

# Apunte de Regresión Lineal

María Eugenia Szretter Noste  
Carrera de Especialización en Estadística  
para Ciencias de la Salud  
Facultad de Ciencias Exactas y Naturales,  
Universidad de Buenos Aires

Agosto - Octubre de 2017

Tabla 1: Observaciones a nuestra disposición. Aquí  $X_1$  quiere decir, la variable  $X$  medida en el individuo 1, etc.

Individuo	Variable $X$	Variable $Y$
1	$X_1$	$Y_1$
2	$X_2$	$Y_2$
$\vdots$	$\vdots$	$\vdots$
$n$	$X_n$	$Y_n$

En estas notas, estamos pensando en que medimos ambas variables en la misma unidad: puede tratarse de un individuo, un país, un animal, una escuela, etc. Comencemos con un ejemplo.

duo. El más utilizado de todos es el que se conoce como *coeficiente de correlación*, que se simboliza con una letra griega *rho*:  $\rho$  ó  $\rho_{XY}$  y se define por

$$\begin{aligned}\rho_{XY} &= E \left[ \left( \frac{X - \mu_X}{\sigma_X} \right) \left( \frac{Y - \mu_Y}{\sigma_Y} \right) \right] \\ &= \frac{cov(X, Y)}{\sigma_X \sigma_Y},\end{aligned}$$

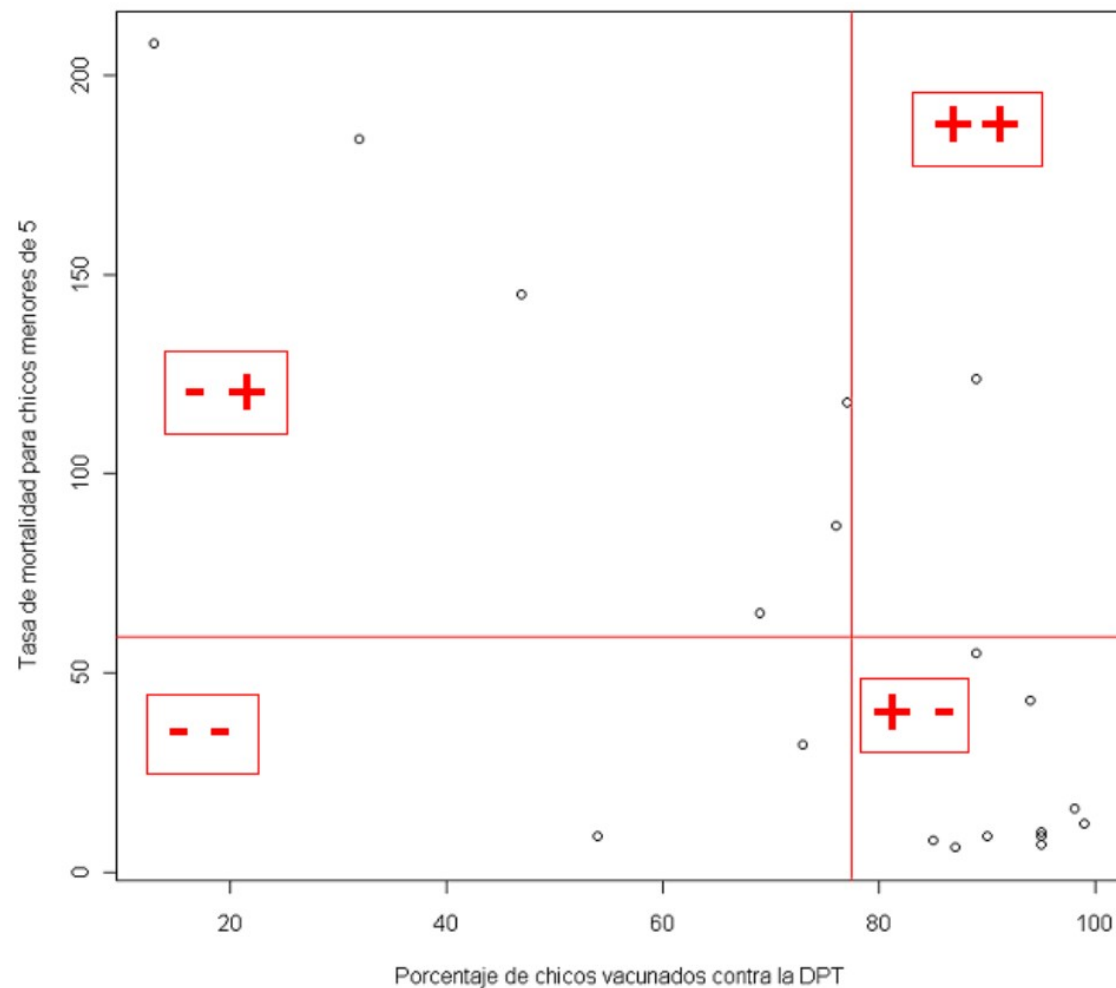
o sea, el número promedio a nivel población del pro

$$\hat{\mu}_X = \overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i,$$

ducto de los desvíos estándares. Como  
por  $\hat{\sigma}_X^2 = S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X})^2$ , de correlación de Pearson, o coeficiente  
esti  
de

$$r = \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X}) (Y_i - \overline{Y})}{S_X \cdot S_Y}.$$

Scatter plot de tasa de mortalidad versus  
porcentaje inmunizado



**Test para  $\rho = 0$**  Los supuestos para llevar a cabo el test son que los pares de observaciones  $(X_1, Y_1), \dots, (X_n, Y_n)$  sean independientes entre sí, idénticamente distribuidos, y tengan distribución (conjunta) normal bivariada (ver la definición de esto en la Observación 1.1). En particular, esto implica que cada una de las muestras  $X_1, \dots, X_n$  e  $Y_1, \dots, Y_n$  tengan distribución normal. Si la hipótesis nula es verdadera, entonces el estadístico

$$T = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

que no es más que  $\hat{\rho}$  dividido por un estimador de su desvío estándar, tiene distribución  $t$  de Student con  $n - 2$  grados de libertad, lo cual notaremos

$$T \sim t_{n-2} \text{ bajo } H_0.$$

### 1.3. Coeficiente de correlación de Spearman

1. Se ordena cada muestra por separado, de menor a mayor. A cada observación se le calcula el ranking que tiene (o rango, o número de observación en la muestra ordenada). De este modo, la observación más pequeña de las  $X$ 's recibe el número 1 como rango, la segunda recibe el número 2, etcétera, la más grande de todas las  $X$ 's recibirá el rango  $n$ . Si hubiera dos o más observaciones empatadas en algún puesto (por ejemplo, si las dos observaciones más pequeñas tomaran el mismo valor de  $X$ , entonces se promedian los rangos que les tocarían: cada una tendrá rango 1,5, en este ejemplo, ya que  $\frac{1+2}{2} = 1,5$ . En el caso en el que las tres primeras observaciones fueran empatadas, a las tres les tocaría el promedio entre 1, 2 y 3, que resultará ser  $\frac{1+2+3}{3} = 2$ ). A este proceso se lo denomina *ranquear las observaciones  $X$* . Llamemos  $R(X_i)$  al rango obtenido por la  $i$ -ésima observación  $X$ .
2. Se reemplaza a cada observación  $X_i$  por su rango  $R(X_i)$ .
3. Se ranquean las observaciones  $Y$ , obteniéndose  $R(Y_i)$  de la misma forma en que se hizo en el ítem 1 para las  $X$ 's.
4. Se reemplaza a cada observación  $Y_i$  por su rango  $R(Y_i)$ . Observemos que conocemos la suma de todos los rangos de ambas muestras (es la suma de  $1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2}$ ).
5. Se calcula el coeficiente de correlación de Pearson entre los pares  $(R(X_i), R(Y_i))$ . El valor obtenido es el coeficiente de correlación de Spearman, que denotaremos  $r_S$ .

## 2.2. Modelo lineal simple

El modelo de regresión lineal es un modelo para el vínculo de dos variables aleatorias que denominaremos  $X = \textit{variable predictora o covariable}$  e  $Y = \textit{variable dependiente o de respuesta}$ . El modelo lineal (simple pues sólo vincula una variable predictora con  $Y$ ) propone que

$$Y = \beta_0 + \beta_1 X + \varepsilon, \tag{2}$$

donde  $\varepsilon$  es el término del error. Esto es que para cada valor de  $X$ , la correspondiente observación  $Y$  consiste en el valor  $\beta_0 + \beta_1 X$  más una cantidad  $\varepsilon$ , que puede ser positiva o negativa, y que da cuenta de que la relación entre  $X$  e  $Y$  no es exactamente lineal, sino que está expuesta a variaciones individuales que hacen que el

