

ANALISIS DE SENSIBILIDAD EN REGRESION

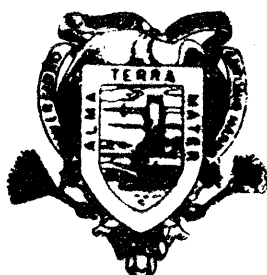
OSCAR ALEJANDRO MARTINEZ JAIME

TESIS

**PRESENTADA COMO REQUISITO PARCIAL
PARA OBTENER EL GRADO DE
MAESTRO EN CIENCIAS
EN ESTADISTICA EXPERIMENTAL**



**BIBLIOTECA
EGIDIO G. REBONATO
BANCO DE TESIS
U.A.A.**



Universidad Autónoma Agraria

"Antonio Narro"

PROGRAMA DE GRADUADOS

Buenvista, Saltillo, Coah.

NOVIEMBRE DE 2001

13516

Universidad Autónoma Agraria Antonio Narro
Subdirección de Postgrado

ANÁLISIS DE SENSIBILIDAD EN REGRESIÓN
TESIS

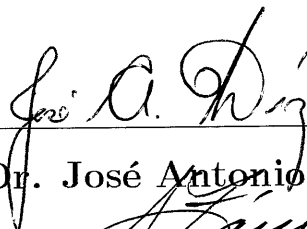
Por:

OSCAR ALEJANDRO MARTÍNEZ JAIME

Elaborada bajo la supervisión del comité particular
de asesoría y aprobada como requisito parcial, para optar al grado de:

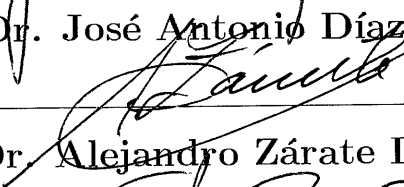
MAESTRO EN CIENCIAS
EN ESTADÍSTICA EXPERIMENTAL
Comité Particular

Asesor principal:



Dr. José Antonio Díaz García

Asesor:

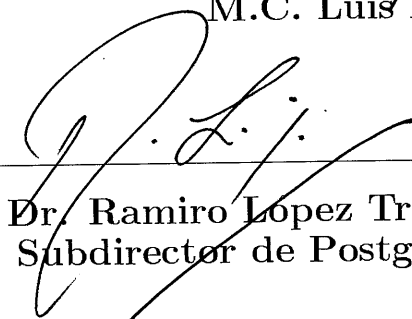


Dr. Alejandro Zárate Lupercio

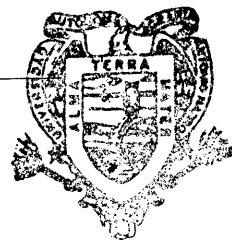
Asesor:



M.C. Luis Rodríguez Gutiérrez



Dr. Ramiro López Trujillo
Subdirector de Postgrado



Buenavista, Saltillo, Coahuila. Noviembre de 2001

BIBLIOTECA
EGIDIO G. REYNOLATO
BANCO DE TIEMPOS
U.A.A.N.

AGRADECIMIENTOS

A mi asesor principal, el Dr. José Antonio Díaz García, por sus consejos, sugerencias e invaluable ayuda, sin la cual, no hubiera sido posible la culminación de este trabajo de tesis.

A mis asesores, el Dr. Alejandro Zárate Lupercio y al M.C. Luis Rodríguez Gutiérrez, por su participación en la revisión del presente trabajo de tesis.

A mis profesores, el Dr. Mario Cantú Sifuentes, y a los M.C. Emilio Padrón Corral, Roberto Coronado Niño y Félix de Jesús Sánchez Pérez, por contribuir en mi formación académica.

DEDICATORIAS

*A mi esposa, Erika Elizabeth,
por su gran amor
y apoyo incondicional en todo momento.*

y

*A mi hija, Alison Elizabeth,
por ser el regalo más preciado
que Dios me ha dado.*

COMPENDIO

Análisis de Sensibilidad en Regresión

POR

OSCAR ALEJANDRO MARTÍNEZ JAIME

MAESTRÍA

ESTADÍSTICA EXPERIMENTAL

UNIVERSIDAD AUTÓNOMA AGRARIA "ANTONIO NARRO"

BUENAVISTA, SALTILLO, COAHUILA. NOVIEMBRE DE 2001

Dr. José Antonio Díaz García -Asesor-

Palabras clave: Medidas de Diagnóstico, Distancia de Mahalanobis Generalizada, Outliers, Puntos Palanca, Puntos Influyentes.

En el presente trabajo, se describen algunas herramientas básicas de Álgebra Lineal, y varias medidas de diagnóstico que permiten detectar observaciones influyentes, y se propone una modificación de la distancia de Cook, basándose en la distancia de Mahalanobis generalizada, en el contexto del modelo de regresión lineal con distribución normal multivariada. Se establece además, la distribución exacta del estadístico basado en esta distancia de Mahalanobis generalizada, la cual proporciona puntos críticos para identificar "outliers" en un conjunto de datos. Este procedimiento, se ilustra con un ejemplo, en el caso de la regresión lineal múltiple y multivariada.

ABSTRACT

Sensitivity Analysis in Regression

BY

OSCAR ALEJANDRO MARTÍNEZ JAIME

MASTER OF SCIENCE

EXPERIMENTAL STATISTICS

UNIVERSIDAD AUTÓNOMA AGRARIA "ANTONIO NARRO"

BUENAVISTA, SALTILLO, COAHUILA. NOVEMBER 2001

Dr. José Antonio Díaz García -Advisor-

Key Words: Diagnostic Measures, Generalized Mahalanobis Distance, Outliers,
High Leverage Points; Influential Points.

In this work, it describe some basic tools of Linear Algebra, and several diagnostics measures that it permit to detect influential observations, and a modification of the classical Cook's distance is proposed, providing us with a generalized Mahalanobis distance in the context of multivariate normal linear regression models. It establish the exact distribution of a pivotal type statistics based on this generalized Mahalanobis distance, which provides critical points for the identification of data points. It illustrate the procedure with an example, in the context of multiple and multivariate linear regression.

ÍNDICE DE CONTENIDO

	Página
ÍNDICE DE FIGURAS	ix
INTRODUCCIÓN	1
LA MATRIZ DE PREDICCIÓN	6
Notación	6
Resultados	7
Supuestos	9
Las Matrices \mathbf{P} y $(\mathbf{I} - \mathbf{P})$	10
Propiedades de la Matriz de Predicción	13
Omitiendo o Adicionando una Observación	18
Condiciones para Valores Grandes de p_{ii}	21
Omitiendo Múltiples Filas de \mathbf{X}	22
Eigenvalores de \mathbf{P} y $(\mathbf{I} - \mathbf{P})$	23
Distribución de p_{ii}	27
EFFECTOS DE UNA OBSERVACIÓN SOBRE LA ECUACIÓN DE REGRESIÓN	29
Método de la Eliminación de una Observación	30
Medidas Basadas en Residuales	30
Outliers, Puntos Palanca y Puntos Influyentes ..	36
Medidas Basadas en la Curva de Influencia.....	40

Medidas Basadas en el Volumen de Elipsoides de Confianza	52
Medidas Basadas en la Función de Verosimilitud	58
Medidas Basadas en un Subconjunto de Coeficientes de Regresión	60
Método Basado en la Diferenciación	64
EFFECTOS DE MÚLTIPLES OBSERVACIONES SOBRE LA ECUACIÓN DE REGRESIÓN	67
Método de la Eliminación de Múltiples Observaciones	68
Medidas Basadas en Residuales	68
Medidas Basadas en la Curva de Influencia	71
Medidas Basadas en el Volumen de Elipsoides de Confianza	75
Medidas Basadas en la Función de Verosimilitud	76
Medidas Basadas en un Subconjunto de Coeficientes de Regresión	78
ARTÍCULO	
“Una modificación de la distancia de Cook”	79
LITERATURA CITADA	96

ÍNDICE DE FIGURAS

Figura No.		Página
1	Datos originales para el ajuste del <i>SCORE</i> y un gráfico Q-Q para los residuales en la regresión simple de <i>SCORE</i> sobre <i>AGE</i>	91
2	Identificación de la influencia y “outliers”, basados en a) Distancia de Cook, b) Distancia de Cook Modificada de acuerdo con (5.20), c) Distancia de Draper y John. $n = 21$, $q = 2$, $p = 1$ y $s^2(i)$. La varianza de los residuales sin la i -ésima observación se usa en todos los casos, de acuerdo con la Nota 2	92
3	Identificación de “outliers” basada en la Distancia de Mahalanobis sobre la matriz de residuales, y detección de “outliers” basada en la Distancia de Cook Modificada dada en el Teorema 3	93

INTRODUCCIÓN

La *regresión* es una herramienta analítica inferencial que establece la relación funcional entre variables, y constituye uno de los métodos estadísticos más ampliamente utilizado. Al conjunto de técnicas que permiten construir y evaluar modelos que describan la relación entre variables, y formular inferencias basadas en los modelos obtenidos, se les conoce como *técnicas de regresión*; y al análisis estadístico que resulta de aplicarlas, se le denomina *análisis de regresión*.

Considerando el modelo de regresión general

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (1.1)$$

donde:

- \mathbf{Y} , es un vector de tamaño $n \times 1$ de variables dependientes (variables respuesta),
- \mathbf{X} , es una matriz de tamaño $n \times k$ (con $n > k$) de variables independientes (variables predictoras o explicatorias),
- $\boldsymbol{\beta}$, es un vector de tamaño $k \times 1$ de coeficientes o parámetros desconocidos, y
- $\boldsymbol{\epsilon}$, es un vector de tamaño $n \times 1$ de errores aleatorios.

Si el modelo dado en la ecuación (1.1) es lineal, entonces se le llama *modelo de regresión lineal*, en el cual, si se involucra sólo una variable independiente, la regresión lineal se dice ser simple, pero si existen dos o más variables independientes en el modelo, entonces se dice ser múltiple. Uno de los propósitos del análisis de regresión lineal, es el de utilizar las variables explicatorias, cuyos valores son

conocidos, para predecir el comportamiento de la variable respuesta, aunque se puede decir que el objetivo general de esta técnica, es hacer mejores inferencias sobre la variable de mayor interés.

Los elementos que determinan una ecuación de regresión lineal, son las observaciones, las variables y los supuestos asociados al modelo, estos últimos se dan en la tercer sección del segundo capítulo. El ajuste de esta ecuación por el método de *cuadrados mínimos*, es el procedimiento de modelación que será utilizado.

