

Modelos Lineales y Lineales Generalizados

Manuel Gijón Agudo

Octubre 2018 -

Índice

1. Regresión lineal simple	3
1.1. Estimación de parámetros	5
1.1.1. Mínimos cuadrados	5
1.1.2. Máxima verosimilitud	6
1.2. Valores predichos y residuos	7
1.3. Estimación de la varianza	7
1.3.1. Método de los momentos	7
1.3.2. Método de máxima verosimilitud	7
1.3.3. Interpretación de los parámetros	9
1.4. Inferencia en los parámetros del modelo	9
1.4.1. Distribución de $\hat{\beta}$	9
1.4.2. Inferencia	10
1.4.3. Tabla ANOVA	11
1.4.4. Test OMNIBUS	12
1.5. Inferencia en los valores preedichos	12
1.5.1. Distribución de los valores	12
1.5.2. Intervalo de predicción	13
1.5.3. Intervalo de confianza	13
1.6. Más sobre los residuos	14
1.6.1. Distribución del vector de residuos	14
1.6.2. <i>Standardized and Studentized residuals</i>	15
1.7. Importante tener en cuenta	15
1.7.1. Multicolinearidad	15
1.7.2. Apalancamiento (<i>Leverage</i>)	16
1.7.3. Valores influyentes: Distancia de Cook	16
1.7.4. Valores influyentes: DFBETA, DFBETAS	17
1.8. Bondad de ajuste del modelo	17
1.8.1. R^2	17
1.8.2. R^2 ajustado	18

<i>Modelos Lineales y Lineales Generalizados</i>	2
1.8.3. Análisis de los residuos	18
1.9. Resumen	18
1.10. Ejemplo completo	20
2. Errores comunes	21
3. ANOVA	21
3.0.1. Ejemplos	21
3.0.2. Ejercicios	21
3.1. ANCOVA	21
3.1.1. Ejemplos	21
3.1.2. Ejercicios	21
4. Regresión lineal generalizada	22
4.1. Genealidades	22
4.2. Binomial Response Models	23
4.3. Poisson Response Models	24

1. Regresión lineal simple

Objetivo: Nuestro objetivo será siempre explicar el comportamiento de una variable aleatoria Y en función de unos ciertos valores X_1, \dots, X_p .

Dado un $n \in \mathbb{Z}^+$ denominaremos Y_i a la muestra de Y obtenida cuando $X_j = x_{ij} \in \mathbb{R} \forall i, j$.

Definición: denominamos el **Modelo Lineal** como:

$$\forall i, \quad Y_i = \beta_0 + x_{i1}\beta_1 + \dots + x_{i(p-1)}\beta_{p-1} + e_i = \mu_i + e_i$$

donde β_0 es denominado **intercepto** (*intercept*) y los términos e_i los **errores**.

Hipótesis:

- $\forall i \in \{1, 2, \dots, n\}, e_i \sim N(0, \sigma_i^2)$
- **Homeodasticidad** (*homeodasticity*): $\forall i \in \{1, 2, \dots, n\}, \sigma_i^2 = \sigma^2$
- $\forall i, y \in \{1, 2, \dots, n\}, i \neq j$ e_i es **independiente** de e_j
- Los valores de X son fijos o variables aleatorias **independientes** de los errores.

En **forma matricial** escribiremos el modelo de la siguiente manera:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & x_{13} & \cdots & x_{1(p-1)} \\ 1 & x_{21} & x_{22} & x_{23} & \cdots & x_{2(p-1)} \\ \vdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \cdots & x_{n(p-1)} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

Siendo así, definimos $Y_{n \cdot 1} = (Y_1, Y_2, \dots, Y_n)^t$, $X_{n \cdot p} = (x_{ij})$, $\beta_{p \cdot 1} = (\beta_0, \beta_1, \dots, \beta_{p-1})^t$, $e_{n \cdot 1} = (e_1, e_2, \dots, e_n)^t$ y escribimos el modelo como:

$$Y = X\beta + e \iff \mu = E(Y|X) = X\beta$$

$$Y|X \sim N(X\beta, \sigma^2 \cdot Id_n)$$

De la última línea se desprende lo siguiente:

$$E((Y - X\beta)(Y - X\beta)^t) = E\left(\begin{pmatrix} e_1^2 & e_1e_2 & e_1e_3 & \cdots & e_1e_n \\ e_2e_3 & e_2^2 & e_2e_3 & \cdots & e_2e_n \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ e_ne_1 & e_ne_2 & e_ne_3 & \cdots & e_n^2 \end{pmatrix}\right) = \sigma^2 \cdot Id_n$$

Observación: Las variables X_1, X_2, \dots, X_n pueden ser función de otro conjunto de variables. Por ejemplo, podríamos tener un conjunto $\{Z_1, Z_2, \dots, Z_m\}$ tales que, para $m \in \mathbb{N}$:

$$X_i = g_i(Z_1, Z_2, \dots, Z_m), \quad i \in \{1, 2, \dots, n\}$$

Ejemplos de modelos lineales:

- Comparamos la presión sanguínea (Y) en dos tipos de individuos, unos que han tomado cierta medicación y un grupo de control:

$$Y_{ij} = \mu_i + e_{ij}, \quad \forall i \in \{1, 2\}, \quad j \in \{1, 2, \dots, n\}$$

En forma matricial lo podemos expresar de la siguiente manera:

$$\begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ \vdots \\ Y_{2n_2} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} + \begin{pmatrix} e_{11} \\ \vdots \\ e_{1n_1} \\ \vdots \\ e_{2n_2} \end{pmatrix}$$

Los modelos conocidos como **modelos de regresión** son un caso particular de modelos lineales en los que las covariables son continuas o discretas, en ningún caso categóricas.