1. La esperanza condicional de  $Y$  depende de  $X$  de manera lineal, es decir

$$E(Y | X) = \beta_0 + \beta_1 X \quad (6)$$

o, escrito de otro modo

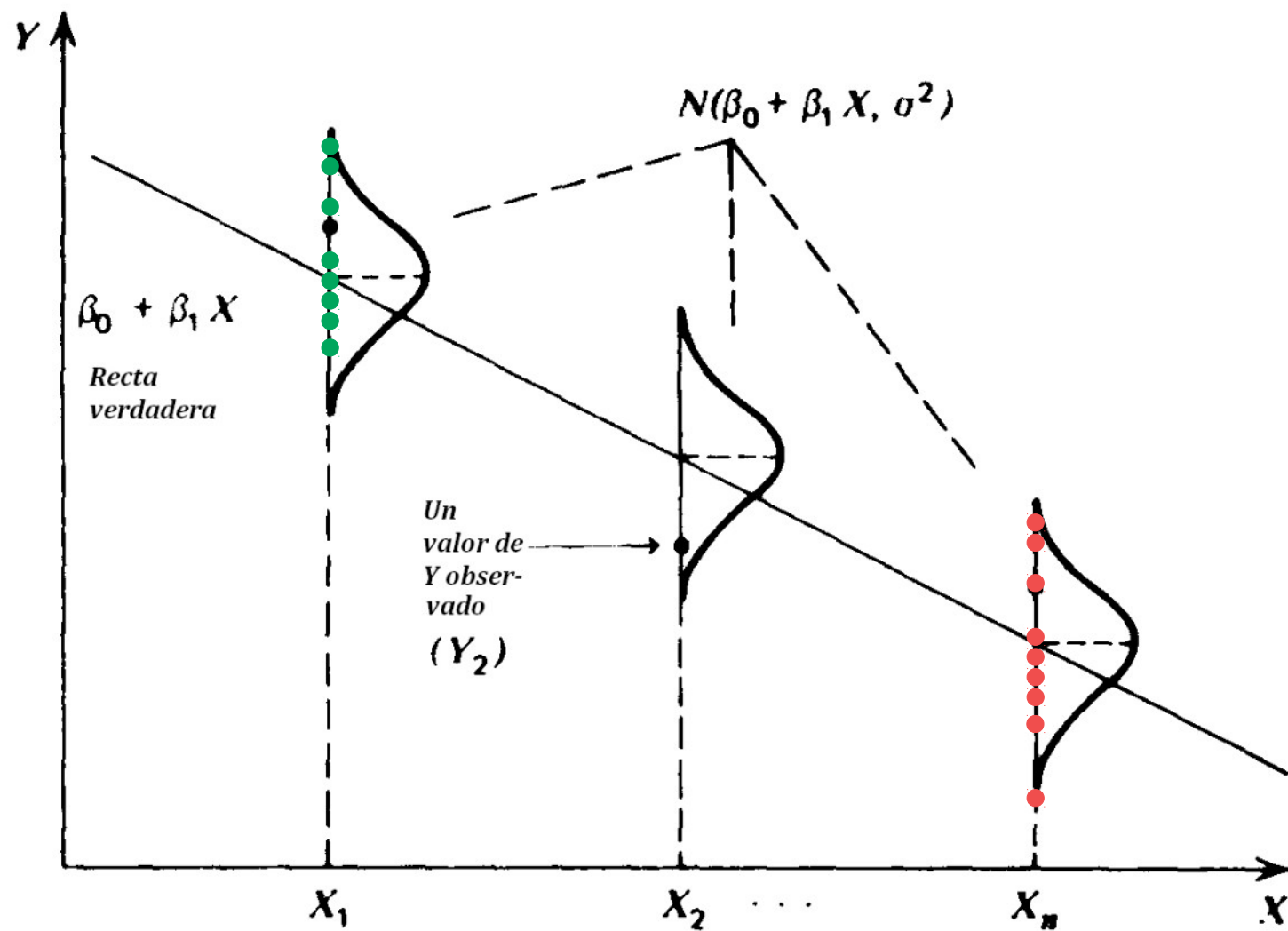
$$E(Y | X = x_i) = \beta_0 + \beta_1 x_i \quad (7)$$

donde  $\beta_0, \beta_1$  son los parámetros del modelo, o coeficientes de la ecuación. A la ecuación (6) se la suele llamar **función de respuesta**, es una recta.

2. La varianza de la variable respuesta  $Y$  dado que la predictora está fijada en  $X = x$  la denotaremos por  $Var(Y | X = x)$ . Asumimos que satisface

$$Var(Y | X = x_i) = \sigma^2,$$





$$Y = \beta_0 + \beta_1 X + \varepsilon,$$

## 2.5. Estimación de los parámetros $\beta_0$ y $\beta_1$

Los coeficientes del modelo se estiman a partir de la muestra aleatoria de  $n$  observaciones  $(X_i, Y_i)$  con  $1 \leq i \leq n$ . Llamaremos  $\hat{\beta}_0$  y  $\hat{\beta}_1$  a los estimadores de  $\beta_0$  y  $\beta_1$ . Los valores  $\hat{\beta}_0$  y  $\hat{\beta}_1$  corresponderán a la recta de ordenada al origen  $\hat{\beta}_0$  y pendiente  $\hat{\beta}_1$  que “mejor ajuste” a los datos  $(X_1, Y_1), \dots, (X_n, Y_n)$  observados. Para encontrarlos, debemos dar una noción de bondad de ajuste de una recta cualquiera con ordenada al origen  $a$  y pendiente  $b$  a nuestros datos. Tomemos las distancias verticales entre los puntos observados  $(X_i, Y_i)$  y los puntos que están sobre la recta  $y = a + bx$ , que están dados por los pares  $(X_i, a + bX_i)$ . La distancia entre ambos es  $Y_i - (a + bX_i)$ . Tomamos como función que mide el desajuste de la recta a los datos a

$$g(a, b) = \sum_{i=1}^n (Y_i - (a + bX_i))^2, \quad (8)$$

$$\frac{\partial g(a, b)}{\partial a} = \sum_{i=1}^n 2(Y_i - (a + bX_i))(-1)$$

$$\frac{\partial g(a, b)}{\partial b} = \sum_{i=1}^n 2(Y_i - (a + bX_i))(-X_i)$$

Las igualamos a cero para encontrar  $\hat{\beta}_0$  y  $\hat{\beta}_1$ , sus puntos críticos. Obtenemos

$$\sum_{i=1}^n \left( Y_i - \left( \hat{\beta}_0 + \hat{\beta}_1 X_i \right) \right) = 0 \quad (9)$$

$$\sum_{i=1}^n \left( Y_i - \left( \hat{\beta}_0 + \hat{\beta}_1 X_i \right) \right) X_i = 0. \quad (10)$$

Las dos ecuaciones anteriores se denominan las *ecuaciones normales* para regresión lineal. Despejamos de ellas las estimaciones de los parámetros que resultan ser

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}, \quad (11)$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}. \quad (12)$$

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + e_i,$$

Modelo ajustado.

```
> ajuste<-lm(headcirc~gestage)
```

```
> summary(ajuste)
```

Coefficients:

	Estimate	Std. Error
(Intercept)	3.91426	1.82915
gestage	0.78005	0.06307

---

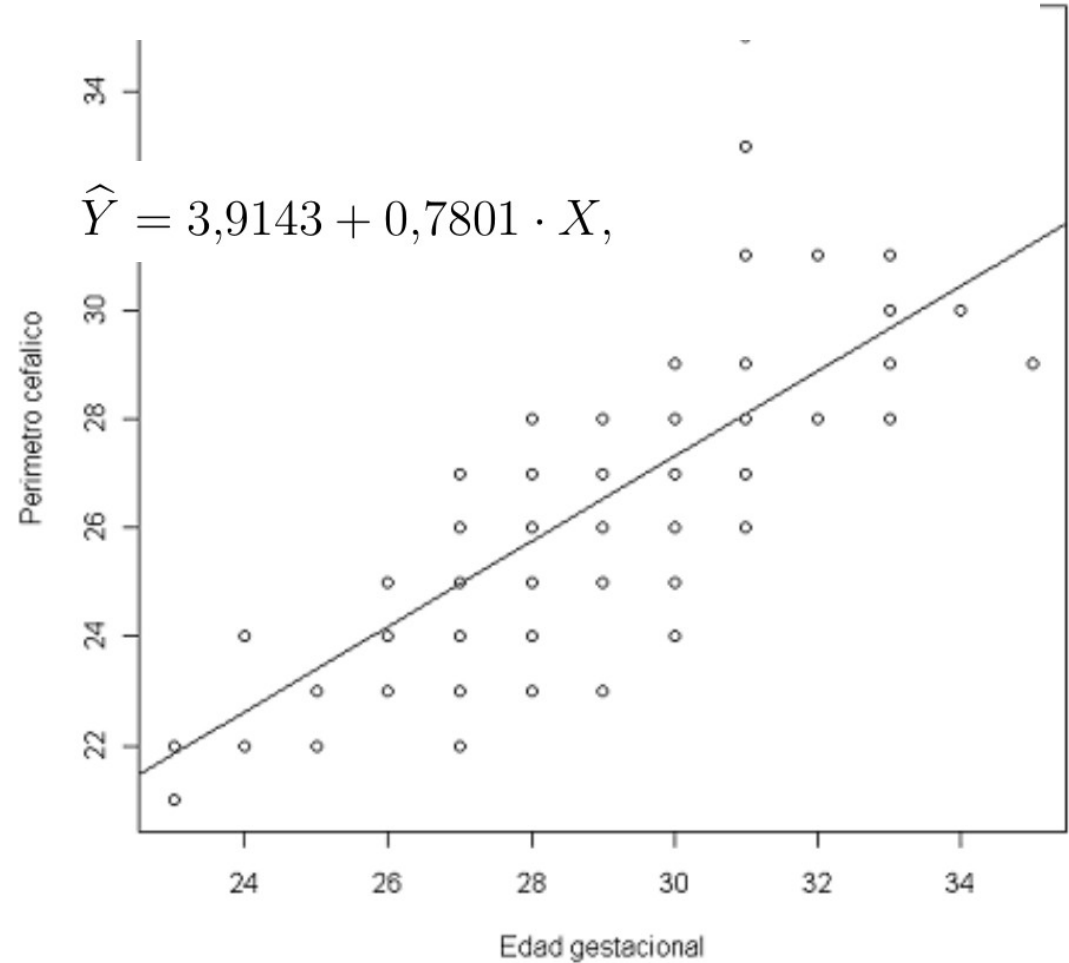
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.59 on

