

CAPÍTULO 2

MODELO DE REGRESIÓN LOGÍSTICA

2.1 INTRODUCCIÓN

La Regresión Logística es una técnica estadística multivariante que nos permite estimar la relación existente entre una variable dependiente no métrica, en particular dicotómica y un conjunto de variables independientes métricas o no métricas.

El Análisis de Regresión Logística tiene la misma estrategia que el Análisis de Regresión Lineal Múltiple, el cual se diferencia esencialmente del Análisis de Regresión Logística por que la variable dependiente es métrica; en la práctica el uso de ambas técnicas tienen mucha semejanza, aunque sus enfoques matemáticos son diferentes.

La variable dependiente o respuesta no es continua, sino discreta (generalmente toma valores 1,0). Las variables explicativas pueden ser cuantitativas o cualitativas; y la ecuación del modelo no es una función lineal de partida, sino exponencial; si bien, por sencilla transformación logarítmica, puede finalmente presentarse como una función lineal.

Así pues el modelo será útil en frecuentes situaciones prácticas de investigación en que la respuesta puede tomar únicamente dos valores: 1, presencia (con probabilidad p); y 0, ausencia (con probabilidad $1-p$).

El modelo será de utilidad puesto que, muchas veces, el perfil de variables puede estar formado por caracteres cuantitativos y cualitativos; y se pretende hacer participar a todos ellos en una única ecuación conjunta.

El modelo puede acercarse mas a la realidad ya que muchos fenómenos, como los del campo epidemiológico, se asemejan más a una curva que a una recta. Además la curva exponencial elegida como mejor ajuste, puede ser transformada logarítmicamente en una ecuación lineal de todas las

variables, siendo así que el aparato matemático estudiado para la regresión lineal múltiple será aplicable; aunque el investigador tenga, al final, que deshacer la transformación para interpretar sus conclusiones.

Si para el Modelo de Regresión Logística una variable regresora de tipo categórica tiene c niveles habrá que generar $c-1$ variables ficticias (dummy) a fin que todas las posibilidades de la variable queden bien representadas en el modelo logístico.

Cuando todas las variables regresoras son categóricas entonces se usa el modelo Log lineal, ver Mc Cullagh (1983).

2.2 OBJETIVOS DE LA REGRESIÓN LOGÍSTICA

El objetivo primordial de esta técnica es el de modelar como influyen las variables regresoras en la probabilidad de ocurrencia de un suceso particular.

Sistemáticamente tiene dos objetivos:

1. Investigar como influye en la probabilidad de ocurrencia de un suceso, la presencia o no de diversos factores y el valor o nivel de los mismos.
2. Determinar el modelo más parsimonioso y mejor ajustado que siendo razonable describa la relación entre la variable respuesta y un conjunto de variables regresoras.

2.3 REGRESIÓN LOGÍSTICA Y OTROS MÉTODOS RELACIONADOS

El objetivo general de la Regresión Logística es predecir la probabilidad de un evento de interés en una investigación, así como identificar las variables predictoras útiles para tal predicción.

Se pueden usar varios métodos multivariantes para predecir una variable respuesta de naturaleza dicotómica a partir de un grupo de variables regresoras.

El Análisis de Regresión Lineal Múltiple y el Análisis Discriminante son dos métodos eficaces pero plantean problemas cuando la variable respuesta es binaria.

En el Análisis de Regresión Lineal Múltiple cuando la variable respuesta toma solo dos valores, se violan los supuestos de necesarios para efectuar inferencias, los problemas que se plantean son:

1. La distribución de los errores aleatorios no es normal.
2. Los valores predictados no pueden ser interpretados como probabilidades como en la Regresión Logística, porque no toman valores dentro del intervalo $[0,1]$.

El Análisis Discriminante permite la predicción de pertenencia de la unidad de análisis a uno de los dos grupos pre-establecidos, pero se requiere que se cumplan los supuestos de multinormalidad de las variables regresoras y la igualdad de matrices de covarianzas de los dos grupos, pueden ser diferentes también; para que la regla de predicción sea óptima, Johnson (1982).

La Regresión Logística requiere mucho menos supuestos que el AD, por ello cuando satisfacen los supuestos requeridos para el AD, la Regresión Logística trabaja bien.

A continuación se describirá un paralelo entre la Regresión Lineal Múltiple y la Regresión Logística, debido a que ambos tienen el mismo objetivo, predecir la variable respuesta a partir de las variables regresoras.

2.4 REVISIÓN DEL MODELO DE REGRESIÓN LINEAL MÚLTIPLE

La diferencia básica entre los modelos del Análisis de Regresión Lineal Múltiple y de la Regresión Logística es naturaleza de la relación entre la variable respuesta y las variables regresoras.

Para el Análisis de Regresión Lineal Múltiple, consideremos y una variable respuesta cuantitativa y x_1, x_2, \dots, x_k variables regresoras o llamadas también explicativas; y se desea describir la relación que hay entre la variable respuesta y las variables explicativas, si entre la variable respuesta y las regresoras hay una relación lineal se espera que:

$$E(y_i) = \mathbf{b}_0 + \mathbf{b}_1 x_{i1} + \mathbf{b}_2 x_{i2} + \dots + \mathbf{b}_k x_{ik} \quad , (2.1)$$

para $i=1, 2, \dots, n$

donde:

y_i es el valor de la variable respuesta cuantitativa para el i -ésimo objeto.

$\mathbf{b}_j ; j = 0, 1, 2, \dots, k$ son los parámetros.

Siendo n el número de objetos u observaciones.

Aunque (2.1) no de valores exactos, se espera que varíe linealmente con las variables regresoras, esto es:

$$E(y_i | \mathbf{x}_i^0) = \mathbf{b}_0 + \mathbf{b}_1 x_{i1} + \mathbf{b}_2 x_{i2} + \dots + \mathbf{b}_k x_{ik} ,$$

para $i=1, 2, \dots, n$

(2.2)

siendo $\mathbf{x}_i^0 = (x_{i0}, x_{i1}, x_{i2}, \dots, x_{in})$ la i -ésima observación, con $x_{i0} = 1$,

(2.2) toma valores reales y en forma vectorial es:

$$E(y_i | \mathbf{x}_i^0) = \mathbf{x}_i^0 T \mathbf{b}^0 \quad (2.3)$$

donde $\mathbf{b}^0 = (\mathbf{b}_0, \mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k)$ es el vector de parámetros.

Pero en (2.3) hay otras variables regresoras que pueden influir linealmente sobre y_i , por tanto cada valor de y_i está variando alrededor de $E(y_i)$ a esa variación lo denotamos con e_i , esto es:

$$\begin{aligned} e_i &= y_i - E(y_i | x_i) \\ &= y_i - x_i^T \beta \end{aligned} \quad (2.4)$$

de (2.4):

$$y_i = x_i^T \beta + e_i \quad (2.5)$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + e_i, \quad \text{para } i=1, 2, \dots, n \quad (2.6)$$

llamado Modelo de Regresión Lineal Múltiple poblacional, a e_i se le llama error aleatorio y tiene las siguientes propiedades:

$$\begin{aligned} E(e_i) &= 0 \\ V(e_i) &= \sigma^2 \\ Cov(e_i, e_j) &= 0 \quad " i \neq j \quad (2.7) \\ Cov(e_i, X_j) &= 0 \end{aligned}$$

las variables regresoras no son variables aleatorias y el comportamiento de y es la respuestas a aquellas, así mismo e_i es una variable aleatoria no observable.

Generalizando el Modelo de Regresión Lineal Múltiple, (2.6), mediante el álgebra matricial está dada por:

$$\hat{y} = X\beta + e \quad (2.8)$$

donde:

$$\begin{aligned} \hat{y}^T &= (y_1, y_2, \dots, y_n), \text{ vector de variables respuestas observadas} \\ X &= (1, x_1, x_2, \dots, x_k) \text{ matriz de rango completo y con} \end{aligned}$$

$$x_i^T = (1, x_{i1}, x_{i2}, \dots, x_{ik})$$

$$\mathbf{b}^T = (\mathbf{b}_0, \mathbf{b}_1, \dots, \mathbf{b}_k)$$

$$\mathbf{e}^T = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n)$$

así mismo (2.1) en forma matricial es:

$$E(\hat{y}) = X\mathbf{b} \quad (2.9)$$

el objetivo es estimar los parámetros del modelo (2.6), los mismos que son estimados mediante el método de mínimos cuadrados.

