

Tarea # 1
Teoría de la probabilidad y estadística
Fecha de entrega: 13 de diciembre de 2021, hasta 11 pm.

El trabajo puede hacerse en grupos de 2 personas. Es permitida solo la discusión de los problemas de la tarea entre los integrantes del grupo; sin embargo, cada estudiante debe montar una copia de la respuesta al classroom y especificar en la solución el compañero con el que trabajó.

Problema 1. Hace cerca de dos años a nivel mundial se desató la pandemia generada por el virus SARS-CoV-2. En Colombia se han llevado registros de los casos que han resultado positivos como portadores del virus, lo cual ha permitido a los científicos desarrollar estrategias que permitan disminuir los efectos negativos de la pandemia. En el archivo 2021-11-20.zip podrá encontrar la base de datos donde se han registrado los casos positivos, reportados por el instituto nacional de salud. Con el fin de que usted ponga en práctica los nuevos conocimientos adquiridos en el curso, se proponen las siguientes actividades:

a. Cree una función en R que permita construir (Dibujar) un histograma con un conjunto cualquiera de datos. Esta función debe permitir al usuario seleccionar que tipo de histograma desea (Frecuencias, Frecuencias relativas o densidades). También debe permitir al usuario indicar cuantas clases desea en el histograma, pero si el usuario no lo indica, la función debe seleccionar la cantidad de clases con el criterio visto en clase.

Para realizar la tarea se importaron los datos del archivo de Excel original, en el que estaban todos los datos, este fue el primer archivo subido por el profesor y decidí hacerlo con este ya que creo que tenía un poco más de complejidad y era más interesante. Se añaden las filas del documento a una lista llamada *datos*.

```
1 datos<-read.csv("C:\\Users\\santi\\Desktop\\Tarea_1\\Salida_Datos_Abiertos.csv")
2 ~ for (i in length(datos)){
3   datos[i]<-c(datos[i])
4 ~ }
```

Debido a que todos los datos del documento no eran numéricos y R solo permite realizar histogramas con valores numéricos, se crea una función llamada ***convertir*** que se encarga de convertir los valores de las columnas a valores numéricos para así tener mayor facilidad a la hora de hacer los histogramas. El argumento de entrada para esta función es la lista con los datos a los que se les desea realizar el cambio por valores numéricos.

La función retorna una lista de listas, la primera lista es la que tiene los valores numéricos y la segunda lista es la que contiene los valores en letras o fechas según sea el caso.

```
6 ▾ convertir <- function(columna) {
7   j=1
8   casos_unicos<-c()
9 ▾   for (i in 1:length(datos[[columna]])){
10 ▾     if ((datos[[columna]][[i]] %in% casos_unicos)==FALSE){
11       casos_unicos[[j]]=datos[[columna]][[i]]
12       j=j+1
13 ▾     }
14 ▾   }
15   lista<-c()
16 ▾   for (k in casos_unicos){
17     pos=which(datos[[columna]]==k)
18 ▾     for (l in pos){
19       lista[l]=pos[1]
20 ▾     }
21 ▾   }
22   resultados<-c()
23   resultados[[1]]=lista
24   resultados[[2]]=casos_unicos
25   return(resultados)
26 ▾ }
```

Luego se utiliza la función **convertir** para cada una de las columnas necesarias, se separan por columnas con datos no numéricos y columnas con datos numéricos.

```
28 #No numericos
29 fecha_hoy_casos=convertir(1)[[1]]
30 Fecha_Not=convertir(3)[[1]]
31 Departamento_nom=convertir(5)[[1]]
32 Ciudad_municipio_nom=convertir(7)[[1]]
33 Sexo=convertir(10)[[1]]
34 Fuente_tipo_contagio=convertir(11)[[1]]
35 Ubicacion=convertir(12)[[1]]
36 Estado=convertir(13)[[1]]
37 Pais_viajo_1_cod=convertir(14)[[1]]
38 Pais_viajo_1_nom=convertir(15)[[1]]
39 Recuperado=convertir(16)[[1]]
40 Fecha_inicio_sintomas=convertir(17)[[1]]
41 Fecha_muerte=convertir(18)[[1]]
42 Fecha_diagnostico=convertir(19)[[1]]
43 Fecha_recuperado=convertir(20)[[1]]
44 Tipo_recuperacion=convertir(21)[[1]]
45 nom_grupo=convertir(23)[[1]]
46
47 #Numericos:
48 Caso=datos[[2]]
49 Departamento=datos[[4]]
50 Ciudad_municipio=datos[[6]]
51 Edad=datos[[8]]
52 unidad_medida=datos[[9]]
53 per_etn_=datos[[22]]
```