Multiple R-squared: 0.6095,

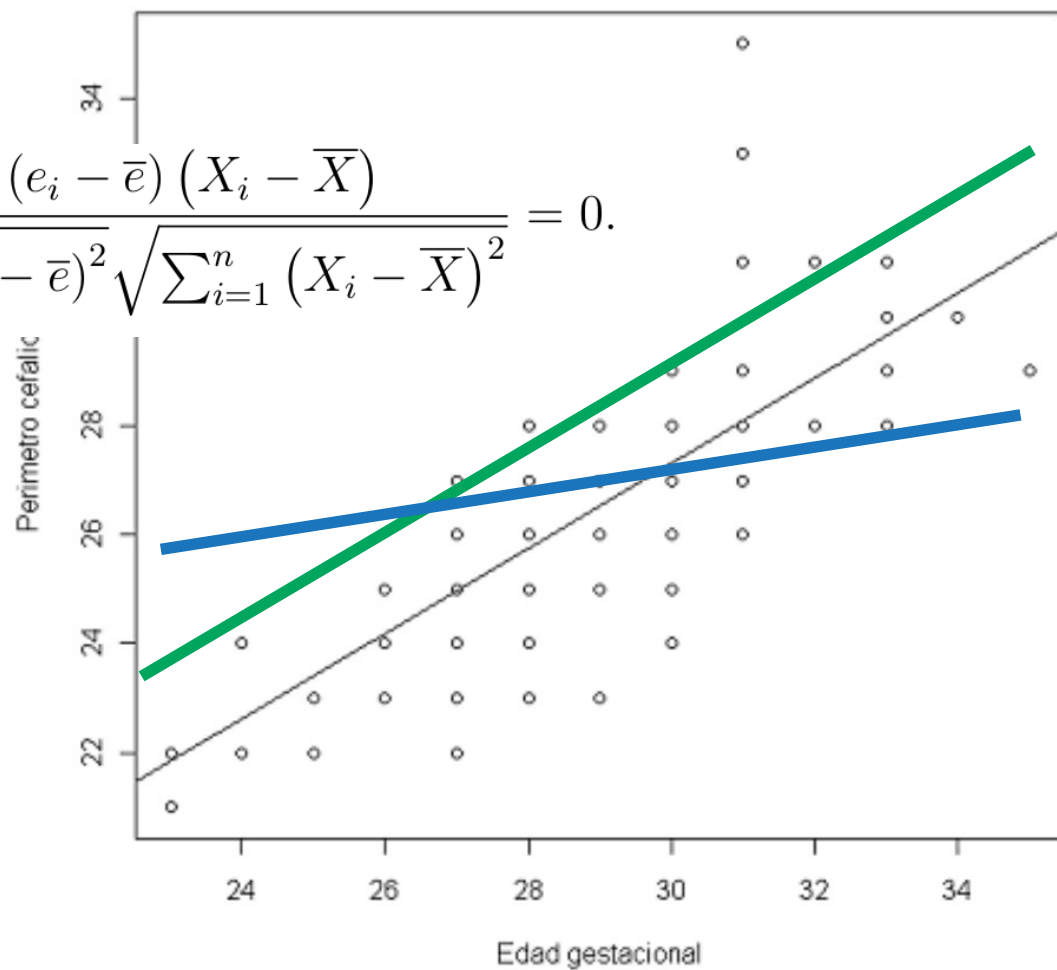
F-statistic: 152.9 on 1 and 98 Df

$$0 = \sum_{i=1}^n \left( Y_i - \left( \hat{\beta}_0 + \hat{\beta}_1 X_i \right) \right) = \sum_{i=1}^n e_i.$$



$$0 = \sum_{i=1}^n \left( Y_i - \left( \hat{\beta}_0 + \hat{\beta}_1 X_i \right) \right) = \sum_{i=1}^n e_i.$$

$$r = r((X_1, e_1), \dots, (X_n, e_n)) = \frac{\sum_{i=1}^n (e_i - \bar{e})(X_i - \bar{X})}{\sqrt{\sum_{i=1}^n (e_i - \bar{e})^2} \sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}} = 0.$$



## 2.4. Supuestos del modelo lineal

Los supuestos bajo los cuales serán válidas las inferencias que haremos más adelante sobre el modelo

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad (4)$$

son los siguientes:

$$\varepsilon_i \sim N(0, \sigma^2), \quad 1 \leq i \leq n, \quad \text{independientes entre sí.}$$

1. los  $\varepsilon_i$  tiene media cero,  $E(\varepsilon_i) = 0$ .
2. los  $\varepsilon_i$  tienen todos la misma varianza desconocida que llamaremos  $\sigma^2$  y que es el otro parámetro del modelo,  $Var(\varepsilon_i) = \sigma^2$ . A este requisito se lo suele llamar *homoscedasticidad*.
3. los  $\varepsilon_i$  tienen distribución normal
4. los  $\varepsilon_i$  son independientes entre sí, y son no correlacionados con las  $X_i$ .

$$Y \mid X = x_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{j=1}^n (X_j - \bar{X})^2} = \frac{1}{\sum_{j=1}^n (X_j - \bar{X})^2} \sum_{i=1}^n (X_i - \bar{X}) Y_i \\ &= \sum_{i=1}^n \frac{(X_i - \bar{X})}{\sum_{j=1}^n (X_j - \bar{X})^2} Y_i = \sum_{i=1}^n c_i Y_i\end{aligned}\quad (20)$$

donde

$$c_i = \frac{(X_i - \bar{X})}{\sum_{j=1}^n (X_j - \bar{X})^2} = \frac{(X_i - \bar{X})}{S_{XX}}, \quad (21)$$

$$S_{XX} = \sum_{j=1}^n (X_j - \bar{X})^2.$$

$$E(\hat{\beta}_1) = \beta_1$$

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

## 2.9. Inferencia sobre $\beta_1$

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}\right).$$

## 2.10. Inferencia sobre $\beta_0$

Esta inferencia despierta menos interés que la de  $\beta_1$ . Aunque los paquetes estadísticos la calculan es infrecuente encontrarla en aplicaciones. Bajo los supuestos del modelo lineal, puede calcularse la esperanza y varianza del estimador de  $\beta_0$ , que resultan ser

$$\begin{aligned}E\left(\widehat{\beta}_0\right) &= \beta_0 \\Var\left(\widehat{\beta}_0\right) &= \sigma^2 \left( \frac{1}{n} + \frac{\overline{X}^2}{\sum_{j=1}^n (X_j - \overline{X})^2} \right).\end{aligned}$$

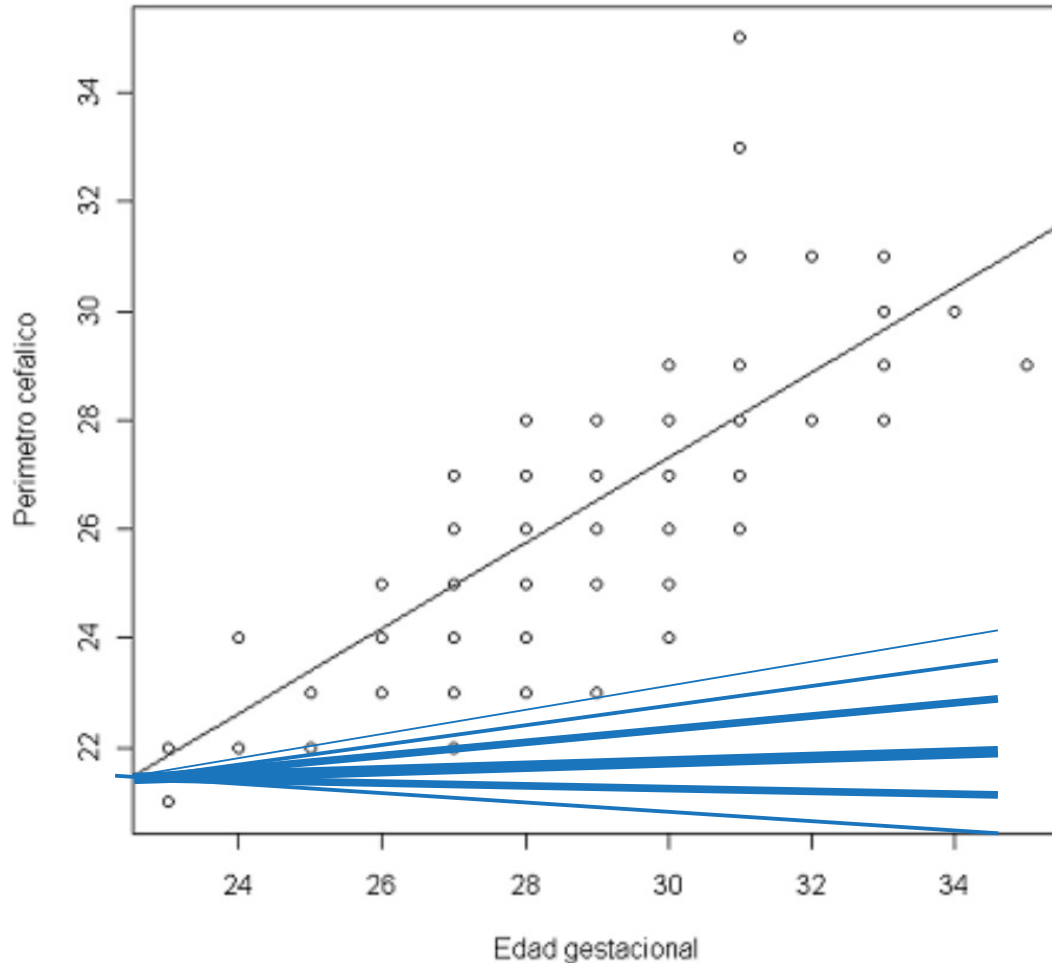
Nuevamente, las conclusiones son condicionales a los valores de los  $X$ 's observados. La varianza puede estimarse por

$$\widehat{Var}\left(\widehat{\beta}_0\right) = \widehat{\sigma}^2 \left( \frac{1}{n} + \frac{\overline{X}^2}{\sum_{j=1}^n (X_j - \overline{X})^2} \right)$$