En la práctica, es común encontrar algunos problemas relacionados con la información disponible (observaciones), tal como ocurre cuando se tiene la presencia de una o varias observaciones discrepantes (también llamados puntos anómalos), las cuales son inherentes a un conjunto de datos determinado. En este caso, el diagrama de dispersión de los datos, constituye un procedimiento informal que permite detectar con facilidad dichos puntos anómalos, pues aparecen en forma aislada, es decir, se alejan de la nube (o patrón) que siguen los demás puntos.

A las observaciones extremas, que individualmente o en conjunto, afectan la estimación mínimo cuadrática de los parámetros de la regresión, es decir, β y σ^2 , al ajustar el modelo de regresión, se les llama *observaciones influyentes*, de las cuales, Chatterjee y Hadi (1988), distinguen tres tipos, a saber: “*outliers*”, *puntos palanca* y *puntos influyentes*, cuyas definiciones se exponen en la primera sección del tercer capítulo.

Se han desarrollado algunos métodos para detectar *observaciones influyentes*, entre los que se encuentran, el *método de la eliminación de observaciones* y el *método basado en la diferenciación*. El primero, estudia cómo el ajuste del modelo de regresión se ve perturbado al omitir una observación en particular, o bien, un conjunto de observaciones; mientras que en el segundo, la perturbación del ajuste es con respecto a ciertos parámetros del modelo.

Dentro del gran número de medidas de diagnóstico correspondientes al primer método, cabe destacar aquellas que miden la influencia de la i -ésima observación (potencialmente influyente), en la estimación de β exclusivamente, cuando dicha observación es omitida. Esto es, si se suprime la i -ésima observación, y se compara el vector que no considera esta observación $\hat{\beta}_{(i)}$, con el vector que considera el conjunto de datos completo $\hat{\beta}$, a través de su diferencia $(\hat{\beta} - \hat{\beta}_{(i)})$, se establece el impacto que tiene la observación i sobre la estimación.

En el diagnóstico del análisis de regresión lineal, comúnmente, se utilizan medidas basadas en residuales, en las que la versión estudentizada, juega un papel fundamental, tal como lo sugieren Atkinson (1981 y 1982) y Velleman y Welsch (1981); estas medidas se describen en la primera sección del tercer capítulo.

Por su parte Hampel (1968 y 1974), propone las medidas basadas en la curva de influencia para $\hat{\beta}$, dentro de las cuales se encuentra la *distancia de Cook*, cuya definición requiere considerar, que bajo normalidad, la región de confianza $100(1 - \alpha)\%$ para β , se obtiene a partir de

$$\frac{(\beta - \hat{\beta})^T (\mathbf{X}^T \mathbf{X}) (\beta - \hat{\beta})}{k \hat{\sigma}^2} \leq F(\alpha; k, n - k),$$

donde $F(\alpha; k, n - k)$ es el punto superior α de la distribución F centrada con k y $(n - k)$ grados de libertad. Esta desigualdad define una región elipsoidal centrada en $\hat{\beta}$, y permite medir la influencia de la i -ésima observación por el cambio en el centro del elipsoide de confianza, cuando la i -ésima observación es omitida. En forma análoga, Cook (1977), propone la medida

$$C_i = \frac{(\beta - \hat{\beta}_{(i)})^T (\mathbf{X}^T \mathbf{X}) (\beta - \hat{\beta}_{(i)})}{k \hat{\sigma}^2}, \quad i = 1, 2, \dots, n,$$

para considerar la influencia de la i -ésima observación sobre el centro del elipsoide de confianza, es decir, sobre $\hat{\beta}$. Asimismo, Welsch y Kuh (1977), presentan la distancia de Welsch y la de Welsch-Kuh, mientras que Atkinson (1981), realiza una

modificación de la distancia de Cook. Estas medidas se muestran en la primera sección del tercer capítulo.

Una alternativa a las medidas anteriores, la constituyen, medidas de la influencia de la i -ésima observación, que se basan en el cambio en el volumen del elipsoide de confianza para β , cuando esta observación es omitida, las cuales se muestran en la primera sección del tercer capítulo, y son: el estadístico Andrews-Pregibon sugerido por Andrews y Pregibon (1978), la relación de varianzas señalada en Belsley *et al.* (1980) y el estadístico Cook-Weisberg propuesto por Cook y Weisberg (1980).

Se concluye esta sección con las definiciones de las medidas basadas en la función de verosimilitud propuestas por Cook y Weisberg (1982).

La mayor parte de las aportaciones son recopiladas en Chatterjee y Hadi (1988) y en Belsley *et al.* (1980), en donde se observa que para la mayoría de las medidas de diagnóstico, el estadístico de prueba no tiene una distribución exacta, que mediante una regla de decisión permita establecer el punto crítico (o punto de corte) que facilite la detección de observaciones influyentes.

En el presente trabajo se plantean las bases y la importancia del análisis de sensibilidad en regresión. Como tema central, se propone una modificación de la distancia de Cook, basándose en la distancia de Mahalanobis generalizada, en el contexto de la regresión lineal con distribución normal multivariada. Se establecen además, las distribuciones exactas de los estadísticos basados en esta distancia modificada, a través de la cual, se obtiene un punto crítico que permita detectar una o varias observaciones influyentes.

En el segundo capítulo, se da la notación, algunos resultados y supuestos asociados con el modelo dado en la ecuación (1.1), que son requeridos en los capítulos posteriores; asimismo, se describen las propiedades de las matrices de predicción P

y de residuales $(I - P)$, y se mencionan los efectos sobre P , que resultan al omitir una o varias observaciones, y por último, se muestra la distribución que siguen los elementos de la diagonal de la matriz de predicción, es decir, los p_{ii} .

Las medidas de diagnóstico, correspondientes al método de la eliminación de una observación, y el método basado en la diferenciación son descritas en el tercer capítulo.

Por su parte, en el cuarto capítulo, se generaliza algunas de las medidas descritas en el tercer capítulo, al caso de la influencia de múltiples observaciones sobre el ajuste de la ecuación de regresión, limitándose sólo a aquellas que corresponden al método de la eliminación de múltiples observaciones.

En el quinto capítulo, se presenta el artículo “**Una modificación de la distancia de Cook**”, donde considerando el modelo de regresión lineal multivariado con distribución normal, se propone una modificación de la distancia de Cook para el caso de una observación, basándose en la distancia generalizada de Mahalanobis, vea el Teorema 1. En el Teorema 2 se define cómo calcular el cuadrado de la distancia de Cook modificada para detectar k observaciones influyentes. En los teoremas 3 y 4 se encuentran las distribuciones de las distancias de Cook modificadas propuestas en los teoremas 1 y 2, respectivamente. Se concluye el artículo con una aplicación de las técnicas descritas a un ejemplo de regresión lineal simple y otro de regresión lineal múltiple multivariado.

LA MATRIZ DE PREDICCIÓN

Notación

A continuación se propone la notación que será utilizada durante el desarrollo del presente trabajo.

En primera instancia, se define \mathbf{I} como la matriz identidad, $\mathbf{0}$ como la matriz nula y $\mathbf{1}$ como el vector de unos.

Sea \mathbf{u}_i , el i -ésimo vector unitario, es decir, un vector con uno en la posición i y ceros en el resto de las posiciones.

Asimismo, se utiliza \mathbf{M}^T , \mathbf{M}^{-1} , \mathbf{M}^{-T} , para denotar la transpuesta, la inversa y la inversa de la transpuesta de la matriz \mathbf{M} , respectivamente. Mientras que el rango, la traza, el determinante y la norma de la matriz \mathbf{M} son dados por $r(\mathbf{M})$, $tr(\mathbf{M})$, $det(\mathbf{M})$ y $\|\mathbf{M}\|$, respectivamente.

Sean \mathbf{y}_i , $i = 1, 2, \dots, n$, el i -ésimo elemento de \mathbf{Y} ; \mathbf{x}_i^T , la i -ésima fila de \mathbf{X} ; y sea \mathbf{X}_j , $j = 1, 2, \dots, k$, la j -ésima columna de \mathbf{X} .

Cuando se hace referencia a la i -ésima observación, se considera el vector fila $(\mathbf{x}_i : \mathbf{y}_i)$, esto es, la i -ésima fila de la matriz aumentada $\mathbf{Z} = (\mathbf{X} : \mathbf{Y})$.

La matriz \mathbf{X} con la i -ésima fila eliminada se denota por $\mathbf{X}_{(i)}$. Del mismo modo, $\hat{\boldsymbol{\beta}}_{(i)}$ es conocido como el vector de parámetros estimado cuando la i -ésima observación ha sido eliminada.

Se define a $\hat{\epsilon}$ o $\hat{\epsilon}_{Y.X}$, como el vector de residuales, cuando se hace una regresión de la variable \mathbf{Y} con la(s) variable(s) \mathbf{X} .

Las expresiones $E(\cdot)$, $Var(\cdot)$, $Cov(\cdot, \cdot)$ y $Corr(\cdot, \cdot)$, son usadas para definir el valor esperado, la varianza, la covarianza y el coeficiente de correlación de la(s) variable(s) aleatoria(s) que se indican por puntos en los paréntesis, respectivamente.

Resultados

Algunos resultados importantes que serán requeridos, son dados a continuación.

1. El método de estimación más comúnmente utilizado es el de Cuadrados Mínimos, el cual requiere minimizar la expresión

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}),$$

para lo cual se deriva con respecto a $\boldsymbol{\beta}$, y resulta el sistema de ecuaciones normales

$$(\mathbf{X}^T \mathbf{X})\boldsymbol{\beta} = \mathbf{X}^T \mathbf{Y},$$

el cual tiene solución única, si y solo si $(\mathbf{X}^T \mathbf{X})^{-1}$ existe, dando lugar al estimador de $\boldsymbol{\beta}$, es decir,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y},$$

el cual tiene las siguientes propiedades:

- (a) $\hat{\boldsymbol{\beta}}$ es un estimador insesgado para $\boldsymbol{\beta}$, esto es, $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$,

- (b) $\hat{\beta}$ es el mejor estimador linealmente insesgado para β , lo cual indica que la varianza del estimador, es decir, $\mathbf{Var}(\hat{\beta}) = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ es la más pequeña, de entre la clase de estimadores linealmente insesgados, y
- (c) Asumiendo que los errores siguen una distribución normal, como se verá en la siguiente sección, por (a) y (b), se tiene que

$$\hat{\beta} \sim \mathcal{N}_k(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}).$$

2. El vector de valores ajustados (predichos), es decir, $\hat{\mathbf{Y}}$, se obtiene así,

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{P}\mathbf{Y}, \quad (2.1)$$

donde

$$\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \quad (2.2)$$

este vector $\hat{\mathbf{Y}}$, cuenta con las siguientes propiedades:

- (a) $\mathbf{E}(\hat{\mathbf{Y}}) = \mathbf{X}\beta$,
- (b) $\mathbf{Var}(\hat{\mathbf{Y}}) = \sigma^2 \mathbf{P}$, y
- (c) Bajo normalidad, por (a) y (b), se tiene que

$$\hat{\mathbf{Y}} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{P}),$$

la cual es una distribución singular de rango k , pues \mathbf{P} es singular.

