

# Lectura de datos

Maximiliano Cacace

24/5/2021

## Leer los datos

```
library(readxl)
Becarios_2020 <- read_excel("C:/Users/hp/Desktop/TODO/UNSAM/Doctorado Ciencias Humanas/Métodos cuantitativos y análisis de grandes datos/BIG_DATA_CACACE/Becarios 2020.csv")
#View(Becarios_2020)

Becarios_2021 <- read_excel("C:/Users/hp/Desktop/TODO/UNSAM/Doctorado Ciencias Humanas/Métodos cuantitativos y análisis de grandes datos/BIG_DATA_CACACE/Becarios 2021.csv")
#View(Becarios_2021)

#“{r CSV} library(readr) Becarios_2020 <- read_delim("C:/Users/hp/Desktop/TODO/UNSAM/Doctorado Ciencias Humanas/Métodos cuantitativos y análisis de grandes datos/BIG_DATA_CACACE/Becarios 2020.csv", ";", escape_double = FALSE, trim_ws = TRUE) View(Becarios_2020)

Becarios_2021 <- read_delim("C:/Users/hp/Desktop/TODO/UNSAM/Doctorado Ciencias Humanas/Métodos cuantitativos y análisis de grandes datos/BIG_DATA_CACACE/Becarios 2021.csv", ";", escape_double = FALSE, trim_ws = TRUE) View(Becarios_2021) #“

str(Becarios_2020)

## tibble [165 x 30] (S3: tbl_df/tbl/data.frame)
##  $ Edad                                : num [1:165] 22 37 22 22 22
##  $ Nacionalidad                       : chr [1:165] "Argentina" "Argentina" "Argentina" "Argentina" "Argentina"
##  $ DNI                                 : num [1:165] 40738524 29740000 40738524 29740000 40738524
##  $ Carrera_que_cursa                  : chr [1:165] "Ingeniería de Sistemas" "Ingeniería de Sistemas" "Ingeniería de Sistemas" "Ingeniería de Sistemas" "Ingeniería de Sistemas"
##  $ Tiempo_de_beca                     : chr [1:165] "Este es el tiempo de beca" "Este es el tiempo de beca" "Este es el tiempo de beca" "Este es el tiempo de beca" "Este es el tiempo de beca"
##  $ Costeo_de_estudios                 : chr [1:165] "Aportes familiares" "Aportes familiares" "Aportes familiares" "Aportes familiares" "Aportes familiares"
##  $ Con_quien_vive                     : chr [1:165] "Solo" "Familiar" "Solo" "Familiar" "Solo"
##  $ Por_que_eligio_la_UNSAM             : chr [1:165] "Por el programa de beca" "Por el programa de beca" "Por el programa de beca" "Por el programa de beca" "Por el programa de beca"
##  $ Educacion Primaria                 : chr [1:165] "Estatad" "Estatad" "Estatad" "Estatad" "Estatad"
##  $ Educacion Secundaria               : chr [1:165] "Estatad" "Estatad" "Estatad" "Estatad" "Estatad"
##  $ Relacion_con_la_carrera             : chr [1:165] "SI" "NO" "SI" "NO" "SI"
##  $ Repitencia_o_abandono_temporal_de_estudios : chr [1:165] "No" "No" "No" "No" "No"
##  $ Tiempo_transcurrido_desde_el_secundario_hasta_ingresar_a_la_universidad : chr [1:165] "Ingresé al secundario" "Ingresé al secundario" "Ingresé al secundario" "Ingresé al secundario" "Ingresé al secundario"
##  $ Que_hizo_en_ese_tiempo             : chr [1:165] "Terminé el secundario" "Terminé el secundario" "Terminé el secundario" "Terminé el secundario" "Terminé el secundario"
##  $ Trabajo                           : chr [1:165] "No, por decisión propia" "No, por decisión propia" "No, por decisión propia" "No, por decisión propia" "No, por decisión propia"
##  $ Que_actividad_desempeña            : chr [1:165] NA "Docente particular" NA "Docente particular" NA "Docente particular"
##  $ Horas_semanales_de_trabajo          : chr [1:165] NA "Entre 10 y 20 horas" NA "Entre 10 y 20 horas" NA "Entre 10 y 20 horas"
##  $ Aprobacion_del_CPU                 : chr [1:165] "Si" "Si" "Si" "Si" "Si"
##  $ Año_de_inicio_del_CPU              : num [1:165] 2017 2016 2017 2016 2017
##  $ Cuatrimestre_de_inicio_CPU         : chr [1:165] "1er Cuatrimestre" "1er Cuatrimestre" "1er Cuatrimestre" "1er Cuatrimestre" "1er Cuatrimestre"
##  $ Alguna_dificultad_en_el_CPU_describir_en_caso_afirmativo : chr [1:165] "No" "No" "No" "No" "No"
##  $ Cuantos_cuatrimestres_le_llevo_aprobar_el_CPU : num [1:165] 1 1 1 1 NA 1
##  $ Año_de_inicio_de_carrera           : num [1:165] 2017 2016 2017 2016 2017
##  $ Cuatrimestre_de_inicio_de_carrera : chr [1:165] "2do Cuatrimestre" "2do Cuatrimestre" "2do Cuatrimestre" "2do Cuatrimestre" "2do Cuatrimestre"
```

```
## $ Cantidad_de_materias_que cursa : num [1:165] 4 4 3 4 NA 0
## $ Finales_pendientes : num [1:165] 3 1 0 0 NA 1
## $ Seguimiento_del_plan_de_estudios : chr [1:165] "Parcialment
## $ Dificultades_observadas_en_el_seguimiento_del_plan_de_estudios : chr [1:165] "Acumulación
## $ Horas_semanales_destinadas_a_estudiar : chr [1:165] "Más de 20 h
## $ Porque_decidio_estudiar_en_la_Universidad : chr [1:165] "Vocación"
```

## Cambiar factor y nombre de respuesta de variable

```
Becarios_2020$Tiempo_de_beca <-
  factor(as.character(Becarios_2020$Tiempo_de_beca),
        labels = c("Este es el primer año" = "0.5",
                    "1 año" = "1",
                    "2 años" = "2",
                    "3 años" = "3",
                    "4 años" = "4")
  )
```

```
Becarios_2020$Tiempo_de_beca <-
  as.numeric(Becarios_2020$Tiempo_de_beca)
```

```
Becarios_2020$Tiempo_de_beca %>%
  class()
```

```
## [1] "numeric"
```

```
table(Becarios_2020$Tiempo_de_beca)
```

```
##
## 1 2 3 4 5
## 31 17 14 6 97
```

```
Becarios_2020$Cuatrimestre_de_inicio_CPU <-
  factor(as.character(Becarios_2020$Cuatrimestre_de_inicio_CPU),
        labels = c("1er Cuatrimestre" = "1",
                    "2do Cuatrimestre" = "2")
  )
```

```
Becarios_2020$Cuatrimestre_de_inicio_CPU <-
  factor(as.numeric(Becarios_2020$Cuatrimestre_de_inicio_CPU),
        labels = c("1" = "1",
                    "2" = "2")
  )
```

```
Becarios_2020$Cuatrimestre_de_inicio_CPU %>%
  class()
```

```
## [1] "factor"
```

```
table(Becarios_2020$Cuatrimestre_de_inicio_CPU)
```

```
##
## 1 2
## 132 33
```

```
Becarios_2020$Cuatrimestre_de_inicio_CPU %>%
  class()
```

```

## [1] "factor"
table(Becarios_2020$Cuatrimestre_de_inicio_CPU)

##
## 1 2
## 132 33
nrow(Becarios_2020)

## [1] 165
ncol(Becarios_2020)

## [1] 30
nrow(Becarios_2021)

## [1] 151
ncol(Becarios_2021)

## [1] 30
#Ratio = 8
#Nominal = 22
Becarios_2020$Edad %>% #Ratio
  class()

## [1] "numeric"
Becarios_2020$Nacionalidad %>% #Nominal
  class()

## [1] "character"
Becarios_2020$DNI %>% #Nominal
  class()

## [1] "numeric"
Becarios_2020$Carrera_que_cursa %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Tiempo_de_beca %>% #Ratio
  class()

## [1] "numeric"
Becarios_2020$Costeo_de_estudios %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Con_quien_vive %>% #Nominal
  class()

## [1] "character"
Becarios_2020$`Por_que_eligio_la_UNSAM` %>% #Nominal
  class()

