

# Inferencia Estadística: Proyecto Parte B

6 de julio de 2022

*Universidad Nacional de Colombia*

Profesor: Mario Enrique Arrieta Prieto

Ander Steven Cristancho Sánchez      John Anderson Guarín López  
Juan David Carrascal Ibañez

## Problema 1:

**A. Comparación de varios intervalos de confianza para una proporción en una muestra aleatoria Bernoulli.**

**Solución:**

1. Para llegar al intervalo de confianza de la posibilidad 1, se siguió el siguiente procedimiento.  
Tenemos que:

$$p \left( -z_{1-\frac{\alpha}{2}} \leq \frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{p(1-p)}} \leq z_{1-\frac{\alpha}{2}} \right) \approx 1 - \alpha$$

De aquí, haremos lo siguiente:

$$\begin{aligned} p \left( -z_{1-\frac{\alpha}{2}} \leq \frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{p(1-p)}} \leq z_{1-\frac{\alpha}{2}} \right) &= p \left( \left| \frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{p(1-p)}} \right| \leq z_{1-\frac{\alpha}{2}} \right) \\ &= p \left( \left| \frac{\sqrt{n}(\hat{p}_n - p)}{\sqrt{p(1-p)}} \right|^2 \leq z_{1-\frac{\alpha}{2}}^2 \right) \\ &= p \left( n(\hat{p}_n - p)^2 \leq z_{1-\frac{\alpha}{2}}^2 (p(1-p)) \right) \\ &= p \left( n\hat{p}_n^2 - 2n\hat{p}_n p + np^2 \leq z_{1-\frac{\alpha}{2}}^2 p - z_{1-\frac{\alpha}{2}}^2 p^2 \right) \\ &= p \left( n\hat{p}_n^2 - 2n\hat{p}_n p + np^2 - z_{1-\frac{\alpha}{2}}^2 p + z_{1-\frac{\alpha}{2}}^2 p^2 \leq 0 \right) \\ &= p \left( p^2(n + z_{1-\frac{\alpha}{2}}^2) - p(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) + \hat{p}_n^2 \leq 0 \right) \\ &\approx 1 - \alpha \end{aligned}$$

Ahora, resolvemos la desigualdad para p, usando la fórmula cuadrática, con lo que obtenemos que:

$$\begin{aligned}
p &= \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) \pm \sqrt{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2)^2 - 4(n + z_{1-\frac{\alpha}{2}}^2)(\hat{p}_n^2)}}{2(n + z_{1-\frac{\alpha}{2}}^2)} \\
&= \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) \pm \sqrt{4n^2\hat{p}_n^2 + 4n\hat{p}_nz_{1-\frac{\alpha}{2}}^2 + z_{1-\frac{\alpha}{2}}^4 - 4n^2\hat{p}_n^2 - 4n\hat{p}_nz_{1-\frac{\alpha}{2}}^2}}{2(n + z_{1-\frac{\alpha}{2}}^2)} \\
&= \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) \pm \sqrt{4n\hat{p}_nz_{1-\frac{\alpha}{2}}^2 + z_{1-\frac{\alpha}{2}}^4 - 4n\hat{p}_n^2z_{1-\frac{\alpha}{2}}^2}}{2(n + z_{1-\frac{\alpha}{2}}^2)} \\
&= \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) \pm \sqrt{z_{1-\frac{\alpha}{2}}^2(4n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2 - 4n\hat{p}_n^2)}}{2(n + z_{1-\frac{\alpha}{2}}^2)} \\
&= \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) \pm z_{1-\frac{\alpha}{2}}\sqrt{4n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2 - 4n\hat{p}_n^2}}{2(n + z_{1-\frac{\alpha}{2}}^2)}
\end{aligned}$$

Ahora, como el coeficiente de  $p^2$  es positivo y queremos los valores que hacen que esa inecuación sea menor o igual a 0, la desigualdad se cumple para cualquier  $p$  entre el intervalo conformado por las dos soluciones, a las que llamaremos  $L(\hat{p}_n)$  siendo el límite inferior del intervalo, y  $U(\hat{p}_n)$  siendo el límite superior del intervalo. Por tanto,  $(L(\hat{p}_n), U(\hat{p}_n))$  es un intervalo de confianza aproximado de nivel  $1 - \alpha$  para  $p$ , el cual es dado por:

$$ICA_{100(1-\alpha)\%}(p) = (L(\hat{p}_n), U(\hat{p}_n))$$

Siendo  $L(\hat{p}_n)$  :

$$L(\hat{p}_n) = \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) - z_{1-\frac{\alpha}{2}}\sqrt{4n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2 - 4n\hat{p}_n^2}}{2(n + z_{1-\frac{\alpha}{2}}^2)}$$

Y siendo  $U(\hat{p}_n)$  :

$$U(\hat{p}_n) = \frac{(2n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2) + z_{1-\frac{\alpha}{2}}\sqrt{4n\hat{p}_n + z_{1-\frac{\alpha}{2}}^2 - 4n\hat{p}_n^2}}{2(n + z_{1-\frac{\alpha}{2}}^2)}$$

2. **Implementación del algoritmo** Presentaremos el algoritmo completo en el informe pues consideramos que fue de vital importancia en el desarrollo del proyecto. En el código correspondiente a este punto describimos los pasos más relevantes.

```

1 alpha <- 0.05
2 a<-seq(from = 0.05,to = 0.95,by = (0.05))
3 b<-c(5,10,50,100,200,500,1000)
4 z_0.975 <- qnorm(0.975)

```

```

5
6 p_n <- function(x){
7   mean(x)
8 }
9 se2 <- function(x){
10   sqrt((p_n(x)*(1-p_n(x)))/n)
11 }
12 error2 <- function(x){
13   z_0.975*(se2(x))
14 }
15
16 l1 <- function(x){
17   ((2*n*p_n(x)+(z_0.975)^2)-(z_0.975)*sqrt(4*n*p_n(x)+(z_0.975)^2-4*n*(p_n(x))
18     ^2))/(2*(n+(z_0.975)^2))
19 }
20 u1 <- function(x){
21   ((2*n*p_n(x)+(z_0.975)^2)+(z_0.975)*sqrt(4*n*p_n(x)+(z_0.975)^2-4*n*(p_n(x))
22     ^2))/(2*(n+(z_0.975)^2))
23 }
24 l2 <- function(x){
25   p_n(x)-error2(x)
26 }
27 u2 <- function(x){
28   p_n(x)+error2(x)
29 }
30
31 l3 <- function(x){
32   (sin(asin(sqrt(p_n(x)))-(z_0.975)/(2*sqrt(n))))^2
33 }
34 u3 <- function(x){
35   (sin(asin(sqrt(p_n(x)))+(z_0.975)/(2*sqrt(n))))^2
36 }
37
38 long1 <- function (x){
39   u1(x)-l1(x)
40 }
41 long2 <- function (x){
42   u2(x)-l2(x)
43 }
44 long3 <- function (x){
45   u3(x)-l3(x)
46 }
47
48 for (n in b) {
49   Cobertura1<-paste("Cob", "P1",n, sep="_")
50   assign(Cobertura1,numeric())
51   Cobertura2<-paste("Cob", "P2",n, sep="_")
52   assign(Cobertura2,numeric())
53   Cobertura3<-paste("Cob", "P3",n, sep="_")
54   assign(Cobertura3,numeric())
55 }

