(ifelse(edad %in% (26:40), 1, 0)),
          miembros_edad_4160 = sum(ifelse(edad %in% (41:64), 1, 0)),
          miembros_edad_65mas= sum(ifelse(edad>=65, 1, 0)),

          rat_miembros_cuartos = miembros/habita,

          aedu2=aedu^2,
          edad2=edad^2 )
```

Las líneas siguientes contienen la sentencia que estima, para los jefes de hogar, el modelo probit para la probabilidad de ser pobre. Para ello empleamos el comando `glm` en el cual definimos la variable independiente y todo el conjunto de regresores, indicamos el data frame referido y la familia de modelos que buscamos estimar. Esta estimación la guardamos en el objeto “probit” que luego visualizamos con un `summary`

Para computar los efectos marginales para el rango 0 a 22 años de educación del jefe de hogar, empleamos el comando `margins`. Para ello indicamos donde almacenamos nuestra estimación (objeto probit), para qué variables deseamos calcular los efectos (aedu) y sobre qué valores evaluarlos (0:22).

```
store <- summary(margins(probit, variables = "aedu", at = list(aedu = 0:22)))
store
```

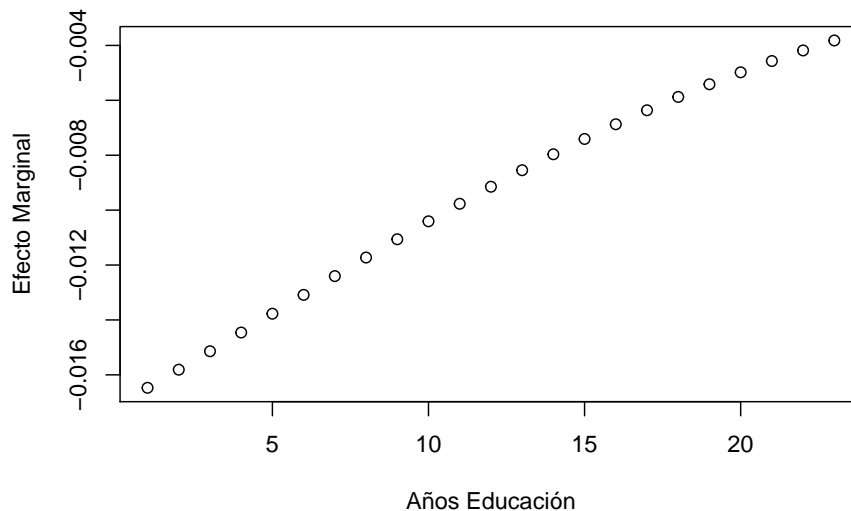
##	factor	aedu	AME	SE	z	p	lower	upper
##	aedu	0.0000	-0.0165	0.0029	-5.6272	0.0000	-0.0222	-0.0107
##	aedu	1.0000	-0.0158	0.0028	-5.7352	0.0000	-0.0212	-0.0104
##	aedu	2.0000	-0.0151	0.0026	-5.8741	0.0000	-0.0202	-0.0101
##	aedu	3.0000	-0.0145	0.0024	-6.0478	0.0000	-0.0191	-0.0098
##	aedu	4.0000	-0.0138	0.0022	-6.2616	0.0000	-0.0181	-0.0095
##	aedu	5.0000	-0.0131	0.0020	-6.5223	0.0000	-0.0170	-0.0092
##	aedu	6.0000	-0.0124	0.0018	-6.8395	0.0000	-0.0160	-0.0088
##	aedu	7.0000	-0.0117	0.0016	-7.2263	0.0000	-0.0149	-0.0085
##	aedu	8.0000	-0.0111	0.0014	-7.7006	0.0000	-0.0139	-0.0082
##	aedu	9.0000	-0.0104	0.0013	-8.2882	0.0000	-0.0129	-0.0079
##	aedu	10.0000	-0.0098	0.0011	-9.0270	0.0000	-0.0119	-0.0076

```
##      aedu 11.0000 -0.0091 0.0009 -9.9743 0.0000 -0.0109 -0.0074
##      aedu 12.0000 -0.0085 0.0008 -11.2219 0.0000 -0.0100 -0.0071
##      aedu 13.0000 -0.0080 0.0006 -12.9242 0.0000 -0.0092 -0.0068
##      aedu 14.0000 -0.0074 0.0005 -15.3614 0.0000 -0.0084 -0.0065
##      aedu 15.0000 -0.0069 0.0004 -19.0911 0.0000 -0.0076 -0.0062
##      aedu 16.0000 -0.0064 0.0003 -25.3519 0.0000 -0.0069 -0.0059
##      aedu 17.0000 -0.0059 0.0002 -37.0835 0.0000 -0.0062 -0.0056
##      aedu 18.0000 -0.0054 0.0001 -56.6339 0.0000 -0.0056 -0.0052
##      aedu 19.0000 -0.0050 0.0001 -52.3155 0.0000 -0.0052 -0.0048
##      aedu 20.0000 -0.0046 0.0001 -32.5281 0.0000 -0.0048 -0.0043
##      aedu 21.0000 -0.0042 0.0002 -21.7719 0.0000 -0.0046 -0.0038
##      aedu 22.0000 -0.0038 0.0002 -15.9736 0.0000 -0.0043 -0.0034
```

Esta información la guardamos en un objeto llamado “store”, del cual nos interesa recuperar el valor de los coeficientes almacenados bajo el nombre AME (Average Mean Effect). Con ellos generamos un vector “y” que denota los efectos para cada valor de años de educación, los cuales guardamos en el vector “x”. Finalmente graficamos la relación.

```
y = store$AME
x = seq(0:22)

plot(x, y,
      ylab = "Efecto Marginal",
      xlab = "Años Educación")
```



Una forma más directa de graficar los efectos marginales es a partir del comando `cplot` de la familia `margins`, que estima automáticamente estos mismos valores a partir de la estimación del modelo `probit`

```
#cplot(probit, "aedu", what = "effect", main = "Average Marginal Effect of Weight")
```