```

```

## [1] "character"
Becarios_2020$Educacion Primaria %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Educacion_Secundaria %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Relacion_con_la_carrera %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Repitencia_o_abandono_temporal_de_estudios %>% #Nominal
  class()

## [1] "character"
Becarios_2020$`Tiempo_transcurrido_desde_el_secundario_hasta_ingresar_a_la_universidad` %>% #Ratio
  class()

## [1] "character"
Becarios_2020$Que_hizo_en_ese_tiempo %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Trabajo %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Que_actividad_desempeña %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Horas_semanales_de_trabajo %>% #Ratio
  class()

## [1] "character"
Becarios_2020$Aprobacion_del_CPU %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Año_de_inicio_del_CPU %>% #Nominal
  class()

## [1] "numeric"
Becarios_2020$Cuatrimestre_de_inicio_CPU %>% #Nominal
  class()

## [1] "factor"
Becarios_2020$Alguna_dificultad_en_el_CPU_describir_en_caso_afirmativo %>% #Nominal
  class()

```

```

## [1] "character"
Becarios_2020$Cuantos_cuatrimestres_le_llevo_aprobar_el_CPU %>% #Ratio
  class()

## [1] "numeric"
Becarios_2020$Año_de_inicio_de_carrera %>% #Nominal
  class()

## [1] "numeric"
Becarios_2020$Cuatrimestre_de_inicio_de_carrera %>% #Nominal
  class()

## [1] "character"
Becarios_2020$`Cantidad_de_materias_que_cursa` %>% #Ratio
  class()

## [1] "numeric"
Becarios_2020$Finales_pendientes %>% #Ratio
  class()

## [1] "numeric"
Becarios_2020$Seguimiento_del_plan_de_estudios %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Dificultades_observadas_en_el_seguimiento_del_plan_de_estudios %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Horas_semanales_destinadas_a_estudiar %>% #Ratio
  class()

## [1] "character"
Becarios_2020$Porque_decidio_estudiar_en_la_Universidad %>% #Nominal
  class()

## [1] "character"
Becarios_2020$Edad

## [1] 22 37 22 22 26 25 25 21 26 22 20 26 20 36 18 20 17 19 26 34 21 22 19 25 25
## [26] 18 22 59 32 25 22 20 26 21 29 22 20 27 33 23 22 19 22 45 28 32 33 25 60 30
## [51] 27 23 24 24 18 18 24 26 48 27 25 19 20 25 25 29 22 25 18 26 29 32 22 28 21
## [76] 23 22 22 19 22 20 20 23 24 42 24 18 25 29 25 19 33 22 43 27 23 26 35 24 20
## [101] 32 20 22 21 20 20 21 25 23 22 22 19 30 20 17 25 37 19 20 20 21 23 37 18 28
## [126] 31 23 20 19 19 20 32 23 20 26 22 25 27 26 21 27 21 25 23 23 35 24 20 20 22
## [151] 20 33 34 36 19 19 25 23 26 20 25 18 19 19 24

Becarios_2021$Edad

## [1] 26 23 22 22 26 55 21 27 18 21 37 21 20 44 27 24 24 23 29 26 36 23 21 48 28
## [26] 23 28 23 17 40 19 42 46 29 33 34 26 19 30 36 31 20 24 25 25 20 49 28 18 30
## [51] 23 20 26 30 26 23 19 21 27 30 32 23 25 38 22 23 20 26 26 21 20 21 25 27 43
## [76] 26 20 22 26 34 45 24 27 33 36 26 25 21 21 23 21 22 26 42 24 19 23 35 31 22

```

```
## [101] 23 25 38 20 23 28 21 22 25 38 19 25 22 34 23 19 24 51 21 26 41 23 33 22 27
## [126] 39 23 18 28 22 20 22 26 23 36 21 25 25 21 23 20 23 19 29 34 20 21 23 22 24
## [151] 29
```

```
table(Becarios_2020$Edad)
```

```
##
## 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 42 43 45 48 59
## 2 8 14 22 9 21 12 8 18 11 6 3 4 2 1 5 4 2 2 2 3 1 1 1 1 1
## 60
## 1
```

```
table(Becarios_2021$Edad)
```

```
##
## 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42
## 1 3 7 11 15 12 20 7 10 14 6 5 4 4 2 1 3 4 1 4 1 3 1 1 1 2
## 43 44 45 46 48 49 51 55
## 1 1 1 1 1 1 1 1
```

```
table(Becarios_2020$Carrera_que_cursa)
```

```
##
##          Ingeniería Ambiental          Ingeniería Biomédica
##                20                12
##          Ingeniería Electrónica          Ingeniería en Energía
##                10                10
##          Ingeniería en Telecomunicaciones          Ingeniería Espacial
##                4                3
##          Ingeniería Industrial          Licenciatura en Biotecnología
##                8                40
##          Licenciatura en Física Médica  Tecnicatura en Diagnóstico por Imágenes
##                4                38
##          Tecnicatura en Electromedicina  Tecnicatura en Programación Informática
##                2                12
##          Tecnicatura en Redes Informáticas
##                2
```

```
Becarios_2020 %>%
  group_by(Costeo_de_estudios) %>%
  summarize(
    mean(Edad),
    median(Edad)
  )
```

```
## # A tibble: 5 x 3
##   Costeo_de_estudios   `mean(Edad)` `median(Edad)`
##   <chr>              <dbl>         <dbl>
## 1 Aportes familiares    22.4          22
## 2 Beca                  23.8          22
## 3 Con dificultad para costear sus estudios  26.9          25.5
## 4 Plan social           60           60
## 5 Trabajo              27.4          24.5
```

```
mean(Becarios_2020$Edad) -> promedio_Edad2020
promedio_Edad2020
```

```
## [1] 24.86061
```



```
## [127]  51  60  95 138 141  45  74 125  35  66  71  89  50 113 126  29  46  72
## [145] 101 132  39  47  92 152  20 153  98 146  14 154   2 117 123  85  94  44
## [163]  59  28  49
```

## Tests

```
shapiro.test(Becarios_2020$Edad)
```

```
##
##  Shapiro-Wilk normality test
##
## data:  Becarios_2020$Edad
## W = 0.78464, p-value = 2.625e-14
```

```
t.test(Becarios_2020$Edad, Becarios_2021$Edad)
```