- Estudiamos el nivel de un determinado químico en una planta (Y) en función de su presencia en el suelo (X).

$$Y_i = \beta_1 + x_i \beta_2 + e_i, \quad i = 1, 2, \dots, n$$

O en su forma matricial:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} + \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

Si las covariables son todas categóricas estaremos entonces ante un modelo **ANOVA** (*Analysis of Covariance*).

- Queremos estudiar el nivel de un determinado medicamento (Y) en función de su dosis (X_1) y del género del paciente (X_2).

$$Y_{ij} = \beta_{0i} + \beta_{1i}x_{ij} + e_{ij}, \quad i \in \{1, 2\}, \quad j \in \{1, 2, \dots, n\}$$

En forma matricial:

$$\begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & 0 & 0 \\ 1 & x_{12} & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n_1} & 0 & 0 \\ 0 & 0 & 1 & x_{21} \\ 0 & 0 & 1 & x_{22} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 1 & x_{2n_2} \end{pmatrix} \cdot \begin{pmatrix} \beta_{01} \\ \beta_{11} \\ \beta_{02} \\ \beta_{12} \end{pmatrix} + \begin{pmatrix} e_{11} \\ \vdots \\ e_{1n_1} \\ e_{21} \\ \vdots \\ e_{2n_2} \end{pmatrix}$$

Definimos el **Modelo Nulo** (*Null Model*) al más simple, el que tiene un único parámetro.

$$Y_i = \beta_0 + e_i, \quad i = 1, \dots, n$$

Equivalentemente:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \cdot (\beta_0) + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

Este modelo es equivalente a estudiar una muestra de una variable aleatoria.

Denominamos **Intercepto** (*intercept*) al elemento β_0 . En general consideraremos modelos con intercepto, lo que significa que el modelo nulo será un submodelo del susodicho.

Ejemplos de modelos no lineales:

- Queremos estudiar la producción de leche de unas vacas en Litros (Y_i) en función del número de días (x_i) que hace que nacieron.

$$Y_i = e^{\beta_0 + \beta_1 x_i + \log(x_i)} + e_i, \quad e_i \sim N(0, \sigma^2)$$

- Queremos estudiar la calidad de cierto material en función de su proveedor, de cada uno de ellos elegiremos aleatoriamente una muestra de tamaño b del material de las que nos quedaremos con n muestras seleccionadas aleatoriamente.

$$Y_{ijk} = \mu + \alpha_i + \beta_{j(i)} + e_{(ij)k}$$

$$i = 1, \dots, a, \quad j = 1, \dots, b, \quad k = 1, \dots, n, \quad e_{(ij)k} \sim N(0, \sigma^2)$$

1.1. Estimación de parámetros

1.1.1. Mínimos cuadrados

Sea $y = (y_1, y_2, \dots, y_n)^t$ una realización muestral de Y y $\hat{\beta}$ una estimación del parámetro β .

- **Estimación por mínimos cuadrados** (*Minimum least square*): consiste en minimizar

$$S(\beta) = \|y - \hat{y}\|_2^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \left(y_i - \sum_{j=0}^{p-1} x_{ij} \beta_j \right)^2$$

donde $\hat{y} = \hat{\mu} = X\hat{\beta}$ y la **solución** es:

$$\boxed{\hat{\beta} = (X^t X)^{-1} X^t y}$$

Siempre que $X^t X$ sea invertible.

- **Weighted least squares**: consiste en minimizar

$$S(\beta) = \sum_{i=1}^n w_i (y - \hat{y})^2 = \sum_{i=1}^n w_i \left(y_i - \sum_{j=0}^{p-1} x_{ij} \beta_j \right)^2$$

donde $w_i^{-1} = \text{Var}(Y_i)$.

La **solución** es (obviamente, siempre que $X^t X$ sea invertible):

$$\hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} y$$

donde $V = \text{diag}(w_i)$.

Recalquemos que en ningún caso necesitamos conocer la distribución de Y .

