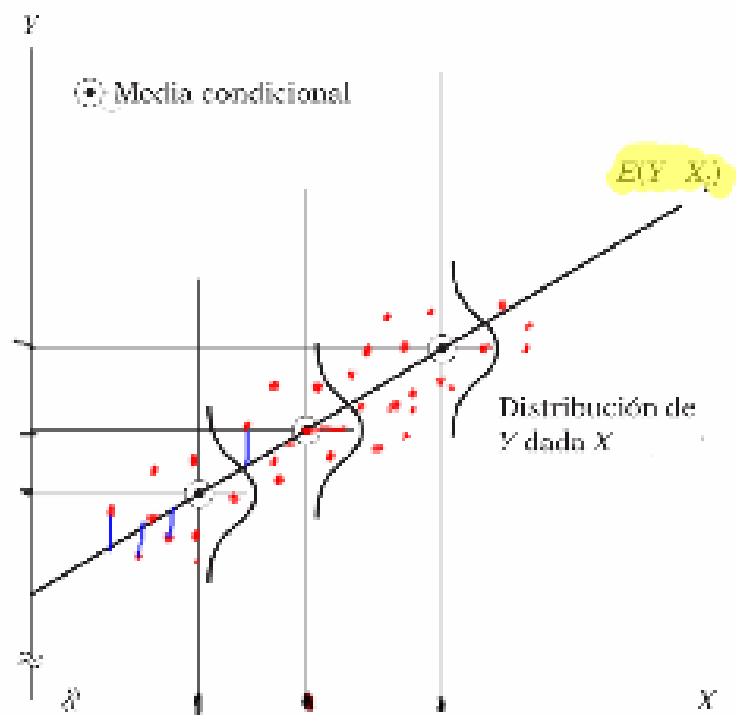


Regresión Lineal

$$Y = f(X)$$

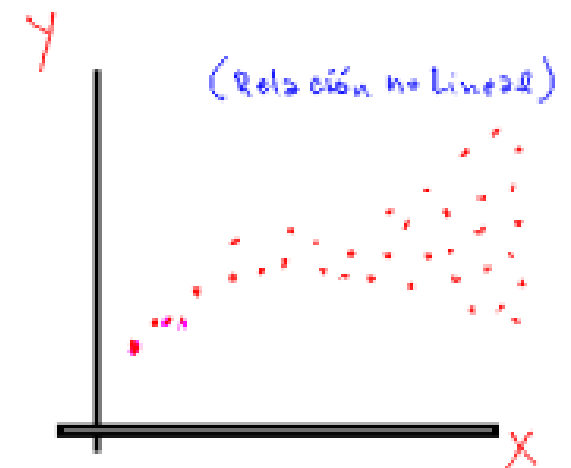
Y : Variable respuesta, predicha, dependiente. Target. → cuantitativa

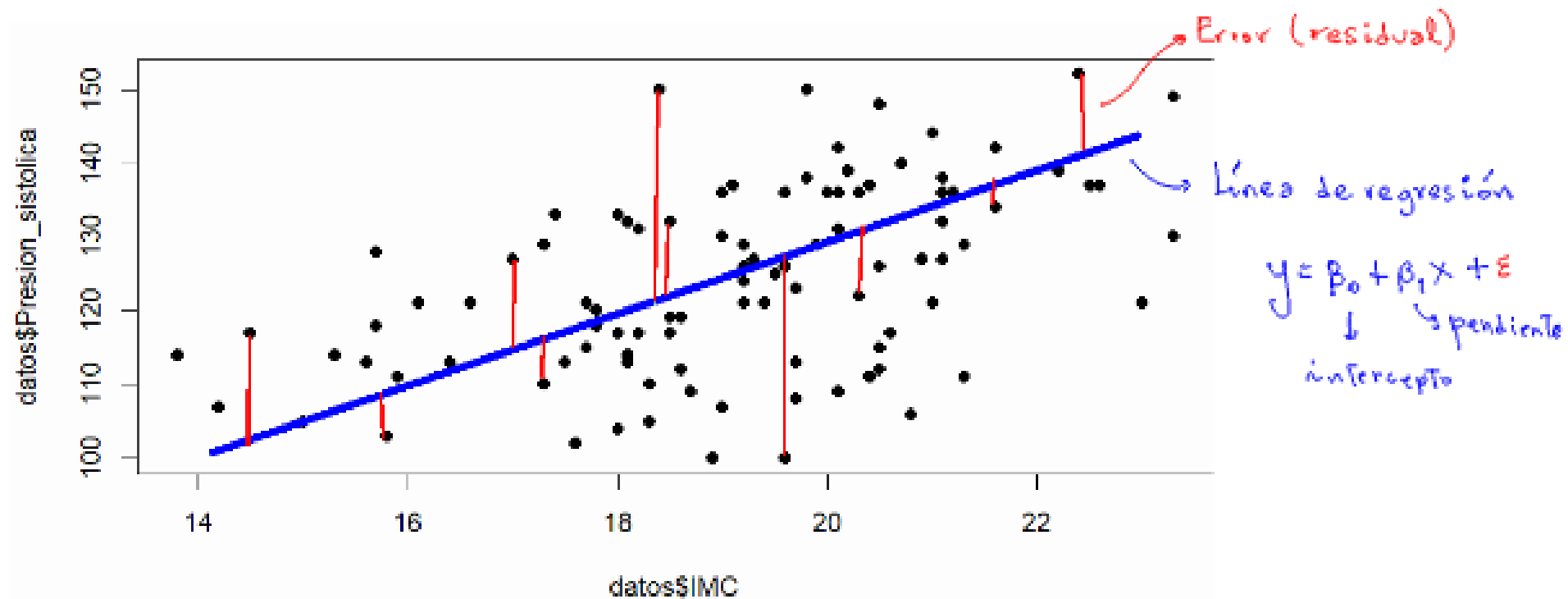
X : Variable explicativa, predictora, independiente. Feature o atributo. → cuantitativa o cualitativa