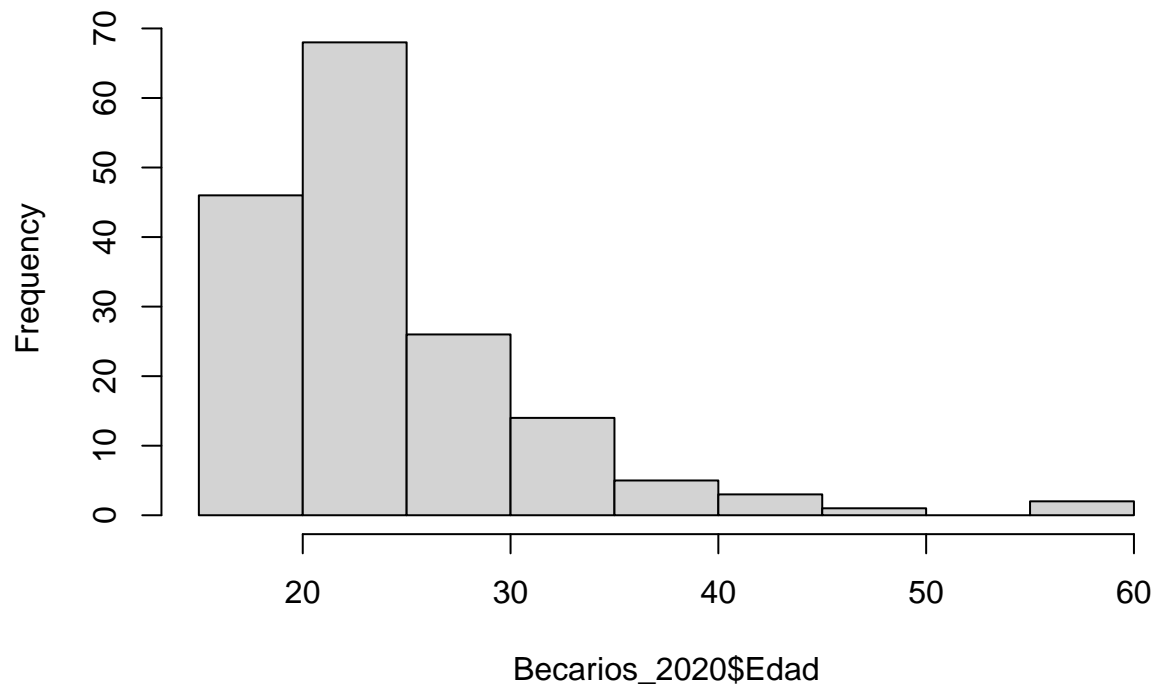
```
##
##  Welch Two Sample t-test
##
## data:  Becarios_2020$Edad and Becarios_2021$Edad
## t = -2.243, df = 304.47, p-value = 0.02562
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -3.3946927 -0.2218436
## sample estimates:
## mean of x mean of y
##  24.86061  26.66887
```

```
##Visualización de datos
```

```
hist(Becarios_2020$Edad)
```