```

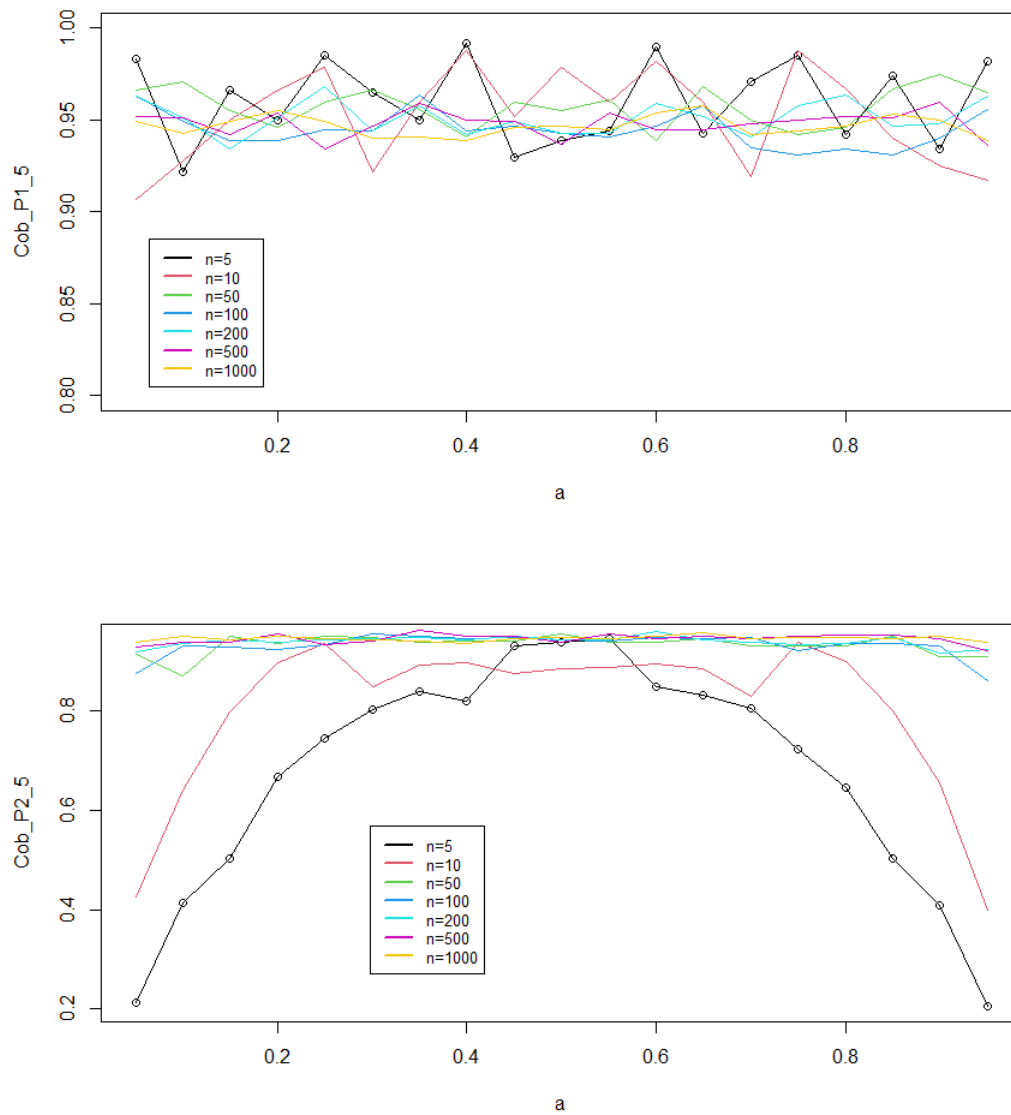
```
56 Longitud1<-paste("Long", "P1", n, sep="_")
57 assign(Longitud1, numeric())
58 Longitud2<-paste("Long", "P2", n, sep="_")
59 assign(Longitud2, numeric())
60 Longitud3<-paste("Long", "P3", n, sep="_")
61 assign(Longitud3, numeric())
62 }
63
64 set.seed(129)
65 for(p in a){
66   for(n in b){
67
68     P1<-paste("P1", n, p, sep="_")
69     assign(P1, numeric())
70
71     P2<-paste("P2", n, p, sep="_")
72     assign(P2, numeric())
73
74     P3<-paste("P3", n, p, sep="_")
75     assign(P3, numeric())
76
77
78     for(i in 1:1000){
79
80       x<-rbinom(n, 1, p)
81
82
83       x_1 <- c(eval(parse(text=P1)), 0)
84       x_2 <- c(eval(parse(text = P1)), long1(x))
85       y_1 <- c(eval(parse(text=P2)), 0)
86       y_2 <- c(eval(parse(text = P2)), long2(x))
87       z_1 <- c(eval(parse(text=P3)), 0)
88       z_2 <- c(eval(parse(text=P3)), long3(x))
89
90       if (l1(x)>=p | u1(x)<=p){
91         assign(P1, x_1)
92       }else{
93         assign(P1, x_2)
94       }
95       if (l2(x)>=p | u2(x)<=p){
96         assign(P2, y_1)
97       }else{
98         assign(P2, y_2)
99       }
100      if (l3(x)>=p | u3(x)<=p){
101        assign(P3, z_1)
102      }else{
103        assign(P3, z_2)
104      }
105    }
106
107
108    c_1 <- (sum(eval(parse(text=P1))!=0)/1000)
```

```

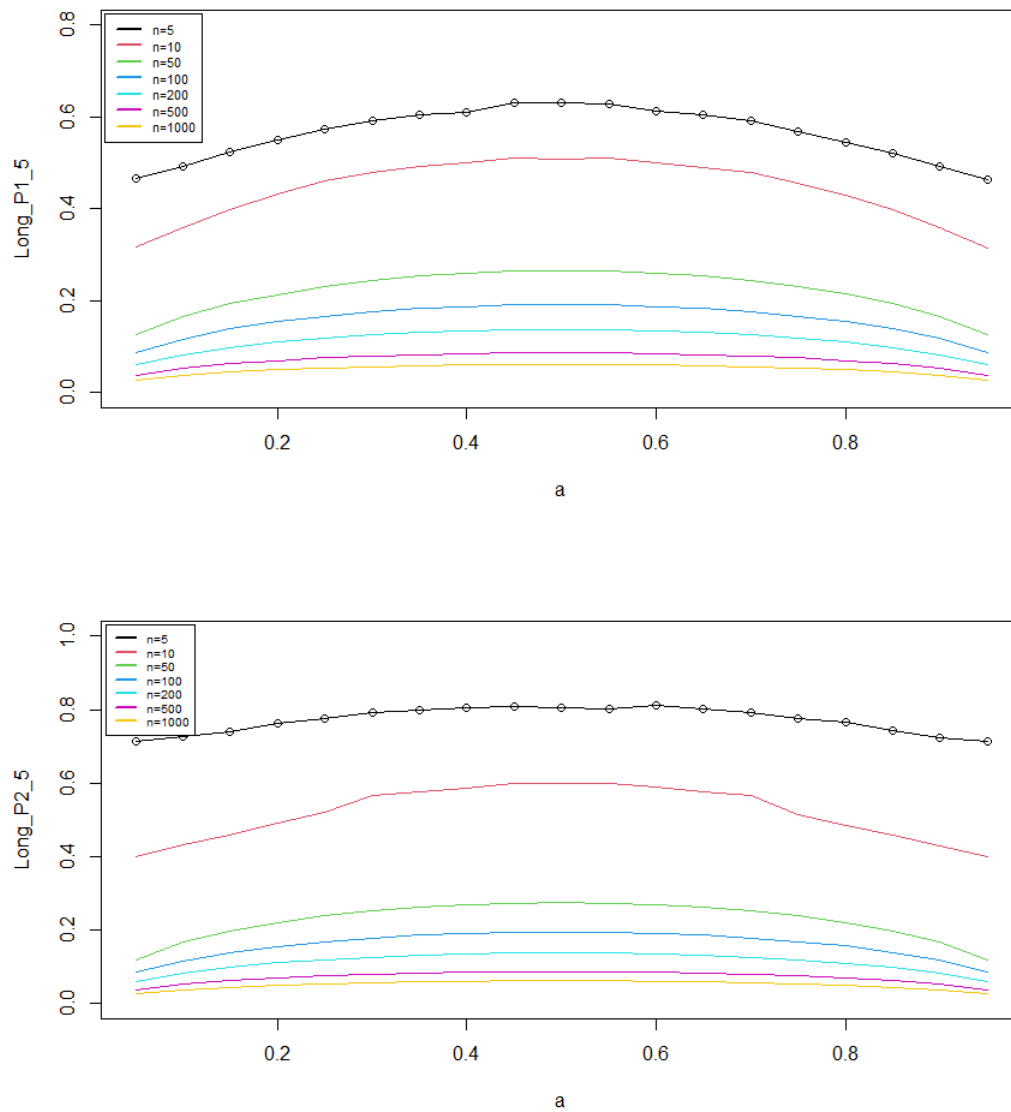
109   c_2 <- (sum(eval(parse(text=P2))!=0)/1000)
110   c_3 <- (sum(eval(parse(text=P3))!=0)/1000)
111   c_4 <- c(eval(parse(text=paste("Cob", "P1", n, sep="_"))), c_1)
112   c_5 <- c(eval(parse(text=paste("Cob", "P2", n, sep="_"))), c_2)
113   c_6 <- c(eval(parse(text=paste("Cob", "P3", n, sep="_"))), c_3)
114
115   l_1 <- sum(eval(parse(text=P1)))/sum(eval(parse(text=P1))!=0)
116   l_2 <- sum(eval(parse(text=P2)))/sum(eval(parse(text=P2))!=0)
117   l_3 <- sum(eval(parse(text=P3)))/sum(eval(parse(text=P3))!=0)
118   l_4 <- c(eval(parse(text=paste("Long", "P1", n, sep="_"))), l_1)
119   l_5 <- c(eval(parse(text=paste("Long", "P2", n, sep="_"))), l_2)
120   l_6 <- c(eval(parse(text=paste("Long", "P3", n, sep="_"))), l_3)
121
122
123   CP_I1<-paste("Cob", "Prom", P1, sep="_")
124   assign(CP_I1, c_1)
125   CP_I2<-paste("Cob", "Prom", P2, sep="_")
126   assign(CP_I2, c_2)
127   CP_I3<-paste("Cob", "Prom", P3, sep="_")
128   assign(CP_I3, c_3)
129
130   LL_I1<-paste("Long", "Prom", P1, sep="_")
131   assign(LL_I1, l_1)
132   LL_I2<-paste("Long", "Prom", P2, sep="_")
133   assign(LL_I2, l_2)
134   LL_I3<-paste("Long", "Prom", P3, sep="_")
135   assign(LL_I3, l_3)
136
137   assign(paste("Cob", "P1", n, sep="_"), c_4)
138   assign(paste("Cob", "P2", n, sep="_"), c_5)
139   assign(paste("Cob", "P3", n, sep="_"), c_6)
140
141
142   assign(paste("Long", "P1", n, sep="_"), l_4)
143   assign(paste("Long", "P2", n, sep="_"), l_5)
144   assign(paste("Long", "P3", n, sep="_"), l_6)
145 }
146 }

```

3. Los plots correspondientes las coberturas promedio, para cada tamaño de muestra y cada p, son los siguientes:

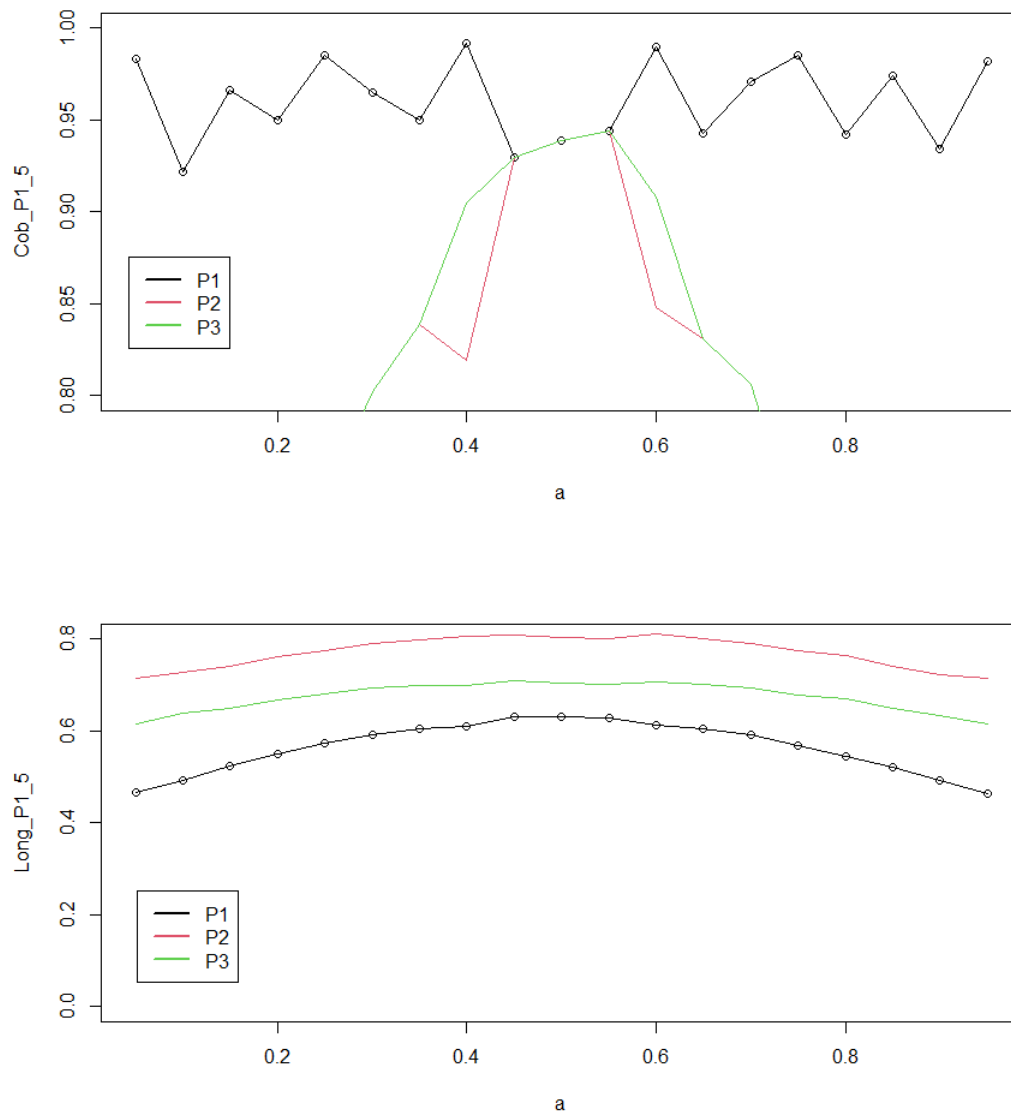


Y los plots correspondientes a las longitudes promedio, para cada tamaño de muestra y cada  $p$ , son los siguientes:

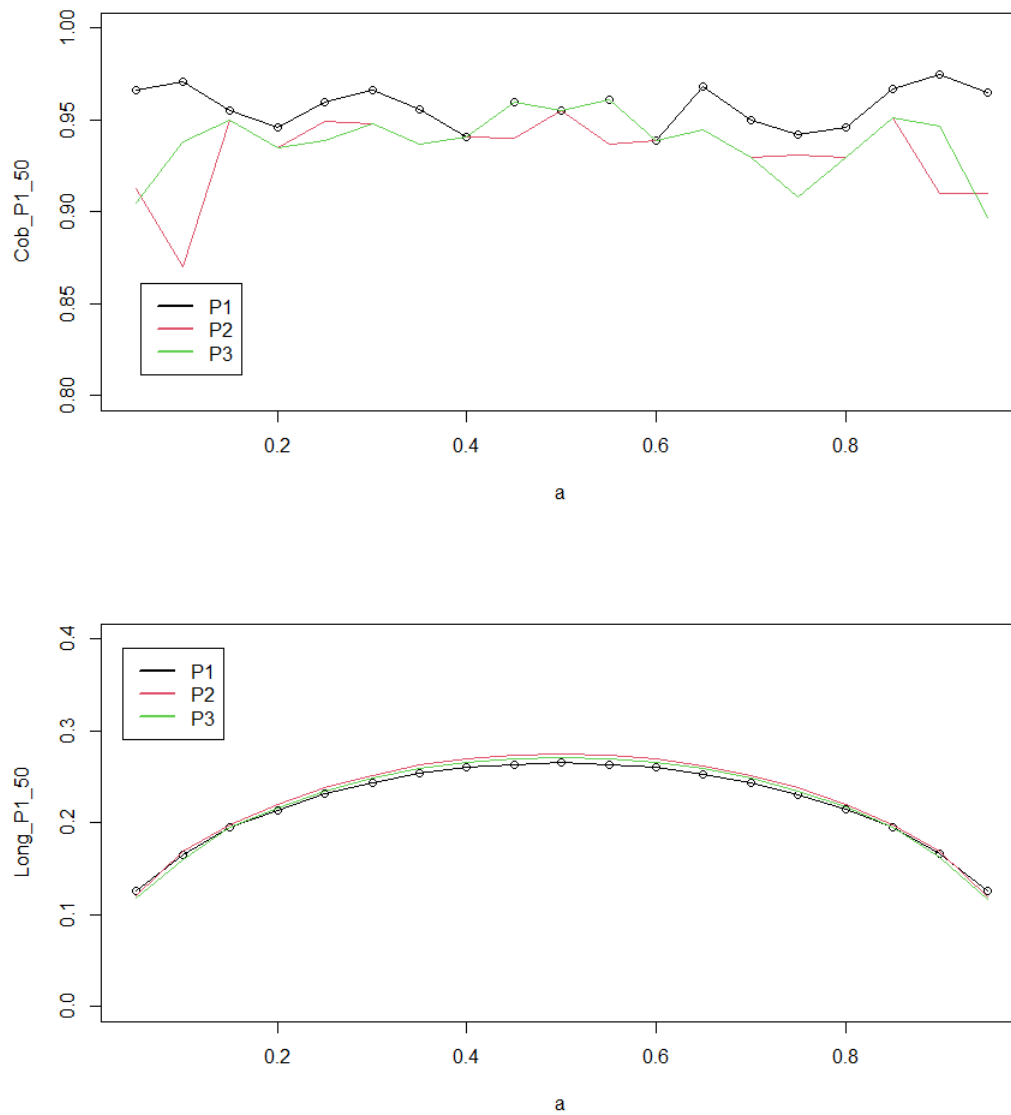


4. Los plots correspondientes a la cobertura y longitud promedio, cuando el tamaño de muestra es  $n = 5$ , para cada posibilidad, son los siguientes:

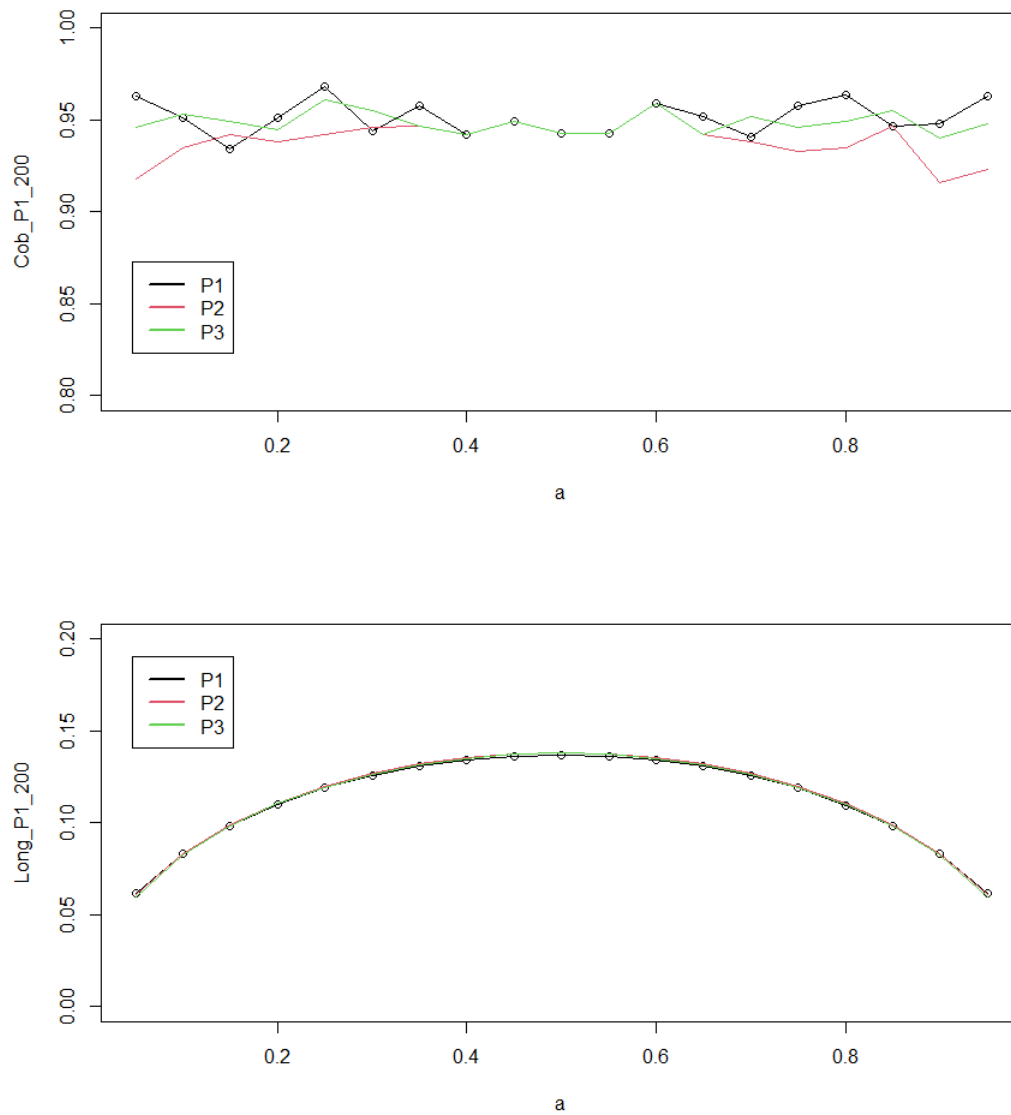




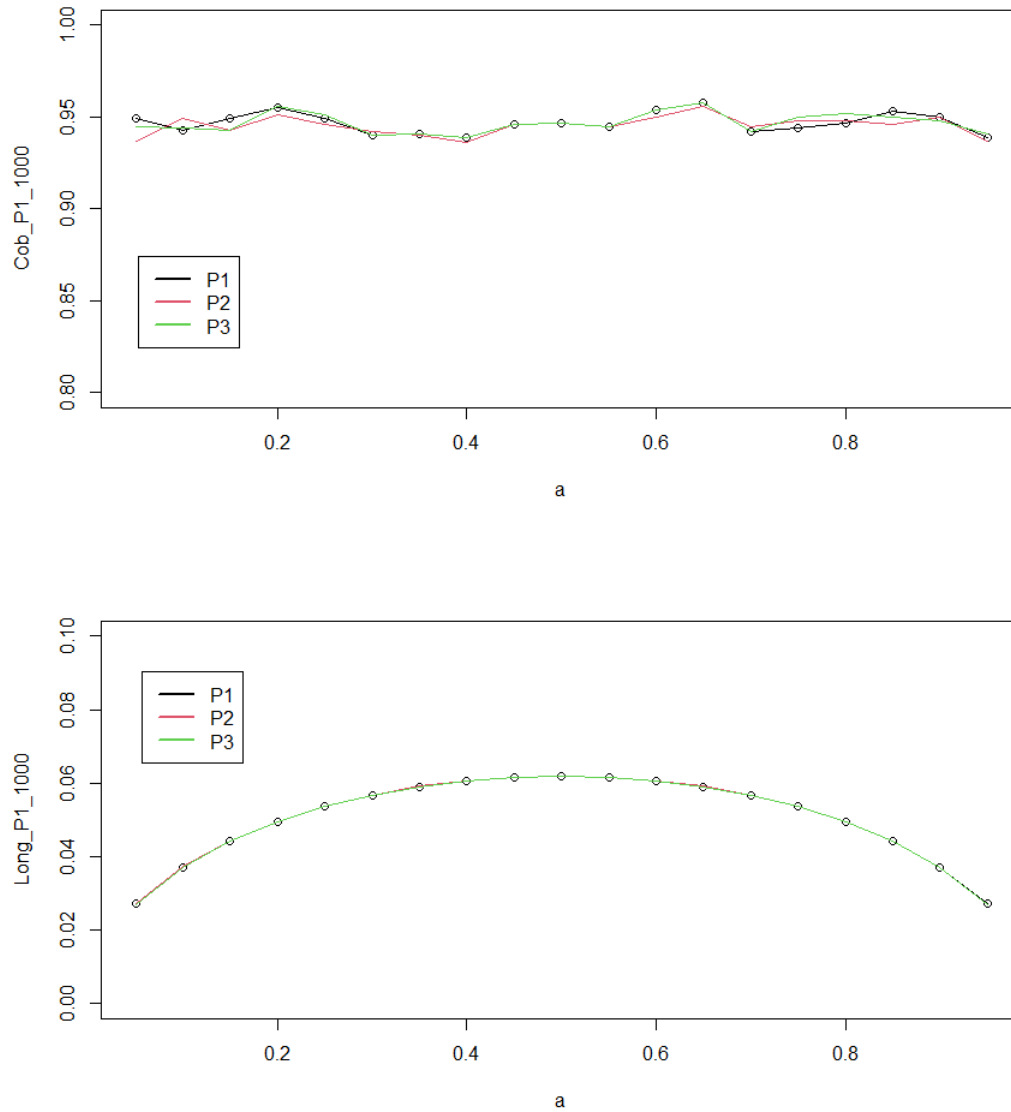
Los plots correspondientes a la cobertura y longitud promedio, cuando el tamaño de muestra es  $n = 50$ , para cada posibilidad, son los siguientes:



Los plots correspondientes a la cobertura y longitud promedio, cuando el tamaño de muestra es  $n = 200$ , para cada posibilidad, son los siguientes:



Y por último, los plots correspondientes a la cobertura y longitud promedio, cuando el tamaño de muestra es  $n = 1000$ , para cada posibilidad, son los siguientes:



5. Podemos concluir que cuando aumentamos el tamaño de muestra, la longitud de los intervalos promedio se reduce mucho, y la cobertura promedio de los intervalos no sólo aumenta sino que además parece estabilizarse.