Sea \hat{y}_i la estimación de y_i , entonces:

$$\hat{y}_i = \hat{\mathbf{b}}_0 + \hat{\mathbf{b}}_1 x_{i1} + \hat{\mathbf{b}}_2 x_{i2} + \dots + \hat{\mathbf{b}}_k x_{ik},$$

$$\text{para } i=1, 2, \dots, n \quad (2.10)$$

o equivalentemente:

$$\hat{y}_i = x_i^T \hat{\mathbf{b}} \quad (2.11)$$

siendo:

$$\hat{\mathbf{b}}^T = (\hat{\mathbf{b}}_0, \hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_k), \text{ el vector de parámetros estimados.}$$

(2.10) en su forma matricial es:

$$\hat{\mathbf{y}} = X\hat{\mathbf{b}} \quad (2.12)$$

los residuos ordinarios r_i es la contraparte muestral de \mathbf{e}_i y está dado por:

$$r_i = y_i - \hat{y}_i \text{ para } i=1, 2, \dots, n \quad (2.13)$$

en forma vectorial es:

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}} \quad (2.14)$$

El estimador de \mathbf{b} se obtiene usando el método de mínimos cuadrados, ver que consiste en minimizar la suma de cuadrados del error y está dada por:

$$SCE = \sum_{i=1}^n e_i^2 = \mathbf{e}^T \mathbf{e} = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})$$

con respecto a \mathbf{b} , esta suma de cuadrados se expresa en forma cuadrática como::

$$(\mathbf{y} - X\mathbf{b})^T (\mathbf{y} - X\mathbf{b}) \quad (2.15)$$

al minimizarlo se obtiene que:

$$(X^T X)\mathbf{b} = X^T \mathbf{y} \quad (2.16)$$

llamada ecuaciones normales.

Como $(X^T X)$ es invertible, por que es simétrica de tamaño $(k+1) \times (k+1)$ y de rango completo, entonces la solución del sistema lineal es:

$$\mathbf{b} = (X^T X)^{-1} X^T \mathbf{y} \quad (2.17)$$

que es sensible a observaciones pobremente ajustados y a los puntos extremos de X , Montgomery y Peck (1992).

El vector (2.12) de valores estimados para el vector de variables respuesta es:

$$\hat{\mathbf{y}} = X\mathbf{b}$$

$$\hat{\mathbf{y}} = X(X^T X)^{-1} X^T \mathbf{y} \quad (2.18)$$

a la matriz $X(X^T X)^{-1} X^T$, se le llama matriz de cambio o de proyección denotada por H , entonces (2.18) es:

$$\hat{\mathbf{y}} = H\mathbf{y} \quad (2.19)$$

El vector de residuos es:

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}}$$

$$\mathbf{r} = \mathbf{y} - H\mathbf{y} \quad (2.20)$$

$$\mathbf{r} = (I - H)\mathbf{y} \quad (2.21)$$

$$\mathbf{v} = M\mathbf{y} \quad (2.22)$$

El vector \mathbf{v} describe las desviaciones de los valores observados de los ajustados y la matriz M es el subespacio en el cual cae \mathbf{v} .

El vector residual es importante para detectar puntos 'extraños'. A la matriz H se le llama matriz sombrero o de proyección, ver Cook y Weisberg (1982). Ahora veamos como queda expresado la suma de cuadrados de los residuos, denotada por SCE :

$$SCE = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) \quad (2.23)$$

Reemplazando $\hat{\mathbf{y}}$ por $X\hat{\mathbf{b}}$: (2.23) es:

$$SCE = (\mathbf{y} - X\hat{\mathbf{b}})^T (\mathbf{y} - X\hat{\mathbf{b}}) \quad (2.24)$$

Y reemplazando $\hat{\mathbf{b}}$ por $(X^T X)^{-1} X^T \mathbf{y}$:

$$SCE = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T X (X^T X)^{-1} X^T \mathbf{y} \quad (2.25)$$

La matriz $H = X (X^T X)^{-1} X^T$, entonces:

$$SCE = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T H \mathbf{y} \quad (2.26)$$

$$SCE = \mathbf{y}^T (I - H) \mathbf{y} \quad (2.27)$$

Sea $M = I - H$, entonces (2.27) es:

$$SCE = \mathbf{y}^T M \mathbf{y} \quad (2.28)$$

Bajo el supuesto que $\mathbf{e} \sim N(\mathbf{0}, \mathbf{S}^2 I_n)$, las observaciones y_1, y_2, \dots, y_n son independientes y distribuidas como una Normal n-variada con $E(\mathbf{y}) = X\mathbf{b}$ con matriz de varianzas y covarianzas $\mathbf{S}^2 I_n$.

En la Regresión Logística, se estima la probabilidad de que un evento ocurra; es decir, el valor esperado de y y dado las variables regresoras, debe tomar valores entre 0 y 1. La relación entre las variables regresoras y la dependiente no es lineal. Las estimaciones de probabilidad estarán siempre entre 0 y 1, así, el valor de la variable respuesta se puede definir como una probabilidad de que ocurra o no un evento sujeto a control.

En la Regresión Logística, se seleccionan los coeficientes, del modelo, que hacen que los resultados sean los más “probables”. Como el modelo de Regresión Logística no es lineal, se requiere de un algoritmo iterativo para estimar los parámetros.

En las secciones siguientes se detallarán los aspectos teóricos y la aplicación de la Regresión Logística.

2.5 REGRESIÓN LOGÍSTICA SIMPLE

Este modelo tiene la forma

$$y_i = b_0 + b_1x_i + e_i \quad \text{para } i = 1, 2, \dots, n \quad (2.29)$$

De esto se deduce que:

$$\text{Si } y = 1, e_i = 1 - b_0 - b_1x_i \quad (2.30)$$

$$\text{Si } y = 0, e_i = -b_0 - b_1x_i \quad (2.31)$$

Por tanto e_i , no puede tener distribución normal debido a que toma valores discretos, el Modelo de Regresión Lineal Simple, no es aplicable para el caso de variable respuesta de tipo dicotómico.

En el Análisis de Regresión Lineal simple, el punto inicial del proceso de estimación del modelo es un gráfico de dispersión de la variable respuesta versus la regresora, pero este gráfico resulta limitado cuando sólo hay dos valores posibles para la variable respuesta, por tanto se debe usar otros gráficos, éstos resultan de la suavización de los valores de la variable respuesta, representando después los valores de la variable respuesta versus la regresora.

La notación que se usará en el presente trabajo para la Regresión Logística es misma que emplea Hosmer y Lemeshow (2000).

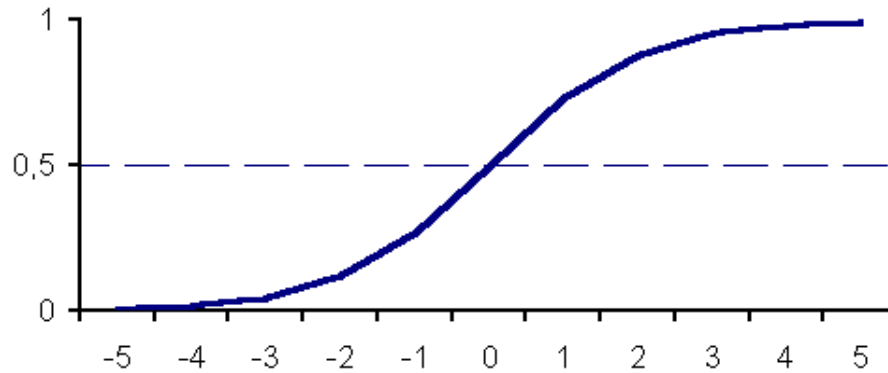
$$\text{Sea } p(x) = E(y|x) \quad (2.32)$$

Que representa la media condicional de $y = 1$ dado x , donde $p(x)$ representa la probabilidad de que ocurra $y = 1$, ciertamente no se espera que (2.32) tenga relación lineal dentro del rango de la variable regresora.