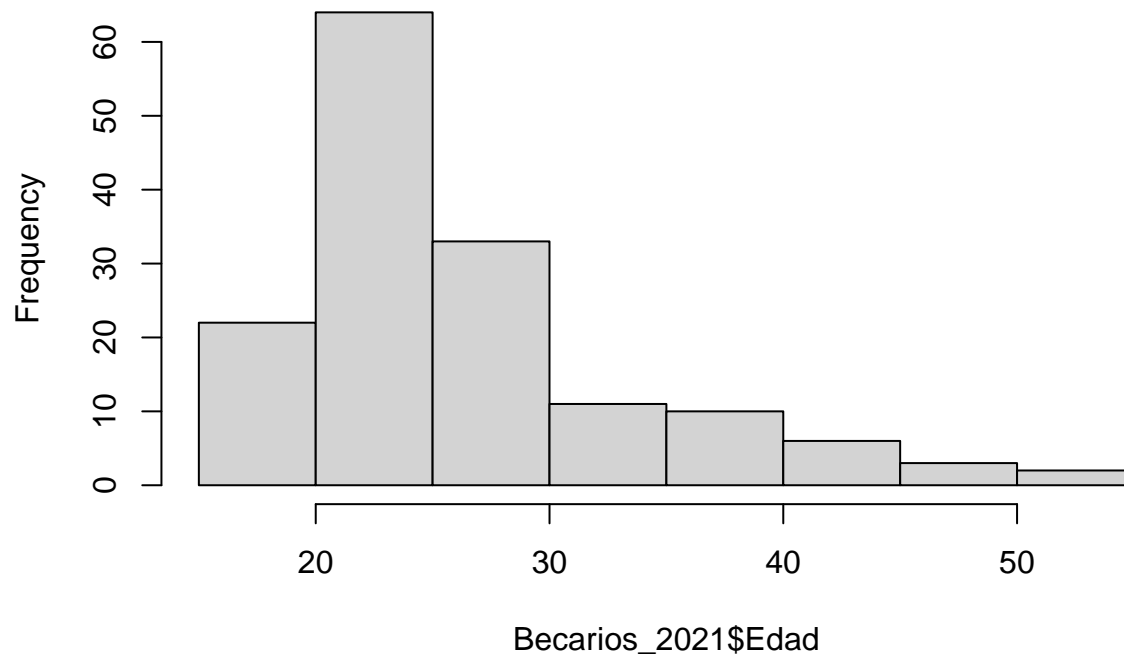


**Histogram of Becarios\_2020\$Edad**



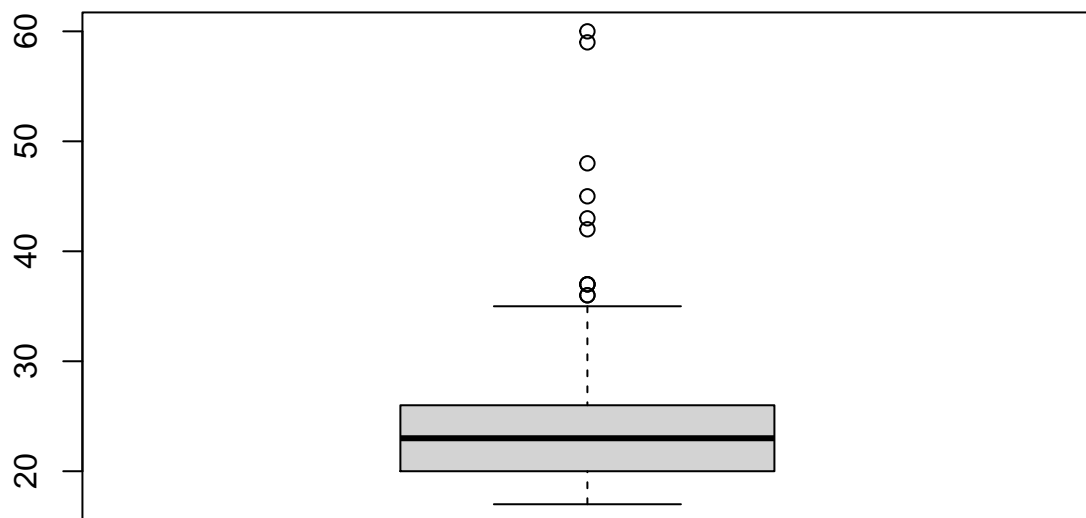
```
hist(Becarios_2021$Edad)
```

**Histogram of Becarios\_2021\$Edad**

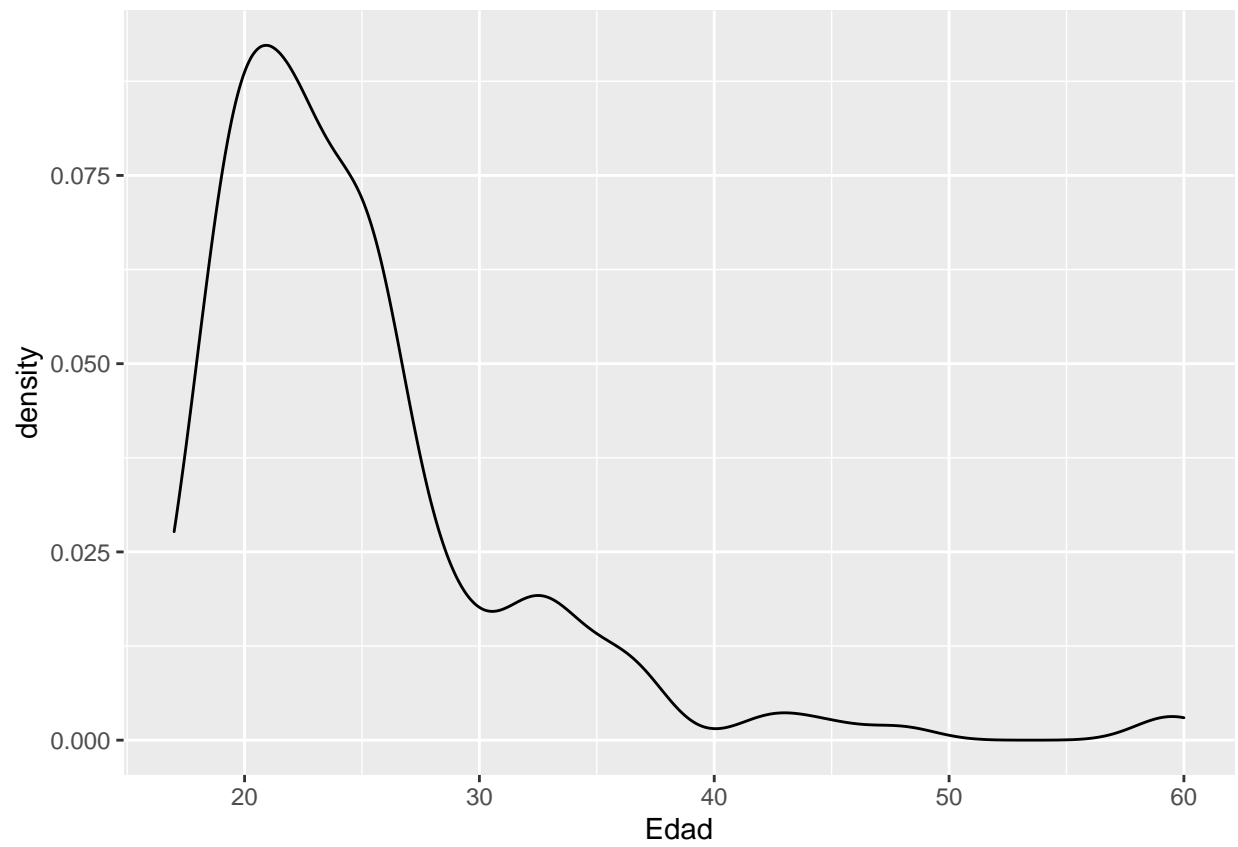


```
#polygon(Becarios_2020$Edad)
```

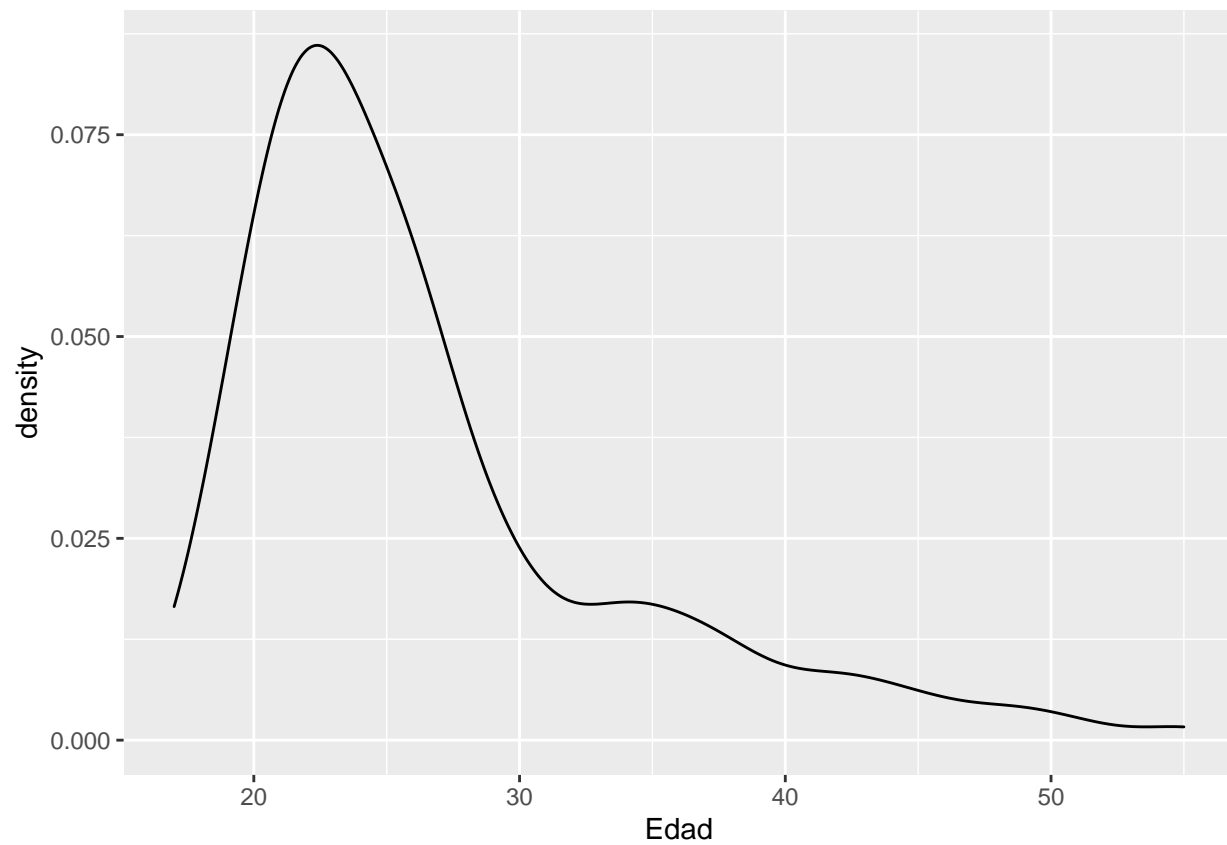
```
boxplot(Becarios_2020$Edad) #visualización de la dispersión
```