Ahora sí se crea la función para graficar cualquiera de las columnas que hay en el documento, la función se llama **histograma** y tiene como argumentos de entrada el conjunto de datos que se desea graficar, cualquiera de las columnas del documento, tipo de histograma que se desea realizar y el número de clases, para este caso si es un argumento vacío (""), la función toma como el numero de clases el valor de 2246 que es el valor aproximado de la raíz cuadrada del total de datos que hay en cada una de las columnas.

```
55 ▾ histograma <- function(conjunto,tipo,clases){
56 ▾   if (clases==""){
57     clases=2246
58 ▾   }
59 ▾   else{
60 ▾   }
61   largo=length(conjunto)
62   conjunto=sort(conjunto)
63   ancho=ceiling((conjunto[[largo]]-conjunto[[1]])/clases)
64   minimo=conjunto[[1]]-1
65   maximo=conjunto[[largo]]+1
66   limites<-c()
67   a=1
68 ▾   while (maximo>minimo){
69     limites[[a]]=minimo+ancho
70     a=a+1
71     minimo=limites[[a-1]]
72 ▾   }
73   frecuencia<-c()
74   b=1
75   conteo=0
76 ▾   for (i in 1:length(conjunto)){
77     lim_sup=limites[[b]]
78 ▾     if (lim_sup>=conjunto[[i]]){
79       conteo=conteo+1
80 ▾     }
81 ▾     else{
82       frecuencia<- c(frecuencia,conteo)
83       b=b+1
84       conteo=0
```

```

55 ▾ histograma <- function(conjunto, tipo, clases){
56 ▾   if (clases==""){
57     clases=2246
58 ▾   }
59 ▾   else{
60 ▾   }
61   conjunto=sort(conjunto)
62   ancho=ceiling((conjunto[[5045412]]-conjunto[[1]])/clases)
63   minimo=conjunto[[1]]-1
64   maximo=conjunto[[5045412]]+1
65   limites<-c()
66   a=1
67 ▾   while (maximo>minimo){
68     limites[[a]]=minimo+ancho
69     a=a+1
70     minimo=limites[[a-1]]
71 ▾   }
72   frecuencia<-c()
73   b=1
74   conteo=0
75 ▾   for (i in 1:length(conjunto)){
76     lim_sup=limites[[b]]
77 ▾     if (lim_sup>=conjunto[[i]]){
78       conteo=conteo+1
79 ▾     }
80 ▾     else{
81       frecuencia<- c(frecuencia, conteo)
82       b=b+1
83       conteo=0
84 ▾     }
85 ▾   }
86   frecuencia<- c(frecuencia, conteo)
87   relativa<-c()
88   densidad<-c()
89 ▾   for (i in 1:length(frecuencia)){
90     f=frecuencia[[i]]
91     relativa<- c(relativa, f/length(conjunto))
92     densidad<- c(densidad, f/(length(conjunto)*ancho))
93 ▾   }
94 ▾   if (tipo="Frecuencia"){
95     barplot(frecuencia)
96 ▾   }
97 ▾   else if (tipo="Relativa"){
98     barplot(relativa)
99 ▾   }
100 ▾   else if (tipo="Densidad"){
101     barplot(densidad)
102 ▾   }
103 ▾   else{
104     print("Seleccione entre Frecuencia, Relativa o Densidad")
105 ▾   }
106 ▾ }

```

b. Utilice la función construida por usted mismo, para realizar un histograma de densidades con las edades de las personas contagiadas en Medellín. Interprete los resultados, mostrando los factores relevantes de su resultado.