$$\mu_{Y|X=1}, \mu_{Y|X=3}$$

$$Y|X \sim N$$





$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i = \mu_i + \epsilon_i \quad i = 1, \dots, 100$$

Presión sistólica \nearrow
 IMC \nearrow
 intercepto \nwarrow β_0
 pendiente \nwarrow β_1

$$\mu_i = \beta_0 + \beta_1 X_i \rightarrow \text{poblacional } (\beta_0, \beta_1 \text{ son parámetros})$$

$$\hat{\mu}_i = \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i \rightarrow \text{muestreal } (\hat{\beta}_0, \hat{\beta}_1 \text{ son estimadores})$$

$$\hat{\mu}_i = \hat{Y}_i = 64.5577 + 3.0957 X_i$$

$\hat{\beta}_0$: Cuando $X_i = 0 \Rightarrow \hat{\mu}_i = 64.5577 \rightarrow$ La presión sistólica media cuando IMC = 0. **No tiene sentido**

$\hat{\beta}_1$: $\hat{Y}_i = 64.5577 + 3.0957 X_i$

$X_i \rightarrow X_{i+1}$

$\hat{Y}_i^* = 64.5577 + 3.0957 (X_i + 1)$

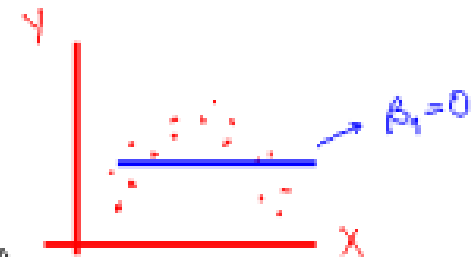
$\hat{Y}_i^* - \hat{Y}_i = 3.0957 = \hat{\beta}_1$

$$\left(\bar{X} - t_{(1-\alpha/2; n-1)} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{(1-\alpha/2; n-1)} \frac{s}{\sqrt{n}} \right)$$

$$IC(\beta_j) = \hat{\beta}_j \pm t_{1-\alpha/2, n-1} s_{\hat{\beta}_j} \rightarrow \mathcal{R}$$

$$Y = \beta_0 + \beta_1 X + \epsilon$$

$$\mu = \beta_0 + \beta_1 X \begin{cases} \beta_1 = 0 \rightarrow X \text{ e } Y \text{ no tienen relación lineal} \\ \beta_1 \neq 0 \rightarrow X \text{ e } Y \text{ sí tienen relación lineal} \end{cases}$$



```

> modelo |> aov() |> summary()

```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
IMC	1	3924	3924	31.86	1.61e-07 ***
Residuals	98	12068	123		

 signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Handwritten annotations:
 - Regression line points to the model command.
 - GL Reg points to the Df=1 row.
 - SC Reg points to the Sum Sq=3924 row.
 - CM Reg points to the Mean Sq=3924 row.
 - F calc points to the F value=31.86 row.
 - GLE points to the Df=98 row.
 - SCE points to the Sum Sq=12068 row.
 - CME points to the Mean Sq=123 row.