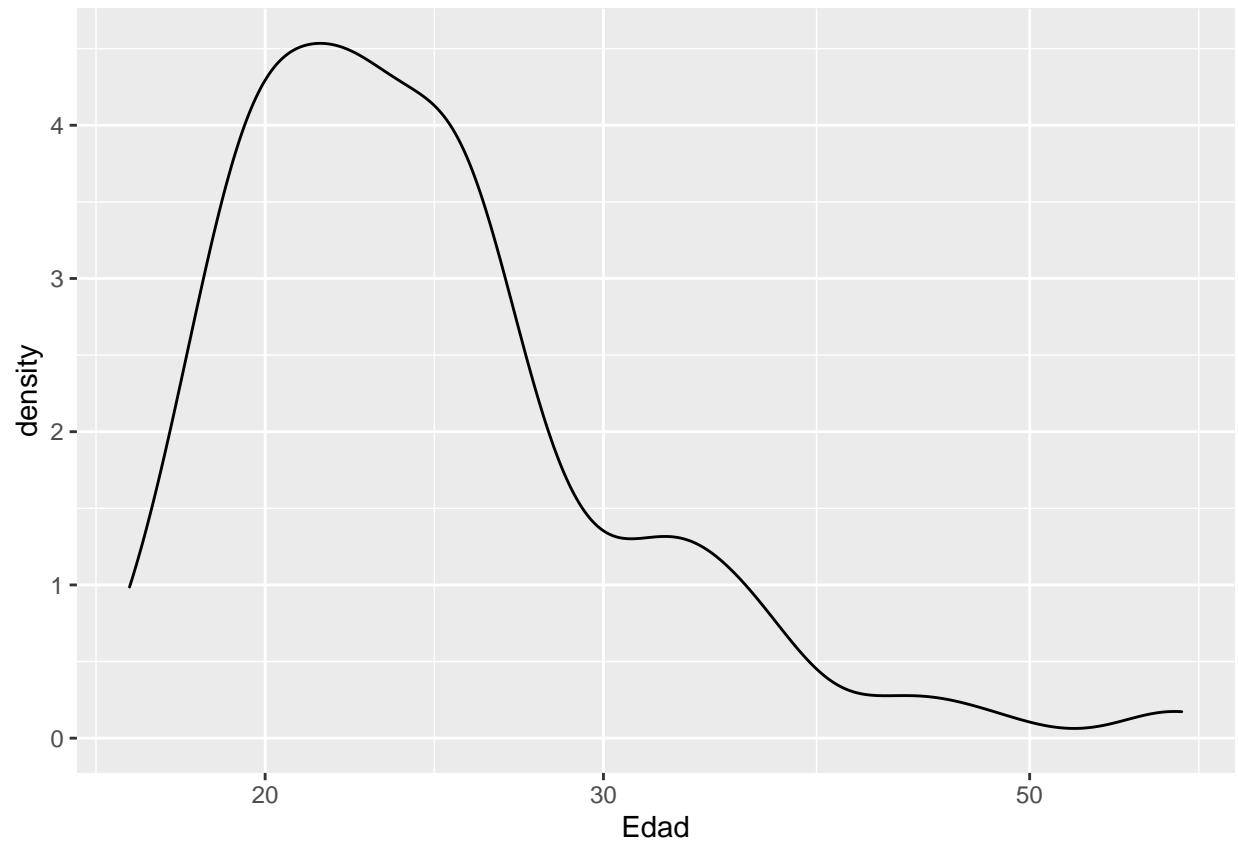
```
library(ggplot2)
Becarios_2020 %>%
  ggplot(aes(Edad))+
  geom_density()
```



```
Becarios_2021 %>%  
  ggplot(aes(Edad))+  
  geom_density()
```

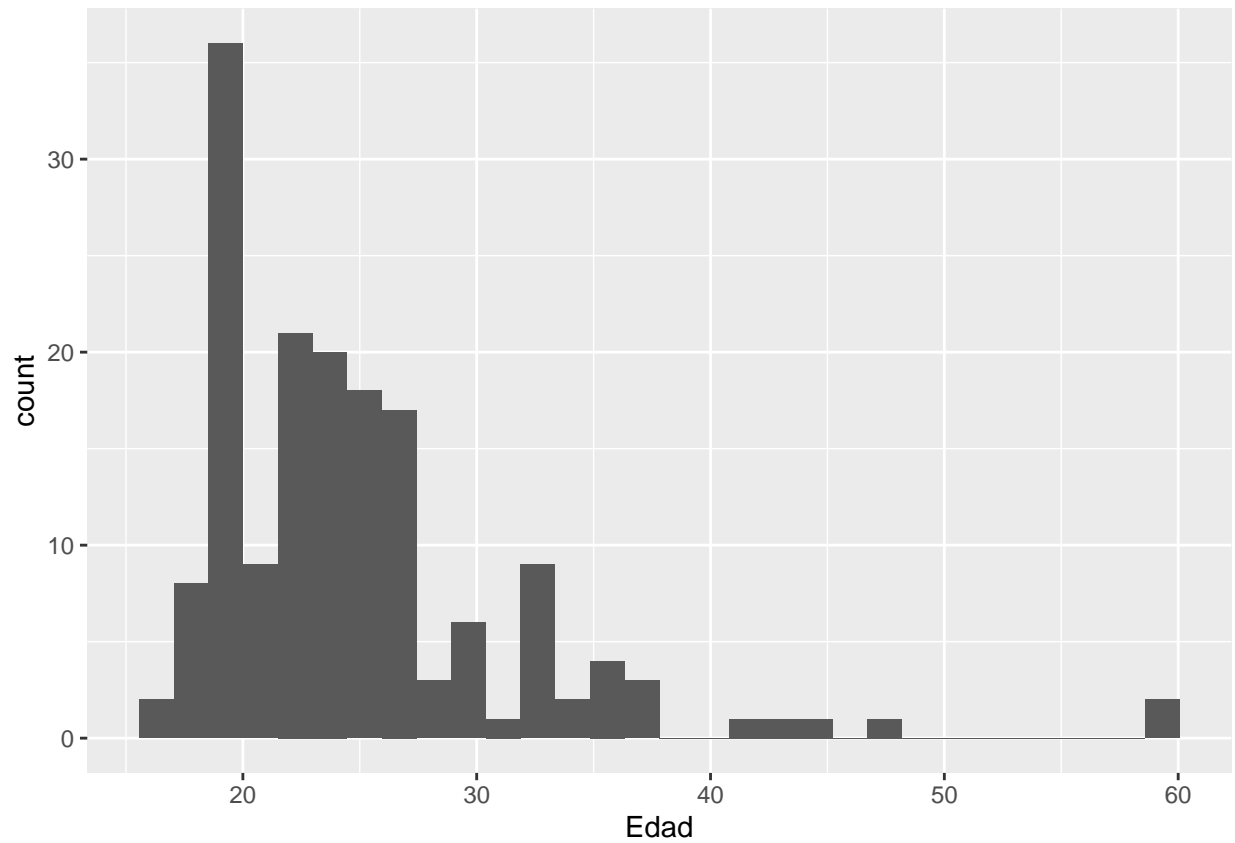


```
library(ggplot2)
Becarios_2020 %>%
  ggplot(aes(Edad))+
  geom_density()+
  scale_x_log10()
```