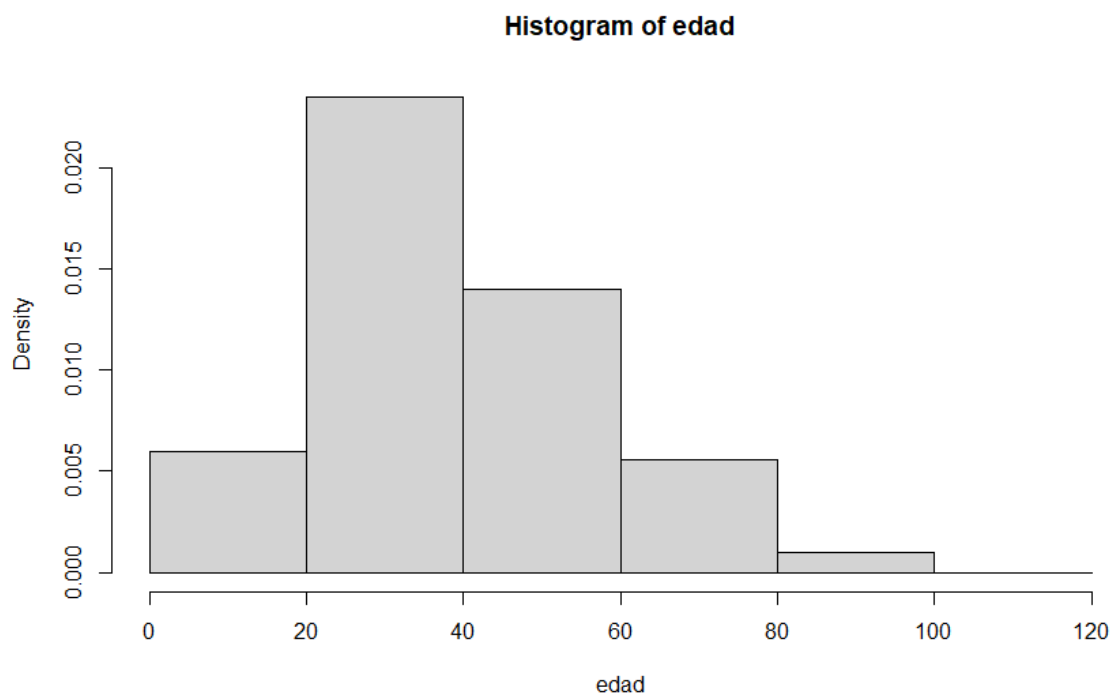
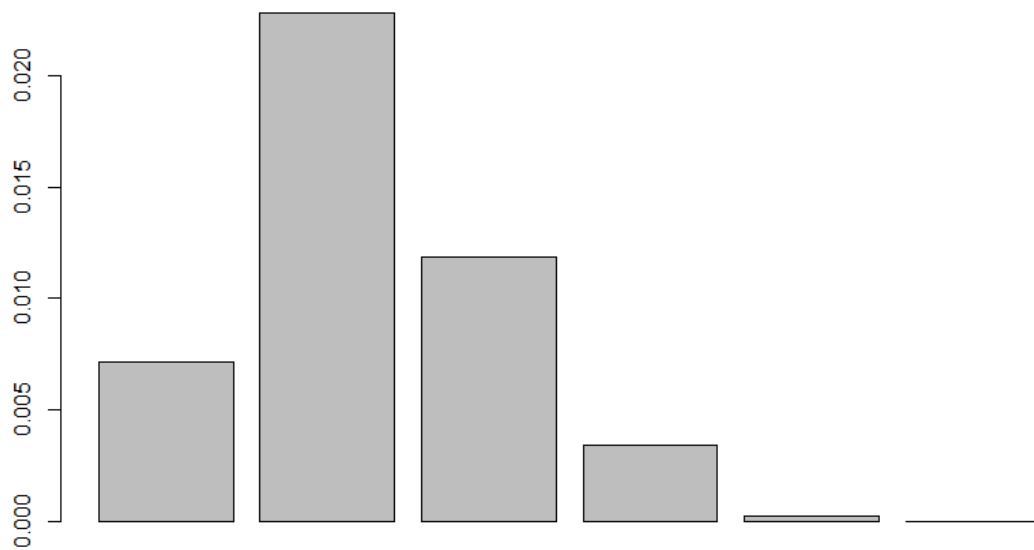
Para este punto se realizó otra función en la cual se filtran los resultados para solo los habitantes de la ciudad de Medellín quienes aparecen en el documento como *MEDELLIN*.