3. El vector de residuales $\hat{\epsilon}$, se calcula así,

$$\hat{\epsilon} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{P}\mathbf{Y} = (\mathbf{I} - \mathbf{P})\mathbf{Y},$$

y presenta las siguientes propiedades:

- (a) $\mathbf{E}(\hat{\epsilon}) = \mathbf{0}$,
- (b) $\mathbf{Var}(\hat{\epsilon}) = \sigma^2(\mathbf{I} - \mathbf{P})$,

(c) Nuevamente asumiendo normalidad, por (a) y (b), se tiene que

$$\hat{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{P})),$$

la cual es una distribución singular de rango $n - k$, pues $(\mathbf{I} - \mathbf{P})$ es singular, y

(d) $\frac{\hat{\epsilon}^T \hat{\epsilon}}{\sigma^2} \sim \chi^2(n - k)$, donde $\hat{\epsilon}^T \hat{\epsilon}$ es la suma de cuadrados de residuales.

Supuestos

El análisis de regresión general se basa en suposiciones, en seguida se muestran aquellas que serán de utilidad en el presente estudio.

1. El supuesto de linealidad implícito en el modelo, el cual establece que cada valor respuesta observado \mathbf{y}_i , puede ser escrito como una función lineal de la i -ésima fila de \mathbf{X} , o sea,

$$\mathbf{y}_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \quad i = 1, 2, \dots, n$$

2. El supuesto que permite encontrar un estimador único de $\boldsymbol{\beta}$, para lo cual es necesario que $(\mathbf{X}^T \mathbf{X})^{-1}$ exista, o equivalentemente, que el $r(\mathbf{X}) = k$.
3. El supuesto distribucional, en el que se requiere que
 - (a) \mathbf{X} sea medida sin errores,
 - (b) ϵ_i no dependa de \mathbf{x}_i^T , para $i = 1, 2, \dots, n$, y
 - (c) $\epsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$.

4. El supuesto que indica que todas las observaciones son igualmente confiables y deben influir de igual manera en la determinación de los estimadores de $\boldsymbol{\beta}$ y σ^2 .

Las Matrices P y $(I - P)$

Recordando que la matriz P dada en la ecuación (2.2), juega un papel muy importante en el análisis de regresión lineal, ya que es necesaria para obtener los estimadores de los parámetros incluidos en el modelo. Chatterjee y Hadi (1986), le llaman *matriz de predicción*, porque es la transformación matricial que aplicada a Y , produce el vector de los valores predichos, es decir, \widehat{Y} . Similarmente, a $(I - P)$ se le denomina *matriz de residuales*, ya que al aplicarse a Y , genera el vector de residuales estimado, esto es, $\widehat{\epsilon}$.

En cuanto a notación se refiere, los elementos de P se definen de la siguiente manera,

1. p_{ij} , $i, j = 1, 2, \dots, n$ es el ij -ésimo elemento de P , y se obtiene así,

$$u_i^T P u_j = u_i^T X (X^T X)^{-1} X^T u_j = x_i^T (X^T X)^{-1} x_j = p_{ij},$$

2. p_{ii} , $i = 1, 2, \dots, n$ es el i -ésimo elemento de la diagonal de P , y similarmente, se encuentra así,

$$u_i^T P u_i = u_i^T X (X^T X)^{-1} X^T u_i = x_i^T (X^T X)^{-1} x_i = p_{ii},$$

Algunos resultados donde están involucrados los elementos de P , son dados a continuación.

1. $Var(\widehat{y}_i) = \sigma^2 p_{ii}$,
2. $Var(\widehat{\epsilon}_i) = \sigma^2 (1 - p_{ii})$,
3. $Cov(\widehat{\epsilon}_i, \widehat{\epsilon}_j) = -\sigma^2 p_{ij}$, y
4. $Corr(\widehat{\epsilon}_i, \widehat{\epsilon}_j) = \frac{-p_{ij}}{\sqrt{(1 - p_{ii})} \sqrt{(1 - p_{jj})}}$,

notando que el coeficiente de correlación es determinado totalmente por los elementos de P .

Una característica importante de las matrices de predicción y de residuales, es que son simétricas e idempotentes, es decir, son proyecciones ortogonales. P proyecta ortogonalmente cualquier vector n -dimensional en el espacio k -dimensional de las columnas de X , esto es, $PX = X$, por lo que se deduce que $(I - P)X = 0$, y el vector de residuales estimado se puede expresar como

$$\hat{\epsilon} = (I - P)Y = (I - P)(X\beta + \epsilon) = (I - P)\epsilon,$$

en donde la relación entre $\hat{\epsilon}$ y ϵ depende sólo de P .

Ahora, si se considera la igualdad dada en la ecuación (2.1), se observa que el i -ésimo valor predicho \hat{y}_i , se puede calcular de la siguiente manera,

$$\hat{y}_i = \sum_{j=1}^n p_{ij}y_j = p_{ii}y_i + \sum_{j \neq i} p_{ij}y_j, \quad i = 1, 2, \dots, n,$$

de donde se sigue que,

$$\frac{\partial \hat{y}_i}{\partial y_i} = p_{ii}, \quad i = 1, 2, \dots, n$$

Por consiguiente, p_{ii} puede interpretarse como el valor “palanca” que cada y_i tiene en determinado \hat{y}_i (Hoaglin y Welsch, 1978). Similarmente, p_{ij} puede ser interpretado como el valor “palanca” que cada y_j tiene en determinado \hat{y}_i .

Por otro lado, el recíproco de p_{ii} se puede considerar como el número de observaciones que determinan \hat{y}_i , Huber (1977 y 1981); esto es, si $p_{ii} = 1$, \hat{y}_i está determinado por una sola observación (solamente por y_i), por lo tanto se requiere un grado de libertad para ajustar y_i , Belsley *et al.* (1980); si $p_{ii} = 0$, y_i no tiene influencia sobre \hat{y}_i ; y si $p_{ii} = 0.5$, \hat{y}_i está determinado por dos observaciones. Una amplia variación en los valores que toma p_{ii} , indican un espaciamiento no homogéneo en las filas de X (Behnken y Draper, 1972).

Se observa que los elementos de P tienen una interesante interpretación geométrica, esto es,

1. Cuando \mathbf{X} contiene una columna constante o cuando las columnas de \mathbf{X} son centradas, la forma cuadrática $\mathbf{V}^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{V} = c$, donde \mathbf{V} es un vector $k \times 1$ y c es una constante, define contornos elípticos k -dimensionales centrados al promedio de las columnas de \mathbf{X} , es decir, $\bar{\mathbf{X}}$. El conjunto convexo más pequeño que contiene a los n puntos en el espacio \mathbf{X} , está contenido en los elipsoides que satisfacen la desigualdad $c \leq \max(\mathbf{p}_{ii})$. En consecuencia, \mathbf{p}_{ii} está determinado por la colocación de \mathbf{x}_i en el espacio \mathbf{X} , es decir, un valor grande (o pequeño) de \mathbf{p}_{ii} , indica que \mathbf{x}_i se encuentra lejos (o cerca) de la nube de puntos del conjunto de datos.
2. El volumen del elipsoide de confianza $100(1-\alpha)\%$ para β es monótono creciente con \mathbf{p}_{ii} , es decir, un \mathbf{p}_{ii} grande, provoca un incremento en el elipsoide de confianza para β cuando la i -ésima observación es eliminada.
3. Considerando la regresión de \mathbf{u}_i sobre \mathbf{X} . El vector de residuales para esta regresión es

$$\hat{\epsilon}_{\mathbf{u}_i, \mathbf{X}} = (\mathbf{I} - \mathbf{P})\mathbf{u}_i,$$

el cual, simplemente es la i -ésima columna (fila) de $(\mathbf{I} - \mathbf{P})$. La norma del vector de residuales es

$$\|\hat{\epsilon}_{\mathbf{u}_i, \mathbf{X}}\| = \sqrt{\mathbf{u}_i^T(\mathbf{I} - \mathbf{P})^T(\mathbf{I} - \mathbf{P})\mathbf{u}_i} = \sqrt{\mathbf{u}_i^T(\mathbf{I} - \mathbf{P})\mathbf{u}_i} = \sqrt{1 - \mathbf{p}_{ii}}$$

El ángulo θ_i entre $\hat{\epsilon}_{\mathbf{Y}, \mathbf{X}}$ y $\hat{\epsilon}_{\mathbf{u}_i, \mathbf{X}}$ se encuentra de la siguiente forma,

$$\begin{aligned} \cos\theta_i &= \frac{\hat{\epsilon}_{\mathbf{Y}, \mathbf{X}}^T \hat{\epsilon}_{\mathbf{u}_i, \mathbf{X}}}{\|\hat{\epsilon}_{\mathbf{Y}, \mathbf{X}}\| \cdot \|\hat{\epsilon}_{\mathbf{u}_i, \mathbf{X}}\|} \\ &= \frac{\mathbf{Y}^T(\mathbf{I} - \mathbf{P})^T(\mathbf{I} - \mathbf{P})\mathbf{u}_i}{\sqrt{\mathbf{Y}^T(\mathbf{I} - \mathbf{P})^T(\mathbf{I} - \mathbf{P})\mathbf{Y}} \sqrt{1 - \mathbf{p}_{ii}}} \\ &= \frac{\mathbf{Y}^T(\mathbf{I} - \mathbf{P})\mathbf{u}_i}{\sqrt{\hat{\epsilon}^T \hat{\epsilon}} \sqrt{1 - \mathbf{p}_{ii}}} \\ &= \frac{\hat{\epsilon}_i}{\sqrt{\hat{\epsilon}^T \hat{\epsilon}} \sqrt{1 - \mathbf{p}_{ii}}} \end{aligned}$$

De este modo, el i -ésimo residual se puede expresar como

$$\hat{\epsilon}_i = \cos\theta_i \sqrt{\hat{\epsilon}^T \hat{\epsilon} (1 - p_{ii})}$$

Ellenberg (1973, 1976), ha demostrado que si $p_{ii} \neq 1$, entonces

$$\cos^2\theta_i \sim \text{Beta}\left(\frac{1}{2}, \frac{1}{2}(n - k - 1)\right)$$

Por su parte, Beckman y Trussel (1974), han mostrado que si $p_{ii} \neq 1$, entonces

$$\cot\theta_i \sqrt{n - k - 1} \sim t(n - k - 1)$$

Propiedades de la Matriz de Predicción

A continuación se muestran algunas propiedades de la matriz P , las cuales se dan en los siguientes teoremas.