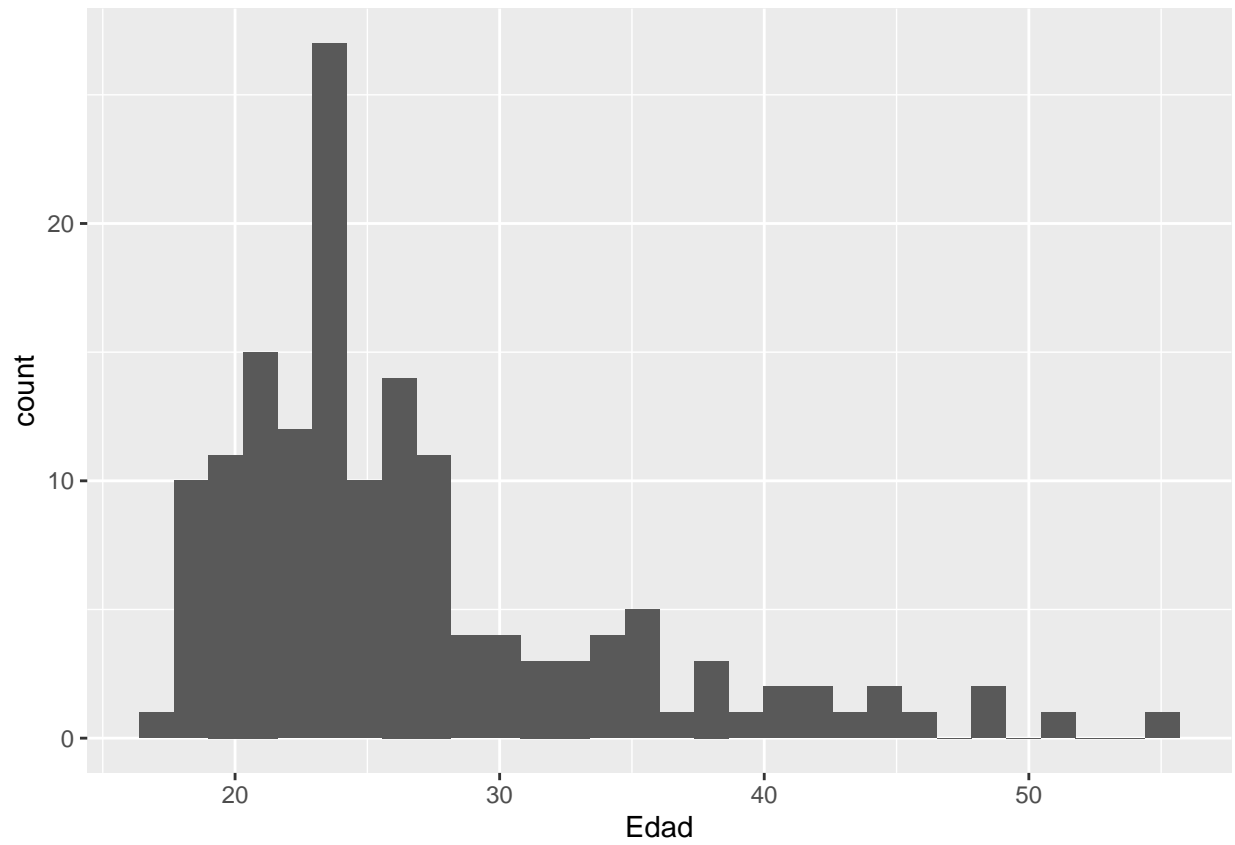
```
library(ggplot2)
options(scipen=100)
Becarios_2020 %>%
  ggplot(aes(x=Edad))+
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



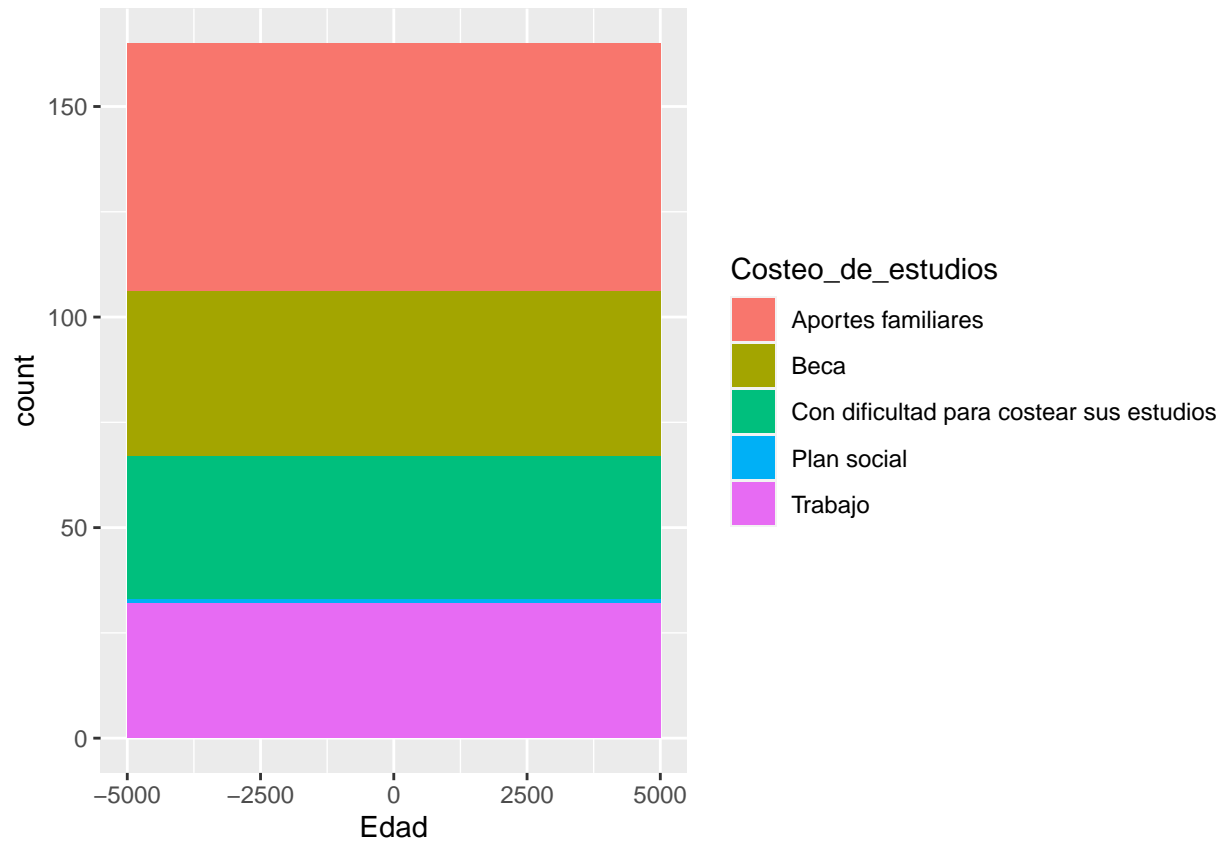
```
Becarios_2021 %>%  
  ggplot(aes(x=Edad))+  
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