Nuevamente, el estadístico  $\widehat{\beta}_0$  tiene distribución normal, su distribución es  $N\left(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\overline{X}^2}{\sum_{j=1}^n (X_j - \overline{X})^2} \right)\right)$ , luego



# ¿ Qué hago con la distribución del estadístico ?



$$E\left(\hat{\beta}_1\right) = \beta_1$$

$$Var\left(\hat{\beta}_1\right) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}\right).$$

$$\hat{\beta}_1 \sim N \left( \beta_1, \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right).$$

Un estimador de la varianza es

$$\widehat{Var} \left( \hat{\beta}_1 \right) = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\text{SSRes}/(n-2)}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

Finalmente, bajo los supuestos del modelo, puede probarse que

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\widehat{Var} \left( \hat{\beta}_1 \right)}} = \frac{\hat{\beta}_1 - \beta_1}{se_{\hat{\beta}_1}}$$

tiene distribución *t de Student con  $n - 2$  grados de libertad* si los datos siguen el modelo lineal, donde

$$se_{\hat{\beta}_1} = \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}} = \sqrt{\frac{\text{SSRes}/(n-2)}{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

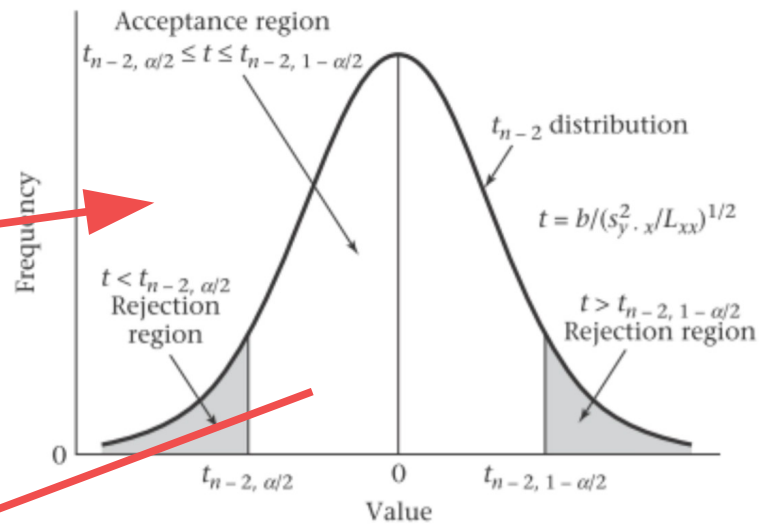
El estimador de  $\sigma^2$  que usaremos será

$$\hat{\sigma}^2 = \frac{1}{n-2} \text{SSRes} = \text{MSRes}. \quad \hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (e_i - \bar{e})^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{1}{n-2} \sum_{i=1}^n \left( Y_i - \hat{Y}_i \right)^2. \quad (17)$$

# Intervalos de Confianza

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\widehat{Var}(\hat{\beta}_1)}} = \frac{\hat{\beta}_1 - \beta_1}{se_{\hat{\beta}_1}}$$

*t de Student con  $n - 2$  grados de libertad*



Con esta distribución podemos construir un intervalo de confianza de nivel  $1 - \alpha$  para  $\beta_1$  que resultará

$$\hat{\beta}_1 \pm t_{n-2; 1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}, \text{ o bien} \quad (18)$$

$$\hat{\beta}_1 \pm t_{n-2; 1-\frac{\alpha}{2}} \cdot se_{\hat{\beta}_1}$$

# Test de Hipótesis

$$H_0 : \beta_1 = b$$

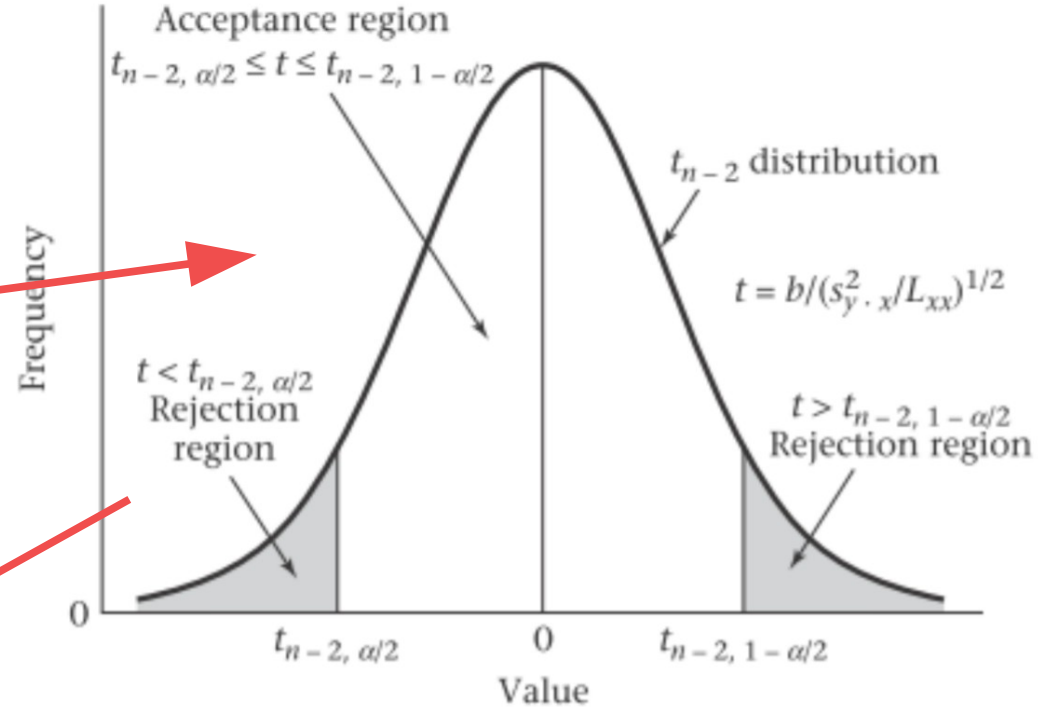
$$H_1 : \beta_1 \neq b.$$

$$T_{obs} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\widehat{Var}(\hat{\beta}_1)}} = \frac{\hat{\beta}_1 - \beta_1}{se_{\hat{\beta}_1}}$$

*t de Student con  $n - 2$  grados de libertad*

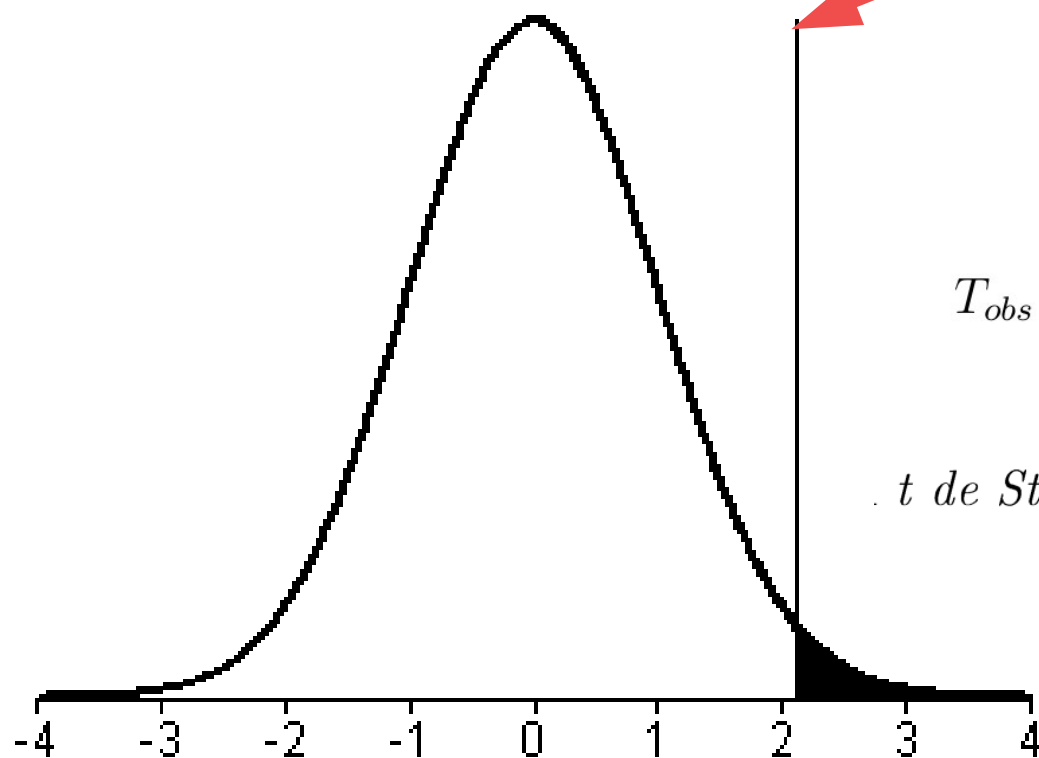
el test rechaza  $H_0$  con nivel  $\alpha$  si

$$T_{obs} \leq t_{n-2, \frac{\alpha}{2}} \quad \text{ó} \quad t_{n-2, 1-\frac{\alpha}{2}} \leq T_{obs},$$