Además de esto, cuando el parámetro  $p$  está muy cerca de los extremos del intervalo  $(0, 1)$ , la cobertura y la longitud se reduce. Esto se debe a que, cuando el parámetro estaba cerca de los extremos, en las simulaciones obteníamos muestras con valores también cerca a los extremos,

y esto afectaba en que, al hallar los intervalos de confianza para cada muestra, se pasaran del soporte del parámetro y además, dieran valores muy parecidos pues la probabilidad de éxito era o muy baja o muy alta, por lo que la longitud del intervalo de confianza era menor pues la proporción muestral no tenía datos muy atípicos. Esto repercute en que, al ser intervalos más pequeños, su cobertura fuera un poco menor.

También podemos concluir que la posibilidad 1 trabaja mucho mejor que la posibilidad 2 y 3, esto debido a que, aunque las tres posibilidades se basan en el teorema del límite central y por tanto requieren una muestra grande para funcionar muy bien, la posibilidad 1 es más confiable pues por ejemplo, a diferencia de la posibilidad 2, se basa únicamente en el teorema del límite central, mientras que la posibilidad 2 también se basa en una estimación consistente del parámetro, para que las cuentas no sean tan engorrosas, sacrificando un poco de precisión con muestras más pequeñas, como se puede ver en los plots del punto 3 y 4.

Pero, como las tres posibilidades están basadas en el teorema del límite central, con muestras muy grandes parecen trabajar muy parecido, y ser igual de confiables. Esto se puede notar en los plots del punto 4, con las longitudes y coberturas de las tres posibilidades cuando  $n = 1000$ .

## Problema 2:

### B. Análisis inferencial sobre datos reales

A continuación se presentan las rutinas usadas para generar las estimaciones bootstrap por intervalo:

```
1 #Bootstrap_1
2 #Funcion generadora de estimadores
3 x_i=NULL
4 var_1<-function(n){
5   for (i in 1:42){
6     x_i<-append(x_i,bootstrap[i,n])
7   }
8   return(var(x_i))
9 }
10 #Funcion que almacena estos estimadores
11 y_i=NULL
12 for (i in 1:1000){
13   y_i=append(y_i,var_1(i))
14 }
15 #Bootstrap_2
16 #Funcion generadora de estimadores
17 x_i=NULL
18 var_2<-function(n){
19   for (i in 1:42){
20     x_i<-append(x_i,bootstrap[i,n])
21   }
22   return(var(x_i)-v_o)
23 }
24 #Funcion que almacena estos estimadores
25 y_i=NULL
26 for (i in 1:1000){
```

```
27 y_i=append(y_i,var_2(i))
28 }
29 #En ambos casos para el punto 2
30 #Limite inferior
31 quantile(y_i,.0275)
32 #Limite superior
33 quantile(y_i,.975)
```

### Solución:

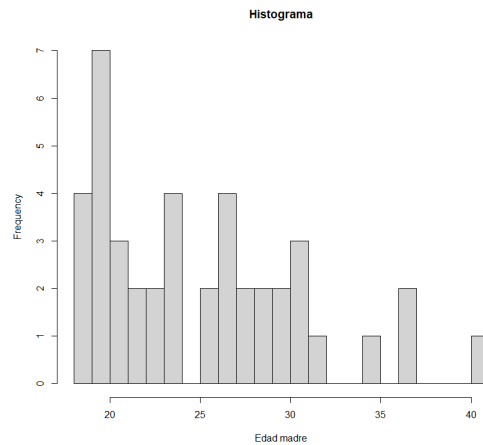
1. Calcule estimaciones puntuales de los siguientes parámetros en la población. No es necesario que use estimadores UMVUE, simplemente justifique qué estimador usaron en cada caso y qué propiedades deseables tiene.

a. ¿Qué porcentaje de nacidos tienen una madre que no supera los 20 años?

### Solución:

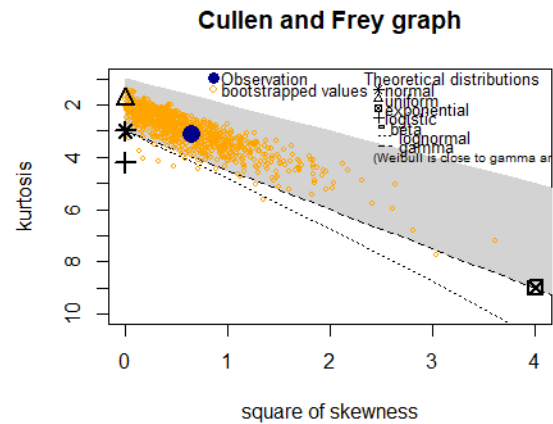
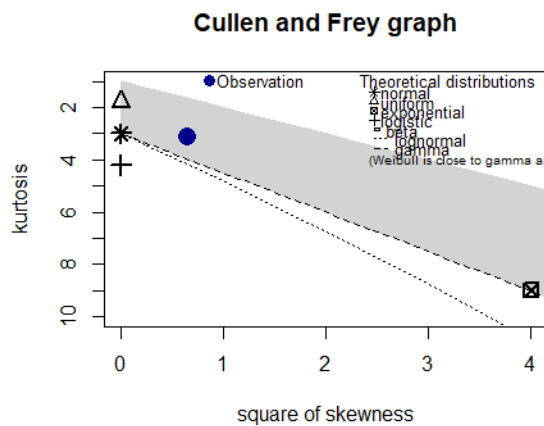
Definimos la variable aleatoria  $X$  : 'Edad de la madre' y podaremos a encontrar su distribución. Estableceremos el data set a trabajar y haremos uso de las siguientes librerías:

```
1 data=read.xlsx("C:/Users/kaido/Downloads/dataset_trabajo1.xlsx",startRow=1,colNames=
  TRUE)
2 #Ver tabla
3 View(data)
4 #Fijar semilla
5 set.seed(233)
6 #Librerias
7 library(openxlsx)
8 library(tidyverse)
9 library(MASS)
10 library(survival)
11 library(readr)
12 library(fitdistrplus)
13 library(datasets)
14 library(dplyr)
```



Haremos uso del diagrama de Cullen y Frey para intentar determinar cuál es la mejor distribución a ser usada:

```
1 # Paso (1) Plot Cullen y Frey
2 descdist( data = data$mage , discrete = FALSE)
3 descdist(data = data$mage, discrete = FALSE, boot=1000)
```



El uso de la muestra bootstrap ver hacia que distribución parecen inclinarse los datos cuando se hace remuestreo.

Ahora bien, si intentamos ajustar los datos a una distribución beta como nos indica el diagrama anterior mediante la implementación:

```
1 fitdlist(data$mage, "beta")
2 max(x)
```

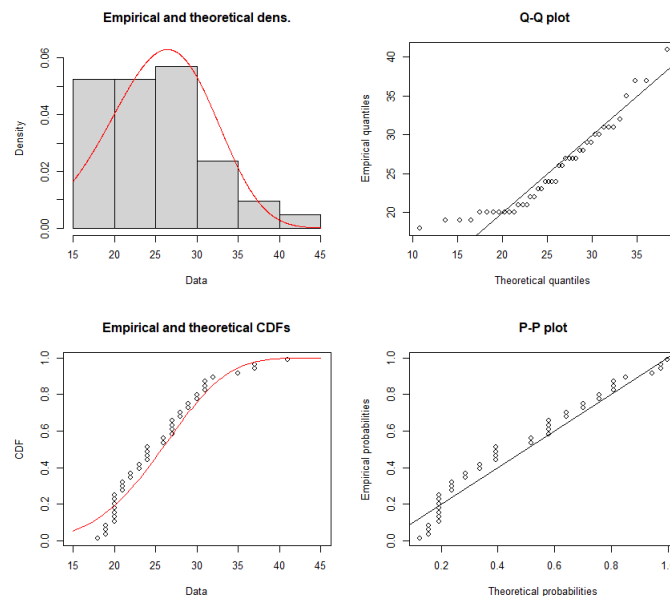
Obtenemos la siguiente salida:

```
1 > fitdist(data$mage, "beta")
2 Error in computing default starting values.
3 Error in manageparam(start.arg = start, fix.arg = fix.arg, obs = data, :
4   Error in start.arg.default(obs, distname) :
5     values must be in [0-1] to fit a beta distribution
6 > max(x)
7 [1] 41
```

Es decir, esto falla precisamente porque hay datos que no pertenecen al soporte de la distribución  $\beta$ . Por otra parte, el diagrama nos indica que los modelos que mejor se ajustan después de  $\beta$  son: normal, Weibull y gamma.

```
1 normal_ = fitdist(data$mage, "norm")
2 weibull_ = fitdist(data$mage, "Weibull")
3 gamma_ = fitdist(data$mage, "gamma")
```

De estos tres el que mejor parece ajustarse es el modelo Weibull



Realizaremos la prueba de bondad de ajuste mediante el test de Kolmogorov- Smirnov y obtendremos los parámetros de la distribución deseada:

```
1 > prueba<-gofstat(weibull_)
2 > prueba$kstest
3 1-mle-weibull
4 "not rejected"
5 > prueba$chisqpvalue
6 [1] 0.3968121 #p-value> 0.05-> no se rechaza la hipotesis nula
7 > print(weibull_)
8 Fitting of the distribution ' weibull ' by maximum likelihood
```



```

9 Parameters :
10     estimate Std. Error
11 shape    4.65238    0.5202477
12 scale    27.83948    0.9810004

```

De esta manera, con un nivel de significancia del 95 % aceptamos que nuestro modelo se aproxima a una distribución Weibull ,  $W(4,65,27,84)$ . Con esto podemos calcular la estadística de interés.

```

1 > pweibull(20, 4.65, scale = 27.83)
2 [1] 0.1936042

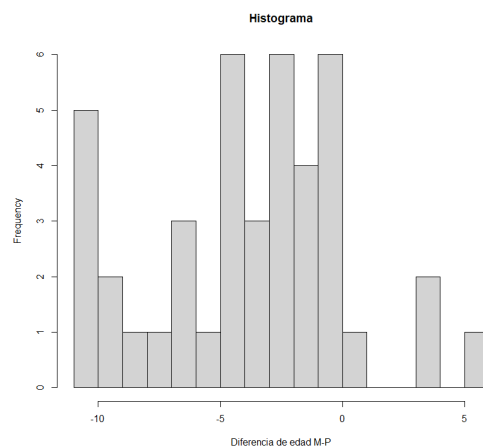
```

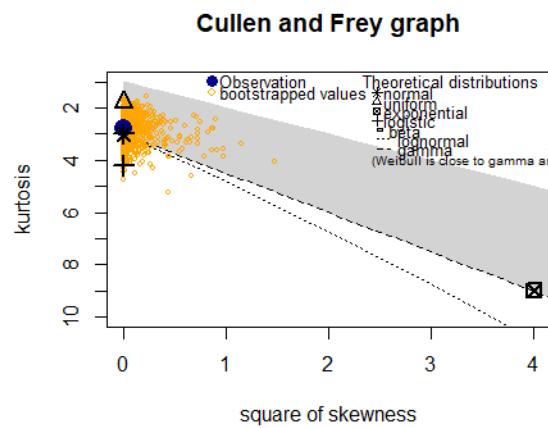
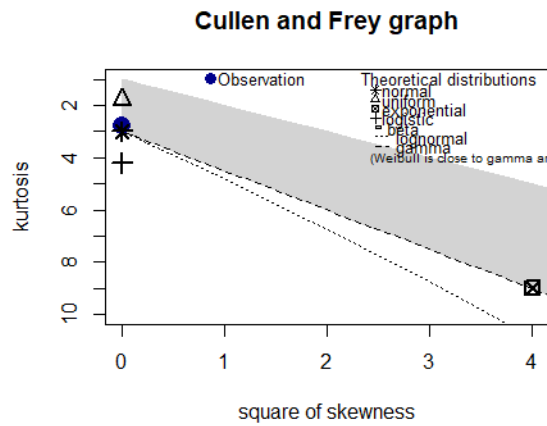
Así, afirmamos que el 19.36 % de nacidos tienen una madre que no supera los 20 años.

b. ¿Qué porcentaje de recién nacidos tienen un padre más joven que la madre?

