

# TP 1 - Muestreo

Gomez Vargas Andrea, Iummato Luciana, Pesce Andrea Gisele

2024-11-11

## Contenido

Paquetes de trabajo	1
Ejercicio I	2
Ejercicio II	7
Ejercicio III	12
Ejercicio IV	19
Ejercicio V (continuación del ejercicio IV)	33
Ejercicio VI	38
Ejercicio VII	39
Ejercicio VIII	40

## Paquetes de trabajo

```
library(tidyverse)
library(survey)
library(readxl)
library(gt)
library(sampling)
library(VIM)
library(binom)
library(openxlsx)

options(scipen = 999)
```

## Ejercicio I

```
datos <- "
Alumno X Y
a 6 14.0
b 9 20.0
c 5 12.0
d 4 10.0
e 2 5.0
f 7 12.0
g 10 24.0
h 4 5.0
i 12 21.0
j 5 9.0
k 8 18.0
l 12 20.0
m 5 8.0
n 9 15.0
o 2 2.5
p 6 11.0
q 11 20.0
r 8 15.0
"

# Convertimos los datos en un data.frame
ejercicio_1 <- read.table(text = datos, header = TRUE)

N= nrow(ejercicio_1)
n=9

#parámetros
parametro_mediaX= sum(ejercicio_1$X)/N
parametro_mediaY=sum(ejercicio_1$Y)/N
parametro_razonX_Y= sum(ejercicio_1$X)/sum(ejercicio_1$Y)

# Generar todas las combinaciones posibles
muestras_posibles <- combn(N, n)

# Ver la cantidad de combinaciones posibles
cantidad_muestras <- ncol(muestras_posibles)
cantidad_muestras

## [1] 48620

#RESOLUCIÓN PUNTOS 2 Y 3 CON UNA MUESTRA SELECCIONADA
# Seleccionamos ahora una muestra aleatoria simple con R de tamaño n
s_mas <- sample(N,n, replace=FALSE)

muestra <- ejercicio_1[s_mas,]

#estimadores para esa muestra
```

```

estimador_mediaX= sum(muestra$X)/n
estimador_mediaY=sum(muestra$Y)/n
estimador_razonX_Y= estimador_mediaX/estimador_mediaY

#varianza de estimadores
s_cuadradoX=var(muestra$X)
s_cuadradoY=var(muestra$Y)

# varianza para la estimacion de la muestra:
# En el MAS
# Var(y_media) = (1-n/N)*S^2/n
Var_X_media = (1-n/N)*s_cuadradoX/n
Var_Y_media = (1-n/N)*s_cuadradoY/n

#coeficiente de variacion para la estimacion de la muestra:
CVMAS_X <- 100*sqrt(Var_X_media)/estimador_mediaX

CVMAS_X

```

```
## [1] 9.10937
```

```

CVMAS_Y <- 100*sqrt(Var_Y_media)/estimador_mediaY

CVMAS_Y

```

```
## [1] 7.246885
```

```

#CALCULO CON SURVEY PARA ESA MUESTRA EN PARTICULAR
muestra$R <- muestra$X/muestra$Y
muestra$pondera <- N/n

muestra$fpc <- N

diseno <- svydesign(id= ~1,weights=~pondera, data=muestra, fpc=~fpc)
diseno

```

```

## Independent Sampling design
## svydesign(id = ~1, weights = ~pondera, data = muestra, fpc = ~fpc)

```

```

# Calcular las medias
media_X <- svymean(~X, diseno)
media_Y <- svymean(~Y, diseno)

# Calcular la razón
razon <- as.numeric(media_X / media_Y)

# Mostrar la razón estimada
razon

```

```
## [1] 0.5114504
```

```

raz <- svyratio(numerator = ~ X, denominator = ~ Y, diseno)

var_razon = SE(raz)^2

cv(raz) * 100

```

```

##           Y
## X 3.743751

```

```

# RESOLUCIÓN DE PUNTOS 2 Y 3 CON EL CALCULO DE TODAS LAS MUESTRAS POSIBLES
# Crear un vector para almacenar los valores de estimadores para cada muestra

```

```

medias_X <- numeric(ncol(muestras_posibles))
medias_Y <- numeric(ncol(muestras_posibles))
razon <- numeric(ncol(muestras_posibles))

```

```

# Calcular estimadores para cada muestra
for (i in 1:ncol(muestras_posibles)) {
  indices <- muestras_posibles[, i]
  medias_X[i] <- mean(ejercicio_1$X[indices])
  medias_Y[i] <- mean(ejercicio_1$Y[indices])
  razon[i] <- medias_X[i] / medias_Y[i]
}

```

```

# Crear el data frame con las tres columnas
muestras_posibles_estimadores <- data.frame(medias_X, medias_Y, razon)

```

```

#esperanza del estimador
esperanza_mediaX=mean(muestras_posibles_estimadores$medias_X) #insesgado
esperanza_mediaY=mean(muestras_posibles_estimadores$medias_Y) #insesgado
esperanza_razon=mean(muestras_posibles_estimadores$razon) #aproximado

esperanza_mediaX

```

```

## [1] 6.944444

```

```

esperanza_mediaY

```

```

## [1] 13.41667

```

```

esperanza_razon

```

```

## [1] 0.5182023

```

```

# Calcular la varianza de los estimadores
varianza_medias_X <- var(muestras_posibles_estimadores$medias_X)
varianza_medias_Y <- var(muestras_posibles_estimadores$medias_Y)
varianza_razon <- var(muestras_posibles_estimadores$razon)

```

```

varianza_medias_X

```

```
## [1] 0.5455813
```

```
varianza_medias_Y
```

```
## [1] 2.147511
```

```
varianza_razon
```

```
## [1] 0.0004556671
```

```
# Calcular la varianza de los estimadores
```

```
CVestimador_X <- 100*sqrt(varianza_medias_X)/esperanza_mediaX
```

```
CVestimador_y <- 100*sqrt(varianza_medias_Y)/esperanza_mediaY
```

```
CVestimador_r <- 100*sqrt(varianza_razon)/esperanza_razon
```

```
CVestimador_X
```

```
## [1] 10.63634
```

```
CVestimador_y
```

```
## [1] 10.92253
```

```
CVestimador_r
```

```
## [1] 4.11931
```

```
#SELECCIÓN DE 10000 MUESTRAS
```

```
# Creo una funcion que seleccione una muestra de tamaño x
```

```
# y estime total de poblacion
```

```
estimo_diezmil <- function() {
```

```
  muestras <- ejercicio_1[sample(nrow(ejercicio_1), 9, replace = FALSE), ] # Tamaño de muestra de 9
```

```
  estim_X <- mean(muestras$X) # Calcular la media de X
```

```
  estim_Y <- mean(muestras$Y) # Calcular la media de Y
```

```
  a <- c(estim_X, estim_Y) # Crear un vector con las estimaciones
```

```
  return(a) # Devolver el vector
```

```
}
```

```
# Crear una lista para almacenar las estimaciones
```

```
lista_estim <- lapply(1:10000, function(x) estimo_diezmil())
```

```
# Convertir la lista en un data frame
```

```
df_estimdiezmil <- data.frame(matrix(unlist(lista_estim),  
                                     nrow = length(lista_estim), byrow = TRUE))
```

```
# Asignar nombres a las columnas
```

```
colnames(df_estimdiezmil) <- c("media_X", "media_Y")
```

```
df_estimdiezmil$razon <- df_estimdiezmil$media_X/df_estimdiezmil$media_Y
```

```
# Calcular la varianza de los estimadores
var_medias_X <- var(df_estimdiezmil$media_X)
var_medias_Y <- var(df_estimdiezmil$media_Y)
var_razon <- var(df_estimdiezmil$razon)
```

```
var_medias_X
```

```
## [1] 0.5536014
```

```
var_medias_Y
```

```
## [1] 2.175541
```

```
var_razon
```

```
## [1] 0.0004490736
```

```
# Calcular la varianza de los estimadores
CV_X <- 100*sqrt(var_medias_X)/mean(df_estimdiezmil$media_X)
CV_y <- 100*sqrt(var_medias_Y)/mean(df_estimdiezmil$media_Y)
CV_r <- 100*sqrt(var_razon)/mean(df_estimdiezmil$razon)
```

```
CV_X
```

```
## [1] 10.72857
```