¿Qué hay de la relación entre $p(x)$ y x para valores intermedios de x ? Se espera una relación curvilínea. Para cualquier valor grande de x , $p(x)$ tomará valores cercanos a 1 y para valores pequeños de x , $p(x)$ tomará

valores cercanos a cero. El gráfico que muestra el comportamiento de $p(x)$ versus x es:

FIGURA N° 2.1



curva en forma de S o sigmoide que tiene las propiedades requeridas para $p(x)$ y que tiene las propiedades de una función de distribución de probabilidad acumulada, para esta probabilidad se usa la función de distribución acumulada de la distribución logística dada por:

$$p(x) = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}} \quad (2.33)$$

(2.33) tiene un gráfico similar a la Figura N° 2.1, cuando $b_0 < 0$ y $b_1 > 0$, además este modelo toma valores en el intervalo $[0,1]$.

Cuando $P[y = 1] = 0.5$ el valor de x es: $-\frac{b_0}{b_1}$, que brinda información

muy útil.

Una transformación de $p(x)$ que es lo central del estudio de la Regresión Logística es la transformación logit, transformación que se define en términos $p(x)$ como:

$$g(x) = \ln \frac{p(x)}{1-p(x)} \quad (2.34)$$

$$= b_0 + b_1 x$$

Lo importante de esta transformación es que tiene muchas propiedades semejantes al Modelo de Regresión Lineal simple, por ejemplo es lineal en sus parámetros, puede ser continua y puede tomar cualquier valor real dependiendo de x .

Para el Modelo de Regresión Lineal simple, la variable respuesta, de (2.4) se expresa como:

$$y = E(y|x) + e \quad (2.35)$$

para la variable respuesta dicotómica lo expresamos como:

$$y = p(x) + e \quad (2.36)$$

veamos que ocurre con este modelo:

Si $y = 1$, $e_i = 1 - p(x)$ y tiene probabilidad $p(x)$

Si $y = 0$, $e_i = -p(x)$ y tiene probabilidad $1 - p(x)$

Entonces e_i tiene distribución binomial con media cero y varianza $p(x)[1 - p(x)]$. Por tanto la distribución condicional de la variable respuesta tiene distribución de probabilidad binomial con media $p(x)$.

El lado izquierdo de (2.34) se llama también logaritmo de ODDS RATIO o razón de probabilidades de $y = 1$ contra $y = 0$, específicamente:

$$ODDS \ RATIO = \frac{p(x)}{1 - p(x)} \quad (2.37)$$

o también llamado razón de ventaja a favor de éxito.

2.6 REGRESIÓN LOGÍSTICA MÚLTIPLE

En esta sección se generaliza el Modelo de Regresión Logística Simple tratado en la sección anterior, es decir consideraremos más de una variable regresora, en donde por lo menos una es de tipo cuantitativo.

2.6.1 MODELO DE REGRESIÓN LOGÍSTICA MÚLTIPLE

Sea el vector de variables regresoras $x^T = (x_1, x_2, \dots, x_k)$ por el momento asumiremos que están medidas por lo menos bajo escala intervalar. Sea la probabilidad condicional para que la variable respuesta sea igual a 1, denotado por:

$$P(y = 1 | x) = p(x) \quad (2.39)$$

el logaritmo del Modelo de Regresión Logística Múltiple está dado por:

$$g(x_i) = b_0 + b_1 x_{i1} + b_2 x_{i2} + \dots + b_k x_{ik}, \text{ para } i=1, 2, \dots, n \quad (2.40)$$

entonces el Modelo de Regresión Logística Múltiple es:

$$p(x) = \frac{e^{g(x)}}{1 + e^{g(x)}} \quad (2.41)$$

Al igual que en el caso de Regresión Lineal Múltiple si es necesario usar variables regresoras categóricas, por ejemplo si una tiene c niveles será necesario incorporar c-1 variables ficticias o llamadas dummy., así entonces el logit para un modelo con k variables regresoras y una variable categórica, por ejemplo la j-ésima es:

$$g(x) = b_0 + b_1 x_{i1} + \dots + \sum_{l=1}^{c-1} b_{jl} D_{jl} + b_k x_{ik},$$

para $i=1, 2, \dots, n$

(2.42)

2.6.2 AJUSTE DEL MODELO DE REGRESIÓN LOGÍSTICA MÚLTIPLE

El ajuste se efectúa a través del uso de los métodos de máxima verosimilitud, los mismos que se encuentran en los softwares estadísticos que permiten analizar datos mediante este método.

Asumiremos que disponemos de una muestra n observaciones independientes

$(x_i, y_i), i=1,2, \dots, n$; donde y_i toma valores 0 ó 1, para estimar

$\mathbf{b}^T = (b_0, b_1, \dots, b_k)$ que es el vector de parámetros desconocidos.

Para el Modelo de Regresión Lineal Múltiple se usa el método de Mínimos Cuadrados para estimar \mathbf{b} , el cual minimiza la suma de cuadrados del error, pero cuando la variable respuesta es binaria aplicar este método no provee las mismas propiedades cuando es usado en variables respuestas continuas.

Por ello se usará el método de Máxima Verosimilitud, ya que obtendremos parámetros estimados que maximizan la probabilidad de obtener un conjunto de datos observados.

La función de verosimilitud expresa la probabilidad de los datos observados como una función de parámetros desconocidos. Los Estimadores de Máxima Verosimilitud de esos parámetros son aquellos que están en concordancia con los datos observados.

Consideremos el Modelo de Regresión Lineal Múltiple con mayor detalle, supongamos que se dispone de n objetos u observaciones donde para cada uno de ellos existe una respuesta que puede ser:

$$y_i = 0 \quad \text{o} \quad y_i = 1$$

Sea $\mathbf{y}^T = (y_1, y_2, \dots, y_n)$ donde $y_i \sim B(1, p_i)$ y sea

$x_i^T = (1, x_{i1}, \dots, x_{ik})$ la i -ésima observación para las k variables explicativas.

Así el Modelo de Regresión Logística está dada por la expresión (2.40):

$$P[y_i = 1 | x_i] = p(x_i) = \frac{e^{g(x_i)}}{1 + e^{g(x_i)}} \quad (2.43)$$

equivalentemente

$$P[y_i = 1 | x_i] = \frac{\exp\left\{b_0 + \sum_{j=1}^k b_j x_{ij}\right\}}{1 + \exp\left\{b_0 + \sum_{j=1}^k b_j x_{ij}\right\}} \quad (2.44)$$

y la probabilidad de que y_i sea igual a cero es:

$$P[y_i = 0 | x_i] = 1 - P[y_i = 1 | x_i] \text{ entonces :}$$

$$P[y_i = 0 | x_i] = \frac{1}{1 + \exp\left\{b_0 + \sum_{j=1}^k b_j x_{ij}\right\}} \quad (2.45)$$

para facilitar la notación usaremos la variable indicadora $x_{i0} = 1, i = 1, 2, \dots, n$.

Entonces (2.44) y (2.45) son respectivamente:

$$P[y_i = 1 | x_i] = p(x_i) = \frac{e^{x_i^T b}}{1 + e^{x_i^T b}} \quad (2.46)$$

$$P[y_i = 0 | \mathbf{x}_i] = 1 - p(\mathbf{x}_i) = \frac{1}{1 + e^{\mathbf{b}^T \mathbf{x}_i}} \quad (2.47)$$

donde: $\mathbf{x}_i^T = (x_{i0}, x_{i1}, \dots, x_{ik})$, es el vector que contiene los valores de las variables explicativas

$\mathbf{b}^T = (\mathbf{b}_0, \mathbf{b}_1, \dots, \mathbf{b}_k)$ es el vector de parámetros a ser estimado.

El i-ésimo logito es:

$$l_i = \text{Ln} \frac{p_i}{1 - p_i} = \sum_{j=0}^k \mathbf{b}_j x_{ij} \quad (2.48)$$

como vemos, (2.48) es una función lineal simple del vector de observaciones \mathbf{x}_i llamada transformación logística de la probabilidad p_i o simplemente Logit o Logito de la ecuación, a la expresión (2.48) también se le llama Modelo Logístico Lineal.