Teorema 1. *La matriz de predicción P es invariante bajo transformaciones lineales no singulares de la forma $X \rightarrow XE$, donde E es cualquier matriz de tamaño $k \times k$ no singular.*

Demostración. Sea E cualquier matriz de tamaño $k \times k$ no singular y sean las matrices de predicción para X y para XE , denotadas por P_X y P_{XE} , respectivamente, entonces sólo hay que mostrar que $P_{XE} = P_X$, para lo cual, se ve que

$$P_{XE} = (XE) \left((XE)^T (XE) \right)^{-1} (XE)^T = X(X^T X)^{-1} X^T = P_X \blacksquare$$

Este teorema indica que las matrices de predicción para X y XE son las mismas, por lo tanto, cuando se hace una regresión, el vector de residuales estimado para XE es el mismo que para X , es decir,

$$\hat{\epsilon}_{Y \cdot XE} = (I - P_{XE})Y = (I - P_X)Y = \hat{\epsilon}_{Y \cdot X}$$

De este modo, si el modelo de la ecuación (1.1) contiene un término constante, los residuales y sus varianzas estimadas son invariantes bajo transformaciones de parámetros de escala y de localización en \mathbf{X} ; mientras que si el modelo no contiene un término constante, solo son invariantes bajo transformaciones de parámetro de escala en \mathbf{X} .

Este teorema, además, permite cualquier reparametrización no singular del modelo de la ecuación (1.1), es decir, si $\boldsymbol{\alpha} = \mathbf{E}^{-1}\boldsymbol{\beta}$, entonces los modelos

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad \text{y} \quad \mathbf{Y} = \mathbf{X}\mathbf{E}\boldsymbol{\alpha} + \boldsymbol{\epsilon},$$

son equivalentes, en el sentido de que producen el mismo $\widehat{\mathbf{Y}}$.

Teorema 2. \mathbf{P} y $(\mathbf{I} - \mathbf{P})$ son matrices simétricas e idempotentes.

Demostración. Para demostrar simetría, hay que ver que

$$\mathbf{P} = \mathbf{P}^T \quad \text{y} \quad \text{que} \quad (\mathbf{I} - \mathbf{P}) = (\mathbf{I} - \mathbf{P})^T$$

Mientras que para mostrar idempotencia, se ve que

$$\mathbf{P} = \mathbf{P}^2 \quad \text{y} \quad \text{que} \quad (\mathbf{I} - \mathbf{P}) = (\mathbf{I} - \mathbf{P})^2 \blacksquare$$

Teorema 3. Sea \mathbf{X} una matriz de tamaño $n \times k$, entonces

1. $\text{tr}(\mathbf{P}) = r(\mathbf{P}) = k$,
2. $\text{tr}(\mathbf{I} - \mathbf{P}) = n - k$, y
3. $\sum_{i=1}^n \sum_{j=1}^n p_{ij}^2 = k$.

Demostración.

1. Para toda matriz de proyección \mathbf{P} , la $\text{tr}(\mathbf{P}) = r(\mathbf{P}) = k$.
2. Se ve que la $\text{tr}(\mathbf{I} - \mathbf{P}) = \text{tr}(\mathbf{I}) - \text{tr}(\mathbf{P}) = n - k$.

3. Si $P = P^T = P^2$, entonces, la $tr(P) = tr(P^2) = tr(P^T P) = \sum_{i=1}^n \sum_{j=1}^n p_{ij}^2 = k$. ■

Para la demostración del Teorema 4, se requiere citar el siguiente lema.

Lema 1. *La inversa de una matriz particionada.*

Se define

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

1. Si A y A_{11} son no singulares, entonces:

$$A^{-1} = \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1} A_{12} M A_{21} A_{11}^{-1} & -A_{11}^{-1} A_{12} M \\ -M A_{21} A_{11}^{-1} & M \end{bmatrix}$$

$$\text{donde } M = (A_{22} - A_{21} A_{11}^{-1} A_{12})^{-1}$$

2. Si A y A_{22} son no singulares, entonces:

$$A^{-1} = \begin{bmatrix} N & -N A_{12} A_{22}^{-1} \\ -A_{22}^{-1} A_{12} N & A_{22}^{-1} + A_{22}^{-1} A_{21} N A_{12} A_{22}^{-1} \end{bmatrix}$$

$$\text{donde } N = (A_{11} - A_{12} A_{22}^{-1} A_{21})^{-1}$$

Demostración. Si se considera $B = A^{-1}$, entonces sólo se debe mostrar que $BA = I$.

■

Teorema 4. Sea $X = (X_1 : X_2)$, donde X_1 es una matriz de tamaño $n \times r$ de rango r y X_2 es una matriz $n \times (k - r)$ de rango $(k - r)$. Sea $P_1 = X_1 (X_1^T X_1)^{-1} X_1^T$ la matriz de predicción para X_1 , y $W = (I - P_1) X_2$ es la proyección ortogonal de X_2 sobre el complemento de X_1 . Por último, sea $P_2 = W (W^T W)^{-1} W^T$ la matriz de predicción para W ; entonces, P puede ser expresada como

$$P = P_1 + P_2$$

$$\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X}_1(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T + (\mathbf{I} - \mathbf{P}_1) \mathbf{R}_2 (\mathbf{I} - \mathbf{P}_1)$$

donde $\mathbf{R}_2 = \mathbf{X}_2 (\mathbf{X}_2^T (\mathbf{I} - \mathbf{P}_1) \mathbf{X}_2)^{-1} \mathbf{X}_2^T$.

Demostración. Se ve que \mathbf{P} se puede calcular así,

$$\mathbf{P} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}_1^T \\ \mathbf{X}_2^T \end{bmatrix}$$

Al invertir la matriz particionada que se encuentra en la parte central de la expresión anterior, utilizando para ello el Lema 1, se llega a que $\mathbf{P} = \mathbf{P}_1 + \mathbf{P}_2$, quedando probado el teorema. ■

Este teorema permite además, particionar la matriz \mathbf{P} en dos o más matrices de predicción.

Teorema 5. Para $i = 1, 2, \dots, n$ y $j = 1, 2, \dots, n$, se tiene que

1. $0 \leq p_{ii} \leq 1$, para todo i ,
2. $0.5 \leq p_{ij} \leq 0.5$, para todo $j \neq i$,
3. Si \mathbf{X} contiene una columna constante, entonces
 - (a) $p_{ii} \geq n^{-1}$, para todo i ,
 - (b) $\mathbf{P}\mathbf{1} = \mathbf{1}$, y
4. Suponga que \mathbf{x}_i ocurre a veces y que $-\mathbf{x}_i$ ocurre b veces, entonces

$$p_{ii} \leq (a + b)^{-1}.$$

Demostración. Para el inciso 3, se tiene

- (a) Si \mathbf{X} contiene una columna constante, entonces se puede definir como

$$\mathbf{X} = (\mathbf{1} : \mathbf{X}_2),$$

por el Teorema 4, se tiene que

$$\begin{aligned}
P_1 &= \mathbf{1}(\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T = n^{-1} \mathbf{1} \mathbf{1}^T = n^{-1} \mathbf{1} \mathbf{1}^T \\
W &= (I - P_1) X_2 = (I - n^{-1} \mathbf{1} \mathbf{1}^T) X_2 \equiv \tilde{X} \\
P_2 &= \tilde{X} (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T
\end{aligned}$$

De este modo, la matriz de predicción P se puede expresar como

$$P = P_1 + P_2 = n^{-1} \mathbf{1} \mathbf{1}^T + \tilde{X} (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T$$

Se ve que cada uno de los elementos de la diagonal de P_1 es igual a n^{-1} . Como P_2 es una matriz de predicción, por el Teorema 5(1), se deduce que los elementos de la diagonal son no negativos, por lo que $p_{ii} \geq n^{-1}$, para todo i .

(b) Se ve que $\tilde{X}^T \mathbf{1} = \mathbf{0}$, y que $P_2 \mathbf{1} = \mathbf{0}$, entonces $P \mathbf{1} = P_1 \mathbf{1} = \mathbf{1}$.

La demostración completa aparece en Chatterjee y Hadi (1988). ■

Teorema 6. Para $i = 1, 2, \dots, n$ y $j = 1, 2, \dots, n$, se tiene que

1. Si $p_{ii} = 1$ o 0 , entonces $p_{ij} = 0$,
2. $(1 - p_{ii})(1 - p_{jj}) - p_{ij}^2 \geq 0$,
3. $p_{ii} p_{jj} - p_{ij}^2 \geq 0$, y
4. $p_{ii} + \frac{\hat{\epsilon}_i^2}{\hat{\epsilon}^T \hat{\epsilon}} \leq 1$.

Demostración.

1. Se observa que el i -ésimo elemento diagonal de P , puede ser escrito como

$$p_{ii} = \sum_{j=1}^n p_{ij}^2 = p_{ii}^2 + \sum_{j \neq i}^n p_{ij}^2,$$

por lo que, si $p_{ii} = 0$, entonces $p_{ij} = 0$; y si $p_{ii} = 1$, entonces $p_{ij} = 0$.

2. Se puede ver en Chatterjee y Hadi (1988).

3. Simplificando la expresión del Teorema 6(2), se comprueba que

$$p_{ii}p_{jj} - p_{ij}^2 \geq 0.$$

4. Se define $Z = (X : Y)$, $P_X = X(X^T X)^{-1} X^T$ y $P_Z = Z(Z^T Z)^{-1} Z^T$.

Utilizando el Teorema 4, se ve que

$$P_Z = P_X + \frac{(I - P_X) Y Y^T (I - P_X)}{Y^T (I - P_X) Y} = P_X + \frac{\widehat{\epsilon} \widehat{\epsilon}^T}{\widehat{\epsilon}^T \widehat{\epsilon}}$$

De esta manera, los elementos de la diagonal de P_Z son menores o iguales a uno, tomando el i -ésimo elemento, resulta

$$p_{ii} + \frac{\widehat{\epsilon}_i^2}{\widehat{\epsilon}^T \widehat{\epsilon}} \leq 1 \blacksquare$$

El Teorema 6(1), 6(2) y 6(3), indica que si p_{ii} es grande (cercano a 1), o pequeño (cercano a 0), entonces p_{ij} es pequeño para toda $j \neq i$. El Teorema 6(4) indica que para p_{ii} grande, el i -ésimo residual estimado $\widehat{\epsilon}_i$ es pequeño.

Omitiendo o Adicionando una Observación

Para determinar los efectos que resultan de omitir o adicionar observaciones en P , son necesarios los siguientes resultados.