Se creó la función *puntob*, la cual tiene como argumentos de entrada las columnas de ciudad y edad en ese orden, la función nos retorna una lista con las edades de las personas de la ciudad de Medellín.

```
109 > puntob <- function(ciudades,edades){  
110   edad<-c()  
111 >   for (i in 1:length(ciudades)){  
112 >     if (ciudades[[i]]=="MEDELLIN"){  
113       edad<- c(edad,edades[[i]])  
114 >     }  
115 >   }  
116   return(edad)  
117 > }
```

Se grafican ambos histogramas, primero el histograma con la función realizada por el estudiante y el segundo histograma es utilizando la función *hist* de R.

```
119 edad=puntob(datos[[7]],datos[[8]])  
120 clases=5  
121 histograma(edad,"Densidad",clases)  
122 hist(edad,freq=FALSE,breaks=clases)
```



Según los histogramas, (gráficamente no es posible identificar una diferencia notable entre ellos), la mayor cantidad de personas tiene o tenía entre los 20 y los 40 años, siendo muy evidente en la gráfica como sobresale frente a los demás, era de esperarse ya que la mayoría

de los trabajadores se encuentran en esta edad, entre los 100 y los 120 años no se tiene mucha cantidad puesto que llegar a esta edad es difícil.

Podría decirse que ambos gráficos son casi iguales dados sus valores en el eje y correspondientes a la densidad y a la altura de las barras del gráfico.

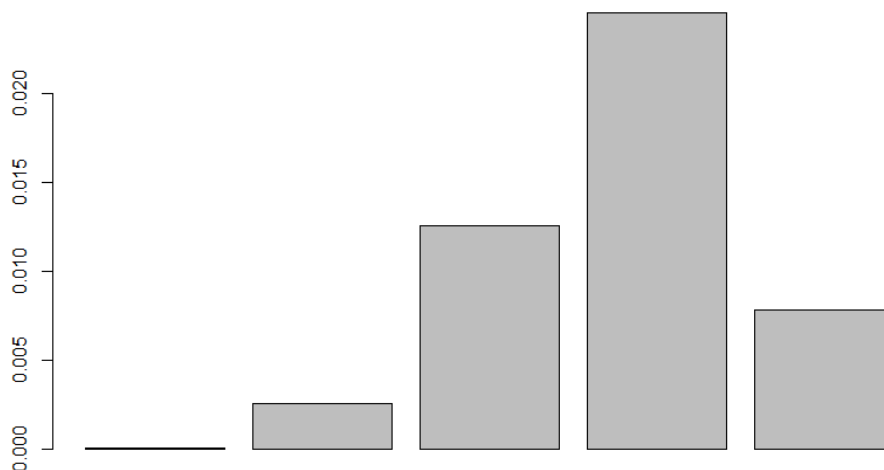
c. Utilice la función construida por usted mismo, para realizar un histograma de densidades con las edades de las personas fallecidas en Medellín. Interprete los resultados y compare el nuevo histograma con el del numeral b.

Se crea una función de manera similar a la anterior, en donde se filtran por la ciudad de Medellín y ahora los fallecidos, obteniendo así la edad de los fallecidos de la ciudad de Medellín.

```
126 ▶ puntoc <- function(ciudades,edades,estado){
127   fallecidos<-c()
128 ▶   for (i in 1:length(ciudades)){
129 ▶     if (ciudades[[i]]=="MEDELLIN"&(estado[[i]]=="Fallecido")){
130       fallecidos<- c(fallecidos,edades[[i]])
131 ▶     }
132 ▶   }
133   return(fallecidos)
134 ▶ }
```