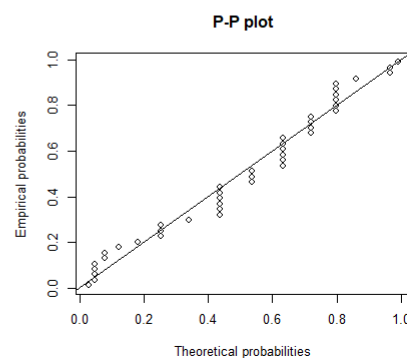
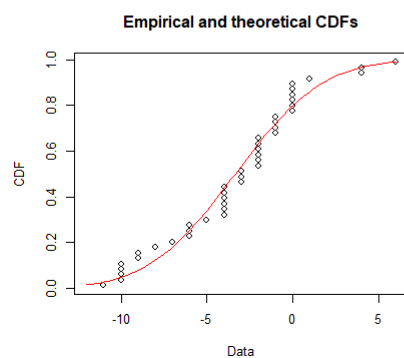
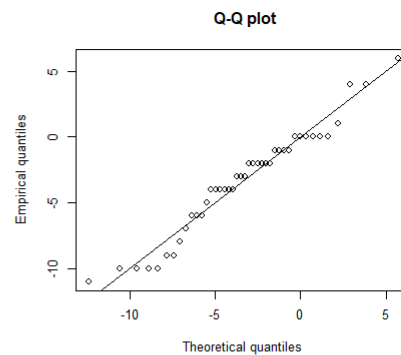
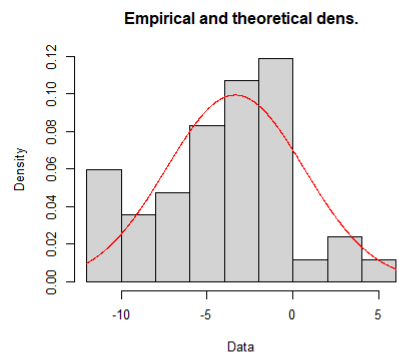
### Solución:

Sean  $X$ : 'Edad de la madre' y  $Y$ : 'Edad del padre', estamos interesados en estudiar la variable aleatoria  $Z = X - Y$ . Implementando las rutinas anteriores a esta nueva variable aleatoria obtenemos los siguientes resultados:





Estos diagramas nos indican que la distribución de estos datos parece ser normal. En efecto, la curva parece ajustarse bien como lo indica el siguiente gráfico:



Haciendo las respectivas pruebas de bondad

```
1 > prueba<-gofstat(normal_)
2 > prueba$kstest
3 1-mle-norm
4 "not rejected"
5 > prueba$chisqpvalue
6 [1] 0.03969815
7 > print(normal_)
8 Fitting of the distribution ' norm ' by maximum likelihood
9 Parameters:
10 estimate Std. Error
11 mean -3.357143 0.6188841
12 sd 4.010827 0.4376170
```

De esta manera, con un nivel de significancia del 95 % aceptamos que nuestro modelo se aproxima a una distribución Normal ,  $N(-3,35, 4,01)$ . Con esto podemos calcular la estadística de interés, para esto debemos calcular  $P(Z \geq 1)$  pues al ser la edad una cantidad discreta no sería correcto usar  $P(Z \geq 0)$ .

```
1 1-pnorm(1, mean=-3.35, sd=4.01)
2 [1] 0.1390078
```

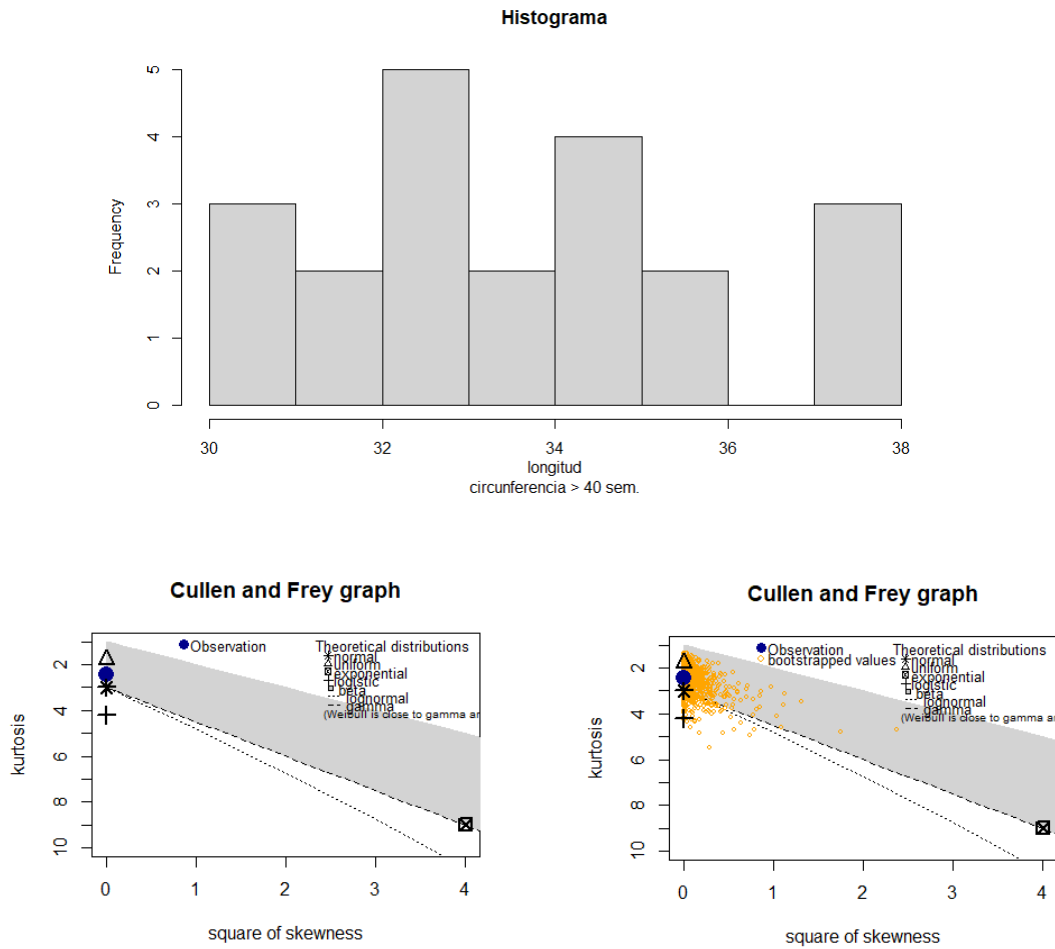
Finalmente, concluimos que el 13,9 % de recién nacidos tienen un padre más joven que la madre.

c. ¿Cuál es el promedio y la desviación estándar de la longitud de la circunferencia de los bebés que no superaron las 40 semanas de gestación?

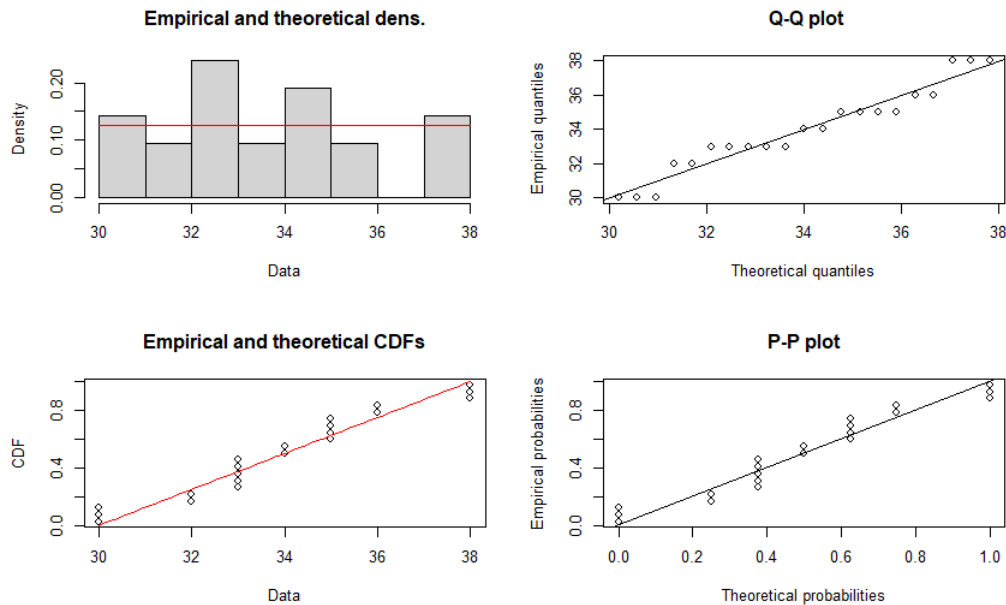
### Solución:

Inicialmente pensamos generar la distribución de estos datos a partir de la definición de distribución condicional, sin embargo esto requería muchos supuestos y debíamos realizar una integral doble. Preferimos tomar una submuestra hecha con los datos de los bebés que no superaron las 40 semanas de gestación y repetir el procedimiento de los puntos anteriores. Esta muestra fue generada con:

```
1 submuestra_c<-data %>%filter(Gestation <40)
2 x=submuestra_c$Headcirc
```



Tenemos evidencia, de que los datos se ajustan a una distribución uniforme. En efecto,



