

# MAT-466: Bondad de ajuste y residuos en GLM

**Felipe Osorio**

[fosorios.mat.utfsm.cl](mailto:fosorios.mat.utfsm.cl)

Departamento de Matemática, UTFSM



Recuerde que, en general,  $D(\mathbf{y}; \hat{\boldsymbol{\mu}})$  **no** sigue una distribución asintótica  $\chi^2(n-p)$ . Por ejemplo,

- ▶ **Poisson:** Cuando  $\mu_i \rightarrow \infty$  para  $i = 1, \dots, n$  tenemos  $D(\mathbf{y}; \hat{\boldsymbol{\mu}}) \sim \chi^2(n-p)$ .
- ▶ **Binomial:**  $Y_i \sim \text{Bin}(n_i, \mu_i)$ , para  $i = 1, \dots, k$ . Así para  $k$  fijo y  $n_i \rightarrow \infty$ ,  $\forall i$  tenemos  $D(\mathbf{y}; \hat{\boldsymbol{\mu}}) \sim \chi^2(n-p)$ .
- ▶ **Normal:** Es bien sabido que para  $\sigma^2$  fijo, tenemos  $D(\mathbf{y}; \hat{\boldsymbol{\mu}})/\sigma^2 \sim \chi^2(n-p)$ .

En general tenemos que (Jørgensen, 1987)

$$D^*(\mathbf{y}; \hat{\boldsymbol{\mu}}) \sim \chi^2(n-p), \quad \text{cuando } \phi \rightarrow \infty$$

Es decir,

- ▶ **Normal:**  $\sum_{i=1}^n (y_i - \hat{\mu}_i)^2 / \sigma^2 \sim \chi^2(n-p)$  si  $\sigma^2 \rightarrow 0$ .
- ▶ **Gamma:**  $D^*(\mathbf{y}; \hat{\boldsymbol{\mu}}) \sim \chi^2(n-p)$  si  $\text{CV} \rightarrow 0$  ( $\phi^{-1/2} \rightarrow \infty$ )



Sabemos que, si  $T \sim \chi^2(n-p)$ . Entonces,  $E(T) = n - p$ . De este modo,

### Regla de trabajo:

Un valor del desvío cercano a  $n - p$  puede indicar que el **modelo está bien ajustado**.

Otra alternativa es usar la estadística chi-cuadrado de Pearson:

$$X^2 = \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)}.$$



Suponga que

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}).$$

El vector de residuos es dado por:

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y},$$

con  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ . En nuestra notación, tenemos  $\hat{\boldsymbol{\mu}} = \mathbf{H}\mathbf{Y} = \hat{\mathbf{E}}(\mathbf{Y})$ .

En particular,

$$h_{ii} = \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i, \quad i = 1, \dots, n,$$

y  $0 \leq h_{ii} \leq 1$ . Además,

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{1}{n} \text{tr}(\mathbf{H}) = \frac{p}{n}.$$



# Residuos en regresión lineal

Bajo el supuesto de normalidad  $Y \sim N_n(X\beta, \sigma^2 I)$ , tenemos

$$e \sim N_n(\mathbf{0}, \sigma^2(I - H)),$$

es decir

$$E(e_i) = 0, \quad \text{var}(e_i) = \sigma^2(1 - h_{ii}), \quad \text{Cov}(e_i, e_j) = -\sigma^2 h_{ij}.$$

De ahí que los residuos tienen varianzas diferentes y son correlacionados.

El residuo estandarizado es definido como:

$$r_i = \frac{e_i}{s\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n.$$

Cook y Weisberg (1982) mostraron que

$$\frac{r_i^2}{n - p} \sim \text{Beta}\left(\frac{1}{2}, \frac{n - p - 1}{2}\right),$$

de este modo

$$E(r_i) = 0, \quad \text{var}(r_i) = 1, \quad \text{Cov}(r_i, r_j) = \frac{-h_{ij}}{\sqrt{(1 - h_{ii})(1 - h_{jj})}}.$$



Considere el residuo studentizado:

$$t_i = \frac{e_i}{s_{(i)}\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n.$$

donde

$$s_{(i)}^2 = \frac{1}{n - p - 1} \sum_{j \neq i}^n (y_j - \hat{\mu}_j)^2 = s^2 \left( \frac{n - p - r_i^2}{n - p - 1} \right),$$

sigue que  $t_i \sim t(n - p - 1)$ .

Una interpretación interesante de  $t_i$  es que corresponde al estadístico  $t$  para probar la hipótesis  $H_0 : \gamma = 0$  en el modelo de salto en la media:

$$Y_j = \mathbf{x}_j^\top \boldsymbol{\beta} + d_j \gamma + \epsilon_j, \quad j = 1, \dots, n,$$

donde  $d_j = 1$  si  $j = i$  y 0 en caso contrario.



## Objetivo:

Evaluar desvios de normalidad de los residuos studentizados  $t_i$ 's.