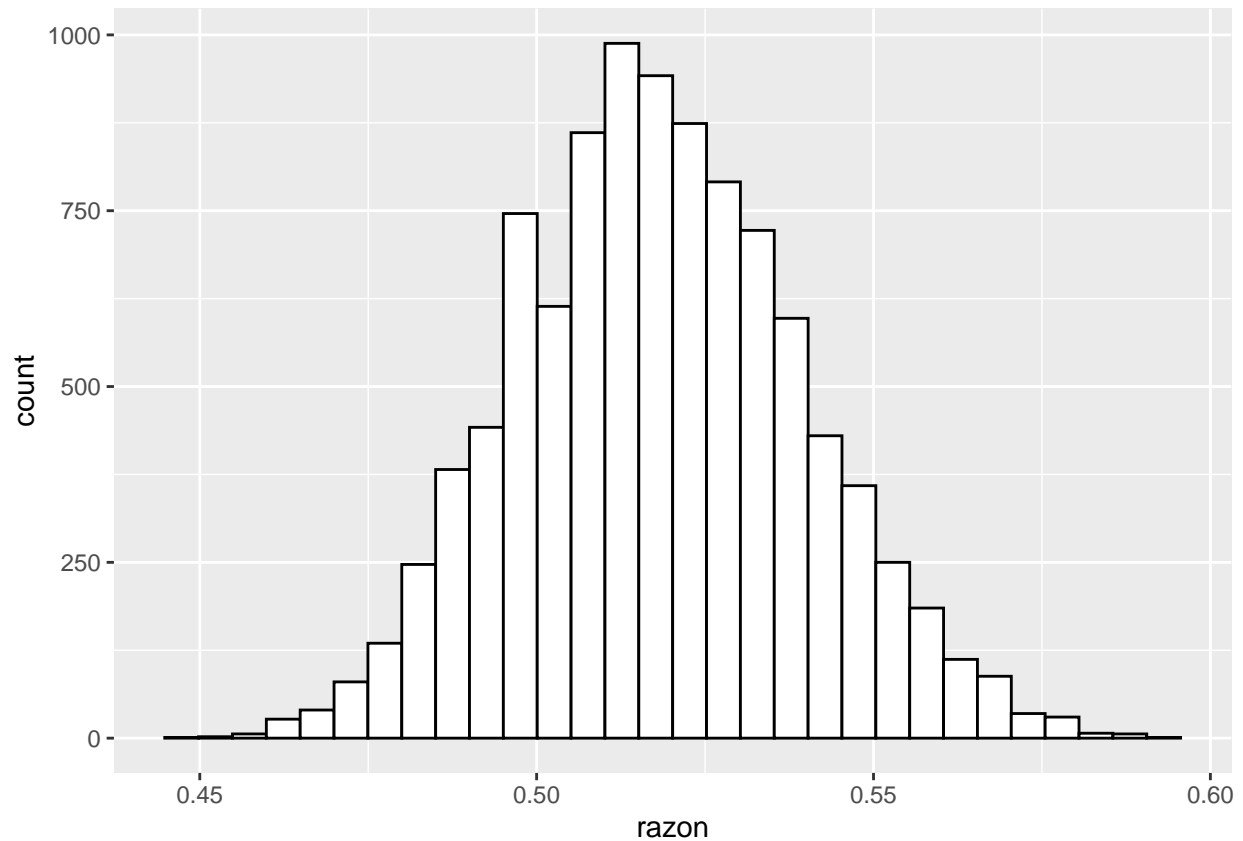
```
CV_y
```

```
## [1] 11.00912
```

```
CV_r
```

```
## [1] 4.089139
```

```
# Grafico la distribucion de las estimaciones
p <-ggplot(df_estimdiezmil, aes(x=razon)) +
  geom_histogram(color="black", fill="white")
p
```



## Ejercicio II

```
df_tabla <- read_excel("tabla_muestras_posibles.xlsx")

#TODAS LAS MUESTRAS POSIBLES
df_muestras <- data.frame(matrix(unlist(combn(df_tabla$Y,10, simplify = FALSE)),
                                ncol=10, byrow=TRUE))

#1 y 2. MEDIA, MEDIANA, MEDIA TRUNCADA DE TODAS LAS MUESTRAS POSIBLES

df_muestras[, 1:10] <- t(apply(df_muestras[, 1:10], 1, sort))

media <- apply(df_muestras[,1:10],1,mean)
mediana <- apply(df_muestras[,1:10],1,median)
media_t <- apply(df_muestras[,2:9],1,mean)

df_muestras$Media <- media
df_muestras$Mediana <- mediana
df_muestras$MediaTrunc <- media_t

#verificación de que la media es un estimador insesgado de la media poblacional
```

```
media_poblacional=mean(df_tabla$Y)
mediana_poblacional=median(df_tabla$Y)
mediaTrunc_poblacional=mean(df_tabla$Y, trim = 0.1)
```

```
media_poblacional
```

```
## [1] 42.435
```

```
mediana_poblacional
```

```
## [1] 28.37879
```

```
mediaTrunc_poblacional
```

```
## [1] 39.34896
```

```
media_muestras=mean(df_muestras$Media)
mediana_muestras=mean(df_muestras$Mediana)
mediaTrunc_muestras=mean(df_muestras$MediaTrunc)
```

```
media_muestras
```

```
## [1] 42.435
```

```
mediana_muestras
```

```
## [1] 30.81752
```

```
mediaTrunc_muestras
```

```
## [1] 40.29746
```

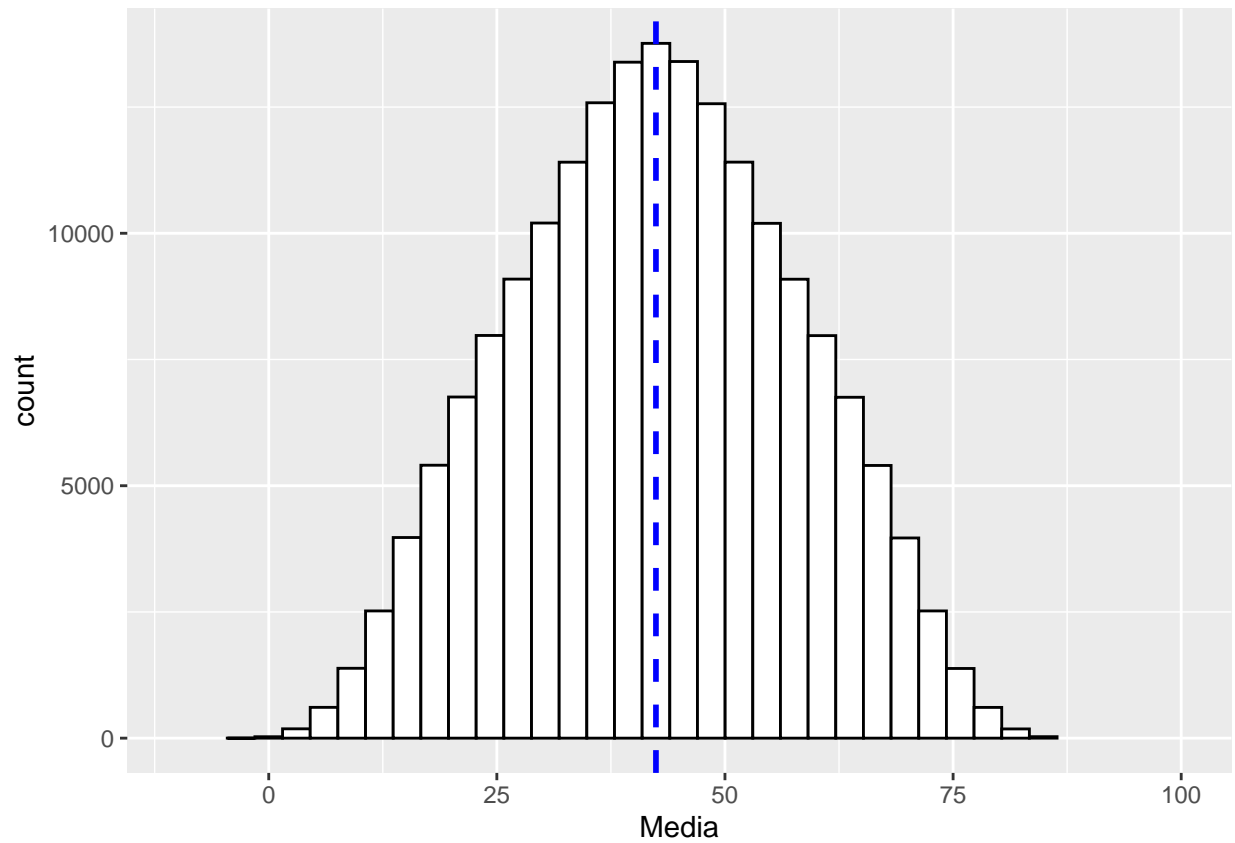
```
#GRÁFICO DE LAS 3 ESTIMACIONES
```

```
#media
```

```
p <- ggplot(df_muestras, aes(x=Media)) +
  geom_histogram(bins=30, color="black", fill="white") +
  coord_cartesian(xlim = c(-10, 100))

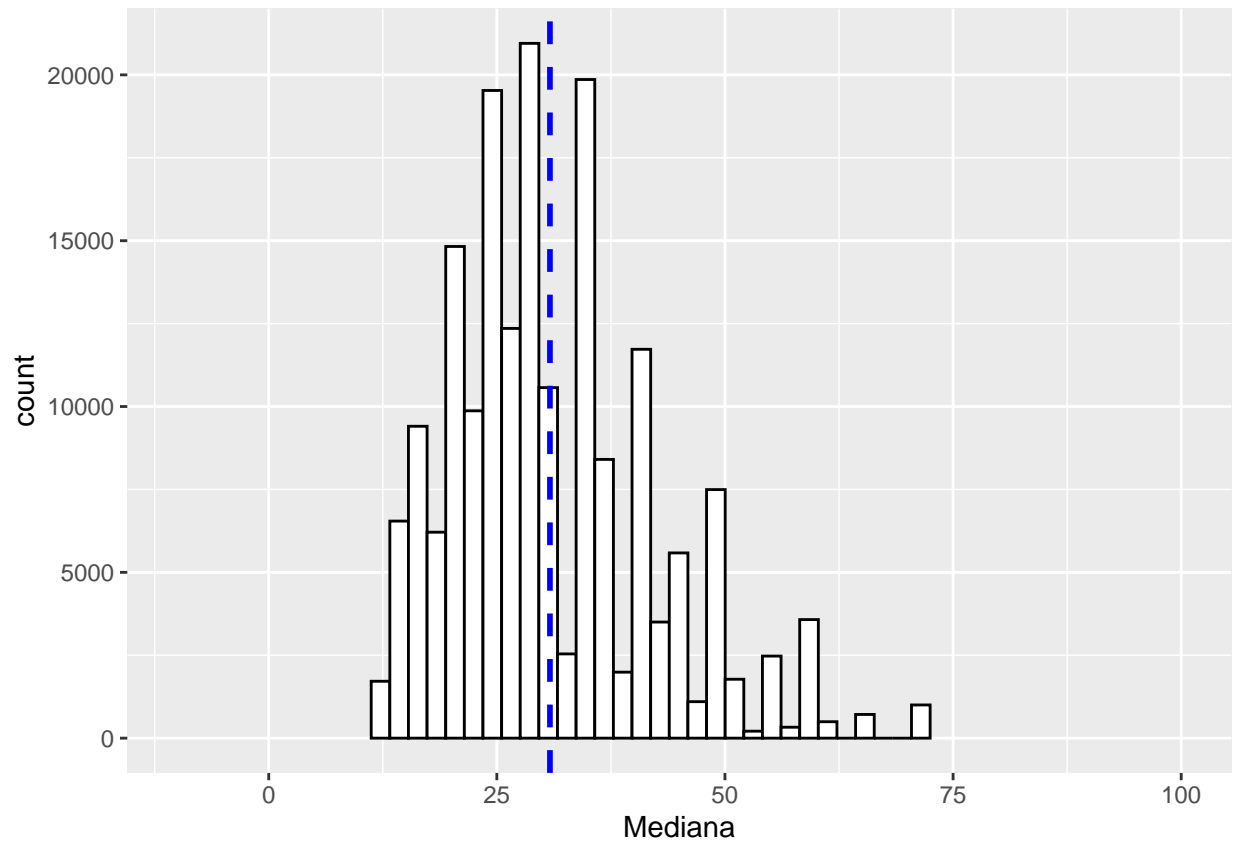
p <- p+ geom_vline(aes(xintercept=mean(Media)),
  color="blue", linetype="dashed", linewidth=1)
p
```