Con los siguientes resultados observamos que falla el calculo del test ks sin embargo para la prueba de chi de Pearson es posible aceptar estos resultados si tenemos una significancia  $\alpha = 0,93$ .

```

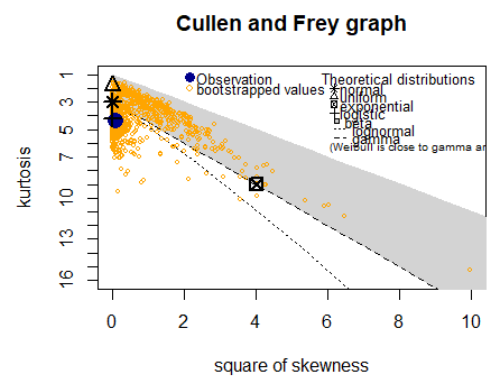
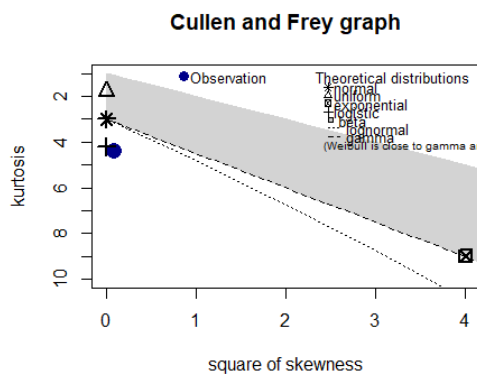
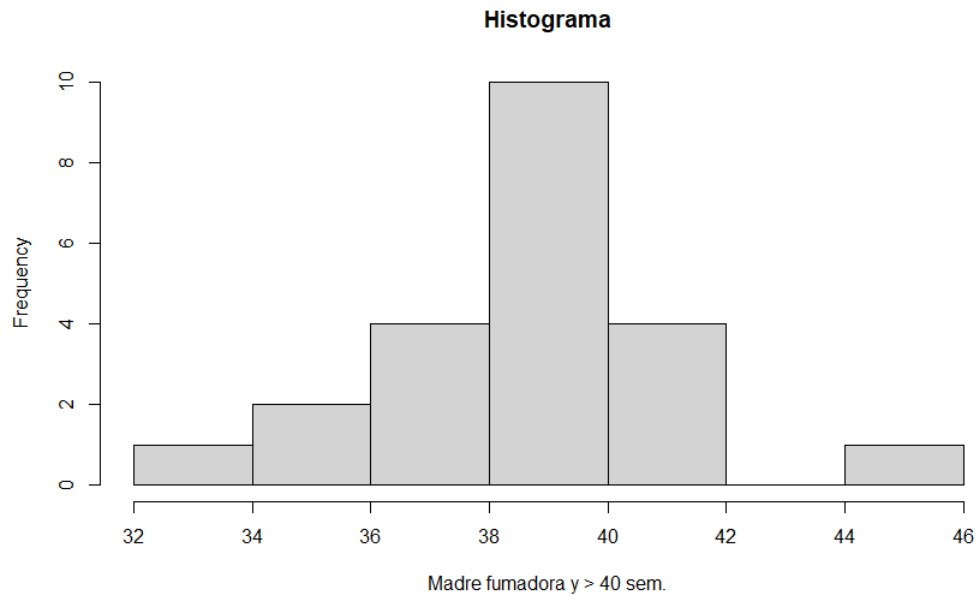
1 > prueba<-gofstat(unif_)
2 > prueba$kstest
3     1-mle-unif
4 "not computed"
5 > prueba$chisqpvalue
6 [1] 0.06854764
7 > print(unif_)
8 Fitting of the distribution ' unif ' by maximum likelihood
9 Parameters:
10      estimate Std. Error
11 min          30         NA
12 max          38         NA
13 > 1/2*(30+38)
14 [1] 34
15 > 1/12*(38-30)^2
16 [1] 5.333333

```

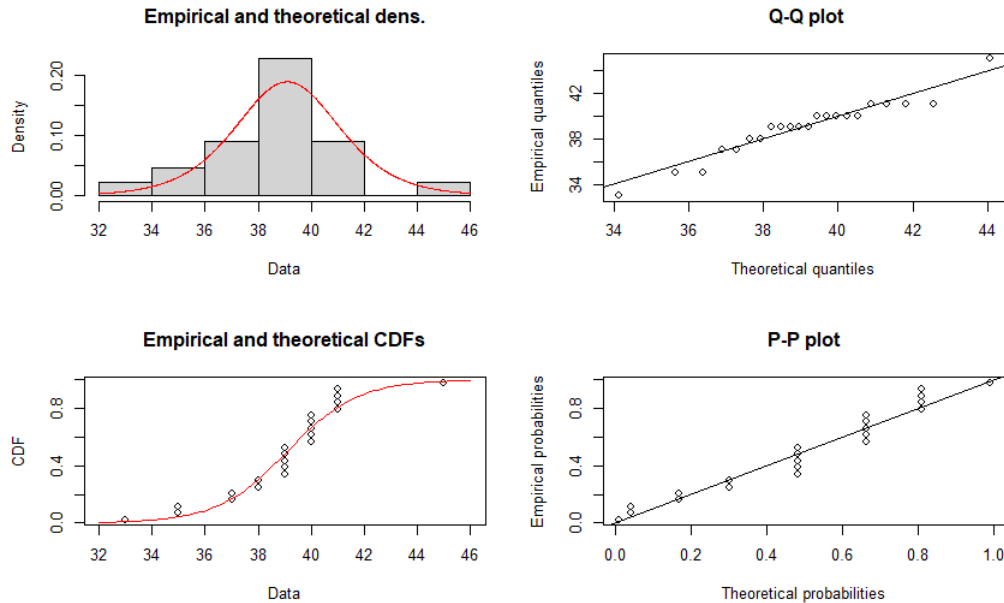
De manera que, las estimaciones de la media y la desviación estándar son 34 y 5.33, respectivamente.

d. ¿Cuál es la proporción de madres fumadoras de los bebés que no superaron las 40 semanas de gestación?

**Solución:**



Los diagramas nos dan bastante evidencia de que estamos trabajando con una distribución logística.



```

1 #Pruebas de bondad de ajuste
2 > logis_=fitdist(x,"logis")
3 > plot(logis_)
4 > prueba<-gofstat(logis_)
5 > prueba$kstest
6     1-mle-logis
7 "not computed"
8 > prueba$chisqpvalue
9 [1] 0.23136
10 > print(logis_)
11 Fitting of the distribution ' logis ' by maximum likelihood
12 Parameters:
13     estimate Std. Error
14 location 39.095971  0.4807186
15 scale    1.320257  0.2411096
16 > plogis(40,location = 39.09, scale = 1.32)
17 [1] 0.6658321

```

De nuevo falla la prueba K-S pero la prueba del p-value no. Así, con un nivel de significancia del 95 % aceptamos que nuestro modelo se aproxima a una distribución logística,  $Logis(39,09,1,32)$ . Y determinamos que la proporción de la pregunta es 66.58 %

2. Calcule e interprete un intervalo confidencial para los siguientes parámetros con un nivel de confianza del 99 %. Primero, verifique si las suposiciones distribucionales requeridas en cada caso se tienen. Si no, use ambos métodos de intervalos de confianza Bootstrap y compare los resultados.

a. Peso promedio del bebé.

### Solución:

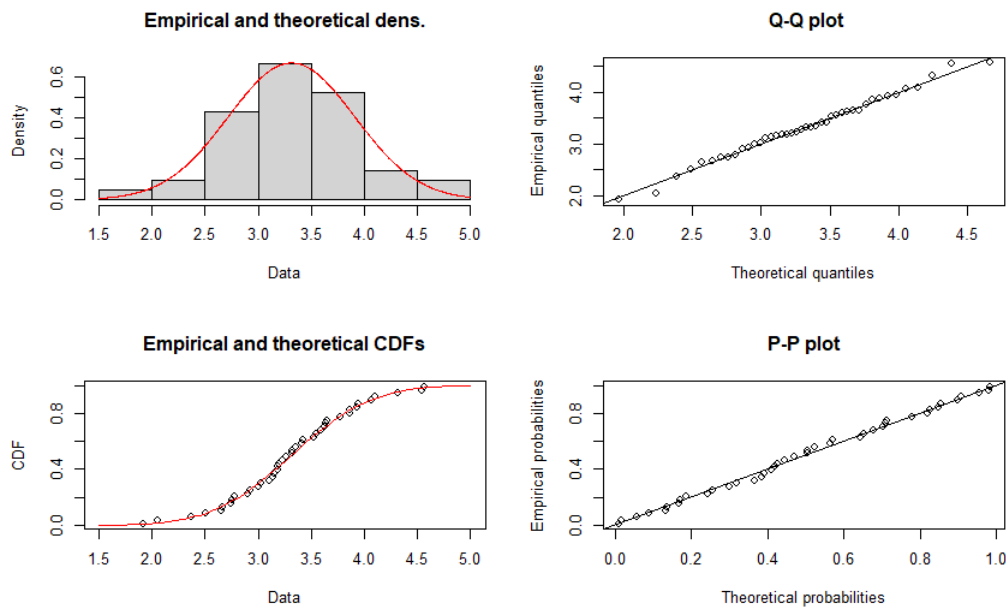
Vamos a verificar que si estos datos vienen de una distribución normal. La siguiente rutina presenta las pruebas de bondad de ajuste y la evidencia gráfica que sustenta la veracidad de esto:

```

1 > #Peso promedio bebes
2 > x=data$Birthweight
3 > #Verificamos si los datos provienen de una dist. normal
4 > normal_ = fitdist(x, "norm")
5 > plot(normal_) #Evidencia grafica
6 > prueba<-gofstat(normal_)
7 > prueba$kstest
8     1-mle-norm
9 "not rejected"
10 > prueba$chisqpvalue
11 [1] 0.7054872

```

Debido a que trabajamos con un nivel de significancia es claro que con este p-value no se rechaza la hipótesis nula de bondad de ajuste.



Con la siguiente rutina calculamos los límites de el intervalo de confianza requerido:

```

1 > #Calculo intervalo de confianza
2 > x_barra=mean(x)
3 > sd=sqrt(var(x))
4 > x_barra-qt(0.995,df=41)*sd/sqrt(42) #Limite inferior
5 [1] 3.061153
6 > x_barra+qt(0.995,df=41)*sd/sqrt(42) #Limite superior

```



```
7 [1] 3.564561
```

De esta manera concluimos que:

$$IC_{99\%}(\mu) = [3,06, 3,56]$$

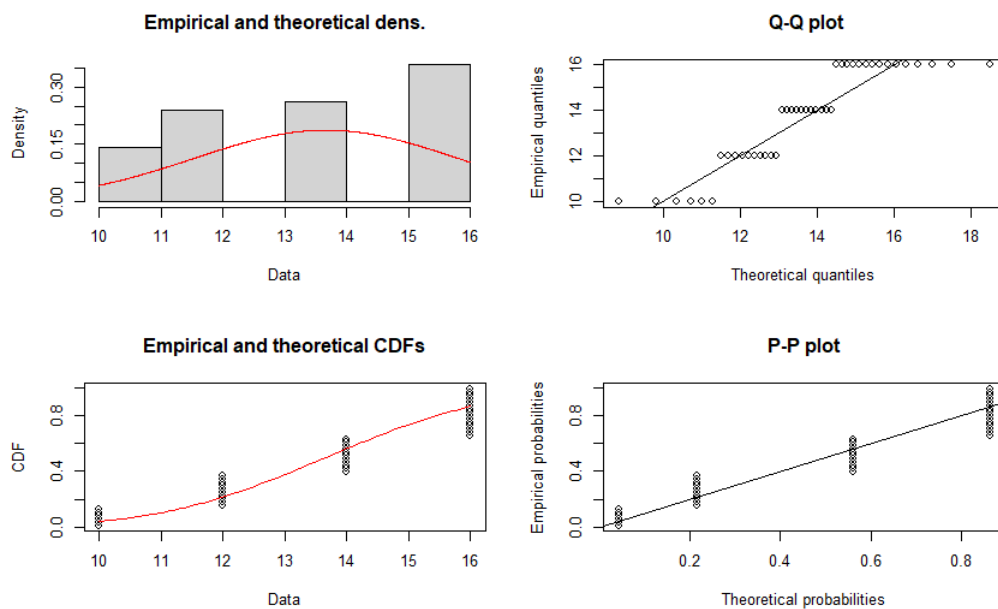
b. Varianza de la cantidad de años de educación del padre.