Se reutiliza la función histograma con el fin de obtener el histograma de densidades, el numero de clases es el mismo que se utilizó en el punto anterior, *clases=5*.

```
135 fallecido=puntoc(datos[[7]],datos[[8]],datos[[13]])
136 histograma(fallecido,"Densidad",clases)
```

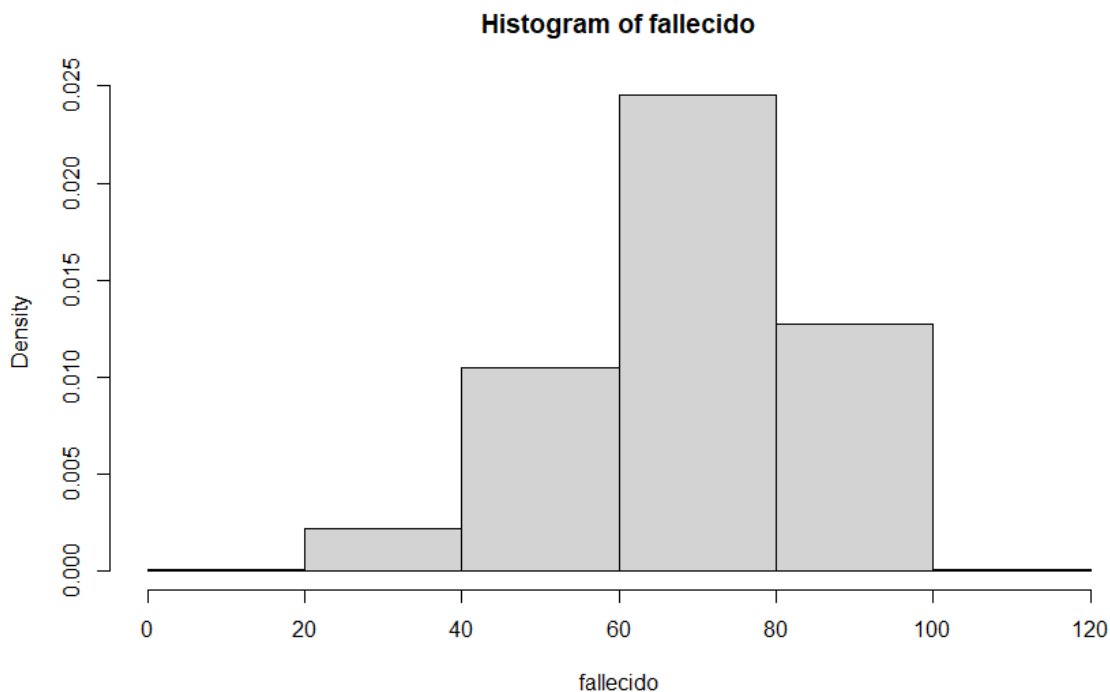


En este histograma podemos observar que aunque los más jóvenes fueron los más afectados en cuanto a contagios, no fue en igual proporción para las muertes, la peor parte la llevaron las personas entre los 60 y 80 años, puesto que como se mencionó en el punto anterior, es muy difícil tener alrededor de los 100 años y por eso no se ven tantas muertes en ese sentido, es claro que los menos afectados en cuanto a muertes fueron las personas entre 0 y 20 años, las personas entre 40 y 60 también sufrieron grandes pérdidas.

d. Realice el histograma de densidades con las edades de las personas fallecidas en Medellín, usando la función hist de R. Compare sus resultados con los de dicha función y explique las diferencias (si las hay).

Por último se tiene el histograma de densidades para las edades de los fallecidos en la ciudad de Medellín, debido a que ya se tenían filtrados los datos se pudo realizar este punto con una sola línea de código.

```
137 hist(fallecido, freq=FALSE, breaks=clases)
```



La diferencia que se tiene respecto al histograma anterior es que el anterior no muestra las muertes entre las personas de 100 y 120 años, esto es debido a que es un numero insignificante para el programa, habría que hacer un análisis mas profundo de los datos para obtener el mismo resultado. Sin embargo la gráfica a simple vista luce igual y los valores están de ambos ejes y están similares.