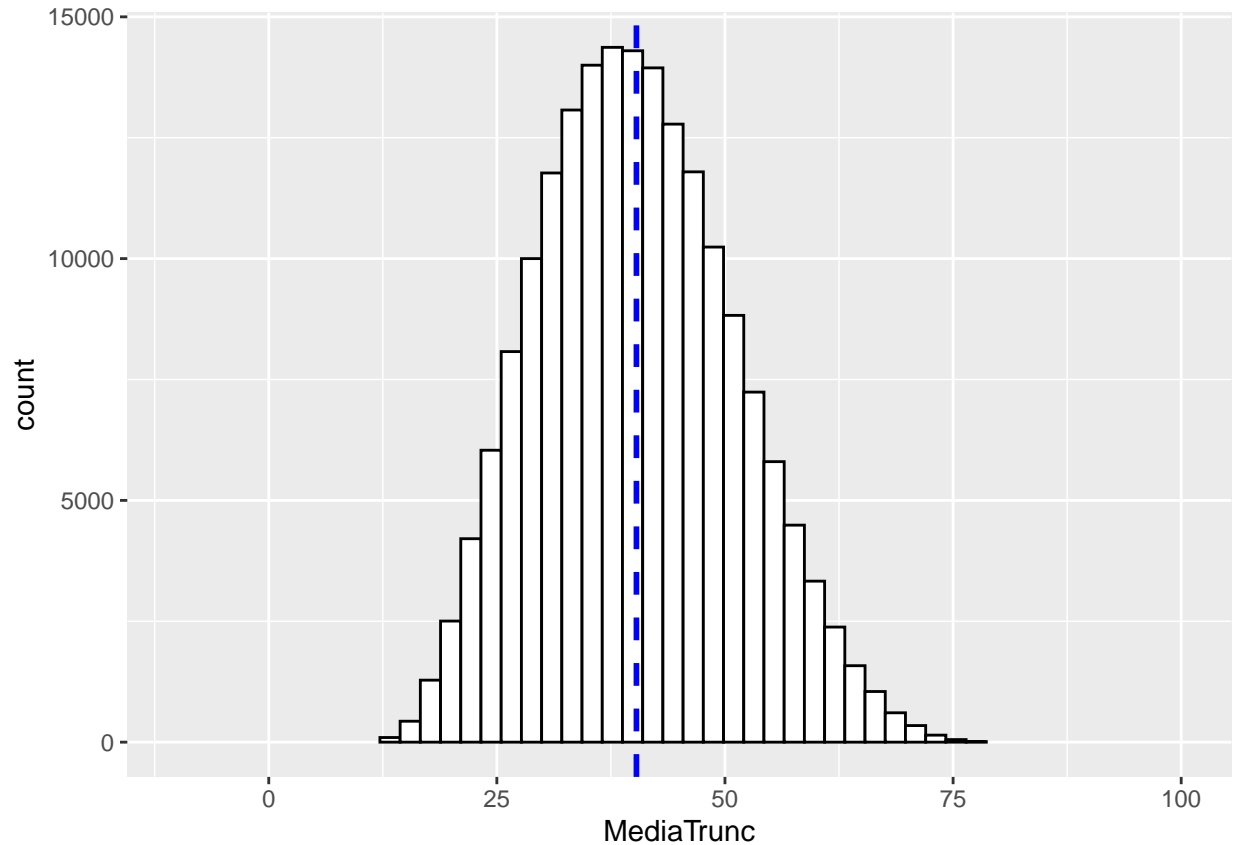
```
#mediana
p_median <- ggplot(df_muestras, aes(x=Mediana)) +
  geom_histogram(bins=30, color="black", fill="white") +
  coord_cartesian(xlim = c(-10, 100))

p_median <- p_median+ geom_vline(aes(xintercept=mean(Mediana)),
  color="blue", linetype="dashed", linewidth=1)
p_median
```



```
#media truncada
p_mediat <- ggplot(df_muestras, aes(x=MediaTrunc)) +
  geom_histogram(bins=30, color="black", fill="white") +
  coord_cartesian(xlim = c(-10, 100))

p_mediat <- p_mediat + geom_vline(aes(xintercept=mean(MediaTrunc)),
  color="blue", linetype="dashed", linewidth=1)
p_mediat
```



*#cv y EMC de los 3 estimadores*

```
cv_media <- sd(df_muestras$Media) / mean(df_muestras$Media)*100
cv_mediana <- sd(df_muestras$Mediana) / mean(df_muestras$Mediana)*100
cv_mediaT <- sd(df_muestras$MediaTrunc) / mean(df_muestras$MediaTrunc)*100

var_media<-var(df_muestras$Media)
var_mediana<-var(df_muestras$Mediana)
var_mediaT<-var(df_muestras$MediaTrunc)

sesgo_media<-media_poblacional - media_muestras
sesgo_mediana<-mediana_poblacional- mediana_muestras
sesgo_mediaT<-mediaTrunc_poblacional - mediaTrunc_muestras

emc_media <- var_media + sesgo_media^2
emc_mediana <- var_mediana + sesgo_mediana^2
emc_mediaT <- var_mediaT + sesgo_mediaT^2

resultados <- data.frame(
  Estimador = c("Media", "Mediana", "Media Truncada"),
  CV = c(cv_media, cv_mediana, cv_mediaT),
  EMC = c(emc_media, emc_mediana, emc_mediaT)
)

gt(resultados)
```

Estimador	CV	EMC
Media	36.24570	236.5710
Mediana	36.73640	134.1180
Media Truncada	26.55862	115.4421

*#parece que el mejor estimador es la media truncada porque tiene menos varianza*

## Ejercicio III

```
# Lectura de las tablas
radios_sexo <- read_excel("cen2010_radios_sexo.xlsx")
radios_tipo <- read_excel("cen2010_radios_tipo.xlsx")
radios_bienes <- read_excel("cen2010_radios_bienes.xlsx")

# Juntamos los archivos para unificarlos en un unico dataframe
radios_2010 = merge(radios_sexo, radios_bienes, by = "Codigo")
radios_2010 = merge(radios_2010, radios_tipo, by = "Codigo")

# Generamos nuevas variables en el dataframe que utilizaremos mas adelante

# Poblacion total en cada radio
radios_2010$Pob_radio <- radios_2010$Varon + radios_2010$Mujer

# Hogares rancho- casilla
radios_2010$Rancho_casilla <- radios_2010$Rancho + radios_2010$Casilla

# Cantidad de viviendas en cada radio
radios_2010$Viv_radio <- radios_2010$Casa +
  radios_2010$Rancho +
  radios_2010$Casilla +
  radios_2010$Departamento +
  radios_2010$Inquilinato +
  radios_2010$Hotel_pension

# Total de hogares
radios_2010$Hogares <- radios_2010$HeladeraSi + radios_2010$HeladeraNo

# Extraemos el codigo de provincia
radios_2010$prov <- floor(radios_2010$Codigo/10000000)

# Creamos etiqueta para las provincia2
radios_2010 <- radios_2010 %>%
  mutate(Provincia = case_when(prov == 2 ~ 'CABA', prov == 6 ~ 'BsAs',
                                prov == 10 ~ 'Catamarca', prov == 14 ~ 'Cordoba',
                                prov == 18 ~ 'Corrientes', prov == 22 ~ 'Chaco',
                                prov == 26 ~ 'Chubut', prov == 30 ~ 'Entre Rios',
```

```

prov == 34 ~ 'Formosa', prov == 38 ~ 'Jujuy',
prov == 42 ~ 'La Pampa', prov == 46 ~ 'La Rioja',
prov == 50 ~ 'Mendoza', prov == 54 ~ 'Misiones',
prov == 58 ~ 'Neuquen', prov == 62 ~ 'Rio Negro',
prov == 66 ~ 'Salta', prov == 70 ~ 'San Juan',
prov == 74 ~ 'San Luis', prov == 78 ~ 'Santa Cruz',
prov == 82 ~ 'Santa Fe', prov == 86 ~ 'Santiago',
prov == 90 ~ 'Tucuman', prov == 94 ~ 'TdFuego'))

```

*# Eliminamos los radios sin viviendas*

```
radios_2010 <- radios_2010[radios_2010$Viv_radio>0,]
```

```
N=nrow(radios_2010)
```

```
n=240
```

*#parámetros*

```
poblacion<- sum(radios_2010$Pob_radio)
```

```
hogares_casa<-sum(radios_2010$Casa)
```

```
hogares_rancho<-sum(radios_2010$Rancho)+sum(radios_2010$Casilla)
```

```
prop_rancho<-(sum(radios_2010$Rancho)+sum(radios_2010$Casilla))/sum(radios_2010$Viv_radio)
```

*#CV*

*# Calculamos la varianza para la estimacion del total:*

*# En el MAS*

*#  $Var(y_{media}) = (1-n/N)*S^2/n$*

*#  $Var(N*y_{media}) = N^2*(1-n/N)*S^2/n$*

```
P=prop_rancho
```

```
Q=1-prop_rancho
```

```
P+Q
```