Lema 2. Sean A y D matrices no singulares de órdenes $k \times k$ y $m \times m$, respectivamente. Sean además B y C , ambas de orden $k \times m$; entonces, existe la siguiente inversa,

$$(A + BDC^T)^{-1} = A^{-1} - A^{-1}B(D^{-1} + C^T A^{-1}B)^{-1}C^T A^{-1}$$

Demostración. Sólo debe verse que

$$(A + BDC^T) \left(A^{-1} - A^{-1}B(D^{-1} + C^T A^{-1}B)^{-1}C^T A^{-1} \right) = I \blacksquare$$

Un caso para el que la inversa de esta matriz no existe, se tiene cuando $A = I$, $B = X$, $C^T = X^T$ y $D = -(X^T X)^{-1}$, así, $(A + BDC^T) = (I - P)$, ésta última, como se sabe, es singular.

Este lema, también es conocido como el Teorema de Sherman-Morrison-Woodbury (Henderson y Searle, 1981).

Sustituyendo $B = x_i$, $C^T = x_i^T$, $D = 1$ y $A = (X_{(i)}^T X_{(i)})$, en la expresión del Lema 2, se obtiene

$$\begin{aligned} (X^T X)^{-1} &= (X_{(i)}^T X_{(i)} + x_i x_i^T)^{-1} \\ &= (X_{(i)}^T X_{(i)})^{-1} - \frac{(X_{(i)}^T X_{(i)})^{-1} x_i x_i^T (X_{(i)}^T X_{(i)})^{-1}}{1 + x_i^T (X_{(i)}^T X_{(i)})^{-1} x_i} \end{aligned} \quad (2.3)$$

Sustituyendo ahora, $B = x_i$, $C^T = x_i^T$, $D = -1$ y $A = X^T X$, resulta

$$\begin{aligned} (X_{(i)}^T X_{(i)})^{-1} &= (X^T X - x_i x_i^T)^{-1} \\ &= (X^T X)^{-1} + \frac{(X^T X)^{-1} x_i x_i^T (X^T X)^{-1}}{1 - x_i^T (X^T X)^{-1} x_i} \end{aligned} \quad (2.4)$$

Para ver los efectos sobre la matriz de predicción, al adicionar una observación x_i^T , se utiliza la ecuación (2.3). A continuación se presentan algunos ejemplos.

1. Si $p_{rc} = x_r^T (X^T X)^{-1} x_c$, entonces

$$p_{rc} = p_{rc(i)} - \frac{p_{ri(i)} p_{ic(i)}}{1 + p_{ii(i)}}, \quad \text{para } r, c \neq i,$$

donde $p_{rc(i)} = x_r^T (X_{(i)}^T X_{(i)})^{-1} x_c$.

2. Si $p_{ri} = p_{ir} = x_r^T (X^T X)^{-1} x_i$, entonces

$$p_{ir} = p_{ir(i)} - \frac{p_{ii(i)} p_{ir(i)}}{1 + p_{ii(i)}} = \frac{p_{ir(i)}}{1 + p_{ii(i)}}, \quad \text{para } r \neq i,$$

donde $p_{ir(i)} = x_r^T (X_{(i)}^T X_{(i)})^{-1} x_i$.

3. Si $p_{ii} = x_i^T (X^T X)^{-1} x_i$, entonces

$$p_{ii} = p_{ii(i)} - \frac{p_{ii(i)}^2}{1 + p_{ii(i)}} = \frac{p_{ii(i)}}{1 + p_{ii(i)}},$$

donde $p_{ii(i)} = x_i^T (X_{(i)}^T X_{(i)})^{-1} x_i$.

Por otra parte, para ver los efectos sobre la matriz de predicción, al omitir una observación x_i^T , se usa la ecuación (2.4). Observando por ejemplo lo que sucede con $p_{cr(i)}$, esto es,

$$p_{cr(i)} = p_{cr} + \frac{p_{ci} p_{ir}}{1 - p_{ii}}, \quad \text{para } r, c \neq i$$

y cuando $r = c$, entonces la expresión anterior se simplifica, resultando

$$p_{rr(i)} = p_{rr} + \frac{p_{ir}^2}{1 - p_{ii}}.$$

Cuando la i -ésima observación es omitida, la anterior expresión muestra que para $r \neq i$, $p_{rr(i)}$ es grande, si p_{ii} , p_{rr} y/o $|p_{ir}|$ son grandes. Ahora, cuando la i -ésima y r -ésima filas de X son iguales, entonces $p_{rr(i)}$ se reduce a

$$p_{rr(i)} = \frac{p_{ii}}{1 - p_{ii}},$$

observando que si p_{ii} está próximo a 0.5, entonces $p_{rr(i)}$ está próximo a 1, en tal caso la i -ésima (r -ésima) observación puede no ser detectada, cuando ambas observaciones están presentes.

Notando que, implícitamente se supone que existe la inversa de $(X_{(i)}^T X_{(i)})$, o equivalentemente, el $r(X_{(i)}) = k$. No existe $(X_{(i)}^T X_{(i)})^{-1}$, si la eliminación de la i -ésima observación resulta en un modelo de rango deficiente.

Condiciones para Valores Grandes de p_{ii}

Para el caso de una regresión simple de la variable Y con X , se tiene

$$P = \frac{X X^T}{\sum_{r=1}^n x_r^2},$$

de donde se deduce que

$$p_{ii} = \frac{x_i^2}{\sum_{r=1}^n x_r^2}, \quad i = 1, 2, \dots, n$$

Si un término constante es incluido en X , y \bar{X} denota el promedio de X , entonces utilizando la prueba del Teorema 5(3(a)), p_{ii} se puede escribir como

$$p_{ii} = \frac{1}{n} + \frac{(x_i - \bar{X})^2}{\sum_{r=1}^n (x_r - \bar{X})^2}, \quad i = 1, 2, \dots, n$$

De este modo, p_{ii} será grande, si x_i se encuentra alejado de la nube de puntos de un determinado conjunto de datos.

En el caso de una regresión múltiple, considerando λ_r y V_r , $r = 1, 2, \dots, k$, los eigenvalores y sus correspondientes eigenvectores normalizados de $X^T X$, respectivamente. Si θ_{ir} es el ángulo entre x_i y V_r , entonces

$$p_{ij} = \|x_i\| \cdot \|x_j\| \sum_{r=1}^k \lambda_r^{-1} \cos^2 \theta_{ir},$$

y así,

$$p_{ii} = x_i^T x_i \sum_{r=1}^k \lambda_r^{-1} \cos^2 \theta_{ir}.$$

Notando que si x_i y V_r , tienen la misma dirección, entonces el $\cos^2 \theta_{ir} = 1$, y si son perpendiculares, entonces $\cos^2 \theta_{ir} = 0$. De igual forma, p_{ii} será grande si x_i se encuentra alejado de la nube de puntos de un determinado conjunto de datos (Cook y Weisberg, 1982).

Omitiendo Múltiples Filas de \mathbf{X}

Los efectos sobre \mathbf{P} , que se tienen al omitir un subconjunto de filas de \mathbf{X} , se muestran en el siguiente teorema.

Teorema 7. *Suponiendo que \mathbf{X} es de orden $n \times k$ de rango k y que existen $m < k$ elementos de la diagonal de \mathbf{P} iguales a uno, y sea $I = \{i : p_{ii} = 1, i = 1, 2, \dots, n\}$ el conjunto de sus índices; entonces, para cualquier $J \supseteq I$, el $r(\mathbf{X}_{(J)}) \leq k - m$, con la igualdad sólo si $J = I$.*

Demostración. Ajustando las filas de \mathbf{X} , de tal forma que las m filas indexadas por I , sean las últimas m filas de \mathbf{X} , entonces \mathbf{X} se puede escribir así,

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_{(I)} \\ \mathbf{X}_I^T \end{bmatrix} \begin{matrix} (n - m) \times k \\ m \times k \end{matrix}$$

Se define a \mathbf{P}_I como el menor principal de \mathbf{P} correspondiente a \mathbf{X}_I^T . Y como $p_{ii} = 1$, $i \in I$, entonces $\mathbf{P}_I = \mathbf{I}$, donde \mathbf{I} es de dimensión $m \times m$. Utilizando el Teorema 6(1), \mathbf{P} puede expresarse como

$$\mathbf{P} = \begin{bmatrix} \mathbf{X}_{(I)}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_{(I)}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix},$$

en donde puede observarse que $\left(\mathbf{X}_{(I)}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_{(I)}^T\right)$ es idempotente con rango $(k - m)$, por lo tanto,

$$(k - m) = r\left(\mathbf{X}_{(I)}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_{(I)}^T\right) = r(\mathbf{X}_{(I)}),$$

de donde se sigue que, para cualquier $J \supseteq I$, el

$$r(\mathbf{X}_{(J)}) \leq r(\mathbf{X}_{(I)}) = (k - m) \blacksquare$$

Beckman y Trussel (1974), señalan que si se emplean las propiedades de la inversa generalizada para \mathbf{X} , bajo las condiciones del modelo de la ecuación (1.1),

entonces si $p_{ii}=1$, el $r(\mathbf{X}_{(i)})=(k-1)$, y por lo tanto $(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1}$ no existe. Esto representa un caso especial de este teorema, donde $m=1$, de manera que si $I=\{i\}$ y $p_{ii}=1$, entonces se tiene que

$$r(\mathbf{X}_{(J)}) \leq r(\mathbf{X}_{(i)}) = (k-1)$$

Notando que este teorema establece que si para algún $i \in I$, $p_{ii}=1$, entonces el

$$r(\mathbf{X}_{(I)}) \leq k.$$

Eigenvalores de P y $(I - P)$

Las propiedades de los eigenvalores de P y $(I - P)$, son dadas en los tres teoremas siguientes.

Teorema 8. *Los eigenvalores de P y $(I - P)$ son 0's o 1's.*

Demostración. La demostración se sigue directamente del hecho de que los eigenvalores de matrices idempotentes son 0's o 1's.

■

Teorema 9. *Existen $(n - k)$ eigenvalores de P iguales a 0, y los restantes k eigenvalores son iguales a 1. Similarmente, k eigenvalores de $(I - P)$ son iguales a 0 y $(n - k)$ eigenvalores son iguales a 1.*

Demostración. Sean λ_i , $i = 1, 2, \dots, n$, los eigenvalores de P . Usando el Teorema 3(1), se tiene

$$\sum_{i=1}^n \lambda_i = \text{tr}(P) = k,$$

pero el Teorema 8 afirma que $\lambda_i=0$ o 1 para todo i , por lo tanto, k de los eigenvalores de P son iguales a 1 y los restantes $(n - k)$ eigenvalores son iguales a 0. Se aplica un argumento similar para $(I - P)$. ■

Para la prueba del Teorema 10 es necesario ver los siguientes dos lemas.