# P-valor

$$p - valor = 2P(T \geq |T_{obs}|),$$



$$H_0 : \beta_1 = b$$

$$H_1 : \beta_1 \neq b.$$

$$T_{obs} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\widehat{Var}(\hat{\beta}_1)}} = \frac{\hat{\beta}_1 - \beta_1}{se_{\hat{\beta}_1}}$$

*t de Student con  $n - 2$  grados de libertad*

# Test, Intervalo y p-valor en el Ejemplo

Call:

```
lm(formula = headcirc ~ gestage)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.5358	-0.8760	-0.1458	0.9041	6.9041

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	3.91426	1.82915	2.14	0.0348 *
gestage	0.78005	0.06307	12.37	<2e-16 ***

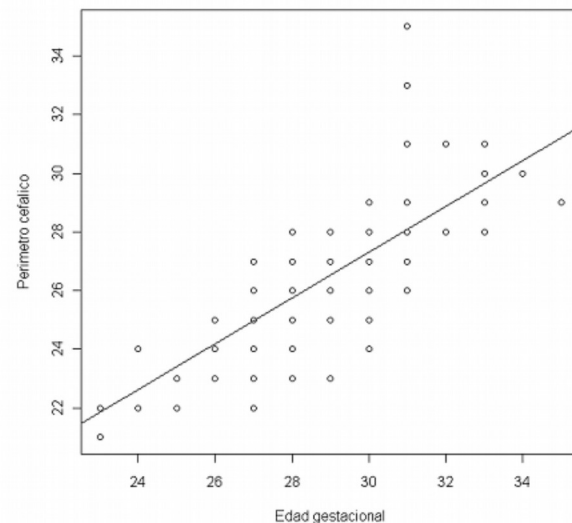
---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.59 on 98 degrees of freedom

Multiple R-squared: 0.6095, Adjusted R-squared: 0.6055

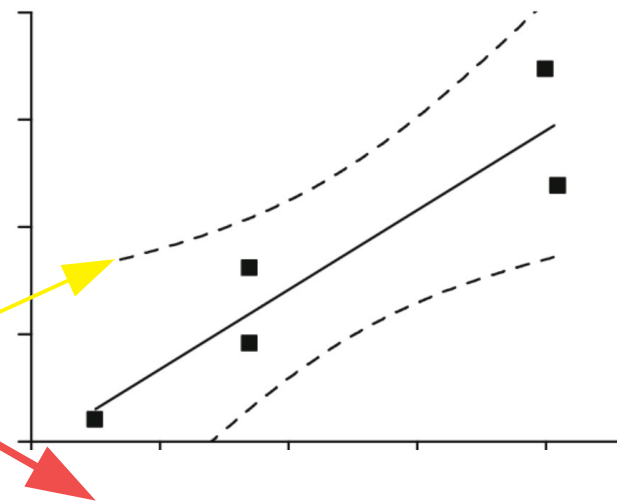
F-statistic: 152.9 on 1 and 98 DF, p-value: < 2.2e-16

$$\hat{\beta}_1 \pm t_{n-2; 1-\frac{\alpha}{2}} \cdot se_{\hat{\beta}_1}$$
$$0,7801 \pm 1,984467 \cdot 0,06307941$$
$$[0,654921, 0,905279]$$



# Inferencia sobre $E(Y|X=x)$

$$E(Y_h | X = x_h) = \beta_0 + \beta_1 x_h.$$



$$E(\hat{Y}_h) = E(Y_h)$$

$$Var(\hat{Y}_h) = \sigma^2 \cdot \left[ \frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right].$$

$$\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h.$$

$$\hat{Y}_h \pm t_{n-2; 1-\frac{\alpha}{2}} \cdot \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}.$$

(24)

# Inferencia sobre Y

$$Y_h = \beta_0 + \beta_1 X_h + \varepsilon_h$$

$$\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h.$$

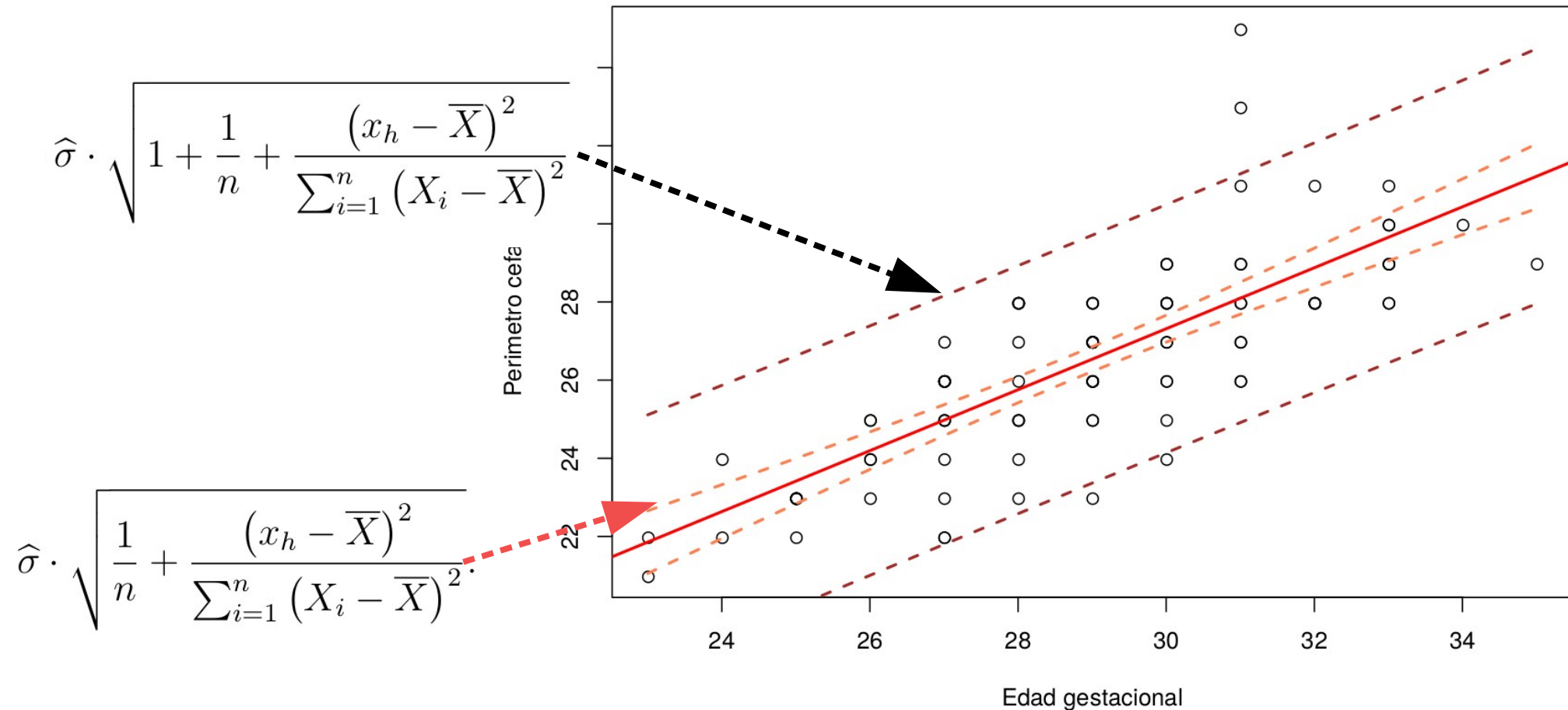
$$E(\hat{Y}_h) = E(Y_h)$$

$$Var(\hat{Y}_h) = \sigma^2 \left[ \frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right].$$

$$\hat{Y}_h \pm t_{n-2; 1-\frac{\alpha}{2}} \cdot \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$



# Predecir Y versus $E(Y|X=x)$



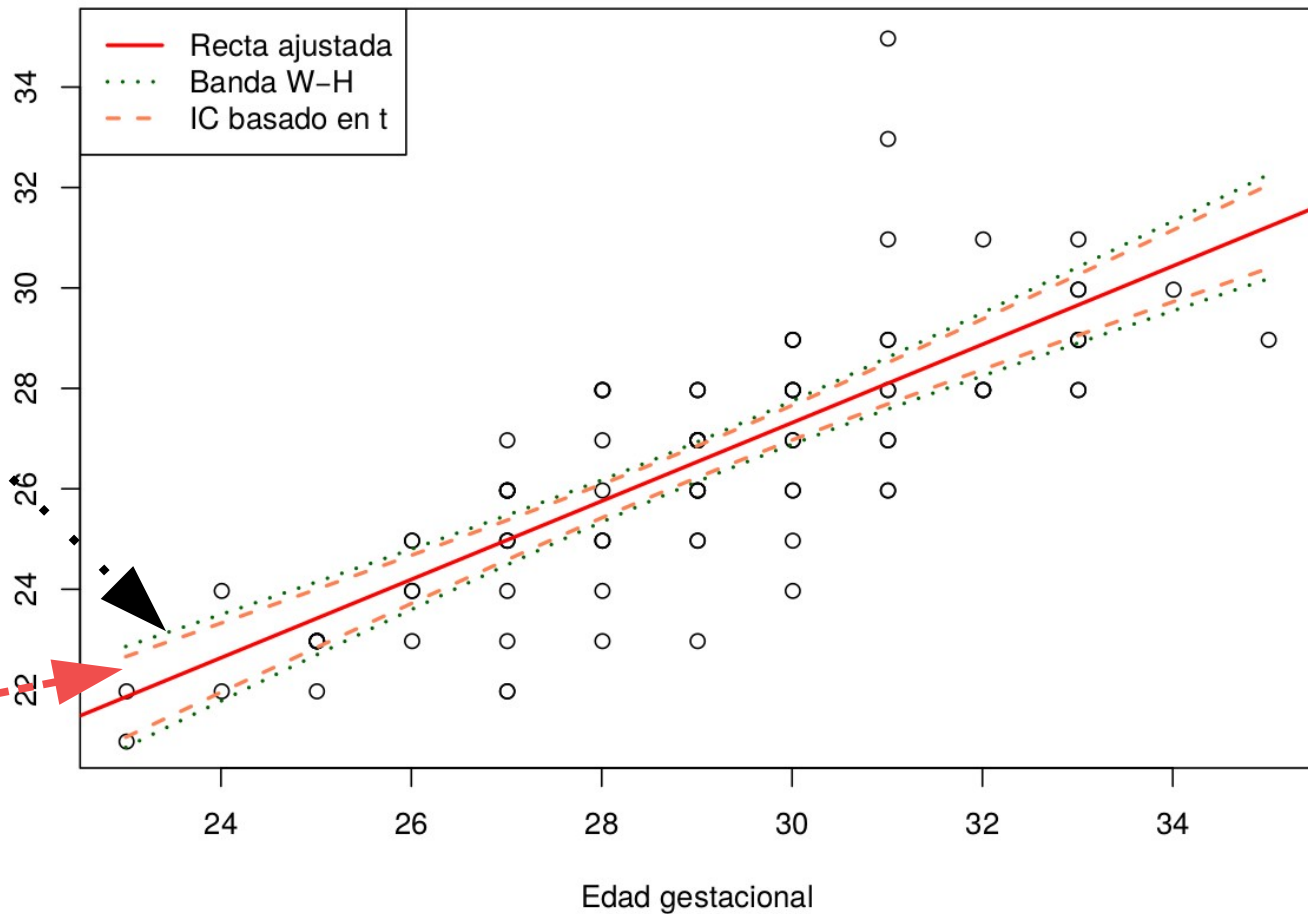
# Nivel Simultaneo

$$\hat{Y}_h \pm W \cdot \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{(x_h - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

$$W = \sqrt{2F_{1-\alpha,2,n-2}},$$

Perimetro cefalico

$$\hat{Y}_h \pm t_{n-2;\frac{\alpha}{2}} \cdot \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{(x_h - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$



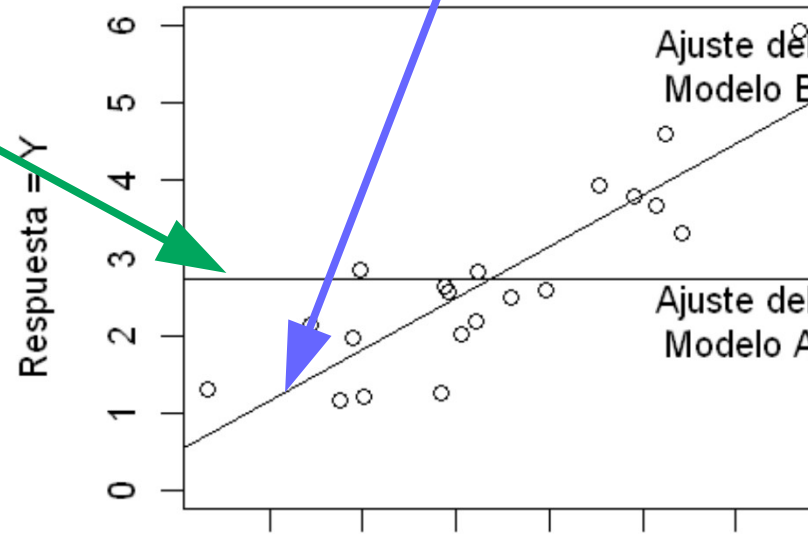
# Análisis de la Varianza

Modelo A:  $E(Y | X) = \mu$ , o escrito de otro modo

Modelo A:  $Y_i = \mu + u_i$  con  $u_i \sim N(0, \sigma_Y^2)$ ,  $1 \leq i \leq n$ , independientes entre sí.

Modelo B:  $E(Y | X) = \beta_0 + \beta_1 X$ , o escrito de otro modo

Modelo B:  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ , con  $\varepsilon_i \sim N(0, \sigma^2)$ ,  $1 \leq i \leq n$ , independientes entre sí.



$$\sum_{i=1}^n \overset{\text{SSTo}}{(Y_i - \bar{Y})^2} = \sum_{i=1}^n \overset{\text{SSRes}}{(Y_i - \hat{Y}_i)^2} + \sum_{i=1}^n \overset{\text{SSReg}}{(\hat{Y}_i - \bar{Y})^2}$$

# Análisis de la Varianza: Demostración

**Proof** [\[edit\]](#)

$$\begin{aligned}
