

Ejemplo clase

```
knitr::opts_chunk$set(warning = FALSE, message = FALSE)
library(tidyverse)
library(tidymodels)
```

Sección 2) Ejercicio

Introducción

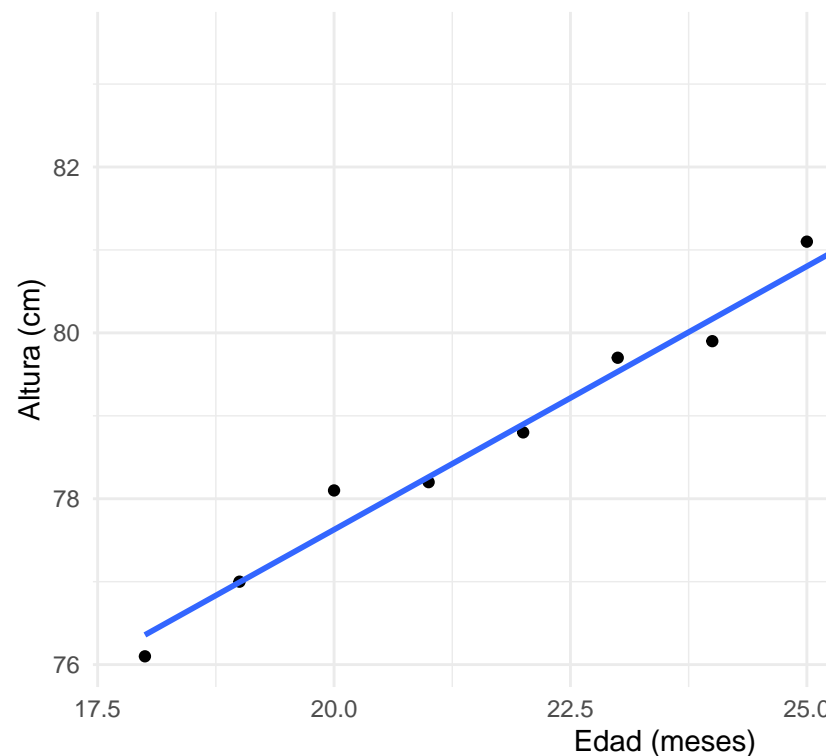
El problema presentado en este ejercicio es determinar si existe una relación lineal entre la edad de un niño y su altura. Se tiene la intuición que a mayor edad, más alto es. A continuación se presentan los datos de edad (meses) y altura (cm) en una muestra de 12 niños

```
data <- read_csv("https://raw.githubusercontent.com/savrgg/class_ITAM_metodos/main/notas_r/ageandheight")
data
```

```
## # A tibble: 12 x 2
##   age height
##   <dbl> <dbl>
## 1    18  76.1
## 2    19   77
## 3    20  78.1
## 4    21  78.2
## 5    22  78.8
## 6    23  79.7
## 7    24  79.9
## 8    25  81.1
## 9    26  81.2
## 10   27  81.8
## 11   28  82.8
## 12   29  83.5
```

```
data %>%
  ggplot(aes(x = age, y = height))+
  geom_point()+
  geom_smooth(method = "lm", se = F)+
  theme_minimal()+
  labs(x = "Edad (meses)",
       y = "Altura (cm)",
       title = "Relación lineal entre altura y edad de niños",
       subtitle = "Muestra de la altura y edad de 12 niños")
```

1. Determine por medio de una gráfica si es razonable pensar en que existe una relación lineal en-
Relación lineal entre altura y edad de niños
Muestra de la altura y edad de 12 niños



entre las variables (variable independiente: edad)

Al comparar gráficamente los datos, podemos observar que los datos si se ajustan por medio de una línea.

2. Determine con una Prueba de Hipótesis si existe una relación lineal entre las variables. Para determinar si existe una relación lineal, se construye una Prueba de Hipótesis del coeficiente de correlación. Recordando, el coeficiente de correlación mide solamente la asociación lineal, mas no la pendiente de la línea

Sea $H_0 : \rho = 0$ y $H_1 : \rho \neq 0$, entonces:

```
n = nrow(data)
ybar = mean(data$height)
xbar = mean(data$age)
sxy = cov(data$age, data$height)
sx2 = var(data$age)
sy2 = var(data$height)
rxy = sxy/(sqrt(sx2)*sqrt(sy2))
T_corr = (rxy*sqrt(n-2))/sqrt(1-rxy**2)
T_corr
```

```
## [1] 29.66465
```

```
2*(1-pt(T_corr, df = n-2))
```

```
## [1] 4.428058e-11
```

El *valor - p* es cercano a 0, entonces rechazamos $H_0 : \rho = 0$, por lo que rechazamos que la correlación sea cero.

```
lm_fit <-
  linear_reg() %>%
  fit(height ~ age, data = data)

tidy(lm_fit)
```