1.1.2. Máxima verosimilitud

El estimador es el siguiente:

$$L(\beta; y) = (\sigma\sqrt{2\pi})^{-n} e^{-\sum_{i=1}^n \frac{(y_i - \sum_{j=0}^{p-1} x_{ij} \beta_j)^2}{2\sigma^2}}$$

Equivalentemente (tomando logaritmos):

$$l(\beta; y) = -n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(y_i - \sum_{j=0}^{p-1} x_{ij} \beta_j)^2}{2\sigma^2}$$

Ahora definimos el siguiente vector:

$$U_j = \frac{\partial l}{\partial \beta_j} = \frac{1}{\sigma^2} \left(X^t (Y - X\beta) \right)_j, \quad \forall j$$

$$U = (U_1, U_2, \dots, U_{p-1})^t$$

que denominaremos **vector de puntuaciones** (*score vector*).

$$U_j = 0 \iff X^t Y = X^t X \beta$$

Luego si $\text{rank}(X^t X) = p$ entonces la solución, el **estimador por máxima verosimilitud** es:

$$\hat{\beta} = (X^t X)^{-1} X^t Y$$

1.2. Valores predichos y residuos

Sea $\hat{Y} = X\hat{\beta} = (\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_n)^t$ nuestro vector de valores predichos.

Si $(y_1, y_2, \dots, y_n)^t$ es una realización muestral de Y , definimos el **raw residual** de la siguiente manera:

$$\text{raw residual} = y_i - \hat{y}_i = \hat{e}_i$$

Notemos que el vector $\hat{e} = Y - X\hat{\beta}$ es ortogonal a las columnas de la matriz X .

$$\begin{aligned} X^t \hat{e} &= X^t (Y - X\hat{\beta}) \\ &= X^t (Y - X((X^t X)^{-1} X^t Y)) \\ &= X^t Y - X^t X((X^t X)^{-1} X^t Y) = X^t Y - X^t Y = 0 \end{aligned}$$

1.3. Estimación de la varianza

1.3.1. Método de los momentos

Asumimos que $p = \text{Rank}(X^t X)$, lo que implica que el rango de $X^t X$ es máximo, luego se cumple:

$$\frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \sim \chi_{n-p}^2$$

Luego inmediatamente se desprende:

$$E\left(\frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right) = n - p \iff E(S^2) = \sigma^2$$

donde:

$$S^2 = \frac{1}{n - p} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Con lo que S^2 es un **estimador insesgado** de σ^2 . Este estimador es conocido como **Error cuadrático medio** (*mean square error*).

1.3.2. Método de máxima verosimilitud

La función de verosimilitud es la siguiente:

$$l(\sigma^2; \mu) = -n \log(\sigma \sqrt{2\pi}) - \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{2\sigma^2}$$

$$\frac{\partial l}{\partial \sigma} = -\frac{n}{\sigma^2} + \frac{1}{(\sigma^2)^2} \sum_{i=1}^n (y_i - \mu_i)^2 = 0$$

hacer los cálculos aquí

$$\Rightarrow \hat{\sigma}^2 = \left(1 - \frac{p}{n}\right) S^2$$

Observación: si n es grande, ambos estimadores son similares pero sin embargo para p y n pequeños ambos estimadores difieren mucho.

Ejemplo: Consideremos el modelo nulo $y_i = \beta_0 + e_i$, $i = 1, \dots, n$. Demostrar que:

$$\widehat{\beta}_0 = \bar{y}$$

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

- Respuesta 1
- Respuesta 2

Ejemplo: Consideremos la regresión lineal simple $y_i = \beta_0 + x_i \beta_1 + e_i$, $i = 1, \dots, n$ y demostremos lo siguiente:

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x}$$

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n (\bar{x})^2} = \frac{Cov(X, Y)}{Var(X)} = r_{XY} \frac{S_Y}{S_X}$$

Antes de comenzar haremos dos observaciones:

- (\bar{x}, \bar{y}) pertenece a la recta de regresión.
- El **coeficiente de correlación** (r_{XY}) mide la relación lineal entre X e Y .
- Respuesta 1
- Respuesta 2

1.3.3. Interpretación de los parámetros

Asumimos el siguiente modelo:

$$\forall i, Y_i = \beta_0 + x_{i1}\beta_1 + x_{i2}\beta_2 + \cdots + x_{ip-1}\beta_{p-1} + e_i = \mu_i + e_i$$

Siendo Y_i la variable respuesta bajo la siguiente condición $X_i = (x_{i1}, x_{i2}, \dots, x_{ip-1})$ y siendo Y_i^* la respuesta bajo las condiciones $X_i^* = (x_{i1}, x_{i2}, \dots, x_{ij} + 1, \dots, x_{ip-1})$ se tiene que:

$$Y_i - Y_i^* = \hat{\beta}_j$$

Entonces nos queda:

- $\hat{\beta}_j$ es el cambio medio obtenido al incrementar en una unidad el valor de x_j y mantener el resto inamovibles.
- Si $\hat{\beta}_0$ es la estimación del intercepto, podemos interpretarlo como la respuesta media en el origen.