Lema 3. *El determinante de una matriz particionada. Se define la partición de S , como sigue,*

$$S = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

1. Si A es no singular, entonces el $\det(S) = \det(A)\det(D - CA^{-1}B)$, y
2. Si D es no singular, entonces el $\det(S) = \det(D)\det(A - BD^{-1}C)$.

Demostración.

1. Sea S definida como antes, y sea además S^* definida de la siguiente manera,

$$S^* = \begin{bmatrix} A^{-1} & 0 \\ -CA^{-1} & I \end{bmatrix},$$

entonces se ve que el $\det(S^*) = \det(A^{-1}) = \frac{1}{\det(A)}$.

Considerando ahora el producto

$$S^*S = \begin{bmatrix} A^{-1} & 0 \\ -CA^{-1} & I \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} I & A^{-1}B \\ 0 & -CA^{-1}B + D \end{bmatrix},$$

por lo que el

$$\det(S^*S) = \det(D - CA^{-1}B),$$

o sea,

$$\det(S^*)\det(S) = \det(D - CA^{-1}B),$$

finalmente se llega a

$$\begin{aligned} \det(S) &= \frac{\det(D - CA^{-1}B)}{\det(S^*)} \\ &= \frac{\det(D - CA^{-1}B)}{\frac{1}{\det(A)}} \\ &= \det(A)\det(D - CA^{-1}B) \end{aligned}$$

2. Se aplica un argumento similar al utilizado en la prueba del Lema 3(1). ■

Lema 4. Sean B y C matrices de orden $k \times m$. Si A es una matriz de orden $k \times k$ no singular, entonces el

$$\det(A - BC^T) = \det(A)\det(I - C^T A^{-1}B).$$

Demostración. Por el Lema 3(1), se tiene que el

$$\det \begin{bmatrix} A & B \\ C^T & I \end{bmatrix} = \det(A)\det(I - C^T A^{-1}B),$$

y por el Lema 3(2), se tiene que el

$$\det \begin{bmatrix} A & B \\ C^T & I \end{bmatrix} = \det(I)\det(A - BI^{-1}C^T) = \det(A - BC^T),$$

Por lo tanto se comprueba que el $\det(A - BC^T) = \det(A)\det(I - C^T A^{-1}B)$.

■

Teorema 10. Sea P_I de orden $m \times m$, un menor de P , dado por las filas y columnas de P indexadas por I . Si $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_m$ son los eigenvalores de P_I , entonces,

1. Los eigenvalores de P_I y de $(I - P_I)$ están entre 0 y 1, esto es, $0 \leq \lambda_j \leq 1$, $j = 1, 2, \dots, m$,
2. $(I - P_I)$ es una matriz positiva definida si $\lambda_m < 1$, si no, es positiva semi-definida, y
3. El $r(X_{(I)}) < k$, si $\lambda_m = 1$.

Demostración.

1. Sea P_{22} de orden $m \times m$ un menor de P , es decir, P_{22} son las últimas m filas y columnas de P . Se considera la partición de P , como

$$P = \begin{bmatrix} P_{11} & P_{12} \\ P_{12}^T & P_{22} \end{bmatrix}$$

Como P_{22} es simétrica, entonces por el Teorema de la descomposición espectral (Searle, 1982), puede escribirse como $P_{22} = V\Lambda V^T$, donde Λ es una matriz diagonal con los eigenvalores de P_{22} en la diagonal, y V es una matriz que contiene como columnas los correspondientes eigenvalores normalizados. Ahora, se ve que $P = PP$, entonces

$$P_{22} = P_{12}^T P_{12} + P_{22} P_{22},$$

y como $P_{12}^T P_{12} \geq 0$, entonces $P_{22} \geq P_{22} P_{22}$. Sustituyendo en ésta última expresión el valor de P_{22} , se tiene

$$V\Lambda V^T \geq V\Lambda V^T V\Lambda V^T = V\Lambda^2 V^T,$$

o sea,

$$V\Lambda V^T - V\Lambda^2 V^T \geq 0,$$

luego,

$$V(\Lambda - \Lambda^2)V^T \geq 0,$$

y finalmente se llega a

$$(\Lambda - \Lambda^2) \geq 0,$$

la cual también se puede poner de la forma $(\lambda - \lambda^2) \geq 0$, de donde se deduce que los eigenvalores de P_{22} están entre 0 y 1.

2. Si $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_m$ son los eigenvalores de P_{22} , entonces $(1 - \lambda_j)$, $j = 1, 2, \dots, m$, son los eigenvalores de $(I - P_{22})$. Ahora, si $\lambda_m = 1$, entonces $(1 - \lambda_m) = 0$ y por lo tanto $(I - P_{22})$ es positiva semidefinida. Pero si $\lambda_m < 1$, entonces $(1 - \lambda_i) > 0$, $i = 1, 2, \dots, m$, y por lo tanto $(I - P_{22})$ es positiva definida.

3. Se ve que la expresión $(\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})$ se puede escribir como $(\mathbf{X}^T \mathbf{X} - \mathbf{X}_I \mathbf{X}_I^T)$, y al aplicar el Lema 4, se obtiene

$$\begin{aligned} \det(\mathbf{X}_{(I)}^T \mathbf{X}_{(I)}) &= \det(\mathbf{X}^T \mathbf{X} - \mathbf{X}_I \mathbf{X}_I^T) \\ &= \det(\mathbf{X}^T \mathbf{X}) \det\left(\mathbf{I} - \mathbf{X}_I^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I\right) \\ &= \det(\mathbf{X}^T \mathbf{X}) \det(\mathbf{I} - \mathbf{P}_I), \end{aligned}$$

de este modo, el

$$\det(\mathbf{I} - \mathbf{P}_I) = \prod_{j=1}^m (1 - \lambda_j),$$

concluyéndose que si $\lambda_m = 1$, entonces, el $\det(\mathbf{X}_{(I)}^T \mathbf{X}_{(I)}) = 0$, o equivalentemente, el $r(\mathbf{X}_{(I)}) < k$, si $\lambda_m = 1$.

■

Distribución de p_{ii}

Recordando que la matriz de predicción \mathbf{P} , depende solamente de \mathbf{X} , la cual se supone es medida sin error. En algunos casos, sin embargo, \mathbf{X} es medida con error, y entonces es razonable suponer que las filas de \mathbf{X} tienen una distribución Normal multivariada con media μ y varianza Σ . En este sentido, la distribución de p_{ii} puede ser encontrada, tal como proponen Belsley, *et al.* (1980), quienes establecen que si las filas de $\tilde{\mathbf{X}} = (\mathbf{I} - n^{-1} \mathbf{1} \mathbf{1}^T) \mathbf{X}$ son independientes e idénticamente distribuídas como una distribución Normal $(k - 1)$ -dimensional, entonces,

$$\frac{(n - k) (p_{ii} - n^{-1})}{(k - 1) (1 - p_{ii})} \sim F(k - 1, n - k),$$

donde $F(k - 1, n - k)$ es la distribución F con $(k - 1, n - k)$ grados de libertad. Pero también, recordando que $\mathbf{1}^T \tilde{\mathbf{X}} = \mathbf{0}$, por lo que no se puede suponer que las filas de $\tilde{\mathbf{X}}$, es decir, $\tilde{\mathbf{x}}_i, i = 1, 2, \dots, n$, sean independientes. Afortunadamente, se da

la dependencia de las filas de $\tilde{\mathbf{X}}$ en cada disminución o incremento que se tiene en n . Lo anterior puede verse al calcular la covarianza entre $\tilde{\mathbf{x}}_j$ y $\tilde{\mathbf{x}}_i$, es decir,

$$\begin{aligned} \mathbf{Cov}(\tilde{\mathbf{x}}_j, \tilde{\mathbf{x}}_i) &= \mathbf{E}(\tilde{\mathbf{x}}_j \tilde{\mathbf{x}}_i^T) - (\mathbf{E}(\tilde{\mathbf{x}}_j))(\mathbf{E}(\tilde{\mathbf{x}}_i))^T, \quad i \neq j \\ &= \mu\mu^T - 2(\mu\mu^T + n^{-1}\Sigma) + \mu\mu^T + n^{-1}\Sigma \\ &= -n^{-1}\Sigma, \end{aligned}$$

la cual converge a cero cuando n se aproxima a infinito.

EFFECTOS DE UNA OBSERVACIÓN SOBRE LA ECUACIÓN DE REGRESIÓN

Cuando un conjunto de datos determinado, se somete a un análisis de regresión, en donde se ajusta el modelo dado en la ecuación (1.1), utilizando el método de Cuadrados Mínimos, entonces, los resultados de este ajuste, pueden ser influenciados por la eliminación o adición de una o varias observaciones. Normalmente, no todas las observaciones presentes en un conjunto de datos, influyen de igual manera en el ajuste de dicho modelo.

Las observaciones, que individual o colectivamente, influyen (en cierta medida) en la ecuación de regresión ajustada, se les llama observaciones influyentes, y éstas pueden ser: "outliers", puntos palanca o puntos influyentes.

Las medidas que se utilizan para detectar observaciones influyentes se han agrupado en dos métodos, que son:

1. **El método de la eliminación de observaciones.** Este procedimiento estudia cómo es que varios resultados del análisis de regresión, cambian cuando alguna o algunas de las observaciones son omitidas.
2. **El método basado en la diferenciación.** Por su parte, este método estudia las modificaciones que sufren varios de los resultados del análisis de regresión con respecto a ciertos parámetros del modelo.

En este capítulo, se analizará la detección de una sola observación influyente.

Método de la Eliminación de una Observación

Existen algunas medidas que permiten detectar una observación influyente, entre las que se pueden mencionar:

1. Medidas basadas en residuales.
2. Outliers, puntos palanca y puntos influyentes.
3. Medidas basadas en la curva de influencia.
4. Medidas basadas en el volumen de elipsoides de confianza.
5. Medidas basadas en la función de verosimilitud.
6. Medidas basadas en un subconjunto de coeficientes de regresión.