3. Utilicé tidymodels (lm) para realizar la regresión lineal, analice los resultados obtenidos

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>         <dbl>     <dbl>    <dbl>    <dbl>
## 1 (Intercept)   64.9       0.508    128.    2.13e-17
## 2 age           0.635     0.0214    29.7  4.43e-11

data.frame(glance(lm_fit))

##   r.squared adj.r.squared   sigma statistic    p.value df    logLik
## 1 0.9887639   0.9876403 0.2559638  879.9915 4.428071e-11  1 0.4192965
##       AIC      BIC deviance df.residual nobs
## 1 5.161407 6.616127 0.6551748         10    12
```

```
B1 <- sxy/sx2
B0 <- ybar - B1*xbar
B1
```

4. Determine analíticamente los coeficientes $\hat{\beta}_0$ y $\hat{\beta}_1$

```
## [1] 0.634965
B0
```

```
## [1] 64.92832
```

5. Calcule analíticamente TSS, ESS, RSS, R^2 y $adjR^2$

```
yhat <- predict(lm_fit, new_data = data %>% select(age))
TSS <- sum((data$height-ybar)**2)
ESS <- sum((yhat-ybar)**2)
RSS <- sum((data$height-yhat)**2)
R2 <- ESS/TSS
adjR2 <- 1- ((1-R2)*(n-1))/(n-1-1)
```

6. Calcule analíticamente $\hat{\sigma}^2$, $V(\hat{\beta}_0)$, $V(\hat{\beta}_1)$

```
sigma2hat <- RSS/(n-2)
sigmaahat <- sqrt(sigma2hat)

sigma2_B0 <- sigma2hat*(1/n + (xbar*xbar)/(sx2*(n-1)))
sigma2_B1 <- sigma2hat/(sx2*(n-1))

sigma_B0 <- sqrt(sigma2_B0)
sigma_B1 <- sqrt(sigma2_B1)

sigma_B0
```

```
## [1] 0.5084102
```

```
sigma_B1
```

```
## [1] 0.02140477
```

7. Determine si el valor de los coeficientes es significativo (realice una prueba para cada uno)

```
T_B1 = (B1-0)/sqrt(sigma2_B1)
```

```
T_B0 = (B0-0)/sqrt(sigma2_B0)
```

```
T_B1
```

```
## [1] 29.66465
```

```
T_B0
```

```
## [1] 127.7085
```

```
2*(1-pt(T_B1, df = n-2))
```

```
## [1] 4.428058e-11
```

```
2*(1-pt(T_B0, df = n-2))
```

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

El valor-p de ambas pruebas es cercano a 0, por lo que podemos rechazar que $H_0 : \beta_1 = 0$ y $H_0 : \beta_0 = 0$

8. Comprueba que:

a) Suma de residuales es igual a 0

```
sum(data$height-yhat)
```

```
## [1] -2.842171e-13
```

b) Residuos no están correlacionados a X

```
cor(data$age, data$height-yhat)
```

```
## .pred
```

```
## [1,] 4.694673e-14
```

c) Residuos no están correlacionados a \hat{Y}

```
cor(yhat, data$height-yhat)
```

```
## .pred
```

```
## .pred 4.75625e-14
```

9. Determine si el de R^2 es significativo

```
Fest = (ESS/1)/(RSS/(n-2))
```

```
Fest
```

```
## [1] 879.9915
```

```
pf(Fest, df1 = 1, df2 = (n-2), lower.tail = F)
```

```
## [1] 4.428071e-11
```

El *valor - p* es cercano a 0, por lo que podemos rechazar $H_0 : \rho^2 = 0$ vs $H_1 : \rho^2 > 0$.

Calcula la predicción a la media para una edad 25 años. Calcule el IC al 95%.

```
yhat_25 <- B0+B1*25
quantil_t = qt(.975, df = n-2)
lim_inf = yhat_25 - quantil_t*sigmahat*sqrt(1/n + ((25-xbar)^2)/(sx2*(n-1)))
lim_sup = yhat_25 + quantil_t*sigmahat*sqrt(1/n + ((25-xbar)^2)/(sx2*(n-1)))
```

IC para predicción media de edad 25: (80.62294, 80.98196)

Calcula la predicción individual para un valor de 30 años. Calcule el IC al 95

```
yhat_30 <- B0+B1*30
quantil_t = qt(.975, df = n-2)
lim_inf = yhat_30 - quantil_t*sigmahat*sqrt(1/n + ((30-xbar)^2)/(sx2*(n-1)))
lim_sup = yhat_30 + quantil_t*sigmahat*sqrt(1/n + ((30-xbar)^2)/(sx2*(n-1)))
```

IC para predicción individual de edad 30: (83.62626, 84.32828)