La **desviación residual estándar** (*residual standard deviation*) es el error asociado a nuestras predicciones, el 95 % de nuestras predicciones tendrán el error en el siguiente intervalo:

$$\left(-t_{n-p, \frac{\alpha}{2}} \hat{\sigma}, t_{n-p, \frac{\alpha}{2}} \hat{\sigma} \right) \simeq (-1,95\hat{\sigma}, 1,95\hat{\sigma})$$

[Explicar el por qué](#)

1.4. Inferencia en los parámetros del modelo

1.4.1. Distribución de $\hat{\beta}$

Si β_0 es el parámetro real, sabemos que:

$$\boxed{\hat{\beta}|X \sim N\left(\beta_0, \sigma^2(X^t X)^{-1}\right)}$$

[Por ser una combinación lineal de variables aleatorias normales.](#)

[Explicar el por qué con más detalle](#)

$$\begin{aligned} E(\hat{\beta}|X) &= (X^t X)^{-1} X^t E(Y|X) \\ &= (X^t X)^{-1} X^t X \beta_0 = \beta_0 \end{aligned}$$

Sabemos también que:

$$\hat{\beta} - \beta_0 = (X^t X)^{-1} X^t (Y - X \beta_0)$$

Luego queda:

$$\begin{aligned} E\left((\hat{\beta}|X - \beta_0) \cdot (\hat{\beta}|X - \beta_0)^t\right) &= (X^t X)^{-1} X^t E\left((Y - X\beta_0) \cdot (Y - X\beta_0)^t\right) X (X^t X)^{-1} \\ &= \sigma^2 (X^t X)^{-1} \end{aligned}$$

Observación: Las componentes de $\hat{\beta}$ no son variables aleatorias independientes.

Denominamos **Matriz de información de Fisher** (*Fisher information matrix*) a la matriz $\mathcal{J} = E(UU^t)$.

Observemos que bajo la hipótesis de la normalidad ocurre lo siguiente:

$$\begin{aligned} \mathcal{J} = E(UU^t) &= E\left(\frac{1}{\sigma^2} X^t (Y - X\beta) \cdot (Y - X\beta)^t X \frac{1}{\sigma^2}\right) \\ &= \frac{1}{\sigma^2} X^t E\left((Y - X\beta) \cdot (Y - X\beta)^t\right) X \frac{1}{\sigma^2} \\ &= \frac{1}{\sigma^2} X^t X \end{aligned}$$

Observación: La matriz de información de Fisher es la inversa de la matriz de varianzas y covarianzas de $\hat{\beta}$.

Ejemplo: Consideremos la regresión lineal simple $y_i = \beta_0 + x_i \beta_1 + e_i$, $i = 1, \dots, n$ y demostremos que, para cada coeficiente, las desviaciones estándar de los estimadores son:

$$\begin{aligned} S_{\hat{\beta}_0} &= S \cdot \left(\frac{1}{n} + \frac{1}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \right)^{1/2} \\ S_{\hat{\beta}_1} &= S \cdot \frac{1}{\left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)^{1/2}} \end{aligned}$$

- Respuesta 1
- Respuesta 2

1.4.2. Inferencia

Sabemos que $\hat{\beta}|X \sim N(\beta_0, \sigma^2 (X^t X)^{-1})$, luego cada parámetro verifica lo siguiente:

$$\boxed{\hat{\beta}_i|X \sim N\left(\beta_{0i}, \sigma^2 [(X^t X)^{-1}]_{ii}\right)}$$

Para cada $a \in \mathbb{R}$ podemos hacer el siguiente test:

$$H_0 : \beta_{0i} = a$$

$$H_1 : \beta_{0i} \neq a$$

Para un nivel de significación α se tiene:

$$\frac{\hat{\beta}_i - a}{\hat{\sigma} \sqrt{[(X^t X)^{-1}]_{ii}}}$$

Luego **rechazaremos la hipótesis nula** si:

$$\left| \frac{\hat{\beta}_i - a}{\hat{\sigma} \sqrt{[(X^t X)^{-1}]_{ii}}} \right| \geq t_{n-p, \alpha/2}$$

En particular, estaremos interesados en el siguiente test:

$$H_0 : \beta_{0i} = 0$$

$$H_1 : \beta_{0i} \neq 0$$

No rechazar H_0 implica que la covariable X_i no tiene una influencia significativa en Y .

Los **Intervalos de confianza** para los parámetros con un nivel de significación α son los siguiente:

$$\boxed{\hat{\beta}_i \pm t_{n-p, \frac{\alpha}{2}} \cdot \hat{\sigma} \sqrt{[(X^t X)^{-1}]_{ii}}}$$

Importante: Si el intervalo contiene el valor cero, la correspondiente X_i **no es estadísticamente significativa**

1.4.3. Tabla ANOVA

Dado $Y_i - \bar{Y} = (Y_i - \hat{Y}_i) + (\hat{Y}_i - \bar{Y})$ y que $\sum_{i=1}^n (Y_i - \hat{Y}_i) \cdot (\hat{Y}_i - \bar{Y}) = 0$ se tiene lo siguiente:

$$\boxed{\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{TSS} = \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{RSS} + \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{RegSS}}$$

[Incluir demostración aquí](#)

Donde cada término es:

- **TSS:** Suma de cuadrados total (*total sum of squares*)
- **RSS:** Suma de cuadrados debida a los residuos (*residual sum of squares*)
- **RegSS:** Suma de cuadrados explicada por la regresión (*regression sum of squares*)

TABLA ANOVA				
Fuente	Grados de libertad.	SS	MSS	F
Regresión	$p - 1$	$\sum_{i=1}^n (Y_i - \bar{Y})^2$	$RegSS/p$	F_0
Residuos	$n - p$	$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	$RSS/(n - p)$	
Total	$n - 1$	$\sum_{i=1}^n (Y_i - \bar{Y})^2$		

Denominamos **tabla ANOVA** a la siguiente construcción:

$$\text{Donde } F_0 = \frac{RegSS/p}{RSS/(n-p)} \text{ y } \hat{\sigma} = S^2 = RSS/(n-p).$$

Ejemplo: INCLUIREMOS AQUÍ UN EJEMPLO COMPLETO DE TODO ESTO

1.4.4. Test OMNIBUS

Si β_1 se corresponde con el intercepto, el **test OMNIBUS** se define como:

$$H_0 : \beta_2 = \beta_3 = \dots = \beta_p = 0$$

$$H_1 : \exists i : \beta_i \neq 0$$

Observación: Este test compara un modelo con el modelo nulo.

Este test se efectúa a partir de la tabla ANOVA de la siguiente manera, si H_0 es cierto:

$$F_0 = \frac{RegSS/p}{RSS/(n-p)} \sim F_{p,n-p}$$

Lo que implica que **rechazamos la hipótesis nula** bajo un nivel de significación α si:

$$F_0 \geq F_{\frac{\alpha}{2}, p, n-p}$$

Observación: El test OMNIBUS no es equivalente a comprobar si cada parámetro es equivalente a cero.

1.5. Inferencia en los valores preedichos

1.5.1. Distribución de los valores

Definimos el **vector de valores preedichos** como $\hat{Y} = X\hat{\beta}$.

$$\hat{Y}|X \sim N\left(X\beta_0, \sigma^2 X(X^t X)^{-1} X^t\right)$$

Por ser combinación lineal de distribuciones normales se tiene:

$$E(\hat{Y}|X) = XE(\hat{\beta}|X) = X\beta_0$$

y también

$$\begin{aligned} E\left((\hat{Y}|X - X\beta_0) \cdot (\hat{Y}|X - X\beta_0)^t \text{ Big}\right) &= XE\left((\hat{\beta}|X - \beta_0) \cdot (\hat{\beta}|X - \beta_0)^t\right)X^t \\ &= X\sigma^2(X^tX)^{-1}X^t \\ &= \sigma^2X(X^tX)^{-1}X^t \end{aligned}$$

Denominamos a la matriz $X(X^tX)^{-1}X^t$ **hat matrix** debido a que:

$$\hat{Y} = X(X^tX)^{-1}X^tY$$

Observación: σ^2 veces la *hat matrix* es la matriz de varianzas y covarianzas de las predicciones.

1.5.2. Intervalo de predicción

Sea $X_0 = (x_{01}, x_{02}, \dots, x_{0p})^t$ un conjunto de condiciones experimentales particulares.

El valor predicho en X_0 es: $\hat{y}_0 = X_0^t\hat{\beta}$.

Recordemos que bajo la normalidad y la hipótesis nula el intervalo de confianza para μ se calcula como:

$$\bar{y} \pm t_{\frac{\alpha}{2}, n-p} S / \sqrt{n}$$

En este caso, para un nivel de significación α , un **intervalo de predicción al 100(1 - α) %** para una observación futura es:

$$\boxed{\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{\hat{\sigma}^2(1 + X_0^t(X^tX)^{-1}X_0)}}$$

Ejemplo: Consideremos la regresión lineal simple $y_i = \beta_0 + x_i\beta_1 + e_i$, $i = 1, \dots, n$ y demostremos lo siguiente:

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{\hat{\sigma}^2 \left(1 + \frac{1}{n} + \frac{(X_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)}$$

1.5.3. Intervalo de confianza

Sea $X_0 = (x_{01}, x_{02}, \dots, x_{0p})^t$ un conjunto de condiciones experimentales particulares.

El valor predicho en X_0 es: $\hat{y}_0 = X_0^t\hat{\beta}$.

Sea $\mu_0 = E(Y|X_0)$, tenemos que $E(\bar{y}_0) = X_0^t \hat{\beta} = E(Y|X_0)$.