Medidas Basadas en Residuales

A pesar de que los errores aleatorios ϵ no pueden ser observados, si pueden ser medidos en cierta forma por los residuales $\hat{\epsilon}$, por lo que estos últimos juegan un papel fundamental en el diagnóstico del análisis de regresión. Recordando que los residuales $\hat{\epsilon}$ pueden ser escritos como $\hat{\epsilon} = (\mathbf{I} - \mathbf{P})\epsilon$, en donde se observa que $\hat{\epsilon}$ mide razonablemente a ϵ , sólo si los elementos de la diagonal de \mathbf{P} son lo suficientemente pequeños. Notando además que, aún cuando los elementos de ϵ sean independientes y con la misma varianza, por la identidad $\hat{\epsilon} = (\mathbf{I} - \mathbf{P})\epsilon$ se ve que los residuales $\hat{\epsilon}$ no son independientes (a menos que \mathbf{P} sea diagonal), y no tienen la misma varianza (a menos que los elementos de la diagonal de \mathbf{P} sean iguales). Por lo que puede mencionarse que, $\hat{\epsilon}$ estima aproximadamente a ϵ , si las filas de \mathbf{X} son homogéneas, así, los elementos de la diagonal de \mathbf{P} son aproximadamente iguales y lo suficientemente pequeños.

Por lo anterior, en vez de usar $\widehat{\epsilon}_i$ en el diagnóstico del análisis de regresión, es preferible utilizar una versión transformada de los residuales, esto es,

$$f(\widehat{\epsilon}_i, \sigma_i) = \frac{\widehat{\epsilon}_i}{\sigma_i}, \quad (3.1)$$

donde σ_i es la desviación estándar del i -ésimo residual.

De la ecuación (3.1), se obtienen los residuales normalizados, estandarizados y dos casos de residuales estudentizados.

Residual Normalizado

El i -ésimo residual normalizado es obtenido al reemplazar σ_i por $\sqrt{\widehat{\epsilon}^T \mathbf{T} \widehat{\epsilon}}$ en la ecuación (3.1), con lo que se obtiene

$$a_i \equiv f\left(\widehat{\epsilon}_i, \sqrt{\widehat{\epsilon}^T \mathbf{T} \widehat{\epsilon}}\right) = \frac{\widehat{\epsilon}_i}{\sqrt{\widehat{\epsilon}^T \mathbf{T} \widehat{\epsilon}}}, \quad i = 1, 2, \dots, n. \quad (3.2)$$

Residual Estandarizado

El i -ésimo residual estandarizado es encontrado al sustituir la expresión

$$\widehat{\sigma} = \sqrt{\frac{\widehat{\epsilon}^T \mathbf{T} \widehat{\epsilon}}{n - k}},$$

en lugar de σ_i en la ecuación (3.1), resultando

$$b_i \equiv f(\widehat{\epsilon}_i, \widehat{\sigma}) = \frac{\widehat{\epsilon}_i}{\widehat{\sigma}}, \quad i = 1, 2, \dots, n. \quad (3.3)$$

Residual Interno Estudentizado

Si en la ecuación (3.1), se considera que

$$\sigma_i = \widehat{\sigma} \sqrt{1 - p_{ii}},$$

entonces resulta el i -ésimo residual interno estudentizado, que es

$$r_i \equiv f\left(\widehat{\epsilon}_i, \widehat{\sigma} \sqrt{1 - p_{ii}}\right) = \frac{\widehat{\epsilon}_i}{\widehat{\sigma} \sqrt{1 - p_{ii}}}, \quad i = 1, 2, \dots, n. \quad (3.4)$$

Residual Externo Estudentizado

Tomando ahora

$$\sigma_i = \hat{\sigma}_{(i)} \sqrt{1 - p_{ii}}$$

en la ecuación (3.1), se obtiene el i -ésimo residual externo estudentizado, dado por

$$r_i^* \equiv f\left(\hat{\epsilon}_i, \hat{\sigma}_{(i)} \sqrt{1 - p_{ii}}\right) = \frac{\hat{\epsilon}_i}{\hat{\sigma}_{(i)} \sqrt{1 - p_{ii}}}, \quad i = 1, 2, \dots, n, \quad (3.5)$$

donde

$$\hat{\sigma}_{(i)}^2 = \frac{Y_{(i)}^T (I - P_{(i)}) Y_{(i)}}{n - k - 1}, \quad i = 1, 2, \dots, n,$$

es el residual que estima el Cuadrado Medio del Error cuando la i -ésima observación es omitida, y

$$P_{(i)} = X_{(i)} (X_{(i)}^T X_{(i)})^{-1} X_{(i)}^T, \quad i = 1, 2, \dots, n,$$

es la matriz de predicción para $X_{(i)}$.

Para encontrar las relaciones entre los cuatro casos de residuales, primero, es conveniente ver que están en función de cuatro cantidades que son: $(n - k)$, $\hat{\epsilon}^T \hat{\epsilon}$, $\hat{\epsilon}_i$ y $(1 - p_{ii})$. Y segundo, es necesario escribir $\hat{\sigma}_{(i)}$ en términos de $\hat{\sigma}$. Esto puede verse al omitir la i -ésima observación, pues es equivalente a ajustar el modelo

$$E(Y) = X\beta + u_i\theta,$$

donde θ es el coeficiente de regresión del i -ésimo vector unitario u_i .

Usando la expresión dada en el Teorema 4, se tiene que

$$\begin{aligned} SSE_{(i)} &= Y_{(i)}^T (I - P_{(i)}) Y_{(i)} \\ &= Y^T \left(I - P - \frac{(I - P) u_i u_i^T (I - P)}{u_i^T (I - P) u_i} \right) Y \\ &= SSE - \frac{\hat{\epsilon}_i^2}{1 - p_{ii}} \end{aligned} \quad (3.6)$$

donde SSE es la Suma de Cuadrados de Residuales y $SSE_{(i)}$ es la Suma de Cuadrados de Residuales cuando la i -ésima observación es omitida. Observando que al dividir ambos lados de la última expresión por $(n - k - 1)$, se obtiene

$$\begin{aligned}\hat{\sigma}_{(i)}^2 &= \frac{(n - k)}{(n - k - 1)}\hat{\sigma}^2 - \frac{\hat{\epsilon}_i^2}{(n - k - 1)(1 - p_{ii})} \\ &= \hat{\sigma}^2 \left(\frac{n - k - r_i^2}{n - k - 1} \right)\end{aligned}\quad (3.7)$$

Además, $b_i = a_i\sqrt{n - k}$, y por lo tanto

$$r_i = \frac{b_i}{\sqrt{1 - p_{ii}}} = a_i\sqrt{\frac{n - k}{1 - p_{ii}}}$$

Si en la ecuación (3.5), se reemplaza $\hat{\sigma}_{(i)}$ por $\sqrt{\frac{SSE_{(i)}}{(n - k - 1)}}$, y posteriormente se utiliza la ecuación (3.6), entonces r_i^* se puede escribir como

$$r_i^* = \frac{\hat{\epsilon}_i}{\sqrt{\frac{(1 - p_{ii})SSE_{(i)}}{n - k - 1}}} = \frac{\hat{\epsilon}_i}{\sqrt{\left(\frac{1 - p_{ii}}{n - k - 1}\right) \left(SSE - \frac{\hat{\epsilon}_i^2}{1 - p_{ii}}\right)}} = \frac{a_i\sqrt{n - k - 1}}{\sqrt{(1 - p_{ii}) - a_i^2}}$$

Asimismo, la sustitución de la ecuación (3.7) en la ecuación (3.5), resulta otra forma de escribir r_i^* , esto es,

$$r_i^* = r_i\sqrt{\frac{n - k - 1}{n - k - r_i^2}}$$

Las distribuciones de r_i y de r_i^* se presentan en los siguientes teoremas, cuyas pruebas pueden verse en Chatterjee y Hadi (1988).

Teorema 11. *Suponiendo que \mathbf{X} en el modelo dado en la ecuación (1.1) es de rango k , entonces:*

1. *Ellenberg (1973), señala que si el $r(\mathbf{X}_{(i)})=k$ y $\epsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2\mathbf{I})$, entonces $\left(\frac{r_i^2}{n - k}\right)$, para $i = 1, 2, \dots, n$, son idénticamente distribuidos como Beta $\left(\frac{1}{2}, \frac{1}{2}(n - k - 1)\right)$, y*
2. *Beckman y Trussel (1974), indican que si el $r(\mathbf{X}_{(i)})=(k - 1)$, entonces r_i no está definido, esto es, $\hat{\epsilon}_i=0$ y $\mathbf{Var}(\hat{\epsilon}_i)=0$.*

Teorema 12. *Beckman y Trussel (1974), afirman que si se supone que \mathbf{X} en el modelo dado en la ecuación (1.1) es de rango k , entonces:*

1. *Si el $r(\mathbf{X}_{(i)})=k$ y $\epsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I})$, entonces r_i^* , para $i = 1, 2, \dots, n$, son idénticamente distribuidos como $t(n - k - 1)$, y*
2. *Si el $r(\mathbf{X}_{(i)})=(k - 1)$, entonces r_i^* no está definido.*

Para probar la suficiencia de los supuestos del modelo de la ecuación (1.1), no se debe preferir alguna de las formas de residuales sobre otra, ya que ninguna es independientemente distribuida, pero para la mayoría de los problemas prácticos, especialmente cuando el tamaño de muestra es grande, se pueden considerar independientes, tal como lo señalan Behnken y Draper (1972), Chatterjee y Price (1977) y Draper y Smith (1981).

Behnken y Draper (1972), mencionan además, que los residuales normalizados y los estandarizados son básicamente equivalentes, pues son múltiplos constantes de $\widehat{\epsilon}_i$, sin embargo, también sugieren que si varían los elementos de la diagonal de \mathbf{P} , y por consiguiente, las varianzas de $\widehat{\epsilon}_i$, para $i = 1, 2, \dots, n$, sería preferible usar r_i .

Algunos otros autores como Atkinson (1981 y 1982) y Velleman y Welsch (1981), prefieren r_i^* sobre r_i , ya que como se ve en el Teorema 12(1),

$$r_i^* \sim t(n - k - 1),$$

para lo cual las tablas de t están disponibles.

Existen dos maneras en que los residuales r_i y r_i^* , permiten detectar "outliers", es decir, dos procedimientos de prueba que ayudan a descubrir observaciones que se alejan de la nube de puntos en un diagrama de dispersión de un conjunto de datos, estos son la prueba formal e informal para detectar un "outlier".

Prueba Formal para Detectar un “Outlier”

En regresión lineal, existen algunos métodos, los cuales suponen que existe máximo un “outlier” en un conjunto de datos dado y generalmente son de dos formas,

1. Para ciertos estadísticos $T(\mathbf{x}_i, \mathbf{y}_i)$, se encuentra $C\alpha$ tal que,

$$P\{T(\mathbf{x}_i, \mathbf{y}_i) > C\alpha : \text{máximo un “outlier” está presente}\} \leq \alpha, \quad \text{y}$$

2. Se declara la i -ésima observación como un “outlier” si $T(\mathbf{x}_i, \mathbf{y}_i) > C\alpha$.

Tietjen *et al.* (1973), sugieren utilizar $T(\mathbf{x}_i, \mathbf{y}_i) = \max|r_i| \equiv r_{\max}$, para el caso de una regresión lineal simple; de manera que los valores críticos aproximados para r_{\max} , están dados por

$$r_{\max} = \sqrt{\frac{(n-k)F}{n-k-1+F}},$$

donde F , es el $100(1 - \frac{\alpha}{n})$ punto de $F(1, n-k-1)$. Lund (1975), proporciona valores críticos para r_{\max} para el caso general de regresión múltiple.