 \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \bar{y} + \hat{y}_i - \hat{y}_i)^2 = \sum_{i=1}^n ((\hat{y}_i - \bar{y}) + \underbrace{(y_i - \hat{y}_i)}_{\hat{\varepsilon}_i})^2 \\
 &= \sum_{i=1}^n ((\hat{y}_i - \bar{y})^2 + 2\hat{\varepsilon}_i(\hat{y}_i - \bar{y}) + \hat{\varepsilon}_i^2) \\
 &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 + 2 \sum_{i=1}^n \hat{\varepsilon}_i(\hat{y}_i - \bar{y}) \\
 &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 + 2 \sum_{i=1}^n \hat{\varepsilon}_i(\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip} - \bar{y}) \\
 &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 + 2(\hat{\beta}_0 - \bar{y}) \underbrace{\sum_{i=1}^n \hat{\varepsilon}_i}_0 + 2\hat{\beta}_1 \underbrace{\sum_{i=1}^n \hat{\varepsilon}_i x_{i1}}_0 + \dots + 2\hat{\beta}_p \underbrace{\sum_{i=1}^n \hat{\varepsilon}_i x_{ip}}_0 \\
 &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 = \text{ESS} + \text{RSS}
 \end{aligned}$$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

SSTo                      SSRes                      SSRreg

$$\begin{aligned}
 \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)) &= 0 \\
 \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)) X_i &= 0.
 \end{aligned}$$

# Tabla de Anova

Tabla 14: Tabla de ANOVA para el modelo de Regresión Lineal Simple

Fuente de variación	Suma de cuadrados	Grados de libertad	Cuadrado medio	F	p-valor
Regresión	SSReg	1	MSReg	$\frac{MSReg}{MSRes}$	$P(F_{1,n-2} \geq F_{obs})$
Residuos	SSRes	$n - 2$	MSRes		
Total	SSTo	$n - 1$			

$$MSReg = \frac{SSReg}{1}$$

$$MSRes = \frac{SSRes}{n-2}$$

$$F = \frac{MSReg}{MSRes} = \frac{SSReg(n-2)}{SSRes}$$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

$SSTo$ 
 $SSRes$ 
 $SSReg$

# Coeficiente de Determinación

100 % de variabilidad ——— SSTo

% de variabilidad explicada ——— SSTo – SSRes

Luego el porcentaje de variabilidad explicada es

$$\frac{SSTo - SSRes}{SSTo} \times 100 \%$$

A la cantidad

$$\frac{SSTo - SSRes}{SSTo} = \frac{SSReg}{SSTo}$$

se la denomina  $R^2$ , o **coeficiente de determinación**

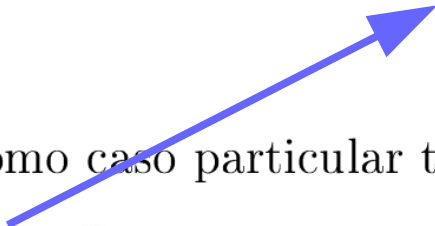
- $0 \leq R^2 \leq 1$
- No depende de las unidades de medición.
- Es el cuadrado del coeficiente de correlación de Pearson para la muestra  $\{(X_i, Y_i)\}_{1 \leq i \leq n}$ . También es el cuadrado del coeficiente de correlación de Pearson para los pares  $\{(\hat{Y}_i, Y_i)\}_{1 \leq i \leq n}$ , es decir, entre los valores de la covariable observados y los predichos por el modelo lineal.
- Mientras mayor es  $R^2$  mayor es la fuerza de la variable regresora ( $X$ ) para predecir a la variable respuesta ( $Y$ ).
- Mientras mayor sea  $R^2$  menor es la SSRes y por lo tanto, más cercanos están los puntos a la recta.
- Toma el mismo valor cuando usamos a  $X$  para predecir a  $Y$  o cuando usamos a  $Y$  para predecir a  $X$ .

# Leverage

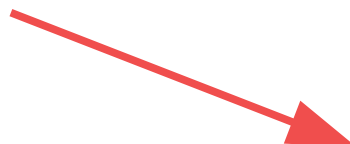
El valor predicho de un dato puede escribirse como combinación lineal de las observaciones