En consecuencia, \hat{y}_0 es un **estimador insesgado** de $E(Y|X_0)$. Es más:

$$Var(\hat{y}_0) = \sigma^2 X_0^t (X^t X)^{-1} X_0$$

Siendo entonces, para μ_0 el **intervalo de confianza** para un nivel de significación α :

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{\hat{\sigma}^2 (X_0^t (X^t X)^{-1} X_0)}$$

Ejemplo: Consideremos la regresión lineal simple $y_i = \beta_0 + x_i \beta_1 + e_i$, $i = 1, \dots, n$ y demostremos lo siguiente:

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{(X_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}$$

1.6. Más sobre los residuos

1.6.1. Distribución del vector de residuos

Definimos el **residuo i-ésimo** como $\hat{e}_i = y_i - \hat{y}_i$ y el **vector de residuos** como:

$$\hat{e} = (Y_1 - \hat{Y}_1, Y_2 - \hat{Y}_2, \dots, Y_n - \hat{Y}_n)^t$$

y verifica que:

$$\hat{e}|X \sim N\left(0, \sigma^2 (Id - X(X^t X)^{-1} X^t)\right)$$

Como tenemos una combinación lineal de variables aleatorias normalmente distribuidas:

$$E(\hat{e}) = X\beta_0 - X\beta_0 = 0$$

Como consecuencia nos encontramos con:

$$E(YY^t) = E(\hat{Y}\hat{Y}^t) + E((Y - \hat{Y}) \cdot (Y - \hat{Y})^t)$$

donde tenemos:

$$\begin{aligned} E((Y - \hat{Y}) \cdot (Y - \hat{Y})^t) &= \sigma^2 Id - \sigma^2 X(X^t X)^{-1} X^t \\ &= \sigma^2 (Id - X(X^t X)^{-1} X^t) \end{aligned}$$

1.6.2. Standarized and Studentized residuals

Obtenemos una estimación de la varianza de $(y_i - \hat{y}_i)$ con la siguiente expresión:

$$S^2 \cdot (1 - [X(X^t X)^{-1} X^t]_{ii})$$

Denotamos $h_{ii} = [X(X^t X)^{-1} X^t]_{ii}$ definimos:

- **Residuos estandarizados** (*Standarized residuals*):

$$\text{Stand. Res} = \frac{y_i - \hat{y}_i}{s\sqrt{1 - h_{ii}}}$$

- **Studentized residuals**:

$$\text{Student. Res} = \frac{y_i - \hat{y}_i}{s_{-i}\sqrt{1 - h_{ii}}}$$

Donde S^2_{-i} es la estimación de la varianza cuando el término i-ésimo es suprimido del modelo.

1.7. Importante tener en cuenta

1.7.1. Multicolinearidad

Definimos la **multicolinearidad** como el hecho de que haya alguna de las variables explicativas estén relacionadas entre sí por medio de una función lineal.

El fenómeno implica que $\det(X^t \cdot X)$ es muy pequeño o incluso cero. Si el fenómeno está presente:

- La interpretación del modelo es muy difícil.
- la varianza del parámetro $\hat{\beta}_i$ es grande.
- La matriz $X^t \cdot X$ puede ser singular, lo que imposibilita calcular las estimaciones de los parámetros.

Observación: incluso con multicolinearidad, las predicciones son correctas si el modelo lo es.

Para **detectar** el fenómeno:

1. Para realizar un *multiple scattered plot* para todos los pares de variables explicativas y calcular el **coeficiente de correlación entre pares de predictores**:

$$r_{x_1 x_2} = \frac{\sum_i x_{1i} x_{2i} - n \bar{x}_1 \bar{x}_2}{\left(\sum_i x_{1i}^2 - n \bar{x}_1^2\right)^{1/2} \left(\sum_i x_{2i}^2 - n \bar{x}_2^2\right)^{1/2}}$$

Observación: Si el término de correlación es relativamente alto, una de las variables debe ser eliminada y repetido el análisis.

2. Es conveniente calcular el **Factor de Infalcción de la Varianza** (*Variance Inflation Factor* (*VIF*)) de cada variable. Cuando calculamos una regresión podemos calcular:

$$R^2 = \frac{\text{RegSS}}{\text{TSS}} \cdot 100$$

Se define el VIF de la siguiente manera:

$$\boxed{\text{VIF}(X_j) = \frac{1}{1 - R^2(X_j)}}$$

donde $R^2(X_j)$ es el coeficiente R^2 obtenido *regressing* la variable X_j con respecto al resto de variables explicativas.

- $VIF = 1 \Rightarrow$ las variables no están correladas.
- $1 < VIF < 5 \Rightarrow$ correlación moderada.
- $VIF > 5$ altamente correladas.

Observación: denominamos **tolerancia** al valor $1 - R^2(x_j)$.

Puede ser probado que:

$$\text{Var}(\hat{\beta}_j) = \text{VIF}(X_j) \cdot \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$$

luego a mayor VIF mayor varianza y menos podemos confiar en ese coeficiente.

1.7.2. Apalancamiento (*Leverage*)

Dada una observación Y_i , su apalancamiento se define como:

$$h_{ii} = [X(X^t X)^{-1} X^t]_{ii}$$

es una medida de la distancia entre $(x_{i1}, x_{i2}, \dots, x_{ip})$ y el centroide $(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)$.

Se verifica que, $\forall i, \frac{1}{n} \leq h_{ii} \leq 1$.

Si tenemos algo similar a cualquiera de los dos siguientes supuestos:

- $h_{ii} > 3 \frac{p}{n}$
- $h_{ii} > 0,9$

Entonces decimos que el apalancamiento es grande y procederemos con cuidado.