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

```
S2_pob <- var(radios_2010$Pob_radio)
```

```
S2_casa <- var(radios_2010$Casa)
```

```
S2_rancho <- var(radios_2010$Rancho + radios_2010$Casilla)
```

```
S2_prop <- P*Q
```

```
VarMAS_MediaPob <- (1-n/N)*S2_pob/n
```

```
VarMAS_MediaCasa <- (1-n/N)*S2_casa/n
```

```
VarMAS_MediaRancho <- (1-n/N)*S2_rancho/n
```

```
VarMAS_pob <- N^2*VarMAS_MediaPob #Recordar: Var(k*X)= k^2*Var(X)
```

```
VarMAS_casa <- N^2*VarMAS_MediaCasa
```

```
VarMAS_rancho <- N^2*VarMAS_MediaRancho
```

*# Calculamos el coeficiente de variacion*

```
CVMAS_Pob <- 100*sqrt(VarMAS_pob)/poblacion
```

```
CVMAS_Casa <- 100*sqrt(VarMAS_casa)/hogares_casa
```

```
CVMAS_Rancho <- 100*sqrt(VarMAS_rancho)/hogares_rancho
```

```
ds_prop<-sqrt(P*Q)
CVMAS_Prop <- 100 *(ds_prop/prop_rancho) #revisar

CVMAS_Pob
```

```
## [1] 4.162589
```

```
CVMAS_Casa
```

```
## [1] 4.127404
```

```
CVMAS_Rancho
```

```
## [1] 14.6853
```

```
CVMAS_Prop #revisar
```

```
## [1] 537.2547
```

El cv de la estimación de rancho o casilla es grande porque el N del universo es menor que en el caso de la población general y los hogares tipo casa. El cv de la proporción de rancho o casilla es grande porque el valor de P es muy pequeño.

```
#n=240*43
```

Para que el cv de la estimación de la población sea aproximadamente 2%, n debe ser 1.000. Para que el cv de la estimación de hogares tipo rancho o casilla sea aproximadamente 2%, n debe ser 10.000

```
# Seleccionamos ahora una muestra aleatoria simple con R
s_mas <- sample(N,n, replace=FALSE)

muestra_radios <- radios_2010[s_mas,]

#estimaciones muestrales N*y_media

muestra_pobla<-N*mean(muestra_radios$Pob_radio)
muestra_casa<-N*mean(muestra_radios$Casa)
muestra_rancho<-N*mean(muestra_radios$Rancho+muestra_radios$Casilla)
muestra_prop<-(sum(muestra_radios$Rancho)+sum(muestra_radios$Casilla))/sum(muestra_radios$Viv_radio)

# Agregamos al data frame el factor de expansion
# (recordar que seleccione una muestra aleatoria simple de radios)
muestra_radios$pondera <- N/n

# Cantidad total de unidades en el marco de muestreo
# lo necesitare luego para survey
muestra_radios$fpc <- N

#junto rancho y casilla
muestra_radios$rancho_casilla <- muestra_radios$Rancho + muestra_radios$Casilla
```

```

# El objeto 'diseno' contiene toda la informacion que sera empleada
# para realizar las estimaciones.

diseno <- svydesign(id= ~1,weights=~pondera, data=muestra_radios, fpc=~fpc)
diseno

## Independent Sampling design
## svydesign(id = ~1, weights = ~pondera, data = muestra_radios,
##      fpc = ~fpc)

# (como es una muestra aleatoria simple ponemos fpc)

# Por ejemplo, si queremos extraer los pesos de un diseno podemos utilizar
pesos <- weights(diseno)

# total población
EstTotalPob <- survey :: svytotal(~Pob_radio, diseno, deff=TRUE, cv=TRUE, ci=TRUE)
EstTotalPob

##              total      SE DEff
## Pob_radio 41776858 1882878      1

# Puedo ahora extraer diferentes valores:
survey :: cv(EstTotalPob)      # -> coeficiente de variacion

##              Pob_radio
## Pob_radio 0.04506989

deff(EstTotalPob)      # -> efecto de diseno

## Pob_radio
##              1

SE(EstTotalPob)      # -> desvio estandar

##              Pob_radio
## Pob_radio 1882878

confint(EstTotalPob) # -> intervalor de confianza (por defecto 95%)

##              2.5 %   97.5 %
## Pob_radio 38086484 45467232

cv(EstTotalPob)

##              Pob_radio
## Pob_radio 0.04506989

```

```

# O pasar los resultados a un data frame
df_EstTotalPob <- as.data.frame(EstTotalPob)

# Quiero cambiar el nombre de las columnas
colnames(df_EstTotalPob) <- c("Estimacion", "SE", "deff")

# Calculemos ahora el intervalo de confianza con un 90% de confianza
# (suponemos que el estimador es aprox normal)
df_EstTotalPob$Li <- df_EstTotalPob$Estimacion-1.64*df_EstTotalPob$SE
df_EstTotalPob$Ls <- df_EstTotalPob$Estimacion+1.64*df_EstTotalPob$SE

# Ahora calculo el CV del estimador
df_EstTotalPob$CV <- 100*df_EstTotalPob$SE/df_EstTotalPob$Estimacion

# total casas
EstTotalcasa <- survey :: svytotal(~Casa, diseno, deff=TRUE, cv=TRUE, ci=TRUE)
EstTotalcasa

```

```

##           total      SE Deff
## Casa 11031522  464576    1

```

```

# Puedo ahora extraer diferentes valores:
survey :: cv(EstTotalcasa)      # -> coeficiente de variacion

```

```

##           Casa
## Casa 0.04211351

```

```

deff(EstTotalcasa)      # -> efecto de diseno

```

```

## Casa
##      1

```

```

SE(EstTotalcasa)      # -> desvio estandar

```

```

##           Casa
## Casa 464576.2

```

```

confint(EstTotalcasa) # -> intervalo de confianza (por defecto 95%)

```

```

##           2.5 %   97.5 %
## Casa 10120970 11942075

```

```

cv(EstTotalcasa)

```

```

##           Casa
## Casa 0.04211351

```



```

# O pasar los resultados a un data frame
df_EstTotalcasa <- as.data.frame(EstTotalcasa)

# Quiero cambiar el nombre de las columnas
colnames(df_EstTotalcasa) <- c("Estimacion", "SE", "deff")

# Calculemos ahora el intervalo de confianza con un 90% de confianza
# (suponemos que el estimador es aprox normal)
df_EstTotalcasa$Li <- df_EstTotalcasa$Estimacion-1.64*df_EstTotalPob$SE
df_EstTotalcasa$Ls <- df_EstTotalcasa$Estimacion+1.64*df_EstTotalPob$SE

# Ahora calculo el CV del estimador
df_EstTotalcasa$CV <- 100*df_EstTotalcasa$SE/df_EstTotalPob$Estimacion

# total rancho y casilla
EstTotalrancho <- survey :: svytotal(~rancho_casilla, diseno, deff=TRUE, cv=TRUE, ci=TRUE)
EstTotalrancho

##              total      SE DEff
## rancho_casilla 578689  82284    1

# Puedo ahora extraer diferentes valores:
survey::cv(EstTotalrancho)      # -> coeficiente de variacion

##              rancho_casilla
## rancho_casilla    0.1421903

deff(EstTotalrancho)      # -> efecto de diseno

## rancho_casilla
##              1

SE(EstTotalrancho)      # -> desvio estandar

##              rancho_casilla
## rancho_casilla    82283.92

confint(EstTotalrancho) # -> intervalor de confianza (por defecto 95%)

##              2.5 % 97.5 %
## rancho_casilla 417415 739962

cv(EstTotalrancho)

##              rancho_casilla
## rancho_casilla    0.1421903

```

```

# O pasar los resultados a un data frame
df_EstTotalrancho <- as.data.frame(EstTotalrancho)

# Quiero cambiar el nombre de las columnas
colnames(df_EstTotalrancho) <- c("Estimacion", "SE", "deff")

# Calculemos ahora el intervalo de confianza con un 90% de confianza
# (suponemos que el estimador es aprox normal)
df_EstTotalrancho$Li <- df_EstTotalrancho$Estimacion-1.64*df_EstTotalPob$SE
df_EstTotalrancho$Ls <- df_EstTotalrancho$Estimacion+1.64*df_EstTotalPob$SE

# Ahora calculo el CV del estimador
df_EstTotalrancho$CV <- 100*df_EstTotalrancho$SE/df_EstTotalrancho$Estimacion

# proporcion
Estproporcion <- survey::svyratio(~rancho_casilla,~Viv_radio, design = diseno, deff = TRUE, cv = TRUE,
Estproporcion

## Ratio estimator: svyratio.survey.design2(~rancho_casilla, ~Viv_radio, design = diseno,
##      deff = TRUE, cv = TRUE, ci = TRUE)
## Ratios=
##              Viv_radio
## rancho_casilla 0.04160196
## SEs=
##              Viv_radio
## rancho_casilla 0.005687017

# Puedo ahora extraer diferentes valores:
estimador <- coef(Estproporcion) # Estimador de la proporción
error_estandar <- SE(Estproporcion) # Error estándar
cv_ <- cv(Estproporcion) # Coeficiente de variación
intervalo_confianza <- confint(Estproporcion) # Intervalo de confianza

#IC 90%

IC90Li <- estimador-1.64*error_estandar
IC90Ls <- estimador+1.64*error_estandar

# Ahora calculo el CV del estimador
CVe <- 100*error_estandar/estimador
CVe

## rancho_casilla/Viv_radio
##              13.67007

```

## Ejercicio IV

Estimación, encuestando en su totalidad una Muestra Sistemática de  $n=240$  radios censales para Total de población, Total de hogares que habitan en viviendas tipo Casa y Total de hogares que habitan en viviendas rancho/ casilla

### Estimación para Población

*#Creo una nueva tabla para este ejercicio con la base que ya armamos en el Ejercicio 3.*

```
radios_2010E4 <- radios_2010
```

```
N <- nrow(radios_2010E4)
```

*# Tamano de la muestra*

```
n <- 240
```

*# Intervalo de selección*

```
I <- floor(N/n)
```

```
print(I)
```

```
## [1] 218
```

*# Parametro poblacional*

```
ParametroPobE4 <- sum(radios_2010E4$Pob_radio)
```

```
print(ParametroPobE4)
```

```
## [1] 40115211
```