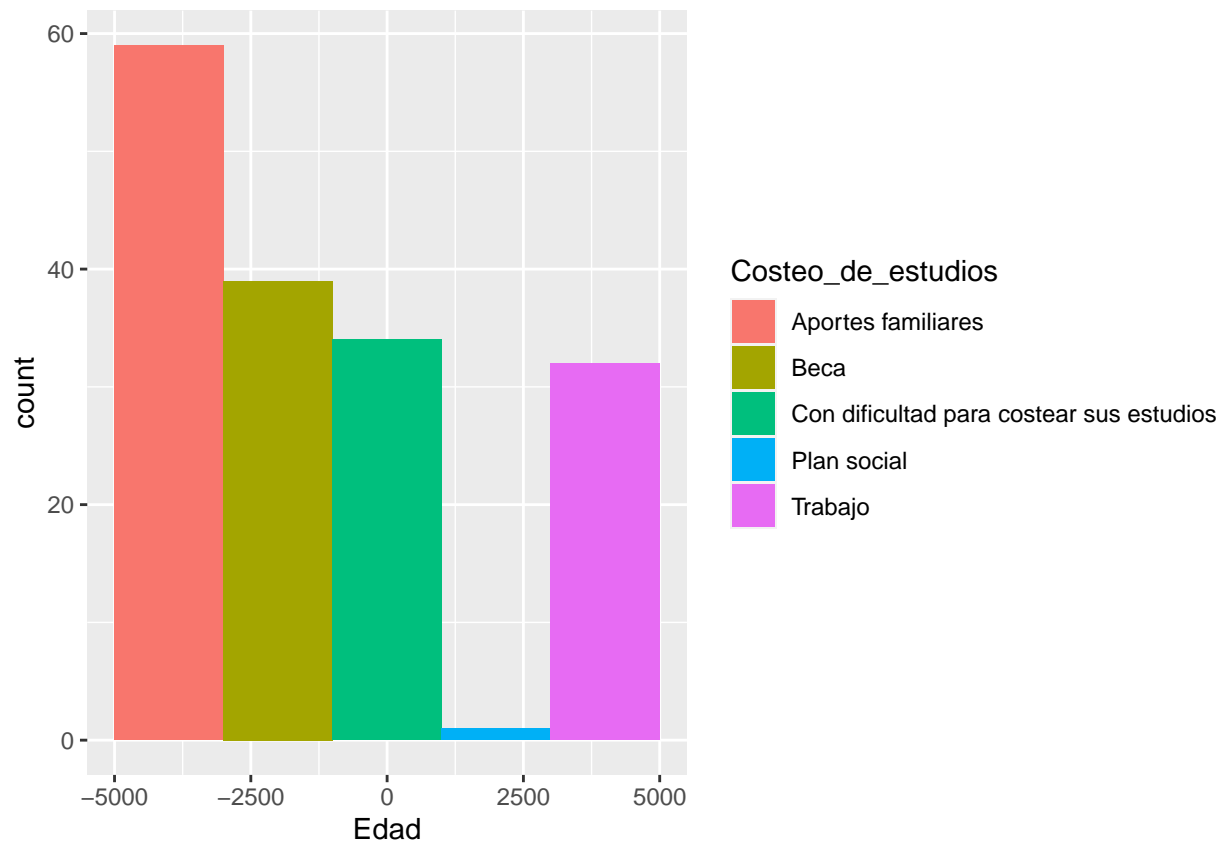


```
library(ggplot2)
Becarios_2020 %>%
  ggplot(aes(x=Edad, fill=Costeo_de_estudios))+
  geom_histogram(binwidth = 10000)
```

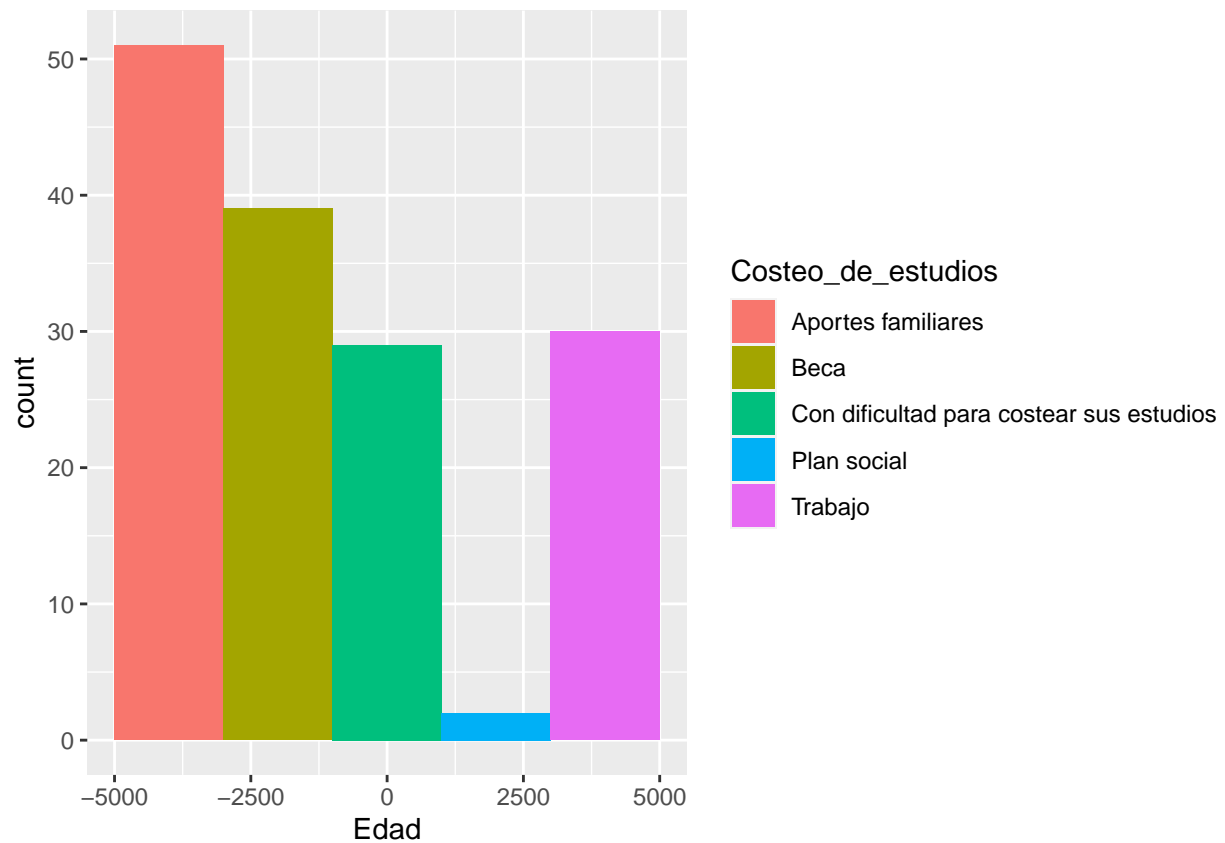




```
Becarios_2020 %>%
  ggplot(aes(x=Edad, fill=Costeo_de_estudios))+
  geom_histogram(binwidth = 10000, position = "dodge")
```