## Notación:

Considere  $Z_i = t_i$ , para  $i = 1, \dots, n$ <sup>1</sup>

## Idea:

Comparar la CDF muestral para los  $Z_i$ 's contra la CDF de la  $N(0, 1)$ .

Asuma que los **residuos**  $Z_i$  están ordenados

$$Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(n)},$$

los  $Z_{(i)}$  son los cuantiles de la CDF muestral, definida como

$$\text{Proportion}(Z \leq Z_{(i)}) = \frac{i}{n}$$

---

<sup>1</sup>La descripción es válida otras medidas de interés.

## QQ-plot en regresión lineal

Asuma que los **residuos**  $Z_i$  están ordenados

$$Z_{(1)} \leq Z_{(2)} \leq \cdots \leq Z_{(n)},$$

los  $Z_{(i)}$  son los cuantiles de la CDF muestral, definida como

$$\text{Proportion}(Z \leq Z_{(i)}) = \frac{i}{n}.$$

Los cuantiles de la distribución teórica, son dados por:

$$q_i^* = \Phi^{-1}\left(\frac{i}{n}\right).$$

Si los errores son aproximadamente normales, se debe tener que el gráfico de los pares  $(q_1^*, Z_{(1)}), \dots, (q_n^*, Z_{(n)})$  sea a recta identidad.





Se ha sugerido la siguiente aproximación para la esperanza de los estadísticos de orden desde  $N(0, 1)$  como:

$$q_i = \Phi^{-1}\left(\frac{i - 3/8}{n + 1/4}\right),$$

de este modo, se utilizará el gráfico cuantil-cuantil (QQ-plot) de los pares  $(q_i, Z_{(i)})$ .

### Observación:

- ▶ Podemos construir QQ-plots para diversas distribuciones.<sup>2</sup>
- ▶ Es difícil chequear visualmente desvios de la distribución de interés.

---

<sup>2</sup>Por ejemplo,  $\chi^2$ ,  $t$  de Student, Poisson, Gama, etc.



## QQ-plot con envelopes en regresión lineal

Envelopes simulados son herramientas gráficas para chequear el ajuste de un modelo. Atkinson (1985) sugirió usar el siguiente procedimiento:

- ▶ Ajustar un modelo de regresión lineal, calcular residuos y estandarizar para obtener varianza unitaria.
- ▶ Generar  $M$  ( $\approx 1000$ ) muestras como respuesta. Para cada muestra ajuste el mismo modelo y calcule los residuos estandarizados
- ▶ Ordenar todos los conjuntos de residuos estandarizados.
- ▶ El envelope consiste de los cuantiles 2.5% inferior y superior de los residuos estandarizados generados en cada posición.



# Herencia de la estatura (Weisberg, 2005)

## *Ejemplo (Herencia de la estatura):*

Se recolectó la altura de  $n = 1375$  madres en UK (bajo 65 años) y una de sus hijas adultas (sobre 18 años).

Exploramos el conjunto de datos por medio del gráfico:

```
> plot(dheight ~ mheight, data = Heights)
```

Ajuste de un modelo de regresión lineal simple

```
> fm <- lm(dheight ~ mheight, data = Heights)
> fm
```

Call:

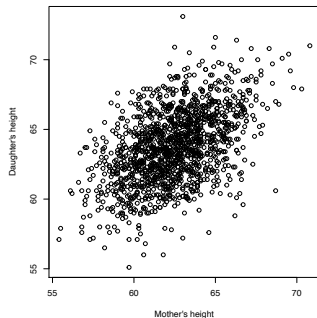
```
lm(formula = dheight ~ mheight, data = Heights)
```

Coefficients:

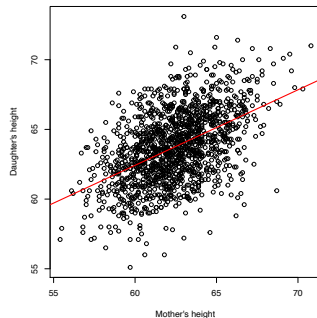
(Intercept)	mheight
29.9174	0.5417



# Herencia de la estatura

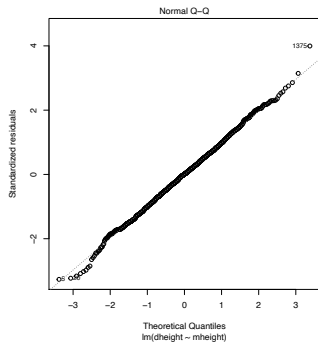


(a) datos estatura

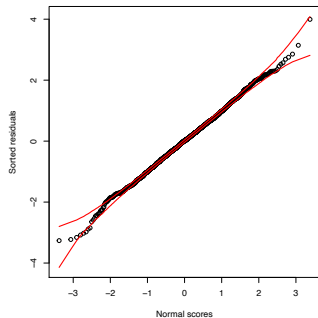


(b) recta ajustada

# Herencia de la estatura



(a) datos estatura



(b) recta ajustada

Note que a la convergencia del proceso iterativo,

$$\hat{\beta} = (\mathbf{X}^\top \widehat{\mathbf{W}} \mathbf{X})^{-1} \mathbf{X}^\top \widehat{\mathbf{W}} \widehat{\mathbf{Z}},$$