1.7.3. Valores influyentes: Distancia de Cook

Una **observación influyente** (*influential observation*) es una observación que tienen un valor desproporcionado sobre los valores de los coeficientes de regresión.

Calculamos la **distancia de Cook** (*Cook's distance*):

Observación: cuanto mayores sean los componentes de c , mayor será la correlación entre las variables explicativas y la variable respuesta.

Una expresión alternativa para el cálculo de R^2 es la siguiente:

$$R^2 = \frac{\text{RegSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

y es interpretado como la proporción de variabilidad en los datos explicado por las covariables.

Ejemplo: Consideremos la regresión lineal simple $y_i = \beta_0 + x_i\beta_i + e_i$, $i = 1, \dots, n$ y demostremos lo siguiente:

$$R^2 = (r_{xy})^2$$

1.8.2. R^2 ajustado

Es obvio que cuantas más covariables tengamos en el modelo más grande será el valor de R^2 . Para penalizar estos modelos y **poder comparar modelos con diferente número de variables** nace este coeficiente:

$$R_{adj}^2 = 1 - \frac{\text{RSS}/(n-p)}{\text{TSS}/(n-1)} = 1 - (1 - R^2) \frac{n-1}{n-p}$$

Observación: Si p aumenta, el numerador de la fracción aumenta y entonces R_{adj}^2 decrece respecto a R^2 .

1.8.3. Análisis de los residuos

Comprobando las hipótesis del modelo:

1. *Scatter plot* de y contra x , debemos observar **linearidad**.
2. *Scatter plot* de \hat{e}_i contra \hat{y}_i , no debemos observar **ninguna tendencia**.
3. *qq-plot* de \hat{e}_i , debemos observar **linealidad**.
4. En ocasiones puede ser útil \hat{e}_i contra el orden en que se han tomado las observaciones.

1.9. Resumen

- **Estimación de parámetros:**

$$\hat{\beta} = (X^t \cdot X)^{-1} X^t y$$

- Valores predichos:

$$\hat{y}_i = (X\hat{\beta})_i$$

- Suma residual de los cuadrados:

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Estimación de la varianza:

$$\hat{\sigma}^2 = S^2 = MSE = \frac{RSS}{n-p}$$

- Desciación estándar para los β 's:

$$S_{\hat{\beta}_j} = S \cdot \sqrt{c_{jj}}$$

donde los coeficientes c_{jj} son los elementos en la diagonal de la matriz $(X^t \cdot X)^{-1}$.

- Intervalos de confianza para los parámetros

$$\hat{\beta}_i \pm t_{\frac{1-\alpha}{2}, n-p} \cdot S_{\hat{\beta}_j}$$

- Coeficiente de determinación:

$$R^2 = \frac{\text{RegSS}}{\text{TSS}} = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

1.10. Ejemplo completo

2. Errores comunes

Llegados a este punto es importante hablar de los diferentes tipos de modelos lineales con los que nos encontraremos en función del tipo de dato que sean las variables explicativas.

3. ANOVA

En este caso alguna de las covariables si no todas serán **factores** (también llamadas **variables categóricas**). Estos modelos se denominan **análisis de la varianza (ANOVA)**.

3.0.1. Ejemplos

3.0.2. Ejercicios

3.1. ANCOVA

En el caso de que los coeficientes

3.1.1. Ejemplos

Example 3.1 *Tenemos datos sobre coche que utilizan diesel o no.*

3.1.2. Ejercicios

4. Regresión lineal generalizada

4.1. Generalidades

Ejemplo: Consideremos la distribución de Poisson $Y \sim \text{Pois}(\mu)$:

Partimos de la función de distribución:

$$f(y; \mu) = e^{-\mu} \frac{\mu^y}{y!}$$

Utilizamos el, a falta de un nombre mejor, truco:

$$f(y; \mu) = e^{\log(f(y; \mu))}$$

$$\log(f(y; \mu)) = -\mu + y \log(\mu) - \log(y!)$$

$$f(y; \mu) = e^{\log(f(y; \mu))} = e^{-\mu + y \log(\mu) - \log(y!)}$$

Igualamos (tomado para cada muestra y_i y cada parámetro μ_i):

$$-\mu_i + y_i \log(\mu_i) - \log(y_i!) = \frac{y \cdot \theta_i - b(\theta_i)}{a(\theta)} + c(y; \theta)$$

Nos basta tomar lo siguiente:

$$\begin{aligned}\theta_i &= \log(\mu_i) \\ \phi &= 1\end{aligned}$$

$$\begin{aligned}a(\phi) &= \phi = 1 \\ b(\theta_i) &= \mu_i = e^{\theta_i} \\ c(y; \phi) &= -\log(y_i!)\end{aligned}$$

Ejemplo: Consideremos la distribución de Normal $Y \sim N(\mu_i, \sigma)$:

Partimos de la función de distribución:

$$f(y; \mu) =$$

Utilizamos la siguiente igualdad:

$$f(y; \mu) = e^{\log(f(y; \mu))}$$

Ejemplo: Consideremos la distribución Binomial (con parámetro n conocido).