### Solución:

Observemos que los datos no provienen de una distribución normal:

```
1 > #Verificamos si los datos provienen de una dist. normal
2 > normal_ = fitdist(x, "norm")
3 > plot(normal_) #Evidencia grafica
4 > prueba<-gofstat(normal_)
5 > prueba$kstest
6 1-mle-norm
7 "rejected"
8 > prueba$chisqpvalue
9 [1] 0.0001371704
```

Con este p-value se rechaza la hipótesis nula.



Sin embargo, usando el diagrama de Cullen y Frey tenemos evidencia para sospechar que la distribución uniforme se ajusta a los datos. De esta manera, la siguiente rutina nos permite verificar la prueba de bondad de ajuste y la estimación de dichos parámetros:

```
1 > #Bondad de ajuste
```

```

2 > unif_=fitdist(x,"unif")
3 > plot(unif_)
4 > prueba<-gofstat(unif_)
5 > prueba$kstest
6 1-mle-unif
7 "rejected"
8 > prueba$chisqpvalue
9 [1] 0

```

A pesar de esto obtenemos peores resultados por lo que concluimos que, entre las distribuciones conocidas, la distribución normal es la que mejor se ajusta a estos datos. Procedemos a hacer el cálculo de este intervalo de confianza mediante los métodos bootstrap:

```

1 #Metodo 1
2 #Limite inferior
3 > quantile(y_i,.0275)
4 2.75%
5 3.317073
6 > #Limite superior
7 > quantile(y_i,.975)
8 97.5%
9 5.865273
10 #Metodo 2
11 #Limite inferior
12 > v_o-quantile(y_i,.975)
13 97.5%
14 3.46806
15 > #Limite superior
16 > v_o-quantile(y_i,.0275)
17 2.75%
18 6.01626

```

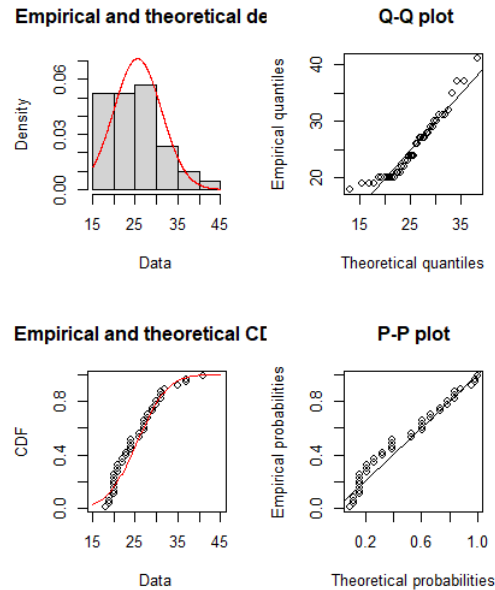
$$IC_{Boot_1}(\sigma^2) = [3,31, 5,86]$$

$$IC_{Boot_2}(\sigma^2) = [3,46, 6,01]$$

c. coeficiente de variación de la edad de la madre.

### Solución:

Veamos que efectivamente, estos datos provienen de una distribución normal



Para las pruebas de bondad de ajuste obtenemos:

```

1 > #Verificamos si los datos provienen de una dist. normal
2 > normal_ = fitdist(x, "norm")
3 > plot(normal_) #Evidencia grafica
4 > prueba<-gofstat(normal_)
5 > prueba$kstest
6      1-mle-norm
7 "not rejected"
8 > prueba$chisqpvalue
9 [1] 0.2328646

```

El p-value calculado nos garantiza que el test no rechaza la hipótesis nula con la significancia dada. Sin embargo, no encontramos una fórmula para este intervalo de confianza, por lo tanto lo estimaremos mediante los métodos bootstrap.

```

1 #Metodo 1
2 > #Limite inferior
3 > quantile(y_i,.0275)
4      2.75%
5 0.1774045
6 > #Limite superior
7 > quantile(y_i,.975)
8      97.5%
9 0.2574538
10 #Metodo 2
11 > #Limite inferior
12 > v_o-quantile(y_i,.975)
13      97.5%
14 0.1861368

```

```

15
16 > #Limite superior
17 > v_o-quantile(y_i,.0275)
18     2.75%
19 0.2661862

```

Con esto calculamos los intervalos de confianza dados por:

$$IC_{Boot_1}(c) = [0,17, 0,25]$$

$$IC_{Boot_2}(c) = [0,18, 0,26]$$

d. proporción de padres fumadores.

### Solución:

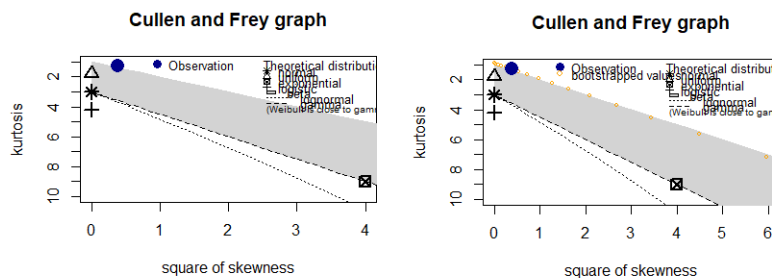
Dado que el data set no contiene explícitamente construiremos esta columna y la añadiremos al dataframe

```

1 fnocig<-as.numeric(data$fnocig>0)
2 data<-cbind(data,fnocig)
3 x=fnocig
4 v_0=length(x[x==1])/42 #Estimacion inicial

```

En este caso, ninguna distribución continua se va ajustar. Es interesante ver como esto representa en el diagrama de Callen



Donde el valor observado sale de la frontera de aceptación y las muestras bootstrap parecen converger puntualmente a un punto de acumulación. Como no proviene de una distribución normal pasamos a hacer uso de los métodos bootstrap.

```

1 #Metodo 1
2 > #Limite inferior
3 > quantile(y_i,.0275)
4 2.75%
5 0.5
6 > #Limite superior
7 > quantile(y_i,.975)
8 97.5%
9 0.7857143

```

```

10 #Metodo 2
11 > #Limite inferior
12 > v_o-quantile(y_i,.975)
13     97.5%
14 0.07893817
15 > #Limite superior
16 > v_o-quantile(y_i,.0275)
17     2.75%
18 0.3646525

```

$$IC_{Boot_1}(c) = [0,5, 0,78]$$

$$IC_{Boot_2}(c) = [0,07, 0,36]$$

**BONUS [+0.2]:** Para las variables de este punto que resultaron no ser normales en la verificación de supuestos, ajuste otra distribución que usted considere y estime los parámetros correspondientes.

**Solución:** Para visualizar el bono, remitirse a los puntos anteriores; en donde se realizó el ajuste de la distribución y la respectiva estimación de los parámetros.

3. Se considera que el peso de un recién nacido sano oscila alrededor de los 3.4 kg. ¿Considera que los datos muestran evidencia para concluir que, en promedio, los bebés de esta población se encuentran sanos?

**Solución:**

Queremos determinar si la media del peso de bebés sanos se aproxima a  $\hat{p}_n$ :

$$\begin{cases} H_0 : \mu = 3,4 \\ H_1 : \mu \neq 3,4 \end{cases}$$

Para ello, recurrimos a la función t-test de R. Por lo que el código queda de la siguiente manera:

```

1 ##### Punto 3 #####
2 ##### Importar Base de Datos #####
3 library(readxl)
4 data <- read_excel("C:/Users/ander/OneDrive - Universidad Nacional de Colombia/
  Documentos/(2022-01) Cuarto Semestre/Inferencia Estadística/Proyecto/Proyecto
  parte B/dataset_trabajo1.xlsx")
5 View(data)
6 ##### Punto 3 #####
7 t.test(x=data$Birthweight, alternative = 'two.sided', mu = 3.4, conf.level = 0.95)

```

Esto nos arroja como resultado lo siguiente:

```

1 One Sample t-test
2
3 data: data$Birthweight
4 t = -0.93518, df = 41, p-value = 0.3552
5 alternative hypothesis: true mean is not equal to 3.4
6 95 percent confidence interval:
7 3.124670 3.501044

```

```

8 sample_estimates:
9 mean of x
10 3.312857

```

Por lo tanto, a un nivel de confianza del 95 % podemos concluir que los datos no muestran evidencia suficiente para rechazar  $H_0$ . Es decir, que la media del peso de los bebés sanos se aproxima a 3.4 kg, con un nivel de confianza del 95 %. Asimismo, teniendo en cuenta que la media de la población del peso de los bebés sanos se encuentra entre  $[3,12kg, 3,50kg]$  a un 95 % de confianza.

4. ¿Es posible concluir que menos del 40 % de madres gestantes son fumadoras?

### Solución:

Queremos determinar si la proporción de madres fumadoras es mayor a 0.4:

$$\begin{cases} H_0 : \hat{\rho}_n \leq 0,4 \\ H_1 : \hat{\rho}_n > 0,4 \end{cases}$$

Para ello, recurrimos a la función t-test de R. Por lo que el código queda de la siguiente manera:

```

1 ##### Punto 4 #####
2 z_critico <- (22/42 - 0.40) / sqrt(0.40 * (1 - 0.40) / 200)
3 z_critico

```

Esto nos arroja como resultado un valor  $z_c = 1,58$ . Con un  $\alpha = 0,05$  y la prueba a una cola tenemos un  $Z = 1,64$ , por lo tanto como  $z_c < Z$ , entonces los datos no evidencian información suficiente para rechazar  $H_0$ .

5. ¿Es plausible afirmar que los bebés de madres no fumadoras tienen, en promedio, una longitud mayor a la de los bebés de madres fumadoras?

### Solución:

Queremos determinar si la media de la longitud de los bebés de las madres no fumadoras es mayor a la de las madres fumadoras. Para ello planteamos la siguiente prueba de hipótesis:

$$\begin{cases} H_0 : \mu_1 - \mu_2 \leq 0 \\ H_1 : \mu_1 - \mu_2 > 0 \end{cases}$$

Teniendo en cuenta que  $\mu_1$  representa al promedio de la longitud de los bebés de las madres no fumadoras y  $\mu_2$  al promedio de la longitud de los bebés de las madres fumadoras.