con  $\widehat{\mathbf{Z}} = \widehat{\boldsymbol{\eta}} + \widehat{\mathbf{W}}^{-1/2} \widehat{\mathbf{V}}^{-1/2} (\mathbf{Y} - \widehat{\boldsymbol{\mu}})$ . Una de las primeras alternativas ha sido definir

$$\mathbf{r}_P = \widehat{\mathbf{W}}^{1/2} (\widehat{\mathbf{Z}} - \widehat{\boldsymbol{\eta}}) = \widehat{\mathbf{V}}^{-1/2} (\mathbf{Y} - \widehat{\boldsymbol{\mu}}),$$

cuyo  $i$ -ésimo elemento es el **residuo Pearson**

$$r_{P_i} = \frac{y_i - \widehat{\mu}_i}{\widehat{V}_i^{1/2}}, \quad i = 1, \dots, n,$$

donde  $\widehat{V}_i = V(\widehat{\mu}_i)$ .



Asumiendo que  $\text{Cov}(\mathbf{Z}) \approx \phi^{-1} \widehat{\mathbf{W}}^{-1}$ , tenemos  $\text{Cov}(\mathbf{r}_P) \approx \phi^{-1}(\mathbf{I} - \widehat{\mathbf{H}})$ , con

$$\mathbf{H} = \mathbf{W}^{1/2} \mathbf{X} (\mathbf{X}^\top \widehat{\mathbf{W}} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{1/2},$$

podemos definir la versión estandarizada

$$t_{S_i} = \frac{\phi^{1/2}(y_i - \hat{\mu}_i)}{\hat{V}_i^{1/2}(1 - \hat{h}_{ii})^{1/2}}, \quad i = 1, \dots, n,$$



Considere la función

$$A(\mu) = \int_0^\mu \frac{dt}{V^{1/3}(t)},$$

se ha sugerido el uso de la función  $A(\cdot)$  para definir residuos cuya distribución puede ser más cercana de la normal. Esto lleva al **residuo Anscombe**

$$t_{A_i} = \frac{\phi^{1/2} \{A(y_i) - A(\hat{\mu}_i)\}}{\hat{V}_i^{1/2} A'(\hat{\mu}_i)}.$$

Algunos ejemplos de  $A(\cdot)$ :

- ▶ **Normal:**  $\mu$ .
- ▶ **Binomial:**  $\int_0^\mu t^{-1/3}(1-t)^{-1/3} dt$ .
- ▶ **Poisson:**  $\frac{3}{2}\mu^{2/3}$ .
- ▶ **Gama:**  $3\mu^{1/3}$ .
- ▶ **Normal inversa:**  $\log \mu$ .





Considere

$$D(\mathbf{y}; \hat{\boldsymbol{\mu}}) = \sum_{i=1}^n d(y_i; \hat{\mu}_i),$$

es decir  $d_i = d(y_i; \hat{\mu}_i)$  es el componente  $i$ -ésimo del desvío. Eso lleva a definir el **residuo componente de desvío**

$$\begin{aligned} r_{D_i} &= \text{sign}(y_i - \hat{\mu}_i) \sqrt{d_i} \\ &= \text{sign}(y_i - \mu_i) \sqrt{2} \{y_i(\tilde{\theta}_i - \hat{\theta}_i) + (b(\theta))\}^{1/2} \end{aligned}$$

para  $i = 1, \dots, n$ . Esto lleva a la versión estandarizada

$$t_{D_i} = \frac{\phi^{1/2} \text{sign}(y_i - \hat{\mu}_i) \sqrt{d_i}}{(1 - \hat{h}_{ii})^{1/2}}, \quad i = 1, \dots, n.$$

*Observación:*

El residuo componente de desvío es posiblemente el más usado en GLM.



Un cuarto tipo de residuo fue propuesto por Williams (1987), definido como:

$$t_{W_i} = \text{sign}(y_i - \hat{\mu}_i) \{ (1 - \hat{h}_{ii}) t_{D_i}^2 + \hat{h}_{ii} t_{S_i}^2 \}^{1/2}, \quad i = 1, \dots, n,$$

que puede ser interpretado como un promedio ponderado entre  $t_{D_i}$  y  $t_{S_i}$ .

Más recientemente Dunn y Smith (1996) introdujeron el **residuo cuantil**, definido como:

$$r_{Q_i} = \Phi^{-1} \{ F(y_i; \hat{\mu}_i, \hat{\phi}) \},$$

donde  $F$  es la CDF asociada al modelo estadístico  $FE(\theta_i, \phi)$  y  $\Phi$  es la CDF de la distribución  $N(0, 1)$ .

### *Observación:*

Evidentemente, salvo la aleatoriedad en  $\hat{\mu}_i$  y  $\hat{\phi}$ ,  $r_{Q_i}$  es exactamente normal.