*# Ordenamiento del marco de muestreo*

```
radios_2010E4 <- radios_2010E4[order(radios_2010E4$Codigo),]
```

```
radios_2010E4$aleatorio <- runif(nrow(radios_2010E4),0,1)
```

```
radios_2010E4 <- radios_2010E4[order(radios_2010E4$Viv_radio),]
```

```
radios_2010E4 <- radios_2010E4[order(radios_2010E4$Pob_radio),]
```

```
radios_2010E4 <- radios_2010E4[order(radios_2010E4$aleatorio),]
```

Definiremos el arranque aleatorio, generando un número aleatorio entre 1 y el intervalo de selección I, que en nuestro caso es 218.

```
#aa <- sample(1:I, 1)
```

```
aa = 75
```

```
print(aa)
```

```
## [1] 75
```

El resultado de aa es 75. Le pondremos un # a la función dado que cada vez que se ejecuta cambiará el valor de aa

```
#Seleccionaremos una muestra sistemática de los datos en radios_2010E4 utilizando el arranque aleatorio
```

```
s = radios_2010E4[ seq(aa,N,I), ]
```

Calcularemos una estimación del total poblacional y su error relativo para evaluar la precisión de la muestra sistemática.

```
#Primero realizaremos la estimación del total poblacional usando el método de Horvitz-Thompson, un esti.
```

```
estim <- I*sum(s$Pob_radio)
print(estim)
```

```
## [1] 42690940
```

```
error_rel <- 100*(ParametroPobE4 - estim) /ParametroPobE4
print(error_rel)
```

```
## [1] -6.420829
```

La estimación de la población usando una muestra de 240 radios censales es muy cercana al parámetro poblacional ParametroPobE4 (40115211), lo cual sugiere que la muestra debería ser buena. El error relativo es muy bajo, por lo que la muestra puede proveer una estimación precisa.

### Estimación para viviendas tipo Casa

```
#Parámetro poblacional para viviendas tipo Casa:
```

```
ParametroCasaE4 <- sum(radios_2010E4$Casa)
print(ParametroCasaE4)
```

```
## [1] 10620866
```

```
#Estimación para la muestra:
```

```
estimCasaE4 <- I * sum(s$Casa)
print(estimCasaE4)
```

```
## [1] 11056306
```

```
#Error relativo
```

```
error_relCasaE4 <- 100 * (ParametroCasaE4 - estimCasaE4) / ParametroCasaE4
print(error_relCasaE4)
```

```
## [1] -4.099854
```

El valor de la estimación calculada a partir de la muestra se aproxima bastante al valor real de la población, por lo cual el muestreo es representativo y el tamaño de la muestra se puede interpretar como adecuado. El error relativo es muy bajo, por lo que la estimación es bastante cercana al parámetro poblacional real.

### Estimación para viviendas tipo Rancho o Casilla

```
# Crear la columna Rancho_Casilla
radios_2010E4$Rancho_Casilla <- radios_2010E4$Rancho + radios_2010E4$Casilla

# Seleccionar la muestra sistemática
s = radios_2010E4[seq(aa, N, I), ]

# Parámetro poblacional
ParametroRanchoCasillaE4 <- sum(radios_2010E4$Rancho_Casilla)
print(ParametroRanchoCasillaE4)
```

```
## [1] 461725
```

```
# Estimación para la muestra
estimRanchoCasillaE4 <- I * sum(s$Rancho_Casilla)
print(estimRanchoCasillaE4)
```

```
## [1] 461288
```

```
# Error relativo
error_relRanchoCasillaE4 <- 100 * (ParametroRanchoCasillaE4 - estimRanchoCasillaE4) / ParametroRanchoCa
print(error_relRanchoCasillaE4)
```

```
## [1] 0.09464508
```

La estimación de 426408 es cercano al parametro 461725, lo cual indica que el muestreo sistemático es preciso para estimar el total de viviendas tipo rancho o casilla. El error relativo es muy bajo siendo un buen indicador de precisión en el muestreo

## Compararemos dos estrategias utilizando como estimador la media muestral

### 1. Muestreo sistemático, ordenando la tabla por Provincia-Total de viviendas del radio

Primero indicamos aquí nuevamente los parámetros que obtuvimos en el punto anterior:

```
print(ParametroPobE4)
```

```
## [1] 40115211
```

```
print(ParametroCasaE4)
```

```
## [1] 10620866
```

```
print(ParametroRanchoCasillaE4)
```

```
## [1] 461725
```

Comenzamos con la estrategia 1 que implica Orden por “Provincia-Total de viviendas del radio”

```

# Ordenar por Provincia y Total de viviendas del radio
Estrategia1E4 <- radios_2010E4[order(radios_2010E4$Provincia, radios_2010E4$Viv_radio), ]

# Tamaño de la muestra y cálculo del intervalo
n <- 240
N <- nrow(Estrategia1E4)
I <- floor(N / n)
print(I)

```

```
## [1] 218
```

```

# Arranque aleatorio
# aa <- sample(1:I, 1)
# La función arroja un aa de 191. Le colocamos # dado que sino arrojará un aa diferente en cada ejecución
aa = 191
print(aa)

```

```
## [1] 191
```

```

# Selección sistemática
muestraE1 <- Estrategia1E4[seq(aa, N, by = I), ]

# Estimaciones para la muestra en Estrategia 1
estimPobE1 <- I * mean(muestraE1$Pob_radio)
estimCasaE1 <- I * mean(muestraE1$Casa)
estimRancho_CasillaE1 <- I * mean(muestraE1$Rancho + muestraE1$Casilla)

print(estimPobE1)

```

```
## [1] 169419.6
```

```
print(estimCasaE1)
```

```
## [1] 45506.59
```

```
print(estimRancho_CasillaE1)
```

```
## [1] 1902.05
```

Ahora revisamos el error relativo para Estrategia 1

```

error_relPobE1 <- 100 * abs(ParametroPobE4 - estimPobE1) / ParametroPobE4
error_relCasaE1 <- 100 * abs(ParametroCasaE4 - estimCasaE1) / ParametroCasaE4
error_relRancho_CasillaE1 <- 100 * abs(ParametroRanchoCasillaE4 - estimRancho_CasillaE1) / ParametroRanchoCasillaE4

print(error_relPobE1)

```

```
## [1] 99.57767
```

```
print(error_relCasaE1)
```

```
## [1] 99.57154
```

```
print(error_relRancho_CasillaE1)
```

```
## [1] 99.58806
```

Los errores relativos son muy altos (cercaos al 100%), lo que indica que la estimación está muy alejada del valor poblacional, pudiendo interpretar que el ordenamiento esté afectando la representatividad.

## 2. Muestreo sistemático, ordenando la tabla por un número pseudo aleatorio

```
# Ordenar por el número aleatorio  
Estrategia2E4 <- radios_2010E4 [order(radios_2010E4$aleatorio), ]
```

```
# Selección sistemática  
muestraE2 <- Estrategia2E4 [seq(aa, N, by = I), ]
```

```
# Estimaciones para la muestra en Estrategia 2  
estimPobE2 <- I * mean(muestraE2$Pob_radio)  
estimCasaE2 <- I * mean(muestraE2$Casa)  
estimRancho_CasillaE2 <- I * mean(muestraE2$Rancho + muestraE2$Casilla)
```

```
print(estimPobE2)
```

```
## [1] 173921.3
```

```
print(estimCasaE2)
```

```
## [1] 46018.89
```

```
print(estimRancho_CasillaE2)
```

```
## [1] 2591.475
```

```
error_relPobE2 <- 100 * abs(ParametroPobE4 - estimPobE2) / ParametroPobE4  
error_relCasaE2 <- 100 * abs(ParametroCasaE4 - estimCasaE2) / ParametroCasaE4  
error_relRancho_CasillaE2 <- 100 * abs(ParametroRanchoCasillaE4 - estimRancho_CasillaE2) / ParametroRanchoCasillaE4  
print(error_relPobE2)
```

```
## [1] 99.56645
```

```
print(error_relCasaE2)
```

```
## [1] 99.56671
```

```
print(error_relRancho_CasillaE2)
```

```
## [1] 99.43874
```

Al parecer ninguna de las estrategias estaría dando resultados dado que los errores relativos son muy altos. La estrategia de muestreo sistemático puede que no sea la más adecuada para este caso.

## Hallar CV, deff, sesgo relativo y EMC de cada estrategia, seleccionando todas las muestras posibles

Para evaluar cada estrategia y medir su eficiencia y precisión, podemos calcular el coeficiente de variación (CV), el efecto del diseño (deff), el sesgo relativo y el error medio cuadrático (EMC).

Para ello, calcularemos las I estimaciones posibles, una para cada arranque aleatorio. Definiremos una función que seleccione una muestra sistemática, con el arranque aleatorio como variable independiente y devuelva la estimación.

Primero, definimos las funciones para seleccionar la muestra sistemática y para calcular los estimadores necesarios para todas las posibles muestras.