A fin de obtener la estimación máximo verosímil para el vector \mathbf{b} , escribimos la función de densidad de probabilidad del vector \mathbf{y} el cual es proporcional a n funciones $B(1, p_i)$, esto es:

$$\begin{aligned} f(y_i; \mathbf{p}_i) &= \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1 - y_i} \\ &= \prod_{i=1}^n \left(\frac{p_i}{1 - p_i} \right)^{y_i} (1 - p_i) \\ &= \left\{ \prod_{i=1}^n (1 - p_i) \right\} \left\{ \prod_{i=1}^n \text{Exp} \left[\text{Ln} \left(\frac{p_i}{1 - p_i} \right)^{y_i} \right] \right\} \\ &= \left\{ \prod_{i=1}^n (1 - p_i) \right\} \text{Exp} \left[\sum_{i=1}^n y_i \text{Ln} \left(\frac{p_i}{1 - p_i} \right) \right] \end{aligned} \quad (2.49)$$

Reemplazando (2.48) en (2.49), se obtiene:

$$\begin{aligned}
 f(y_i; \mathbf{p}_i) &= \prod_{i=1}^n (1 - p_i) \exp\left(\sum_{j=0}^k b_j y_i x_{ij}\right) \\
 &= \prod_{i=1}^n (1 - p_i) \exp\left(\sum_{j=0}^k b_j y_i x_{ij}\right)
 \end{aligned} \tag{2.50}$$

El logaritmo natural de la función (2.50), llamado función soporte es:

$$l(\mathbf{p}_i; y_i) = \sum_{j=0}^k \left(\sum_{i=1}^n y_i x_{ij} \right) b_j + \sum_{i=1}^n \ln(1 - p_i) \tag{2.51}$$

pero (2.47) : $1 - p_i = \left[1 + \text{Exp} \left(\mathbf{b}^T \mathbf{x}_i \right) \right]^{-1}$, entonces
 $\text{Ln} (1 - p_i) = -\text{Ln} \left[1 + \text{Exp} \left(\mathbf{b}^T \mathbf{x}_i \right) \right]$

$$\text{Ln} (1 - p_i) = -\text{Ln} \left[1 + \text{Exp} \left(\sum_{j=0}^k b_j x_{ij} \right) \right] \tag{2.52}$$

reemplazando (2.52) en (2.51), se obtiene:

$$l(\mathbf{p}_i; y_i) = \sum_{j=0}^k \left(\sum_{i=1}^n y_i x_{ij} \right) b_j - \sum_{i=1}^n \text{Ln} \left[1 + \text{Exp} \sum_{j=0}^k b_j x_{ij} \right] \tag{2.53}$$

como vemos (2.53) es una función que ya no depende de \mathbf{p}_i sino de \mathbf{b}_j solamente, entonces lo denotamos como:

$$L(\mathbf{b}) = \sum_{j=0}^k \left(\sum_{i=1}^n y_i x_{ij} \right) b_j - \sum_{i=1}^n \text{Ln} \left[1 + \text{Exp} \left(\sum_{j=0}^k b_j x_{ij} \right) \right] \tag{2.54}$$

es una función que depende exclusivamente del vector \mathbf{b} .

Definamos como:

$$t_j = \sum_{i=1}^n y_i x_{ij} \quad (2.55)$$

entonces reemplazando (2.55) en (2.54) se tiene:

$$L(\mathbf{b}) = \sum_{j=0}^k \mathbf{b}_j t_j - \sum_{i=1}^n \text{Ln} \left[1 + \text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right) \right] \quad (2.56)$$

Como (2.56) es una función exclusiva del vector de parámetros \mathbf{b} , por el Teorema de Factorización de Fisher-Neyman, Bickel y Doksum (1976), se tiene que t_j para $j = 0, 1, \dots, k$ son estadísticas suficientes para los parámetros \mathbf{b}_j , para $j = 0, 1, \dots, k$,

La variable aleatoria t_j dada en la expresión (2.56) es la suma de algunos de los términos de la matriz de diseño X , es decir se incluyen en la suma solamente los elementos que corresponden a una respuesta del tipo $y = 1$.

Las ecuaciones de verosimilitud, se obtienen derivando (2.54) con respecto a los elementos de \mathbf{b} e igualando a cero:

$$\frac{\partial L}{\partial \mathbf{b}_j} = \sum_{i=1}^n y_i x_{ij} - \sum_{i=1}^n x_{ij} \left[\frac{\text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right)}{1 + \text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right)} \right] \quad (2.57)$$

las ecuaciones de verosimilitud de (2.57) son:

$$\sum_{i=1}^n y_i x_{ij} - \sum_{i=1}^n x_{ij} \hat{p}_i = 0 \quad j = 0, 1, 2, \dots, k \quad (2.58)$$

siendo $x_{i0} = 1$, equivalentemente (2.58) es:

$$\sum_{i=1}^n x_{ij} (y_i - \hat{p}_i) = 0 \quad j = 0, 1, 2, \dots, k \quad (2.59)$$

donde:

$$\hat{p}_i = \frac{\text{Exp} \left(\sum_{j=0}^k \hat{b}_j x_{ij} \right)}{1 + \text{Exp} \left(\sum_{j=0}^k \hat{b}_j x_{ij} \right)} ; \text{ para } i=1, 2, \dots, n$$

es el estimador máximo verosímil de \mathbf{p}_i y se obtiene mediante \hat{b}_j y el vector \mathbf{x}_i

La expresión (2.58) en su forma matricial es:

$$X^T (\mathbf{y} - \mathbf{p}) = X\mathbf{b} = \mathbf{0} \quad (2.60)$$

Estas ecuaciones son parecidas a las ecuaciones normales obtenidas para estimar el Modelo de Regresión Lineal Múltiple, pero son no lineales en \mathbf{b} , lo cual hace que se use un método iterativo para determinar los valores del vector \mathbf{b} .

La obtención de \hat{b}_j mediante métodos iterativos; para $j = 0, 1, \dots, k$ se tratará en la siguiente sección, ahora obtendremos la varianza y covarianza de \mathbf{b} .

Sea $X_{(n \times p)}$ la matriz de diseño, con $p=k+1$, con elementos:

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}$$

Las ecuaciones de verosimilitud en su forma matricial, de la expresión (2.60):

$$X^T \hat{y} = X^T \hat{\mathbf{p}} \quad (2.61)$$

donde $\hat{\mathbf{p}}^T = (\hat{p}_1, \hat{p}_2, \dots, \hat{p}_n)$

$$\sum_{i=1}^n x_{ij} (y_i - p_i) = 0 \quad (2.62)$$

El método de estimación de las varianzas y covarianzas lo obtendremos de la matriz de segunda derivada parcial de (2.62); y tiene la forma:

$$\frac{\partial^2 L}{\partial \mathbf{b}_j^2} = - \sum_{i=1}^n x_{ij}^2 p_i (1 - p_i) \quad (2.63)$$

para $j=0, 1, 2, \dots, k$

reemplazando: la ecuación para p_i en (2.63)

$$\frac{\partial^2 L}{\partial \mathbf{b}_j^2} = - \sum_{i=1}^n \frac{x_{ij}^2 \text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right)}{\left[1 + \text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right) \right]^2} \quad (2.64)$$

para $j = 0, 1, \dots, k$

$$\frac{\partial^2 L}{\partial \mathbf{b}_j \partial \mathbf{b}_l} = - \sum_{i=1}^n x_{ij} x_{il} \mathbf{p}_i (1 - \mathbf{p}_i) \quad (2.65)$$

para $j, l = 0, 1, 2, \dots, k$

reemplazando:

$$\frac{\partial^2 L}{\partial \mathbf{b}_j \partial \mathbf{b}_l} = - \sum_{i=1}^n x_{ij} x_{il} \frac{\text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right)}{\left[1 + \text{Exp} \left(\sum_{j=0}^k \mathbf{b}_j x_{ij} \right) \right]^2} \quad (2.66)$$

Tanto (2.64) como (2.65) no son funciones de y_i , entonces la matriz de observación y la matriz de segunda derivada esperada son idénticas.