$$p\text{valor} = 1.61 \times 10^{-7} = 0.000000161 < \alpha$$

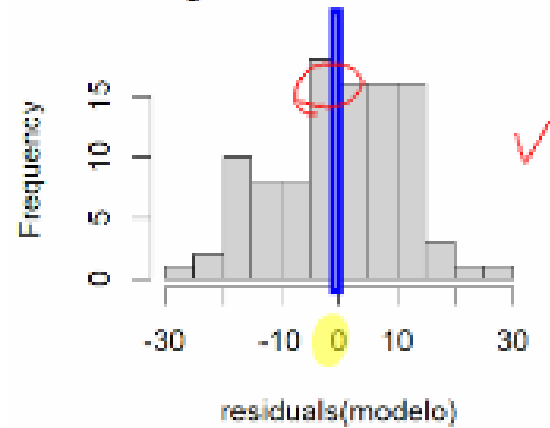
\Downarrow
 Se rechaza H_0
 \Downarrow
 X e Y sí tienen relación lineal

$$Y = \overbrace{\beta_0 + \beta_1 X}^{\text{REGRESIÓN}} + (\varepsilon) \rightarrow \text{Error recoge los efectos de otras variables no observadas}$$

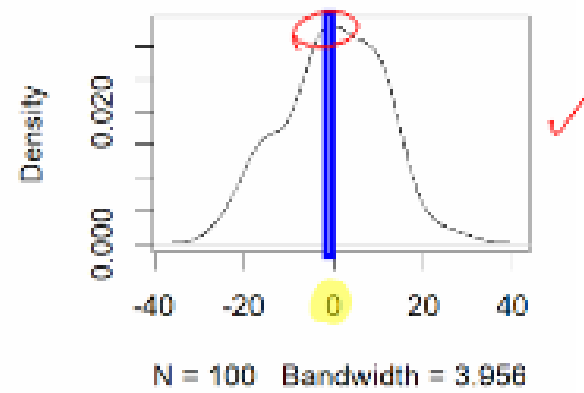
$$E(\varepsilon) = \mu_{\varepsilon} = 0$$