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i = \sum_{k=1}^n h_{ik} Y_k \quad (33)$$

donde


$$h_{ik} = \frac{1}{n} + \frac{(X_i - \bar{X})(X_k - \bar{X})}{S_{XX}}$$

y como caso particular tenemos que


$$S_{XX} = \sum_{k=1}^n (X_k - \bar{X})^2$$

$$\sum_{k=1}^n h_{ik} = 1, \quad \sum_{i=1}^n h_{ik} = 1$$

$$\sum_{i=1}^n h_{ii} = 2$$

$$\frac{1}{n} \leq h_{ii} \leq \frac{1}{s} \leq 1.$$


$$h_{ii} = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{S_{XX}}. \quad (34)$$

# Residuos

$$e_i = Y_i - \hat{Y}_i \approx \varepsilon_i$$

## 3.1.2. Residuos

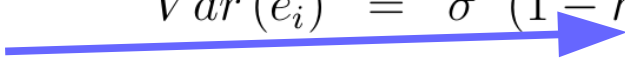
Dijimos en la Sección 2.8 que los residuos son cantidades observables, que representan de alguna manera el correlato empírico de los errores. Para verificar los supuestos del modelo lineal, suelen usarse métodos gráficos que involucran a los residuos. El modelo lineal

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

supone que los errores  $\varepsilon$  tienen media poblacional cero y varianza constante (que denominamos  $\sigma^2$ ), y que son indendientes para distintas observaciones. Sin embargo, ya hemos visto que no ocurre lo mismo con los residuos. Vimos que los residuos no son independientes. Además, puede probarse que

$$\begin{aligned} E(e_i) &= 0 \\ Var(e_i) &= \sigma^2 (1 - h_{ii}) \end{aligned} \quad (37)$$

$h_{ii} = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{S_{XX}}$





# Residuos Estandarizados

## 3.1.3. Residuos estandarizados

Para hacer más comparables a los residuos entre sí, podemos dividir a cada uno de ellos por un estimador de su desvío estándar, obteniendo lo que se denominan *residuos estandarizados*:

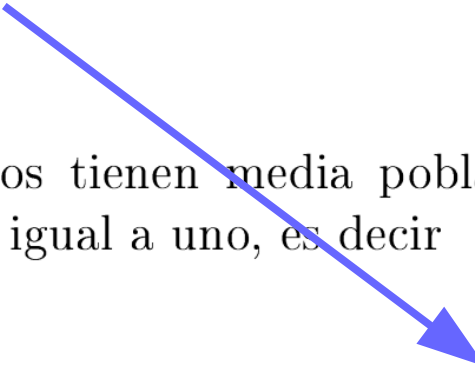
$$rest_i = \frac{e_i}{\sqrt{\hat{\sigma}^2 (1 - h_{ii})}}. \quad (38)$$

Recordemos que el estimador de  $\sigma^2$  bajo el modelo de regresión está dado por

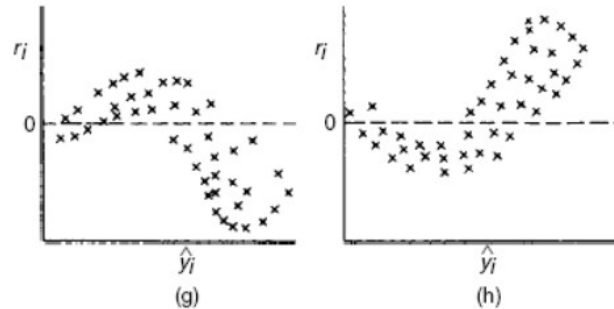
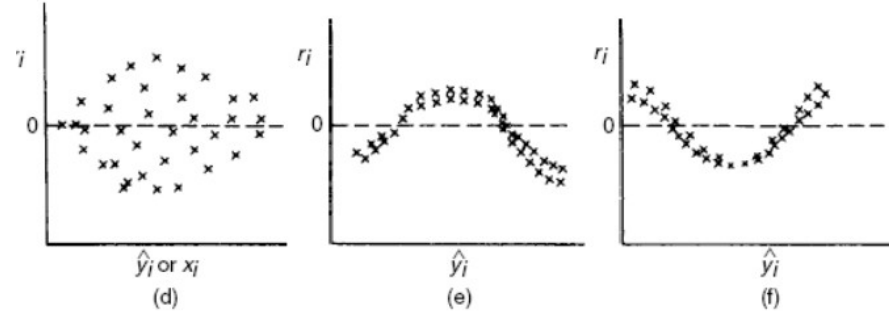
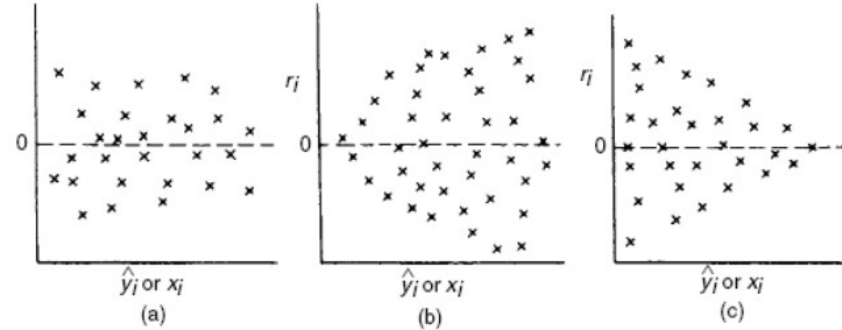
$$\hat{\sigma}^2 = \frac{SSRes}{n - 2}$$

Puede probarse que los residuos estandarizados tienen media poblacional cero (igual que los residuos), y varianza poblacional igual a uno, es decir

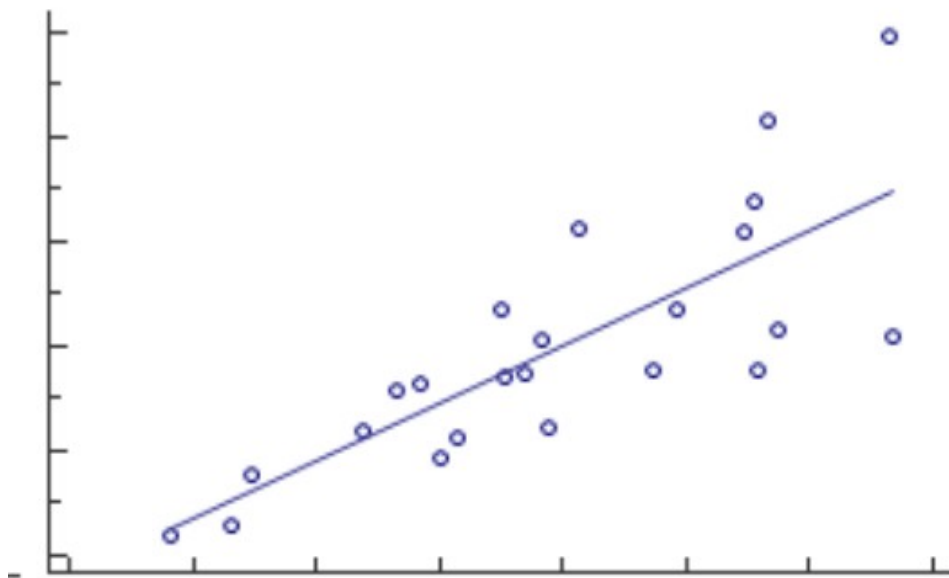
$$\begin{aligned} E(rest_i) &= 0 \\ Var(rest_i) &= 1, \quad \text{para todo } i. \end{aligned}$$


$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

# La Estructura de los Residuos



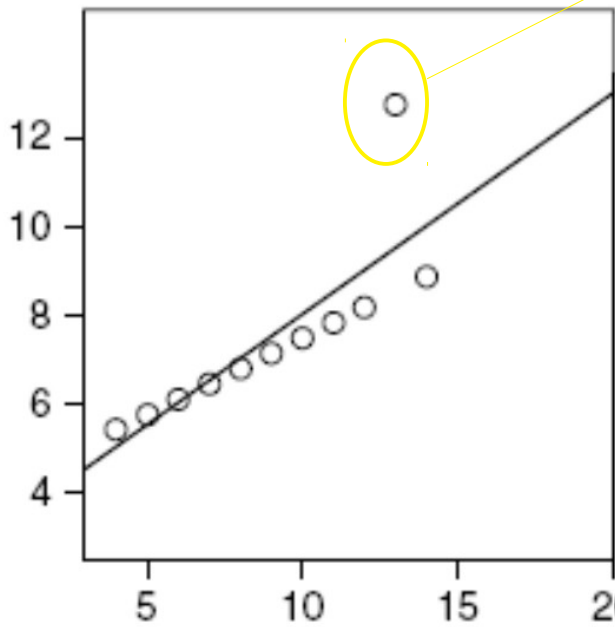
# Heterocedasticidad



$$\text{Var}(Y \mid X = x_i) = \text{Var}(\varepsilon_i) = \frac{\sigma^2}{w_i}$$

$$g_{wls}(a, b) = \sum_{i=1}^n w_i (Y_i - (a + bX_i))^2.$$

# Outliers



1. Eliminamos esa observación de la muestra, de modo que ahora tenemos una muestra con  $n - 1$  casos.
2. Usando el conjunto de datos reducidos volvemos a estimar los parámetros, obteniendo  $\hat{\beta}_{0(i)}$ ,  $\hat{\beta}_{1(i)}$  y  $\hat{\sigma}_{(i)}^2$  donde el subíndice  $(i)$  está escrito para recordarnos que los parámetros fueron estimados sin usar la  $i$ -ésima observación.