Ahora bien la matriz que contiene el negativo de las ecuaciones (2.64) y (2.66) se denota con $I(\hat{\mathbf{b}})$, llamada Matriz de Información; las varianzas y covarianzas de $\hat{\mathbf{b}}_j$ se obtienen tomando la inversa de esta matriz, esto es:

$$\text{Cov}(\hat{\mathbf{b}}) = I^{-1}(\hat{\mathbf{b}}) \quad (2.67)$$

ver Cordeiro (1992).

Los estimadores de la varianza y covarianza, denotada por $\hat{\text{Cov}}(\hat{\mathbf{b}})$, se

obtiene evaluando

$$\text{Cov}(\hat{\mathbf{b}}) \text{ en } \hat{\mathbf{b}}.$$

Entonces la matriz de información estimada, matricialmente tiene la forma:

$$\hat{I}(\hat{\mathbf{b}}) = X' V X \quad (2.68)$$

V es una matriz diagonal, esto es:

$$V = \text{Diag}[\hat{p}_i (1 - \hat{p}_i)]$$

de tamaño $n \times n$, además (2.68) es:

$$\hat{Cov}(\hat{\mathbf{b}}) = (X'VX)^{-1} \quad (2.69)$$

y es de tamaño $(k+1)(k+1)$

escribiremos los elementos de la matriz (2.69)

$$\hat{Cov}(\hat{\mathbf{b}}) = \begin{bmatrix} \hat{s}^2(\hat{\mathbf{b}}_0) & \hat{s}(\hat{\mathbf{b}}_0, \hat{\mathbf{b}}_1) & \dots & \dots & \hat{s}(\hat{\mathbf{b}}_0, \hat{\mathbf{b}}_k) \\ \cdot & \hat{s}(\hat{\mathbf{b}}_1) & \dots & \dots & \hat{s}(\hat{\mathbf{b}}_1, \hat{\mathbf{b}}_k) \\ \cdot & \cdot & \dots & \dots & \cdot \\ \cdot & \cdot & \dots & \dots & \cdot \\ \cdot & \cdot & \dots & \dots & \hat{s}^2(\hat{\mathbf{b}}_k) \end{bmatrix}$$

donde:

$\hat{s}^2(\hat{\mathbf{b}}_j)$ es la varianza estimada de $\hat{\mathbf{b}}_j$

$\hat{s}(\hat{\mathbf{b}}_j, \hat{\mathbf{b}}_l)$ es la covarianza estimada de $\hat{\mathbf{b}}_j$ y $\hat{\mathbf{b}}_l$

$\hat{s}(\hat{\mathbf{b}}_j)$ es el error estandar de $\hat{\mathbf{b}}_j$

La matriz (2.69) será muy útil cuando se discuta el ajuste y la evaluación del Modelo de Regresión Logística.

2.6.3 MÉTODO DE NEWTON – RAPHSON PARA ESTIMAR LOS PARÁMETROS DEL MODELO DE REGRESIÓN LOGÍSTICA.

Este es un método para resolver ecuaciones no lineales, como las obtenidas en (2.57) o equivalentemente en (2.58), y requieren una solución mediante métodos iterativos para hallar la estimación de los parámetros que es el máximo de la función (2.54).

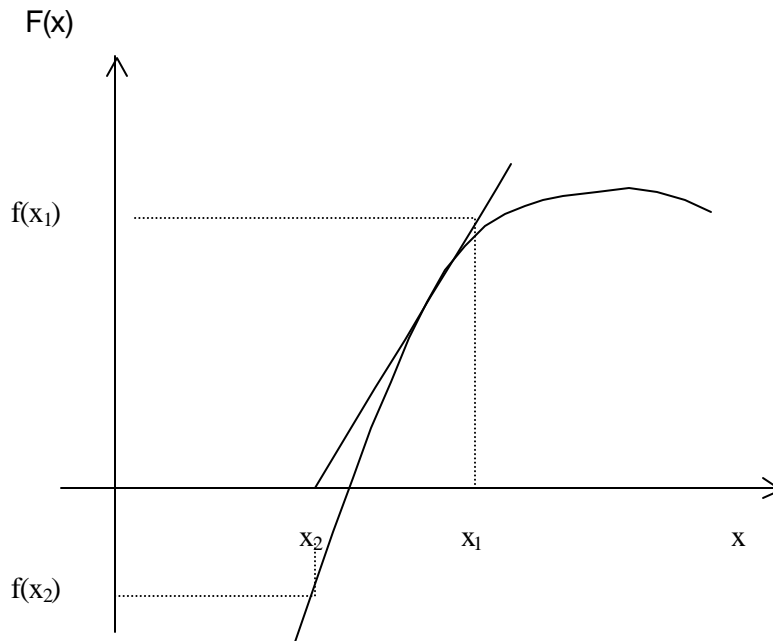
Uno de los métodos más usados para resolver ecuaciones de este tipo, es el de Newton-Raphson, porque converge rápidamente. En la figura Nº 2.2 se ilustra el método.

Tomando como estimación inicial x_1 , se prolonga la tangente a la curva en este punto hasta interceptar con el eje de las abscisas al cual llamaremos punto x_2 , entonces se toma a x_2 como la siguiente aproximación.

Este proceso continua hasta que un valor de x_2 haga que la función sea nula o suficientemente cercana a cero.

Para la estimación del vector \hat{b} se quiere hallar el máximo de una función; entonces usa la primera derivada, porque se anula en el punto máximo y la segunda derivada para calcular las tangentes. En nuestro caso es necesario hallar la segunda derivada para obtener la matriz de varianzas y covarianzas de los parámetros estimados.

Figura Nº 2.2 Interpretación Geométrica del Método Iterativo de Newton Raphson.



Entonces se usa el siguiente esquema iterativo:

$$\hat{\mathbf{b}}^{(t+1)} = \hat{\mathbf{b}}^{(t)} + \left[\mathbf{I}(\hat{\mathbf{b}}^{(t)}) \right]^{-1} \mathbf{S}(\hat{\mathbf{b}}^{(t)}) \quad (2.70)$$

donde:

$\mathbf{S}(\hat{\mathbf{b}})$ y $\mathbf{I}(\hat{\mathbf{b}})$ son las funciones de Score y de Información respectivamente.

La función Score es un vector de tamaño $k+1$, donde el j ésimo elemento de acuerdo a (2.57) es:

$$\frac{\partial L}{\partial \mathbf{b}_j} = \sum_{i=1}^n (y_i - \mathbf{p}_i^{(t)}) x_{ij} \quad (2.71)$$

La cual es similar a la expresión (2.59):

$$\sum_i x_{ij} (y_i - \mathbf{p}_i) = 0 \quad j = 1, 2, \dots, k$$

La Función de información es una matriz de tamaño $(k+1)(k+1)$ donde el ij ésimo elemento (l, j) es:

$$\begin{aligned} \frac{\partial^2 l}{\partial \mathbf{b}_j \partial \mathbf{b}_l} &= -\frac{\partial}{\partial \mathbf{b}_l} \left[\sum_{i=1}^n x_{ij} (y_i - \mathbf{p}_i) \right] \\ &= -\frac{\partial}{\partial \mathbf{b}_l} \left[\sum_{i=1}^n x_{ij} y_i - \sum_{i=1}^n x_{ij} \frac{e^{\mathbf{p}_i^T \mathbf{x}_i}}{1 + e^{\mathbf{p}_i^T \mathbf{x}_i}} \right] \\ &= \sum_{i=1}^n x_{ij} \left[\frac{e^{\mathbf{p}_i^T \mathbf{x}_i} x_{il} (1 + e^{\mathbf{p}_i^T \mathbf{x}_i}) - e^{\mathbf{p}_i^T \mathbf{x}_i} x_{il} e^{\mathbf{p}_i^T \mathbf{x}_i}}{(1 + e^{\mathbf{p}_i^T \mathbf{x}_i})^2} \right] \\ &= \sum_{i=1}^n \frac{x_{ij} x_{il} e^{\mathbf{p}_i^T \mathbf{x}_i}}{(1 + e^{\mathbf{p}_i^T \mathbf{x}_i})^2} \\ &= \sum_{i=1}^n x_{ij} x_{il} \mathbf{p}_i (1 - \mathbf{p}_i) \quad j=0, 1, \dots, k \quad ; \quad l = 0, 1, \dots, k \end{aligned}$$