### Estrategia 1

Total de población

```
# Tamaño de la población y de la muestra
N <- nrow(Estrategia1E4) # total de radios censales
n <- 240 # tamaño de la muestra
I <- floor(N / n) # intervalo de selección

# Valor verdadero del parámetro poblacional
parametro_poblacionalEst1 <- sum(Estrategia1E4$Pob_radio)

# Definimos la función de estimación sistemática
estim_sistemático <- function(aa) {
  s = Estrategia1E4[seq(aa, N, I), ]
  estim <- I * sum(s$Pob_radio)
  return(c(estim))
}

S2 = var(radios_2010E4$Pob_radio)
V_mas = N^2*(1-n/N)*S2/n

estimacionEst1 <- estim_sistemático(2)

print(parametro_poblacionalEst1)
```

```
## [1] 40115211
```

```
print(estimacionEst1)
```

```
## [1] 39672512
```



```

# Calculamos las I estimaciones posibles

lista_arranques <- 1:I

lista_estimaciones <- lapply(lista_arranques, estim_sistematico)

df_estim <- data.frame(matrix(unlist(lista_estimaciones),ncol=1, byrow=TRUE ) )

colnames(df_estim) <- c("Estimacion")

EsperanzaE4 <- mean(df_estim$Estimacion)
SesgoE4 <- EsperanzaE4 - parametro_poblacionalEst1
VarianzaE4 <- var(df_estim$Estimacion)*(N-1)/N
DSE4 <- sqrt(VarianzaE4)
CVestimsE4 <- 100*DSE4/parametro_poblacionalEst1
deff_t1 = VarianzaE4 / V_mas

# EsperanzaE4
# SesgoE4
# VarianzaE4
# DSE4
# CVestimsE4
# deff_t1

pob1 <- data.frame(
  calculo = c("EsperanzaE4", "SesgoE4", "VarianzaE4", "DSE4", "CVestimsE4", "deff_t1"),
  valor = c(EsperanzaE4, SesgoE4, VarianzaE4, DSE4, CVestimsE4, deff_t1)
)

pob1

```

```

##      calculo      valor
## 1 EsperanzaE4 40115211.0000000
## 2      SesgoE4      0.0000000
## 3 VarianzaE4 1070937867919.0131836
## 4         DSE4    1034861.2795535
## 5 CVestimsE4      2.5797229
## 6      deff_t1      0.3840778

```

- La esperanza, es decir, la estimación promedio de la población utilizando todas las muestras posible toma el mismo valor que el parametro poblacional, lo cual sugiere que el estimador es insesgado para esta estrategia.
- El sesgo, que es la diferencia entre la esperanza de la estimación y el valor verdadero, da cero, justo lo que se espera dado que el estimador es insesgado. El valor cero confirma que no hay desviación entre la media de las estimaciones y el valor poblacional.
- La varianza es la medida de la dispersión de las estimaciones.
- La desviación Estándar indica la variabilidad de las estimaciones alrededor de la media.
- Finalmente el Coeficiente de Variación (CV) ayuda a comparar la precisión de la estimación en relación con el valor poblacional real, y el valor que toma (4,05) es bajo CV bajo (menor a 10%) por lo cual indica una buena precisión.

```

# Tamaño de la población y de la muestra
N <- nrow(Estrategia1E4) # total de radios censales
n <- 240                 # tamaño de la muestra
I <- floor(N / n)       # intervalo de selección

# Valor verdadero del parámetro poblacional
parametro_RanchoEst1 <- sum(Estrategia1E4$Rancho_casilla)

# Definimos la función de estimación sistemática
estim_sistematicoR <- function(aa) {
  s = Estrategia1E4[seq(aa, N, I), ]
  estim <- I * sum(s$Rancho_casilla)
  return(c(estim))
}

S2 = var(radios_2010E4$Rancho_casilla)
V_mas = N^2*(1-n/N)*S2/n

estimacionEst1R <- estim_sistematico(2)

print(parametro_RanchoEst1)

```

```
## [1] 461725
```

```
print(estimacionEst1R)
```

```
## [1] 39672512
```

```

# Calculamos las I estimaciones posibles

lista_arranques <- 1:I

lista_estimaciones <- lapply(lista_arranques, estim_sistematicoR)

df_estim <- data.frame(matrix(unlist(lista_estimaciones),ncol=1, byrow=TRUE ))

colnames(df_estim) <- c("Estimacion")

EsperanzaE4 <- mean(df_estim$Estimacion)
SesgoE4 <- EsperanzaE4 - parametro_RanchoEst1
VarianzaE4 <- var(df_estim$Estimacion)*(N-1)/N
DSE4 <- sqrt(VarianzaE4)
CVestimsE4 <- 100*DSE4/parametro_RanchoEst1
deff_r1 = VarianzaE4 / V_mas

# EsperanzaE4
# SesgoE4
# VarianzaE4
# DSE4
# CVestimsE4
# deff_r1

```

```
ranch1 <- data.frame(
  calculo = c("EsperanzaE4", "SesgoE4", "VarianzaE4", "DSE4", "CVestimsE4", "deff_r1"),
  valor = c(EsperanzaE4, SesgoE4, VarianzaE4, DSE4, CVestimsE4, deff_r1)
)
```

```
ranch1
```

```
##      calculo      valor
## 1 EsperanzaE4 461725.0000000
## 2      SesgoE4      0.0000000
## 3 VarianzaE4 3792542292.5885801
## 4          DSE4      61583.6203271
## 5 CVestimsE4      13.3377271
## 6      deff_r1      0.8248942
```

```
# Tamaño de la población y de la muestra
N <- nrow(EstrategialE4) # total de radios censales
n <- 240 # tamaño de la muestra
I <- floor(N / n) # intervalo de selección
```

```
# Valor verdadero del parámetro poblacional
parametro_CasaEst1 <- sum(EstrategialE4$Casa)
```

```
S2c = var(radios_2010E4$Casa)
V_mas_c = N^2*(1-n/N)*S2c/n
```

```
# Definimos la función de estimación sistemática
estim_sistematicoC <- function(aa) {
  s = EstrategialE4[seq(aa, N, I), ]
  estim <- I * sum(s$Casa)
  return(c(estim))
}
```

```
estimacionEst1R <- estim_sistematico(2)
```

```
print(parametro_CasaEst1)
```

```
## [1] 10620866
```

```
print(estimacionEst1R)
```

```
## [1] 39672512
```

```
# Calculamos las I estimaciones posibles
```

```
lista_arranques <- 1:I
```

```
lista_estimaciones <- lapply(lista_arranques, estim_sistematicoC)
```

```
df_estim <- data.frame(matrix(unlist(lista_estimaciones), ncol=1, byrow=TRUE ) )
```

```

colnames(df_estim) <- c("Estimacion")

EsperanzaE4 <- mean(df_estim$Estimacion)
SesgoE4 <- EsperanzaE4 - parametro_CasaEst1
VarianzaE4 <- var(df_estim$Estimacion)*(N-1)/N
DSE4 <- sqrt(VarianzaE4)
CVestimsE4 <- 100*DSE4/parametro_CasaEst1
deff_c1 = VarianzaE4 / V_mas_c

# EsperanzaE4
# SesgoE4
# VarianzaE4
# DSE4
# CVestimsE4
# deff_c1

casa1 <- data.frame(
  calculo = c("EsperanzaE4", "SesgoE4", "VarianzaE4", "DSE4", "CVestimsE4", "deff_c1"),
  valor = c(EsperanzaE4, SesgoE4, VarianzaE4, DSE4, CVestimsE4, deff_c1)
)

casa1

```

```

##      calculo      valor
## 1 EsperanzaE4 10620866.0000000
## 2   SesgoE4      0.0000000
## 3 VarianzaE4 57851369337.9142151
## 4      DSE4    240523.1160157
## 5 CVestimsE4     2.2646281
## 6   deff_c1     0.3010508

```

## Estrategia 2

Total de población

```

# Tamaño de la población y de la muestra
N <- nrow(Estrategia2E4) # total de radios censales
n <- 240                # tamaño de la muestra
I <- floor(N / n)       # intervalo de selección

# Valor verdadero del parámetro poblacional
parametro_poblacionalEst2 <- sum(Estrategia2E4$Pob_radio)

# Definimos la función de estimación sistemática
estim_sistemico2 <- function(aa) {
  s = Estrategia2E4[seq(aa, N, I), ]
  estim2 <- I * sum(s$Pob_radio)
  return(c(estim2))
}

```

```

estimacionEst2 <- estim_sistematico2(2)

print(parametro_poblacionalEst2)

## [1] 40115211

print(estimacionEst2)

## [1] 39552394

S2 = var(Estrategia2E4$Pob_radio)
V_mas = N^2*(1-n/N)*S2/n

# Calculamos las I estimaciones posibles

lista_arranques2 <- 1:I

lista_estimaciones2 <- lapply(lista_arranques2, estim_sistematico2)

df_estim2 <- data.frame(matrix(unlist(lista_estimaciones2),ncol=1, byrow=TRUE ) )

colnames(df_estim2) <- c("Estimacion2")

EsperanzaE4.2 <- mean(df_estim2$Estimacion2)
SesgoE4.2 <- EsperanzaE4.2 - parametro_poblacionalEst2
VarianzaE4.2 <- var(df_estim2$Estimacion2)*(N-1)/N
DSE4.2 <- sqrt(VarianzaE4.2)
CVestimsE4.2 <- 100*DSE4.2/parametro_poblacionalEst2
deff_2 = VarianzaE4.2 / V_mas

# EsperanzaE4.2
# SesgoE4.2
# VarianzaE4.2
# DSE4.2
# CVestimsE4.2
# deff_2

pob2 <- data.frame(
  calculo = c("EsperanzaE4.2", "SesgoE4.2", "VarianzaE4.2", "DSE4.2", "CVestimsE4.2", "deff_2"),
  valor = c(EsperanzaE4.2, SesgoE4.2, VarianzaE4.2, DSE4.2, CVestimsE4.2, deff_2)
)

pob2

##          calculo          valor
## 1 EsperanzaE4.2    40115211.0000000
## 2      SesgoE4.2         0.0000000
## 3 VarianzaE4.2  2453778021947.9614258
## 4        DSE4.2    1566453.9641968
## 5  CVestimsE4.2      3.9048878
## 6        deff_2      0.8800152

```

Rancho - casilla

```
# Tamaño de la población y de la muestra
N <- nrow(Estrategia2E4) # total de radios censales
n <- 240                 # tamaño de la muestra
I <- floor(N / n)       # intervalo de selección

# Valor verdadero del parámetro poblacional
parametro_RanchoEst2 <- sum(Estrategia2E4$Rancho_casilla)

# Definimos la función de estimación sistemática
estim_sistemático2r <- function(aa) {
  s = Estrategia2E4[seq(aa, N, I), ]
  estim2 <- I * sum(s$Rancho_casilla)
  return(c(estim2))
}

S2r = var(radios_2010E4$Rancho_casilla)
V_mas_r = N^2*(1-n/N)*S2r/n

estimacionEst2 <- estim_sistemático2r(2)

print(parametro_RanchoEst2)
```

```
## [1] 461725
```

```
print(estimacionEst2)
```

```
## [1] 370818
```

```
# Calculamos las I estimaciones posibles

lista_arranques2 <- 1:I

lista_estimaciones2 <- lapply(lista_arranques2, estim_sistemático2r)

df_estim2 <- data.frame(matrix(unlist(lista_estimaciones2),ncol=1, byrow=TRUE ) )

colnames(df_estim2) <- c("Estimacion2")

EsperanzaE4.2 <- mean(df_estim2$Estimacion2)
SesgoE4.2 <- EsperanzaE4.2 - parametro_RanchoEst2
VarianzaE4.2 <- var(df_estim2$Estimacion2)*(N-1)/N
DSE4.2 <- sqrt(VarianzaE4.2)
CVestimsE4.2 <- 100*DSE4.2/parametro_RanchoEst2
deff_2r = VarianzaE4.2 / V_mas_r

EsperanzaE4.2
```

```
## [1] 461725
```

```
SesgoE4.2
```