$$V(\varepsilon) = \sigma_{\varepsilon}^2 \text{ constante}$$

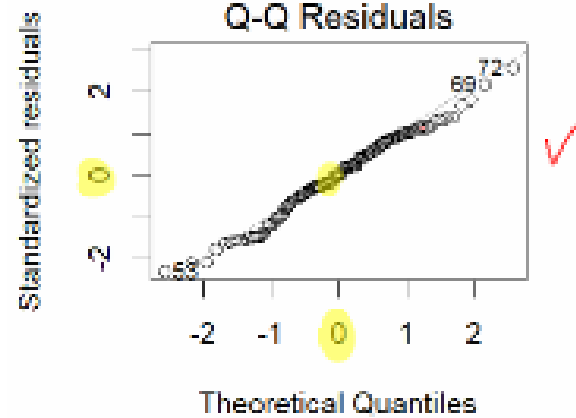
Histograma de los residuales



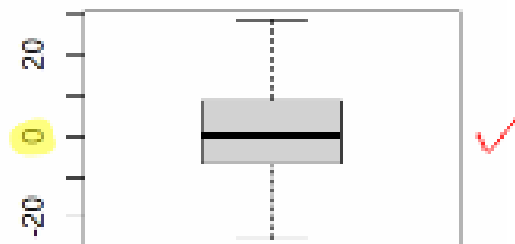
Densidad de los residuales



Q-Q Residuals



Boxplot de los residuales



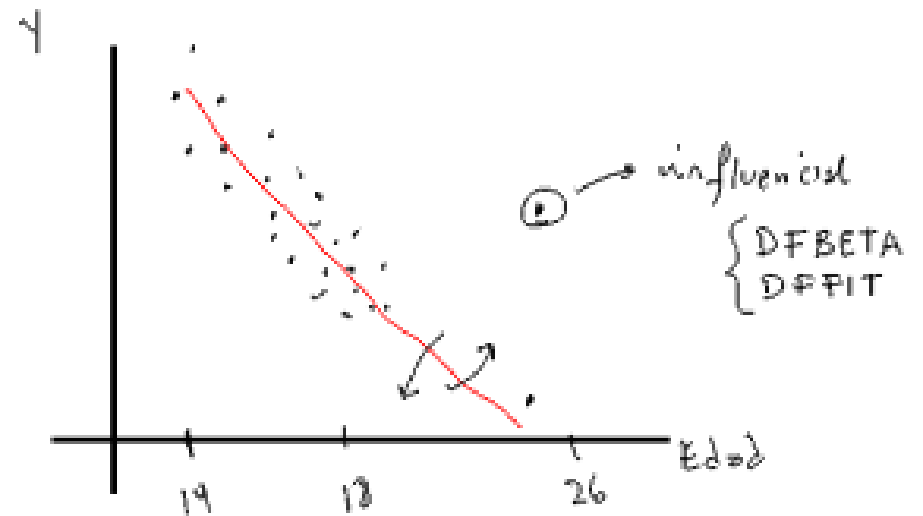
$$E(\varepsilon) = 0$$

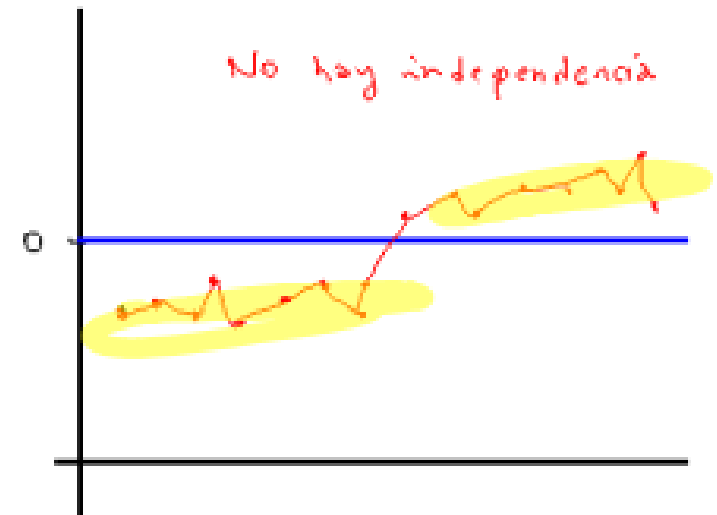
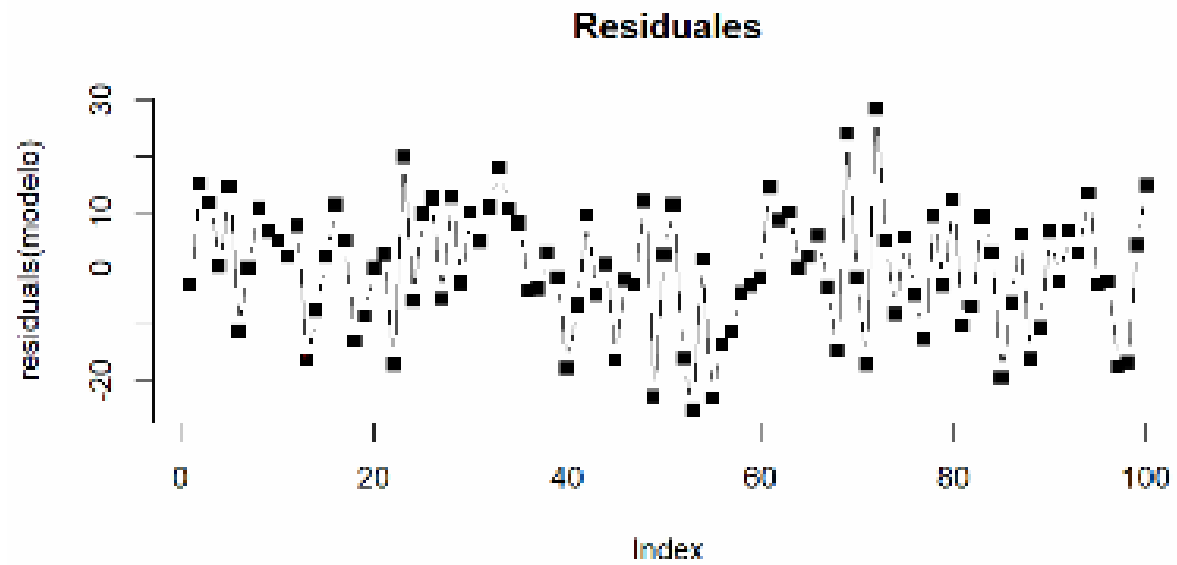
$$E(y - \hat{y}) = 0$$