(2.72)

donde $\mathbf{p}^{(t)}$, es la t-ésima aproximación para \mathbf{p} , obtenida de $\mathbf{b}^{(t)}$ mediante:

$$p_i^{(t)} = \frac{\text{Exp}\left(\sum_{j=0}^k b_j^{(t)} x_{ij}\right)}{\left[1 + \text{Exp}\left(\sum_{j=0}^k b_j^{(t)} x_{ij}\right)\right]} \quad (2.73)$$

Entonces el próximo valor reemplazando en (2.70) es:

$$\mathbf{b}^{(t+1)} = \mathbf{b}^{(t)} - \left\{X^T V^{(t)} X\right\}^{-1} X^T (\mathbf{y} - \mathbf{p}^{(t)}) \quad (2.74)$$

donde $V^{(t)} = \text{Diag}\left[p_i^{(t)}(1 - p_i^{(t)})\right]$

La expresión (2.70) se usa para obtener $\mathbf{p}^{(t+1)}$ y así sucesivamente. Después de dar un valor inicial $\mathbf{b}^{(0)}$, se usa (2.70) para obtener $\mathbf{p}^{(0)}$ y para $t > 0$ las iteraciones siguientes se efectúan usando (2.70) y (2.71).

En el límite, $\mathbf{p}^{(t)}$ y $\mathbf{b}^{(t)}$ converge a los EMV de \mathbf{p} y converge en general en 5 o 6 iteraciones.

Existen software estadísticos como el SAS y el SPSS con programas para estimar una regresión logística usando el método descrito. Una ventaja de este método es que en el paso final del proceso iterativo se obtiene la inversa de la función de información, que es asintóticamente la matriz de varianzas y covarianzas del vector $\mathbf{b}^{(t)}$ y permiten efectuar inferencias sobre los parámetros basado en la teoría normal. Para mayor información se recomienda a Affifi y Clark (1998).

2.6.4 INTERPRETACIÓN DE LOS COEFICIENTES DEL MODELO ESTIMADO

Recordamos del modelo de regresión múltiple que el valor de un coeficiente significaba el cambio en unidades de la variable dependiente por cada unidad de la variable independiente a que se refiere el coeficiente, permaneciendo invariantes los valores del resto de variables independientes del modelo.

A nivel de coeficientes estimados exponencialmente la interpretación es muy similar y la diferencia estriba en que en este caso no se trata del cambio (incremento o disminución) de la probabilidad de la variable dependiente por cada unidad de cambio en las independientes, sino del incremento o disminución que se produce en el cociente entre $P(Y=1) / P(Y=0)$, expresado por:

$$\frac{P(Y = 1)}{P(Y = 0)} = e^{B_0 + B_1 X_1 + B_2 X_2 + \dots + B_K X_K} \quad (2.75)$$

Más aún, están expresados en logaritmos, por lo que sería necesario transformarlos (tomando los valores del antilogaritmo) de tal forma que se evalúe más fácilmente su efecto sobre la probabilidad. Los programas de computador lo hacen automáticamente calculando tanto el coeficiente real como el transformado. Utilizar este procedimiento no cambia en modo alguno la forma de interpretar el signo del coeficiente. Un coeficiente positivo aumenta la probabilidad, mientras que un valor negativo disminuye la probabilidad. Así pues si β es positivo, su transformación (antilog) será mayor a 1, y el odds ratio aumentará. Este aumento se produce cuando la probabilidad prevista de ocurrencia de un suceso aumenta y la probabilidad prevista de su no ocurrencia disminuye. Por lo tanto, el modelo tiene una elevada probabilidad de ocurrencia. De la misma forma, si β es negativo, el antilogaritmo es menor que 1 y el odds ratio disminuye. Un valor de cero equivale a un valor de 1, lo que no produce cambio en el odds. Hair (1999)

2.6.4 PRUEBA DE HIPÓTESIS PARA LOS COEFICIENTES DEL MODELO DE REGRESIÓN LOGÍSTICA.

Usualmente en la estimación del Modelo de Regresión Logística, como en el Modelo de Regresión Lineal Múltiple se efectúan pruebas con objetivos diferentes, siendo estos:

1. Determinar si una variable explicativa tiene coeficiente igual a cero.
2. Determinar si un conjunto de variables explicativas tienen coeficientes igual a cero.
3. Determinar la calidad del ajuste global del modelo.

Veamos para cada objetivo, como se efectúa el análisis.

2.6.5.1 PRUEBA DE WALD

Wald(1943) estudio una prueba asintótica para estimaciones máximos verosímiles, y aseveró que los parámetros estimados en los modelos logísticos tiene una Distribución Normal para muestras grandes.

Esta prueba se usa para evaluar la significancia estadística de cada variable explicativa o regresora.

Sea $\hat{\beta}^{(t)}$ que converge a los EMV de β y y_1, y_2, \dots, y_n variables respuesta binaria independientes cuyas probabilidades satisfacen.

$$\text{Logit}(p_i) = x_i^T \beta$$

donde $p_i = P[y_i = 1/x_i]$

Siendo x_i una observación que contiene los valores de las k variables

explicativas con $x_i^T = (1, x_{i1}, x_{i2}, \dots, x_{ik})$.

Sin pérdida de generalidad, seleccionaremos β_j como el parámetro de interés.

Supóngase que las hipótesis son:

$$\begin{aligned} H_0 : \mathbf{b}_j &= \mathbf{b}_{j0} \\ H_1 : \mathbf{b}_j &\neq \mathbf{b}_{j0} \end{aligned} \quad (2.76)$$

sea $\hat{\mathbf{b}}_j$ un EMV de \mathbf{b}_j y sea:

$I^{-1} = (X^T V X)^{-1}$ la inversa de la matriz de información muestral, entonces la estadística de Wald para dóimar (2.75) es:

$$W = \frac{(\hat{\mathbf{b}}_j - \mathbf{b}_{j0})^2}{\mathbf{s}^2(\hat{\mathbf{b}}_j)} \quad (2.77)$$

donde $\mathbf{s}(\hat{\mathbf{b}}_j)$ es la estimación del error estándar de $\hat{\mathbf{b}}_j$.

Bajo H_0 , $W \sim \mathbf{c}_{(1)}^2$ y para n suficientemente grande se tiene que:

$$z = \frac{\hat{\mathbf{b}}_j - \mathbf{b}_{j0}}{\mathbf{s}(\hat{\mathbf{b}}_j)} \sim N\left(\left(\frac{\hat{\mathbf{b}}_j - \mathbf{b}_{j0}}{\mathbf{s}(\hat{\mathbf{b}}_j)}\right), 1\right) \quad (2.78)$$

por tanto:

$z^2 \sim \mathbf{c}_{(x,1)}^2$, es \mathbf{c}^2 con parámetro de no centralización:

$$\mathbf{x} = \frac{(\mathbf{b}_j - \mathbf{b}_{j0})^2}{\mathbf{s}(\hat{\mathbf{b}}_j)} \quad (2.79)$$

ver Hanck y Donner (1977)

Pero la estadística W , tiene la propiedad que cuando el valor absoluto del coeficiente de Regresión es grande, el error estándar también lo es; esta situación hace que la estadística W sea pequeña y por tanto se puede rechazar \mathbf{b}_j igual a cero, cuando en realidad no debería rechazarse.

Por tanto, cuando se encuentra que un coeficiente es grande, es preferible no usar la estadística de Wald para efectuar dócima individual. Sino se recomienda construir un modelo con y sin esa variable y basarse en la prueba de hipótesis de la diferencia entre los dos modelos, ver Hanck y Donner (1977).