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

```
VarianzaE4.2
```

```
## [1] 5154401262
```

```
DSE4.2
```

```
## [1] 71794.16
```

```
CVestimsE4.2
```

```
## [1] 15.54912
```

```
deff_2r
```

```
## [1] 1.121104
```

```
ranch2 <- data.frame(  
  calculo = c("EsperanzaE4.2", "SesgoE4.2", "VarianzaE4.2", "DSE4.2", "CVestimsE4.2", "deff_2r"),  
  valor = c(EsperanzaE4.2, SesgoE4.2, VarianzaE4.2, DSE4.2, CVestimsE4.2, deff_2r)  
)  
ranch2
```

```
##      calculo      valor  
## 1 EsperanzaE4.2 461725.000000  
## 2      SesgoE4.2      0.000000  
## 3 VarianzaE4.2 5154401262.312323  
## 4      DSE4.2      71794.158971  
## 5 CVestimsE4.2      15.549117  
## 6      deff_2r      1.121104
```

Casa

```
# Tamaño de la población y de la muestra  
N <- nrow(Estrategia2E4) # total de radios censales  
n <- 240 # tamaño de la muestra  
I <- floor(N / n) # intervalo de selección  
  
# Valor verdadero del parámetro poblacional  
parametro_CasaEst2 <- sum(Estrategia2E4$Casa)  
  
# Definimos la función de estimación sistemática  
estim_sistemático2c <- function(aa) {  
  s = Estrategia2E4[seq(aa, N, I), ]  
  estim2 <- I * sum(s$Casa)  
  return(c(estim2))  
}
```

```

S2c = var(radios_2010E4$Casa)
V_mas_c = N^2*(1-n/N)*S2c/n

estimacionEst2 <- estim_sistemico2c(2)

print(parametro_CasaEst2)

## [1] 10620866

print(estimacionEst2)

## [1] 10544878

# Calculamos las I estimaciones posibles

lista_arranques2 <- 1:I

lista_estimaciones2 <- lapply(lista_arranques2, estim_sistemico2c)

df_estim2 <- data.frame(matrix(unlist(lista_estimaciones2),ncol=1, byrow=TRUE ) )

colnames(df_estim2) <- c("Estimacion2")

EsperanzaE4.2 <- mean(df_estim2$Estimacion2)
SesgoE4.2 <- EsperanzaE4.2 - parametro_CasaEst2
VarianzaE4.2 <- var(df_estim2$Estimacion2)*(N-1)/N
DSE4.2 <- sqrt(VarianzaE4.2)
CVestimsE4.2 <- 100*DSE4.2/parametro_CasaEst2
deff_2c = VarianzaE4.2 / V_mas_c

# EsperanzaE4.2
# SesgoE4.2
# VarianzaE4.2
# DSE4.2
# CVestimsE4.2
# deff_2c

casa2 <- data.frame(
  calculo = c("EsperanzaE4.2", "SesgoE4.2", "VarianzaE4.2", "DSE4.2", "CVestimsE4.2", "deff_2c"),
  valor = c(EsperanzaE4.2, SesgoE4.2, VarianzaE4.2, DSE4.2, CVestimsE4.2, deff_2c)
)
casa2

##          calculo          valor
## 1 EsperanzaE4.2    10620866.0000000
## 2      SesgoE4.2         0.0000000
## 3 VarianzaE4.2 167087438689.3264160
## 4         DSE4.2    408763.3039906
## 5  CVestimsE4.2      3.8486815
## 6      deff_2c      0.8695008

```



calculo	Estrategia 1			Estrategia 2		
	total poblacion	casa	rancho	total poblacion	casa	rancho
Esperanza	40.115.211	10.620.866	461.725	40.115.211	10.620.866	461.725
Sesgo	0,00	0,00	0,00	0,00	0,00	0,00
DS	997.448,21	255.208,71	64.954,49	1.779.129,21	455.525,37	67.140,62
CV	2,49	2,40	14,07	4,44	4,29	14,54
deff	0,36	0,34	0,92	1,14	1,08	0,98

## Comparación de estrategias

```
comparativa <- read_xlsx("comparativa.xlsx")

# Mostrar la tabla de comparación
gt(comparativa) %>%
  tab_spanner(columns = ends_with("_1"), label = "Estrategia 1") %>%
  tab_spanner(columns = ends_with("_2"), label = "Estrategia 2") %>%
  cols_label(
    starts_with("pob") ~ "total poblacion",
    starts_with("casa") ~ "casa",
    starts_with("rancho") ~ "rancho"
  ) %>%
  fmt_number(rows = 1, decimals = 0, sep_mark = ".", dec_mark = ",",) %>%
  fmt_number(rows = c(2:5), decimals = 2, sep_mark = ".", dec_mark = ",",)
```

El sesgo es cero en ambas estrategias lo cual es un indicio de que ambas estrategias son estimadores insesgados del total poblacional.

La varianza de la Estrategia 2 es mayor que la de la Estrategia 1, lo que indica que el ordenamiento aleatorio puede introducir más variabilidad en las estimaciones.

El CV de ambas estrategias es igual, lo cual sugiere que las dos estrategias tienen un buen nivel de precisión.

La Estrategia 1, es decir ordenar por Provincia y Total de viviendas, presenta menos variabilidad, por lo cual podría ser la estrategia a elegir si queremos reducir la variabilidad en las estimaciones.

## Ejercicio V (continuación del ejercicio IV)

Probaremos otra estrategia para estimar los parámetros del ejercicio anterior, seleccionando una muestra mediante Madow, con la cantidad de viviendas del radio como variable auxiliar, ordenando la tabla según código de radio (jurisdicción + departamento + fracción + radio)

```
#Creamos la tabla para el ejercicio 5. Eliminamos los radios sin viviendas o poblacion y lo ordenamos p
radios_2010E5 <- radios_2010E4[radios_2010E4$Viv_radio>0,]
radios_2010E5 <- radios_2010E5[!is.na(radios_2010E5$Pob_radio),]

# Hogares rancho- casilla
radios_2010E5$Rancho_casilla <- radios_2010E5$Rancho + radios_2010E5$Casilla
radios_2010E5$Hogares <- radios_2010$HeladeraSi + radios_2010$HeladeraNo
```

```

# N poblacional
N <- nrow(radios_2010E5)

# Tamano de la muestra
n <- 240

# Definimos el vector de probabilidades de seleccion
radios_2010E5$pi_i <- n*radios_2010E5$Viv_radio/sum(radios_2010E5$Viv_radio)
max(radios_2010E5$pi_i)

```

```
## [1] 0.02513294
```

```

# Ordenamos el marco de muestreo por codigo de radio
radios_2010E5 <- radios_2010E5[order(radios_2010E5$Codigo),]

# Controllo que no haya pi_i mayores que 1, nulos ni missing
max(radios_2010E5$pi_i)

```

```
## [1] 0.02513294
```

```
min(radios_2010E5$pi_i)
```

```
## [1] 0.00001740508
```

```
sum( !is.finite(radios_2010E5$pi_i))
```

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

```
summary(radios_2010E5$pi_i)
```

```
##      Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
## 0.00001741 0.00304589 0.00476899 0.00458278 0.00619621 0.02513294
```

Seleccionamos mediante sampling, de la muestra de n=240 radios, mediante Madow, con total de viviendas del radio como variable auxiliar.

```

# Seleccionamos la muestra y definimos el factor de expansion

# Probabilidades de seleccion

pikE5 <- radios_2010E5$pi_i

s= sampling::UPsystematic(pikE5)

muestra_radiosE5 = radios_2010E5[s==1,]
muestra_radiosE5$pondera <- 1/muestra_radiosE5$pi_i

```

```

#Configuramos el diseño de muestra para la estimación en survey

parametro <- sum(radios_2010E5$Rancho_casilla)

DesignE5 <- svydesign(id = ~1, weights = ~pondera, data=muestra_radiosE5)
DesignE5

## Independent Sampling design (with replacement)
## svydesign(id = ~1, weights = ~pondera, data = muestra_radiosE5)

# survey asume ahora muestreo con reposicion

# Estimacion del total de poblacion
EstTotalPobE5 <- svytotal( ~Pob_radio , DesignE5, deff=TRUE, cv=TRUE, ci=TRUE)
EstTotalPobE5

##          total      SE  DEff
## Pob_radio 40758439 1186038 0.3338

100*cv(EstTotalPobE5)

##          Pob_radio
## Pob_radio  2.909921

survey::cv(EstTotalPobE5)

##          Pob_radio
## Pob_radio 0.02909921

df_estim_totalP <- data.frame(EstTotalPobE5)
colnames(df_estim_totalP) <- c("Estimacion", "SE", "deff")

df_estim_totalP$CV <- 100*df_estim_totalP$SE/df_estim_totalP$Estimacion
df_estim_totalP

##          Estimacion      SE  deff      CV
## Pob_radio  40758439 1186038 0.3338293 2.909921

# Estimacion del total de hogares que habitan Casa
EstTotalCasaE5 <- svytotal( ~Casa , DesignE5, deff=TRUE, cv=TRUE, ci=TRUE)
EstTotalCasaE5

##          total      SE  DEff
## Casa 10526596  252718 0.253

100*cv(EstTotalCasaE5)

##          Casa
## Casa 2.400761