Esto vendrá después

Ejemplo: Sabiendo que $E_\theta(Y) = b'(\theta) = \mu = n \frac{e^\theta}{1+e^\theta}$ para una distribución binomial calcular su link canónico ($\eta = g(\mu) = \theta$).

$$y = n \frac{e^\theta}{1+e^\theta}$$

$$\frac{y}{n} = \frac{e^\theta}{1+e^\theta}$$

$$\frac{y}{ne^\theta} = \frac{1}{1+e^\theta}$$

$$\frac{ne^\theta}{y} = 1 + e^\theta \Rightarrow \frac{ne^\theta}{n} - e^\theta = e^\theta \left(\frac{n}{y} - 1 \right) = e^\theta \left(\frac{n-y}{y} \right) = 1$$

$$\ln(1) = 0 = \theta + \ln \left(\frac{n-y}{y} \right)$$

$$\Rightarrow \theta = -\ln \left(\frac{n-y}{y} \right) = \ln \left(\frac{y}{n-y} \right)$$

4.2. Binomial Response Models

Una variable aleatorio es tal que $Y \sim B(p)$ (Bernulli), $0 \leq p \leq 1$, si y solo sí toma valores 1 ó 0 con las siguientes probabilidades:

$$P(Y = 1) = p \text{ and } P(Y = 0) = 1 - p$$

Una variable aleatoria es tal que $Y \sim \text{Bin}(n, p)$ (Binomial), con parámetros.... CONTINUAR AQUÍ

BALA BLA BLA BLA

Definimos los **odds** de una variable aleatoria Binomial como $\text{Odd} = \frac{p}{1-p} \in (0, +\infty)$, tal que verifica:

$$\text{Odd} = \begin{cases} 5 & \text{si } x \leq 2 \\ x^2 - 6x + 10 & \text{si } 2 < x < 5 \\ 4x - 15 & \text{si } x \geq 5 \end{cases}$$

SUSTITUIR APROPIADAMENTE LA MIERDA DE AQUÍ ARRIBA IMPORTANTE!!

Para comparar p_1 con $p_2 \in (-1, 1)$ CONTINUAR A1UÍ BLA BLA BLA BLA

$$\begin{cases} H_0 : p_1 = p_2 & \Longleftrightarrow H_0 : p_1 - p_2 = 0 \\ H_1 : p_1 \neq p_2 & \Longleftrightarrow H_1 : p_1 - p_2 \neq 0 \end{cases}$$

$$Y_i \sim \text{Bin}(m_i, p_i)$$

$$g(\mu) = X\beta \Leftrightarrow g(mp) = X\beta$$

Recordemos el link canónico, el parámetro de ddispersione y la función de varianza son respectivamente:

$$\theta_i = \log \left(\frac{p_i}{1 - p_i} \right)$$

$$\Phi = 1$$

$$V(\mu_i) = \mu_i \left(1 - \frac{\mu_i}{m_i} \right)$$

4.3. Poisson Response Models

La principal característica de estos modelos es que la variable respuesta sigue una distribución de Poisson.

$$Y_i \sim \text{Pois}(\mu_i)$$

$$g(\mu) = X\beta$$

Donde recordemos: $\text{Pois}(Y = k) = e^{-\lambda} \frac{\lambda^k}{k!}$, $E(Y) = \lambda$, $V(Y) = \lambda$.

El **índice de dispersión** será:

$$I(Y) = \frac{V(Y)}{E(Y)} = 1$$

— EXPLICAR MEJOR ESTA MIERDA Ejemplo de las placas de petri overdispersión: tendencia al agrupamiento Poisson: randomness Underdispersion: uno por casilla — Estimamos el parámetro $\Phi = 1$

- $\hat{\Phi} = \frac{X^2}{n-p} >> 1 \Rightarrow \text{overdispersión}$
- $<< 1 \Rightarrow \text{underdispersión}$

- ≈ 1 Poisson

In general $V(Y) = \Phi \cdot \mu$ quasi poisson lo usamos en overdispersion

Estimamos Φ con los datos and like this we will have a more accurate estimation of the true variance.

Una vez tengamos una tabla, si fijamos el parámetro n , lo que nos encontramos es una distribución multinomial. Si no lo fijamos seguimos contando con la distribucúin de poisson

Ejemplo de la tabla $\log(\mu_{ij}) = \beta_0 + \beta_1 \text{Factor 1} + \beta_2 \text{Factor 2}$, en total $1 + a - 1 + b - 1 = a + b - 1$. Tenemos factores, luego puede haber interacciones, sea pues el siguiente modelo a considerar:

$$\beta_0 + \beta_1 \text{Factor 1} + \beta_2 \text{Factor 2} + \beta_3 \text{Interacción}$$

con un total de $a + b - 1 + (a - 1)(b - 1) = \text{RELLENAR} = ab$ luego se corresponde con el modelo completo (*full model*), lo que no tiene sentido considerar.