Para las hipótesis estadísticas:

$$\begin{aligned} H_0 : \mathbf{b}_j &= 0 \\ H_1 : \mathbf{b}_j &\neq 0 \end{aligned} \quad (2.80)$$

La estadística (2.77) es:

$$W = \frac{(\hat{\mathbf{b}}_j)^2}{\mathbf{s}^2(\hat{\mathbf{b}}_j)} \quad (2.81)$$

Bajo H_0 , $W \sim \mathbf{c}_{(1)}^2$ y para n suficientemente grande se tiene que:

$$z = \frac{\hat{\mathbf{b}}_j}{\mathbf{s}(\hat{\mathbf{b}}_j)} \sim N\left(\frac{\hat{\mathbf{b}}_j}{\mathbf{s}(\hat{\mathbf{b}}_j)}, 1\right) \quad (2.82)$$

por tanto:

$$z^2 \sim \mathbf{c}_{(\mathbf{x},1)}^2$$

z^2 se distribuye como una $\mathbf{c}_{(\mathbf{x},1)}^2$ con parámetro de no centralización:

$$\mathbf{x} = \frac{(\mathbf{b}_j)^2}{\mathbf{s}^2(\hat{\mathbf{b}}_j)} \quad (2.83)$$

si la variable explicativa es categórica, los grados de libertad es igual al número de categorías o niveles de la variable menos uno.

2.6.5.2 PRUEBA CHI-CUADRADO

Esta prueba sirve para lograr el segundo objetivo propuesto al iniciarse la sección 2.6.4 y sirve para docimar los coeficientes del modelo logístico.

Para elegir un modelo, se usa la prueba de razón de verosimilitud, Bickel y Docksum (1977), para probar la hipótesis de que los coeficientes \mathbf{b}_j correspondientes a las variables explicativas retiradas, digamos q variables explicativas, del modelo son iguales a cero, siendo la hipótesis estadísticas:

$$H_0 : \mathbf{b}_1 = \mathbf{b}_2 = \dots = \mathbf{b}_q = 0$$

$$H_1 : \mathbf{b}_j \neq 0, \text{ para por lo menos un } j = 1, 2, \dots, q$$

.Esta prueba se basa en la siguiente estadística:

$$\mathbf{c}_q^2 = -2 \left[\text{Ln} L_{p-q} - \text{Ln} L_p \right] \quad (2.84)$$

Bajo la hipótesis de que los coeficientes de las variables retiradas son iguales a cero, la estadística (2.84) tiene una distribución asintótica $\mathbf{c}_{(q)}^2$.

Valores altos para esta estadística indican que una o más de las q variables retiradas tienen coeficiente de regresión distinto de cero.

La estadística \mathbf{c}_q^2 se usa también para probar si una variable explicativa determinada, por ejemplo x_k , muestra una asociación significativa (como factor de riesgo cuando se aplica a casos de enfermedades) para con la variable respuesta en la presencia de las demás variables x_1, x_2, \dots, x_{k-1} .

2.6.5.3 ESTADÍSTICA CHI-CUADRADA DE PEARSON

Esta estadística sirve para lograr el objetivo número 3, es decir evaluar el modelo ajustado en forma global. La estadística se basa en la comparación de los valores observados, y_i ; y sus respectivas probabilidades estimadas,

p_i .

Las hipótesis estadísticas para usar esta estadística son:

$$H_0 : \mathbf{b}_0 = \mathbf{b}_1 = \dots = \mathbf{b}_k = 0$$

$$H_1 : \mathbf{b}_j \neq 0, \text{ para por lo menos un } j = 0, 1, 2, \dots, k$$

esta prueba se basa en la estadística Chi-cuadrado de Pearson, que está dada por:

$$\mathbf{c}^2 = \sum_{i=1}^n \frac{(y_i - \hat{\mathbf{p}}_i)^2}{\hat{\mathbf{p}}_i(1 - \hat{\mathbf{p}}_i)} \quad (2.85)$$

o equivalentemente
$$\mathbf{c}^2 = \sum_{i=1}^n \frac{r_i^2}{v_{ii}} \quad (2.86)$$

donde:

$$r_i = (y_i - \hat{\mathbf{p}}_i)$$

$$v_{ii} = \text{Diag}(\hat{\mathbf{V}}) = \hat{\mathbf{p}}_i(1 - \hat{\mathbf{p}}_i)$$

como observamos la estadística (2.86) es igual a (1.52).

Bajo la hipótesis nula, de que el modelo se ajusta bien a los valores observados, la

estadística (2.86) tiene distribución asintótica Chi-cuadrado $\mathbf{C}_{(n-(k+1))}^2$.

Valores altos de la estadística Chi-cuadrado de Pearson indican discrepancias con el modelo teórico, Cordeiro (1992).

La estadística (2.86) es inestable cuando $\hat{\mathbf{p}}_i$ toma valores cercanos a cero o uno, por ello tomar en cuenta esta observación, cuando se realiza el análisis.

2.6.5.4 DESVIANZA

Otra forma de probar el ajuste global del modelo, es mediante la estadística llamada Desvianza, propuesta por Nelder y Wederburn (1982), es análogo a la suma de cuadrados de los residuales del Modelo de Regresión Lineal Múltiple.

Las hipótesis estadísticas son:

$$H_0 : \mathbf{b}_1 = \dots = \mathbf{b}_k = 0$$

$$H_1 : \mathbf{b}_j \neq 0, \text{ para por lo menos un } j = 1, 2, \dots, k$$

Esta estadística se usa para evitar la inestabilidad de la estadística Chi-cuadrado de Pearson. La Desvianza esta dada por:

$$D_p = \sum_{i=1}^n d_i^2 \quad (2.87)$$

donde :

$$d_i = \begin{cases} \sqrt{-2 \log \hat{p}_i} & \text{si } y_i = 1 \\ \sqrt{-2 \log(1 - \hat{p}_i)} & \text{si } y_i = 0 \end{cases} ; j = 1, 2, \dots, n$$

La Desvianza bajo la hipótesis nula, asintóticamente, es la misma que la distribución Chi-cuadrado de Pearson, es decir se distribuye $\chi^2_{(n-(k+1))}$ y mide la discrepancia o el desvio entre el modelo bajo investigación o actual y el modelo saturado.

La estadística (2.87) para el modelo de regresión logística esta dada por:

$$D = -2 \sum (y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)) \quad (2.88)$$

Cuando el modelo bajo investigación es verdadero se compara el valor D con el valor crítico $\chi^2_{(n-p)}$ de una distribución χ^2 a un nivel de significación igual a α , por tanto :

Si $D > \chi^2_{(n-p)}$ el modelo se rechaza y

Si $D \leq c_{(n-p)}^2$ el modelo no se rechaza.

donde $p = k + 1$

2.6.5.5 TABLA DE CLASIFICACION

También llamada Matriz de Confusión, es una forma sencilla de evaluar el ajuste del Modelo de Regresión Logística, no es tan objetiva pero se usa como indicador de bondad de ajuste.

Es una tabla sencilla de 2x2, en el cual se muestra la distribución de los objetos que pertenecen a las categorías 1 y 2, es decir cuando $y = 0$ y cuando $y = 1$, conjuntamente con la clasificación a cualquiera de las 2 categorías de acuerdo a la probabilidad estimada.

Para interpretar se hace mediante el porcentaje de objetos bien clasificados, esto es, aquellos que mediante la probabilidad estimada permanecen en su respectiva categoría. También se interpreta mediante el porcentaje de objetos mal clasificados, esto es, aquellos que mediante la probabilidad estimada se asignan a categorías diferentes del cual fueron observados.

TABLA DE CLASIFICACION

| GRUPO ACTUAL | GRUPO ESTIMADO | | TOTAL MARGINAL |
|---------------|-------------------|-------------------|-------------------|
| | 0 | 1 | |
| 0 | n_{11} | n_{12} | $n_{11} + n_{12}$ |
| 1 | n_{21} | n_{22} | $n_{21} + n_{22}$ |
| TOTAL MAGINAL | $n_{11} + n_{21}$ | $n_{12} + n_{22}$ | n |

$$\frac{n_{11} + n_{22}}{n} \times 100\% \text{ es el porcentaje de objetos bien clasificados}$$

mediante el Modelo de Regresión Logística estimado.

Por tanto, lo que se debe esperar es que este porcentaje sea lo más alto posible, a fin de concluir que el modelo obtenido clasifica bien a los objetos o individuos.