```

```
survey::cv(EstTotalCasaE5)
```

```
##          Casa
## Casa 0.02400761
```

```
df_estim_totalC <- data.frame(EstTotalCasaE5)
colnames(df_estim_totalC) <- c("Estimacion", "SE", "deff")

df_estim_totalC$CV <- 100*df_estim_totalC$SE/df_estim_totalC$Estimacion
df_estim_totalC
```

```
##      Estimacion      SE      deff      CV
## Casa    10526596 252718.4 0.2529621 2.400761
```

```
# Estimacion del total de hogares que habitan Rancho - casilla
EstTotalRancho_casilla <- svytotal( ~Rancho_casilla, diseno, deff=TRUE, cv=TRUE, ci=TRUE)
EstTotalRancho_casilla
```

```
##          total      SE DEff
## Rancho_casilla 578689 82284    1
```

```
100*cv(EstTotalRancho_casilla)
```

```
##          Rancho_casilla
## Rancho_casilla      14.21903
```

```
survey::cv(EstTotalRancho_casilla)
```

```
##          Rancho_casilla
## Rancho_casilla      0.1421903
```

```
df_estim_totalR <- data.frame(EstTotalRancho_casilla)
colnames(df_estim_totalR) <- c("Estimacion", "SE", "deff")

df_estim_totalR$CV <- 100*df_estim_totalR$SE/df_estim_totalR$Estimacion
df_estim_totalR
```

```
##      Estimacion      SE deff      CV
## Rancho_casilla  578688.5 82283.92    1 14.21903
```

```
# Estimacion de la proporcion de hogares que habitan rancho casilla
EstRatioRancho_casillaE5 <- svyratio(~Rancho_casilla,~Hogares, DesignE5, deff=TRUE, cv=TRUE, ci=TRUE)
EstRatioRancho_casillaE5
```

```
## Ratio estimator: svyratio.survey.design2(~Rancho_casilla, ~Hogares, DesignE5,
##      deff = TRUE, cv = TRUE, ci = TRUE)
## Ratios=
##          Hogares
## Rancho_casilla 0.03375872
## SEs=
##          Hogares
## Rancho_casilla 0.004767134
```

```
survey::cv(EstRatioRancho_casillaE5)
```

```
##                      Hogares  
## Rancho_casilla 0.1412119
```

```
deff <- deff(EstRatioRancho_casillaE5)  
deff
```

```
## [1] 1.001929
```

```
ICE5 <- confint(EstRatioRancho_casillaE5)  
ICE5
```

```
##                      2.5 %      97.5 %  
## Rancho_casilla/Hogares 0.02441531 0.04310213
```

```
# Repetimos diez veces con un for
```

```
estimo <- function(x){  
  s= sampling::UPsystematic(pikeE5)  
  muestra_radiosE5 = radios_2010E5[s==1,]  
  muestra_radiosE5$pondera <- 1/muestra_radiosE5$pi_i  
  EstRatioRancho_casillaE5 <- svyratio(~Rancho_casilla,~Hogares, DesignE5, deff=TRUE, cv=TRUE, ci=TRUE)  
  deff <- deff(EstRatioRancho_casillaE5)  
  cv <- 100*survey::cv(EstRatioRancho_casillaE5)  
  
  caja <- c(EstRatioRancho_casillaE5, deff, cv)  
  return(caja)  
}
```

```
lista_aaE5 <- 1:10
```

```
lista_estimacionesE5 <- lapply(lista_aaE5, estimo)
```

```
df_estimaciones <- data.frame(matrix(unlist(lista_estimaciones),  
                                     nrow=10, byrow=T))
```

```
## Warning in matrix(unlist(lista_estimaciones), nrow = 10, byrow = T): data  
## length [218] is not a sub-multiple or multiple of the number of rows [10]
```

```
df_estim_ratio <- data.frame(Estimacion = EstRatioRancho_casillaE5[[1]],SE= sqrt(EstRatioRancho_casillaE5$SE),  
df_estim_ratio$CV <- 100*df_estim_totalR$SE/df_estim_totalR$Estimacion  
df_estim_ratio$deff <- deff  
colnames(df_estim_ratio) <- c("Estimacion", "SE", "cv", "deff")  
  
df_estim_ratio
```

```
##                      Estimacion      SE      cv      deff  
## Rancho_casilla 0.03375872 0.004767134 14.21903 1.001929
```

Estimadores	Sistematico 1		Sistematico 2		Madow	
	CV	deff	CV	deff	CV	deff
Total población	2.49	0.36	4.44	1.14	2.45	0.26
Total de hogares que habitan en viviendas tipo Casa	2.40	0.34	4.29	1.08	2.31	0.26
Total de hogares que habitan en viviendas rancho/ casilla	14.07	0.92	14.54	0.98	12.22	1.00

```
# Funcion - Método de Madow con variable auxiliar Viv_radio.

estimo_madow <- function(n){
  radios_2010E5$pi_i <- n*radios_2010E5$Viv_radio/sum(radios_2010E5$Viv_radio)
  pikE5 <- radios_2010E5$pi_i
  s = sampling::UPsystematic(pikE5,eps=1e-6)
  muestra_radiosE5 = radios_2010E5[s==1,]
  muestra_radiosE5$pondera <- 1/muestra_radiosE5$pi_i
  estimacionE5 <- sum(muestra_radiosE5$Pob_radio/muestra_radiosE5$pi_i)
  return(estimacionE5)
}

#Madow
lista <- rep(n,1000)
lista_estim <- lapply(lista, estimo_madow)

df_madow <- data.frame(matrix(unlist(lista_estim), nrow=length(lista_estim), byrow=TRUE))
df_madow$diseño <- "Madow"
colnames(df_madow) <- c("Estimacion", "diseño")
```

## Tabla comparativa de resultados

```
tabla <- openxlsx::read.xlsx("tabla_eje4.xlsx")
tabla[, 2:7] <- round(tabla[, 2:7], 2)

gt(tabla) %>%
  tab_spanner(columns = starts_with("sistematico1"), label = "Sistematico 1") %>%
  tab_spanner(columns = starts_with("sistematico2"), label = "Sistematico 2") %>%
  tab_spanner(columns = starts_with("madow"), label = "Madow") %>%
  cols_label(
    ends_with("cv") ~ "CV",
    ends_with("deff") ~ "deff"
  )
```

## Ejercicio VI

### MAS

A partir de la encuesta realizada, el 53% de los electores (212 / 400) planean votar por el candidato X y la utilidad de esta información depende de los supuestos de MAS:

Independencia: Los electores son independientes entre sí, el voto de una persona no afecta el voto de otra.

Aleatoriedad: Los electores son seleccionadas de manera aleatoria en el universo, con la misma probabilidad de ser seleccionados.

Tamaño de la muestra: En este caso, desconocemos el N, pero 400 electores puede ser un buen numero.

De acuerdo a esto, la utilidad de informacion que el estadístico puede brindar depende de si la muestra es representativa de la población. Dado que es una muestra aleatoria simple, y cada individuo tiene la misma probabilidad de ser seleccionado, como hemos mencionado, la muestra sería representativa.

## Coeficiente de Variación

El coeficiente de variación (CV) es una medida que indica la precisión relativa de un estimador. Es la relación entre la desviación estándar del estimador y su valor esperado, multiplicada por 100.

```
n <- 400
p_hat <- 212 / 400

# Desviación estándar de la proporción
sigma_hat <- sqrt((p_hat * (1 - p_hat)) / n)

# Coeficiente de variación (CV)
CVE6 <- (sigma_hat / p_hat) * 100

# Resultados
CVE6
```

```
## [1] 4.708483
```

Dado el tamaño de la muestra y la variabilidad de las respuestas, la estimación de la proporción de votantes de X tiene una variabilidad relativa del 4.71%. Un CV menor al 10% indica que el estimador es bastante preciso y que la proporción muestral obtenida es un buen estimador de la proporción verdadera en la población.

## Conclusión

El estadístico entonces informa que el 53% de los electores en la muestra votarán por X. Esta información es útil, ya que la muestra aleatoria simple es representativa, bajo los supuestos de independencia y aleatoriedad.

El CV del estimador indica que proporción estimada de 53% tiene una variabilidad del 4.71%, indicando que es bastante precisa.

## Ejercicio VII

El candidato quiere saber el tamaño de muestra necesario para estimar el porcentaje de votos con un intervalo de confianza del 95% y una amplitud total de 1%.

Supone que obtendrá un porcentaje cercano al 50%.

La formula para obtener el tamaño de muestra para un intervalo es

$$n = \frac{Z^2 \cdot p \cdot (1 - p)}{E^2}$$

```

Z <- 1.96 # Valor crítico para un intervalo de confianza del 95%
p <- 0.5  # Proporción esperada de votos
E <- 0.01 # Amplitud del intervalo de confianza (1%)

# Calculo
n <- (Z^2 * p * (1 - p)) / E^2

print(n)

```

```
## [1] 9604
```

Para garantizar que el intervalo de confianza tenga una amplitud total de 1% con un nivel de confianza del 95%, se necesitaría una muestra de 9,604 electores.

## Ejercicio VIII

```

# Definir los parámetros
n <- 24 # tamaño de la muestra
x <- 0  # número de hogares con la característica rara
confianza <- 0.90 # Nivel de confianza

# Calcular el intervalo de confianza utilizando el intervalo de Clopper-Pearson
resultado <- binom.confint(x, n, conf.level = confianza, methods = "exact")

# Mostrar el intervalo de confianza
resultado

```

```
## method x n mean lower upper
## 1 exact 0 24 0 0 0.1173462
```

Cuando se calcula un intervalo de confianza para la proporción utilizando el metodo de Clopper Pearson y en ninguna de las unidades de la muestra se observó cierta característica el intervalo toma la forma:

$$\left(0, 1 - \left(\frac{\alpha}{2}\right)^{\frac{1}{n}}\right)$$

En este caso alfa es igual a 0,1 y n = 24, es decir, que el límite superior es igual a  $1 - (0.1/2)^{(1/24)} = 0.1173462$  como muestra la función en *upper*.