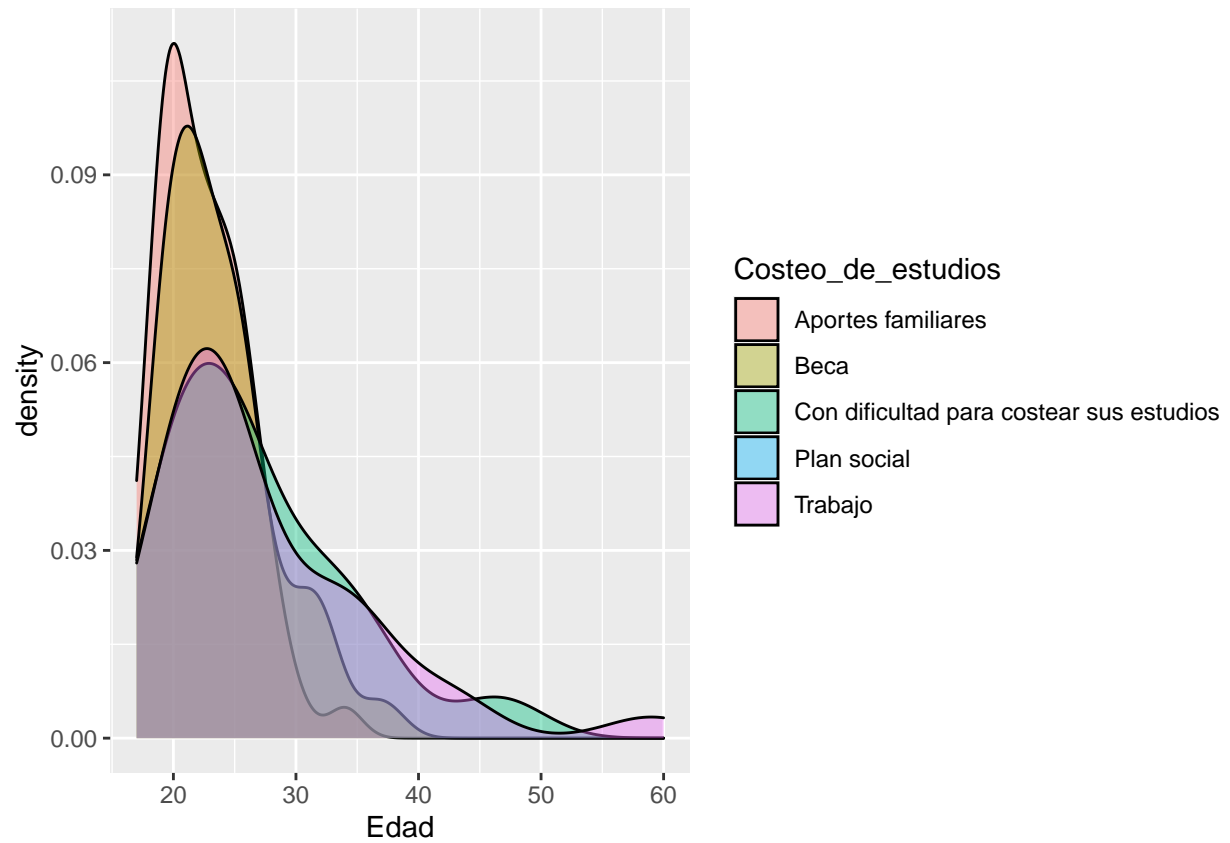


```
Becarios_2021 %>%
  ggplot(aes(x=Edad, fill=Costeo_de_estudios))+
  geom_histogram(binwidth = 10000, position = "dodge")
```

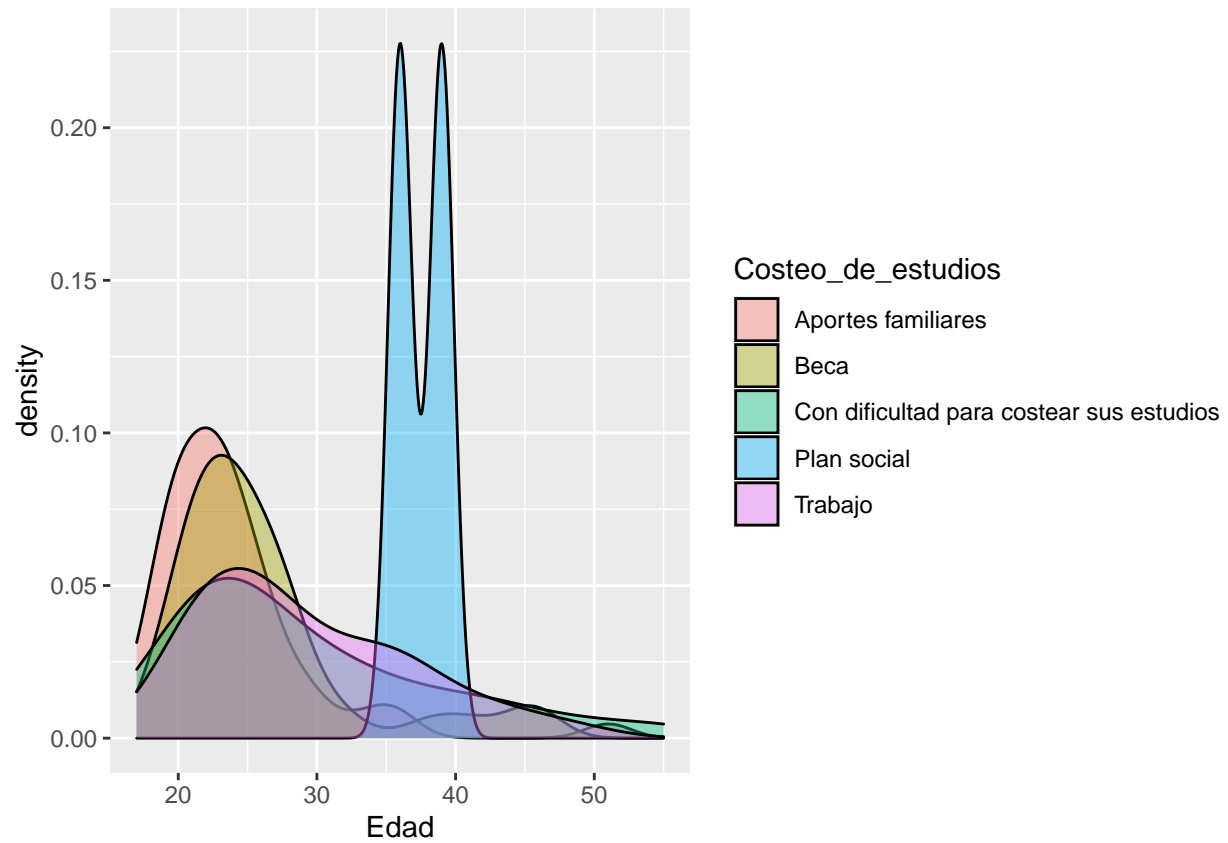


```
Becarios_2020%>%
  ggplot(aes(x=Edad, fill=Costeo_de_estudios))+
  geom_density(alpha=.4)
```

```
## Warning: Groups with fewer than two data points have been dropped.
## Warning in max(ids, na.rm = TRUE): ningun argumento finito para max; retornando
## -Inf
```



```
Becarios_2021%>%
  ggplot(aes(x=Edad, fill=Costeo_de_estudios))+
  geom_density(alpha=.4)
```



```
#library(ggplot2)
#p <- ggplot(Becarios_2020)
#p <- p + aes (x = Edad, y = Costeo_de_estudios)
#p <- p + geom_point()
#p
```