2.6.5.6 CONTRASTE DE BONDAD DE AJUSTE DE HOSMER – LEMESHOW

Este contraste evalúa la bondad de ajuste del modelo, es decir el grado en que la probabilidad predicha coincide con la observada, construyendo una tabla de contingencia a la que aplica un contraste χ^2 . Para ello calcula los deciles de las probabilidades estimadas $(\hat{p}_i; i = 1, 2, \dots, n)$, D_1, D_2, \dots, D_9 y divide los datos observados en 10 categorías dadas por :

$$A_j = \{ \hat{p}_i \in [D_{j-1}, D_j) / i \in \{1, 2, \dots, n\} \} ; j = 1, 2, \dots, 10$$

donde $D_0 = 0$, $D_{10} = 1$.

Sean:

n_j = número de casos en A_j ; $j=1, 2, \dots, 10$

o_j = número de $y_i = 1$ en A_j ; $j=1, 2, \dots, 10$

$$\bar{p}_j = \frac{1}{n_j} \sum_{i \in A_j} \hat{p}_i ; j = 1, 2, \dots, 10$$

El estadístico del contraste viene dado por :

$$T = \sum_{j=1}^{10} \frac{(o_j - n_j \bar{p}_j)^2}{n_j \bar{p}_j (1 - \bar{p}_j)} \quad (2.89)$$

y el p-valor del contraste es $P [\chi^2_8 \geq T_{obs}]$.

2.6.6 DIAGNOSTICO DEL MODELO

Es la evaluación de la bondad de ajuste caso por caso mediante el análisis de los residuos del modelo y de su influencia en la estimación del vector de parámetros del mismo, se realiza usando:

2.6.6.1 RESIDUOS DEL MODELO

Los residuos más utilizados son los siguientes:

Residuos estandarizados.- Son el cociente entre los residuales y una estimación de la desviación estándar.

$$z_i = \frac{y_i - \hat{p}_i}{\sqrt{\hat{p}_i(1 - \hat{p}_i)}} \quad ; \quad i = 1, 2, \dots, n \quad (2.90)$$

Residuos studentizados.- Son el cambio en el valor de la desviación del modelo si el caso es excluido.

$$st_i = \frac{y_i - \hat{p}_{(i)}}{\sqrt{\hat{p}_{(i)}(1 - \hat{p}_{(i)})}} \quad ; \quad i = 1, 2, \dots, n \quad (2.91)$$

donde $\hat{p}_{(i)}$ es la estimación de p_i obtenida eliminando la observación i de la muestra.

Residuos Desviación.- Para cada observación la desviación se calcula :

$$d_i = \begin{cases} \sqrt{-2 \log \hat{p}_i} & \text{si } y_i = 1 \\ \sqrt{-2 \log(1 - \hat{p}_i)} & \text{si } y_i = 0 \end{cases} \quad ; \quad j = 1, 2, \dots, n \quad (2.92)$$

Todos estos residuos se distribuyen aproximadamente como una $N(0,1)$, si el modelo ajustado es correcto.

2.6.6.2 MEDIDAS DE INFLUENCIA

Cuantifican la influencia que cada observación ejerce sobre la estimación del vector de parámetros o sobre las predicciones hechas a partir del mismo, de modo que, cuanto más grande son, mayor es la influencia que ejerce una observación en la estimación del modelo.

Medida de Apalancamiento (Leverage)

Se utiliza para detectar observaciones que tienen un gran impacto en los valores predichos por el modelo.

Se calcula a partir de la matriz $H = W^{1/2} X (X'WX)^{-1} X'W^{1/2}$ donde $W = \text{diag}[\hat{p}_i(1 - \hat{p}_i)]$. El apalancamiento para la observación i -ésima viene dado por el elemento i -ésimo de la diagonal principal de H , h_{ii} , y toma valores entre 0 y 1 con un valor medio de p/n .

Las dos medidas siguientes miden el impacto que tiene una observación en la estimación de \hat{a} .

Distancia de Cook.- Mide la influencia en la estimación de \hat{a} .

$$COOK_i = \frac{1}{p} (\hat{a} - \hat{a}_{(i)})' X'WX (\hat{a} - \hat{a}_{(i)}) \quad (2.93)$$

DFBETA.- Mide la influencia en la estimación de una componente de \hat{a} , \hat{a}_i

$$Dfbeta1_i = \frac{\hat{a}_1 - \hat{a}_{1(i)}}{std(\hat{a}_1)} \quad (2.94)$$

donde $\hat{a}_1, \hat{a}_{1(i)}$ denotan las estimaciones del modelo logístico de \hat{a} y \hat{a}_1 , eliminando la i -ésima observación de la muestra y $std(\hat{a}_1)$ el error estándar en la estimación de \hat{a}_1 .

2.6.7 VARIABLES EXPLICATIVAS CATEGORICAS – VARIABLES DUMMY

Si una de las variables explicativas es categórica, con c valores posibles, se crean $c-1$ variables dicotómicas como variables explicativas también

llamadas variables dummy. Estas variables cuantifican el efecto de un valor de dichas variables con respecto a un valor de referencia.

Estas variables se usan cuando los datos se muestran como categorías, las categorías pueden ser:

Nominales: La variable simplemente indica diferentes categorías, las categorías no pueden ser ordenadas en un orden particular. Ejemplo : Sexo (hombre,mujer).

Ordinales: La variable además de estar agrupada en categorías puede ser ordenada. El que una categoría este en un orden superior que otra implica que su medida representa algo mayor que la otra. Ejemplo: Clase social (baja, media, alta).

Intervalares: La variables no solo puede ser ordenada, sino que su valor mide la distancia entre categorías. Estas tienen estándares de unidades de medida.

Ejemplo: Altura, temperatura, presión sanguínea.

Cuando se tiene variables de este tipo se crean las llamadas variables dummy, si la variable tiene c categorías se usan c-1 variables ficticias o dummy. La variable indica si un dato corresponde a una categoría o no. Veamos un ejemplo de cómo se hace esto:

Supongamos que tenemos una variable clase social, codificada 1: Baja, 2:Media, y 3:Alta, entonces creamos dos variables dummy :

Clase1 : 1 si el dato corresponde a la clase social Baja, 0 si el dato no pertenece a la clase Baja.

Clase2 : 1 si el dato corresponde a la clase social Media, 0 si el dato no pertenece a la clase Media.

Como se puede ver estas nos permiten clasificar cualquier dato en una de las categorías existentes.

Supongamos que creamos las variables dummy Ind1 e Ind2 para una variable X1 de tres categorías, como sigue:

| | | |
|----|------|------|
| X1 | Ind1 | Ind2 |
|----|------|------|

| | | |
|---|---|---|
| 1 | 1 | 0 |
| 2 | 0 | 1 |
| 3 | 0 | 0 |

En este caso la ecuación del modelo ajustado viene dada por:

$$\log\left(\frac{p(\text{Ind1}, \text{Ind2}; \mathbf{b})}{1 - p(\text{Ind1}, \text{Ind2}; \mathbf{b})}\right) = \mathbf{b}_0 + \mathbf{b}_1 \text{Ind1} + \mathbf{b}_2 \text{Ind2} \quad (2.95)$$

Sea $p_i = p[y=1/X1=i]$; $i=1,2,3$. Se tiene:

$$\frac{p_1}{1 - p_1} = e^{\mathbf{b}_0 + \mathbf{b}_1}, \quad \frac{p_2}{1 - p_2} = e^{\mathbf{b}_0 + \mathbf{b}_2}, \quad \frac{p_3}{1 - p_3} = e^{\mathbf{b}_0}$$

Se sigue que:

$$\frac{p_1}{1 - p_1} \bigg/ \frac{p_3}{1 - p_3} = e^{\mathbf{b}_1}, \quad \frac{p_2}{1 - p_2} \bigg/ \frac{p_3}{1 - p_3} = e^{\mathbf{b}_2}$$

Por lo tanto, $e^{\mathbf{b}_i}$, $i = 1,2$ compara los odds ratio correspondientes a $X1=1,2$, frente al de la categoría de referencia $X1=3$.