Por ello, el test realizado en R, queda de la siguiente manera:

```

1 ##### Punto 5 #####
2 x <- data[data$smoker==0, 2]
3 x
4 m_1 <- mean(x$Length)
5 y <- data[data$smoker==1, 2]
6 y
7 m_2 <- mean(y$Length)
8 t.test(x, y, alternative='greater', mu=0, var.equal=TRUE, conf.level=0.95)

```

Al momento de evaluar el test, obtenemos lo siguiente:

```

1 Two Sample t-test
2
3 data:  x and y
4 t = -0.98185, df = 40, p-value = 0.834
5 alternative hypothesis: true difference in means is greater than 0
6 95 percent confidence interval:
7  -2.418802      Inf
8 sample estimates:
9 mean of x mean of y
10 50.90909  51.80000

```

En este caso, asumiendo que las varianzas son iguales, obtenemos un  $t_c = -2,41$ , el cual es menor que el  $t = -0,98$ . Por lo tanto, se concluye que los bebés de las madres no fumadoras, en promedio, tienen una longitud mayor que los bebés de las madres fumadoras (Rechazando  $H_0$ ).

6. ¿Es posible concluir que, en promedio, los padres son mayores en edad a sus correspondientes parejas femeninas?

Dado que la muestras son pareadas trabajemos con el siguiente sistema de hipótesis con cola a izquierda:

**Solución:**

$$\begin{cases} H_0 : \mu_D \leq 0 \\ H_1 : \mu_D > 0 \end{cases}$$

La siguiente rutina nos permite realizar el test para diferencia media poblacional:

```

1 #Punto 4
2 data=read.xlsx("C:/Users/kaido/Downloads/dataset_trabajo1.xlsx",startRow=1,colNames=
   TRUE)
3 #Ver tabla
4 View(data)
5 #Fijar semilla
6 data
7 dif_p_m <- as.numeric( data$fage-data$mage)
8 data <- cbind( data , dif_p_m )
9 x = dif_p_m
10 #Valor critico
11 c=(mean(x)-0)/(sd(x)/sqrt(42))
12 c
13 #Nivel de significancia 95%
14 qnorm(0.975,mean=0,sd=1)

```

Obtenemos como salida:

```

1 >#Valor critico
2 > c=(mean(x)-0)/(sd(x)/sqrt(42))
3 > c
4 [1] 5.359544
5 > #Nivel de significancia 95%

```

```

6 > qnorm(0.05, mean=0, sd=1)
7 [1] -1.644854

```

De manera que como,  $z_c > z_\alpha$  no rechazamos la hipótesis nula. Esto es, tenemos evidencia para pensar que los padres son mayores en edad a sus parejas femeninas.

### Bonus [+0.5]:

**Nota:** Para los puntos 7, 8 y 9 se usó como base el siguiente código, con el fin de importar la base de datos e instalar los paquetes necesarios para la resolución del Bonus. El código en cuestión es el siguiente:

```

1 ##### Bonus Parte B #####
2 ##### Instalacion de Paquetes #####
3 install.packages("readxl")
4 install.packages("dplyr")
5 install.packages("ggplot2")
6 install.packages("GGally")
7 install.packages("Hmisc")
8 install.packages("corrplot")
9 install.packages("PerformanceAnalytics")
10 ##### Carga de Paquetes #####
11 library(readxl)
12 library(dplyr)
13 library(ggplot2)
14 library(GGally)
15 library(Hmisc)
16 library(corrplot)
17 library(PerformanceAnalytics)
18 ##### Importar la Base de Datos #####
19 library(readxl)
20 datos <- read_excel("C:/Users/ander/OneDrive - Universidad Nacional de Colombia/
    Documentos/(2022-01) Cuarto Semestre/Inferencia Estadística/Proyecto/Proyecto
    parte B/dataset_trabajo1.xlsx")
21 View(datos)
22 attach(datos)
23 datos

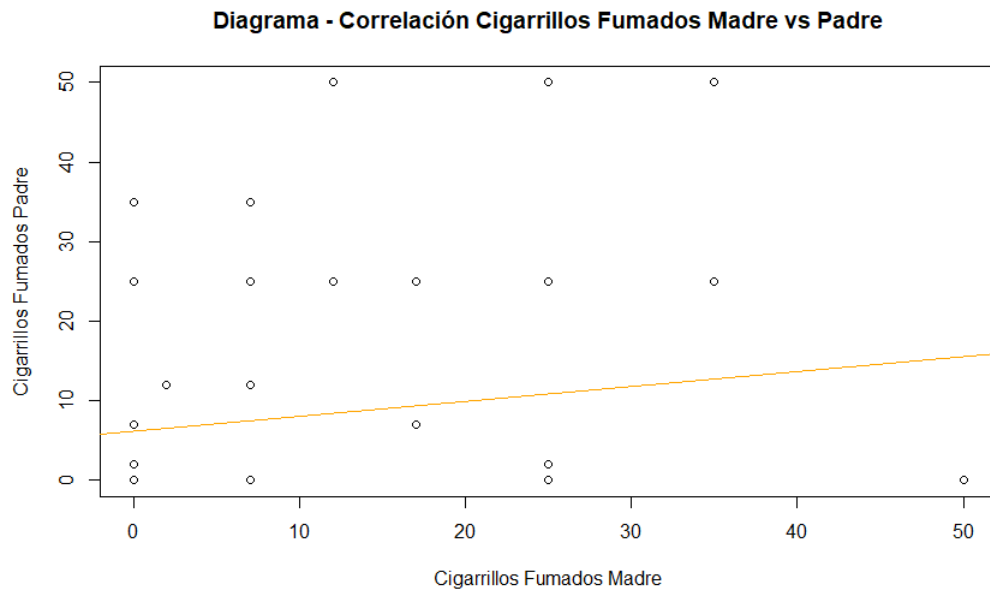
```

Se recomienda, al momento de correr el código, cambiar la ruta de acceso al archivo.

7. ¿Existe una correlación lineal alta entre el consumo de cigarrillo por parte de la madre y del padre? (se considera una correlación alta si supera 0.6). Haga un diagrama de dispersión.

**Solución:** Teniendo en cuenta las variables a considerar. Tanto el consumo de cigarrillo de las madres, como el consumo de cigarrillo de los padres; se puede evidenciar en el siguiente gráfico de dispersión que no hay una relación lineal entre ambas variables. Al calcular numéricamente la correlación se obtiene como resultado 0,26. Lo cual confirma que no existe una relación lineal entre las variables mencionadas. Sin embargo, esto no quiere decir que las variables sean independientes; por lo que solo se puede concluir que no existe una relación lineal entre estas variables.





Asimismo, el código para hallar la correlación entre las variables mencionadas es el siguiente:

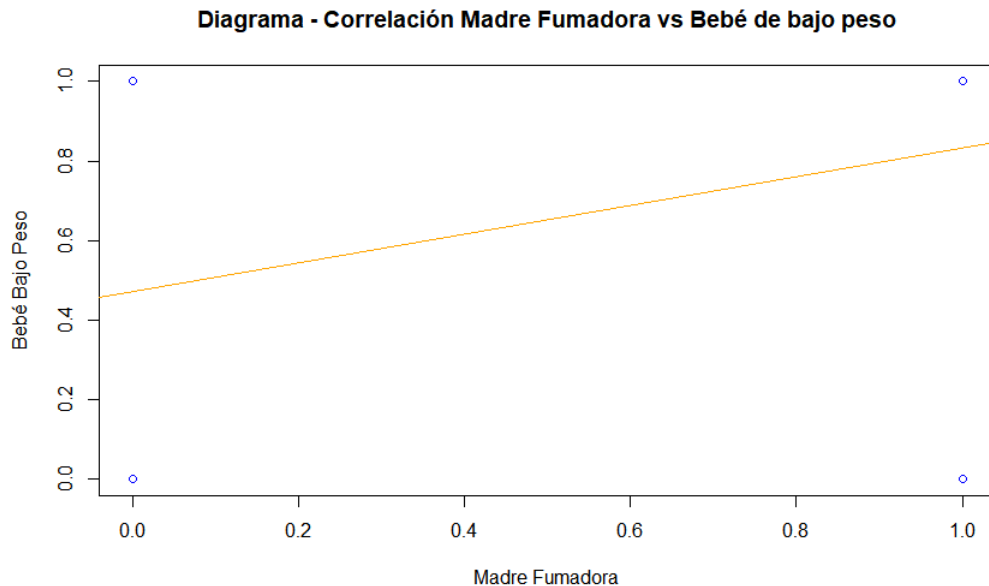
```

1 ##### Punto 7 #####
2 correlacion <- round(cor(datos$mnocig,datos$fnocig),2)
3 correlacion
4 plot(x = datos$mnocig, y = datos$fnocig, col = "blue",main = "Diagrama - Correlacion
   Cigarrillos Fumados Madre vs Padre",xlab="Cigarrillos Fumados Madre",ylab="
   Cigarrillos Fumados Padre")
5 abline(lm(datos$mnocig~datos$fnocig),col = "orange")

```

8. ¿Es posible concluir, estadísticamente hablando, que las variables “madre fumadora” y “bebé de bajo peso” guardan cierta relación o son independientes?

**Solución:** Para responder esta pregunta, hay que tener en cuenta que, dado que las variables en cuestión son variables cualitativas, se debe recurrir al coeficiente de correlación de Spearman. Si queremos ilustrar el diagrama de la correlación, solo nos dará valores de 0 y 1. Tal y como se muestra en la siguiente imagen:



Por ello, para encontrar si las variables "madre fumadora" y "Bebé de bajo peso" tienen cierta relación, solo nos debemos basar en el valor dado en el coeficiente de Spearman; el cual da 0,2530. Por lo tanto se concluye que las variables en cuestión no tienen relación alguna. Sin embargo, esto no garantiza la independencia entre las variables.

Nota: Si las variables son independientes, el coeficiente de correlación es igual a 0. Sin embargo, esta afirmación no es recíproca; dado que si el coeficiente da 0, no garantiza que automáticamente las variables sean independientes.

Finalmente, el código para hallar la correlación entre las variables mencionadas es el siguiente:

```
1 ##### Punto 8 #####
2 correlacion1 <- cor(datos$smoker,datos$lowbwt, method = 'spearman')
3 correlacion1
4 plot(x = datos$smoker, y = datos$lowbwt, col = "blue",main = "Diagrama - Correlacion
   Madre Fumadora vs Bebe de bajo peso",xlab="Madre Fumadora",ylab="Bebe Bajo Peso
   ")
5 abline(lm(datos$smoker~datos$lowbwt),col = "orange")
```

9. ¿Juega la raza un papel en explicar el número de semanas de gestación de las madres embarazadas? Haga un análisis de comparación ANOVA o similar dependiendo del cumplimiento de supuestos.

**Solución:** A continuación, se muestra la tabla que arroja R al momento de realizar un análisis ANOVA entre las variables Raza y Número de Semanas de Gestación.

	Df	Sum Sq	Mean Sq	F Value	Pr(>F)
datos\$Gestation	1	0.04	0.0422	0.038	0.846
Residuals	40	44.43	1.1109		

Tal y como se puede observar, suponiendo un P-Valor de 0,05 se puede interpretar que no existen diferencias estadísticamente significativas entre la raza y las semanas de gestación. Por lo tanto, la raza juega un papel importante en las semanas de gestación de las madres embarazadas.

Finalmente, el código para realizar el análisis ANOVA es el siguiente:

```
1 ##### Punto 9 #####
2 anova <- aov(datos$race ~ datos$Gestation)
3 summary(anova)
```