$$E(y - \hat{\beta}_0 - \hat{\beta}_1 x) = 0$$

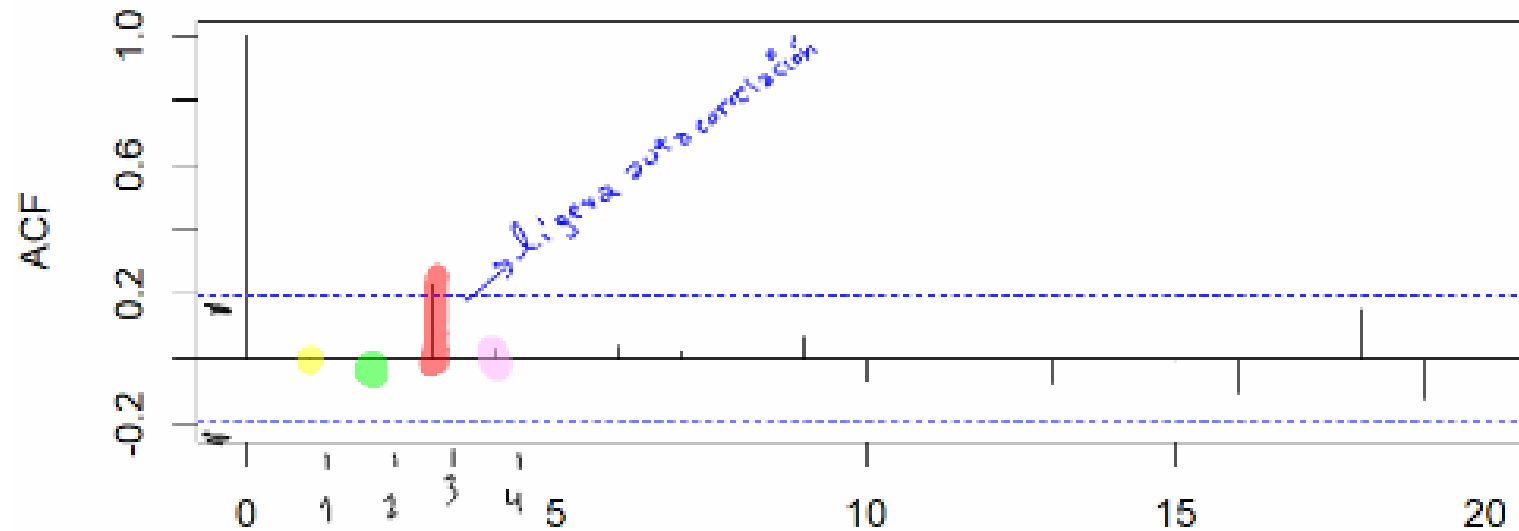
\downarrow \downarrow

$$E(\dots) = 0$$





Series residuals(modelo)



$\hat{\rho}_1 = \text{Cor}(\varepsilon_t, \varepsilon_{t-1})$
 $\hat{\rho}_2 = \text{Cor}(\varepsilon_t, \varepsilon_{t-2})$
 $\hat{\rho}_3 = \text{Cor}(\varepsilon_t, \varepsilon_{t-3})$

Lag

```
> modelo |> durbinwatsonTest(max.lag = 5)
lag Autocorrelation D-W Statistic p-value
1 0.005751365 1.969609 0.934
2 -0.036473991 2.033392 0.790
3 0.228924274 1.467707 0.014
4 0.026582705 1.846881 0.648
5 -0.086195271 2.054406 0.506
Alternative hypothesis: rho[lag] != 0
```

?



Estimación de la media
Predicción individual

Puntuación

$$\hat{\mu} = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\tilde{y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

Intervalo

→ Intervalo de confianza : $IC(\mu)$

→ Intervalo de predicción : $IP(y)$ → más ancho

```
> modelo |> predict(data.frame(IMC = c(18,20,23)))
      1      2      3
120.2800 126.4714 135.7584
> modelo |> predict(data.frame(IMC = c(18,20,23)),
+ level = 0.95, interval = "confidence")
```

	fit	lwr	upr
1	120.2800	117.7909	122.7691
2	126.4714	124.0459	128.8968
3	135.7584	130.9439	140.5729

$IC(\mu)$

```
> modelo |> predict(data.frame(IMC = c(18,20,23)))
      1      2      3
120.2800 126.4714 135.7584
> modelo |> predict(data.frame(IMC = c(18,20,23)),
+ level = 0.95, interval = "prediction")
```

	fit	lwr	upr
1	120.2800	98.11787	142.4421
2	126.4714	104.31629	148.6264
3	135.7584	113.21636	158.3005

$IP(y)$

MÁS ANCHO

$$\hat{V}(\hat{\mu}) = \hat{V}(\hat{y}) = \hat{\sigma}^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{SCX} \right)$$

$$\hat{V}(\hat{y}_0) = \hat{\sigma}^2 \left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{SCX} \right)$$

$$V(\hat{\mu}) < V(\hat{y}_0)$$