3. Para el caso omitido, calculamos el valor ajustado  $\hat{Y}_{i(i)} = \hat{\beta}_{0(i)} + \hat{\beta}_{1(i)}X_i$ . Como el caso  $i$ -ésimo no fue usado en la estimación de los parámetros,  $Y_i$  y  $\hat{Y}_{i(i)}$  son independientes. La varianza de  $Y_i - \hat{Y}_{i(i)}$  puede calcularse y se estima usando  $\hat{\sigma}_{(i)}^2$ .

4. Escribamos

$$t_i = \frac{Y_i - \hat{Y}_{i(i)}}{\sqrt{\widehat{Var}(Y_i - \hat{Y}_{i(i)})}}$$

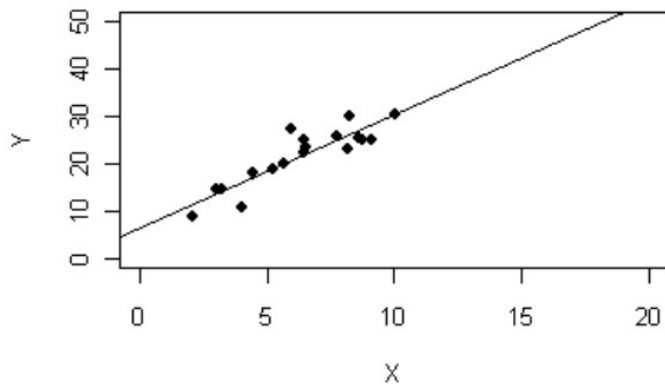
$$t_i = \frac{e_i}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}} = rest_i \sqrt{\frac{n - 3}{n - 2 - rest_i}} \rightarrow t_i \sim t_{n-3}$$

$$rest_i = \frac{e_i}{\sqrt{\hat{\sigma}^2 (1 - h_{ii})}}.$$

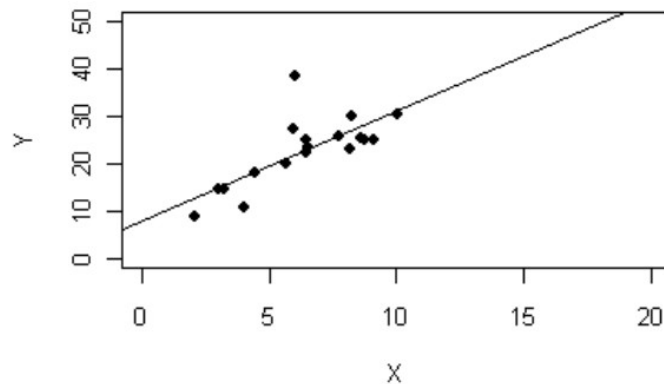
la versión estandarizada del estadístico en consideración. Si la observación  $i$ -ésima sigue el modelo, entonces la esperanza de  $Y_i - \hat{Y}_{i(i)}$  debería ser cero. Si no lo sigue, será un valor no nulo. Luego, si llamamos  $\delta$  a la esperanza poblacional de esa resta,  $\delta = E(Y_i - \hat{Y}_{i(i)})$ , y asumimos normalidad de los errores, puede probarse que la distribución de  $t_i$  bajo la hipótesis  $H_0 : \delta = 0$  es una  $t$  de Student con  $n - 3$  grados de libertad,  $t_i \sim t_{n-3}$  (recordar que

# Observaciones Influyentes

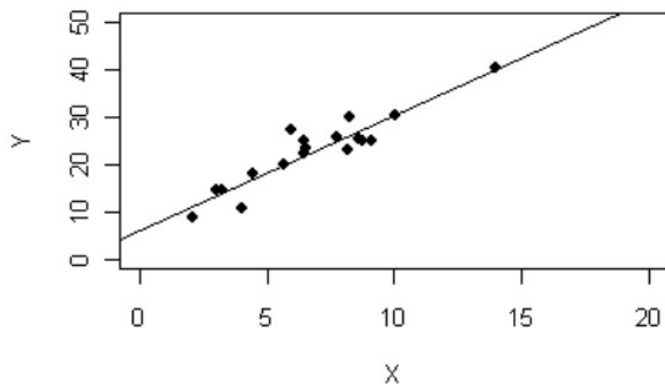
(1) Recta ajustada  $Y = 2.40X + 6.41$



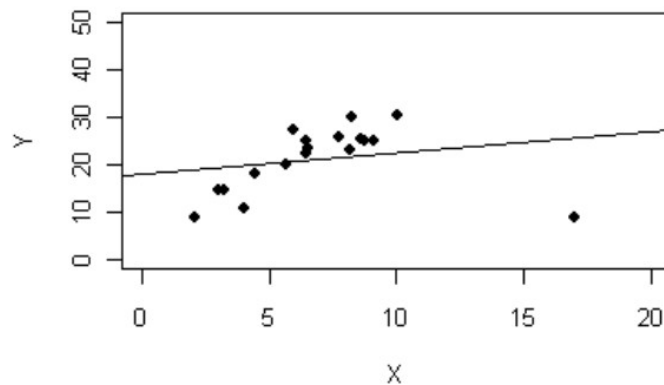
(2) Recta ajustada  $Y = 2.33X + 7.84$



(3) Recta ajustada  $Y = 2.42X + 6.26$



(4) Recta ajustada  $Y = 0.45X + 17.89$



# Distancia de Cook

$$D_i = \frac{\left(\widehat{Y}_{(i)i} - \widehat{Y}_i\right)^2}{2\widehat{\sigma}^2},$$

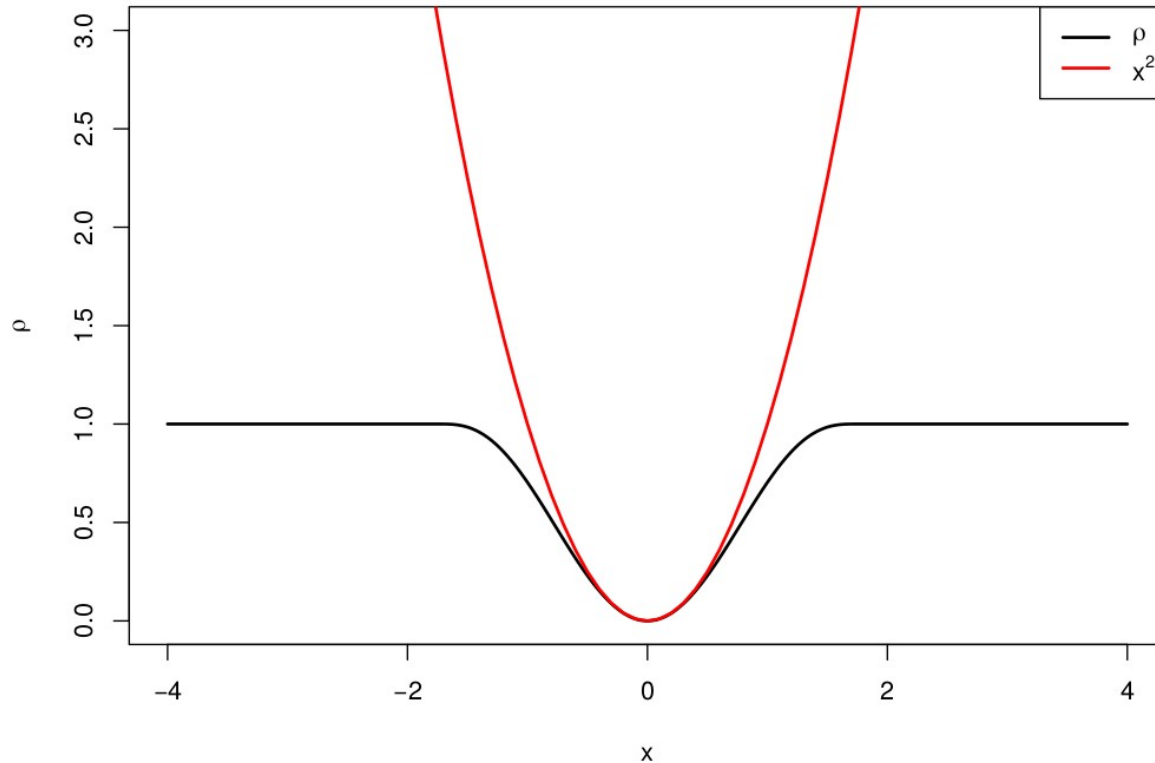
$$D_i = \frac{1}{2} (rest_i)^2 \frac{h_{ii}}{1 - h_{ii}}.$$

- Si  $D_i < \text{percentil } 0,20$  de la distribución  $F_{2,n-2}$  entonces la observación no es influyente.
- Si  $D_i > \text{percentil } 0,50$  de la distribución  $F_{2,n-2}$  entonces la observación es muy influyente y requerirá tomar alguna medida.
- Si  $D_i$  se encuentra entre el percentil 0,20 y el percentil 0,50 de la distribución  $F_{2,n-2}$  se sugiere mirar además otros estadísticos.

# Regresión Robusta

$$g(a, b) = \sum_{i=1}^n (Y_i - (a + bX_i))^2,$$

$$g(a, b) = \sum_{i=1}^n \rho \left( \frac{Y_i - (a + bX_i)}{s_n} \right)$$



# OLS is BLUE



$$y_i = \sum_{j=1}^K \beta_j X_{ij} + \varepsilon_i \quad \forall i = 1, 2, \dots, n$$

$$\sum_{j=1}^K \lambda_j \beta_j$$

- They have mean zero:  $\mathbf{E}[\varepsilon_i] = 0$ .
- They are **homoscedastic**, that is all have the same finite variance:  $\text{Var}(\varepsilon_i) = \sigma^2 < \infty$  for all  $i$  and
- Distinct error terms are uncorrelated:  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \forall i \neq j$ .

$$\hat{\beta}_j = c_{1j}y_1 + \dots + c_{nj}y_n$$

$$\mathbf{E} \left[ \left( \sum_{j=1}^K \lambda_j (\hat{\beta}_j - \beta_j) \right)^2 \right]$$

$$\mathbf{E}[\hat{\beta}_j] = \beta_j$$

minimizes the **sum of squares**

$$\sum_{i=1}^n \left( y_i - \sum_{j=1}^K \hat{\beta}_j X_{ij} \right)^2$$