Prueba Informal para Detectar un “outlier” (Métodos Gráficos)

Los métodos gráficos proporcionan información valiosa acerca de la presencia de “outliers”, de la suficiencia del modelo y de la validación de sus supuestos asociados. Behnken y Draper (1972), sugieren utilizar r_i en estos gráficos, mientras que, Atkinson (1981) prefiere r_i^* , pero como r_i^* es una transformación monótona de r_i , las conclusiones hechas de gráficos basados en estos residuales, usualmente son las mismas.

Algunas de las gráficas de residuales más comúnmente usadas son:

1. Distribución de frecuencias de residuales. Las gráficas de una dimensión, tales como
 - (a) Histogramas,

- (b) Gráficas de puntos,
- (c) Gráficas de tallos y hojas, y
- (d) Gráficas de cajas,

son una vía fácil para observar si la muestra fue extraída de una población Normal.

2. Gráficas de residuales en series de tiempo. Estos gráficos proveen información para validar el supuesto de que los ϵ_i , $i = 1, 2, \dots, n$, son independientes, si no lo son, entonces el problema es conocido como "autocorrelación".
3. Gráficas de probabilidad normal. Este gráfico debe ser aproximadamente una línea recta, donde el intercepto y la pendiente, son estimaciones de μ y σ , respectivamente.
4. Gráficas de residuales contra valores ajustados (predichos). Este gráfico proporciona información acerca de la validación de los supuestos de linealidad y homogeneidad de varianzas.

Outliers, Puntos Palanca y Puntos Influyentes

Los outliers, puntos palanca y puntos influyentes son tres conceptos relacionados entre sí, a continuación se define cada uno de ellos.

"Outliers"

En regresión lineal, se define a un "outlier" como una observación para la cual, el residual estudentizado (ya sea r_i o r_i^*), es grande en magnitud comparado con el que presentan las demás observaciones de un conjunto de datos determinado. El "outlier" se considera como tal, cuando la ecuación de regresión no se ajusta con éxito.

Puntos Palanca

Los puntos palanca, son aquellas observaciones para las cuales el vector \mathbf{x}_i , se encuentra lejos del resto de las observaciones en un conjunto de datos dado. Equivalentemente, un punto palanca es una observación que presenta un p_{ii} grande en comparación con el que se observa en las otras observaciones del conjunto de datos. Las observaciones que se encuentran aisladas en el espacio \mathbf{X} son consideradas como puntos palanca, de manera que en el espacio \mathbf{X} los puntos palanca pueden ser también considerados como “outliers”, por lo tanto, puede decirse que las variables predictoras en \mathbf{X} son las que determinan los puntos palanca, no la variable respuesta en \mathbf{Y} .

Puntos Influyentes

Por su parte, los puntos influyentes, se definen como observaciones que, individual o colectivamente, influyen excesivamente en la ecuación de regresión ajustada, comparadas con las demás observaciones del conjunto de datos.

Sin embargo, una observación puede no tener la misma influencia sobre los resultados del análisis de regresión, por ejemplo, sobre el vector de parámetros $\hat{\beta}$, o sobre el vector de valores ajustados o predichos $\hat{\mathbf{Y}}$. De este modo, en un análisis de regresión, lo importante estriba en saber qué es lo primordial, si la estimación de β o la predicción en \mathbf{Y} , para posteriormente, dedicarse a medir la influencia de las observaciones sobre $\hat{\beta}$ o sobre $\hat{\mathbf{Y}}$, respectivamente.

Es interesante notar que:

1. Los “outliers” no necesariamente son puntos influyentes, y viceversa.
2. Las observaciones con residuales grandes no son deseables, pero se debe tener cuidado, pues un residual pequeño, no necesariamente significa que la observación correspondiente es adecuada, esto generalmente es debido a que en el

ajuste por Cuadrados Mínimos se evitan residuales grandes, lo cual puede hacer que un punto indeseable se acomode a los otros puntos del conjunto de datos. De hecho, existe una tendencia general en los puntos palanca que muestran residuales pequeños pero influyen en el ajuste, ocasionando una ecuación de regresión desproporcionada.

3. Los puntos palanca no necesariamente son puntos influyentes, y viceversa.

Recordando la desigualdad del Teorema 6(4), la cual implica que las observaciones con p_{ii} grande, tienden a tener residuales pequeños, por lo que se consideran puntos palanca. A continuación se mencionan tres cantidades relacionadas con la medición de este tipo de observaciones, que son los elementos de la diagonal de P , la distancia de Mahalanobis y los elementos de la diagonal de P_Z .

Elementos de la Diagonal de P

Los p_{ii} juegan un papel importante en la determinación de los valores ajustados \widehat{Y} , y la magnitud de los residuales, por lo que Hoaglin y Welsch (1978), sugieren la utilización de r_i^* para detectar "outliers" y p_{ii} para detectar puntos palanca que son potencialmente influyentes. Comúnmente se aplican los siguientes tres puntos de corte para p_{ii} :

1. Se puede considerar a $\frac{1}{p_{ii}}$ como el número de observaciones que determinan \widehat{y}_i . Por esto, Huber (1981), señala que puntos con $p_{ii} > 0.2$ se clasifican como puntos palanca. En este sentido, se requiere poner atención especial a aquellas observaciones cuyos valores predichos \widehat{y}_i , estén determinados por cinco o menos observaciones.
2. Por el Teorema 3(1), se puede deducir que el promedio de los elementos de la diagonal de P es $\frac{k}{n}$, esto condujo a Hoaglin y Welsch (1978) a indicar que puntos con $p_{ii} > \frac{2k}{n}$, se consideran puntos palanca.

3. Basándose en las condiciones de la distribución de p_{ii} , descritas en el segundo capítulo, se llega entonces a considerar puntos con

$$p_{ii} \leq \frac{nF(k-1) + (n-k)}{nF(k-1) + n(n-k)}$$

como puntos palanca, donde F es el $100(1-\alpha)$ punto de $F(k-1, n-k)$.

Distancia de Mahalanobis

Se supone que X contiene una columna de unos y que \tilde{X} denota la X centrada excluyendo la columna constante, por lo tanto, un estadístico que mide que tan lejos está \mathbf{x}_i del centro del conjunto de datos, comúnmente se calcula como

$$(n-1)^{-1} \tilde{\mathbf{x}}_i^T (\tilde{X}^T \tilde{X})^{-1} \tilde{\mathbf{x}}_i,$$

donde $\tilde{\mathbf{x}}_i$ es la i -ésima fila de \tilde{X} . Sin embargo, lo que interesa medir es que tan lejos está \mathbf{x}_i del resto de las observaciones, por lo que debe excluirse \mathbf{x}_i cuando se calcula la media y la matriz de varianzas de X , entonces la distancia de Mahalanobis es definida como

$$M_i = (n-2)(\tilde{\mathbf{x}}_i - \bar{\tilde{X}}_{(i)})^T \left(\tilde{X}_{(i)}^T \left(I - (n-1)^{-1} \mathbf{1}\mathbf{1}^T \right) \tilde{X}_{(i)} \right)^{-1} (\tilde{\mathbf{x}}_i - \bar{\tilde{X}}_{(i)}),$$

donde $\bar{\tilde{X}}_{(i)}$ es el promedio de $\tilde{X}_{(i)}$. Usando la ecuación (2.4) y notando además, que

$$\bar{\tilde{X}}_{(i)} = (n-1)^{-1} \tilde{X}_{(i)}^T \mathbf{1} = -(n-1)^{-1} \tilde{\mathbf{x}}_i,$$

puesto que \tilde{X} es centrada, entonces M_i se reduce a

$$M_i = \frac{n(n-2)(p_{ii} - \frac{1}{n})}{(n-1)(1-p_{ii})}, \quad i = 1, 2, \dots, n,$$

donde puede verse que M_i es equivalente a p_{ii} .

Elementos de la Diagonal de P_Z

Una desventaja de utilizar solamente a p_{ii} como medida de diagnóstico para identificar puntos palanca, es que no se dispone de la información contenida en Y ;

en este caso, se utiliza la matriz aumentada $\mathbf{Z} = (\mathbf{X} : \mathbf{Y})$, donde se define al i -ésimo elemento de la diagonal de la matriz de predicción $\mathbf{P}_{\mathbf{Z}}$ como $p_{z_{ii}}$, y recordando el Teorema 6(4), $p_{z_{ii}}$ se puede escribir como

$$p_{z_{ii}} = z_i^T (\mathbf{Z}^T \mathbf{Z})^{-1} z_i = p_{ii} + \frac{\hat{\epsilon}_i^2}{\hat{\epsilon}^T \hat{\epsilon}}$$

En consecuencia, $p_{z_{ii}}$ será grande, siempre que p_{ii} sea grande o $\hat{\epsilon}_i^2$ sea grande, o ambas lo sean; es por esto, que $p_{z_{ii}}$ no debe usarse como medida, pues no distingue entre un punto palanca en el espacio \mathbf{X} y un "outlier" en el espacio \mathbf{Z} .

Medidas Basadas en la Curva de Influencia

Una clase importante de medidas de la influencia de la i -ésima observación sobre los resultados del análisis de regresión, están basadas en la curva de influencia, cuyo concepto fue introducido por Hampel (1968 y 1974), y el cual es dado a continuación.

Curva de Influencia

Esta definición es general, en el sentido de que puede aplicarse a diferentes modelos y métodos de estimación, pero sólo se menciona lo referente a modelos de regresión y estimación mínimo cuadrática. Para estimar algunos parámetros de interés, suponiendo que se tiene un estadístico \mathbf{T} , basado en una muestra aleatoria muy grande, z_1, z_2, \dots, z_n , proveniente de una función de distribución acumulada \mathbf{F} ; al adicionar una observación más z a esta muestra, se observa cómo se modifica al estadístico \mathbf{T} y las conclusiones basadas en \mathbf{T} . Por lo que, la curva de influencia es definida como

$$\psi(z, \mathbf{F}, \mathbf{T}) = \lim_{\epsilon \rightarrow 0} \frac{\mathbf{T}((1 - \epsilon)\mathbf{F} + \epsilon\delta_z) - \mathbf{T}(\mathbf{F})}{\epsilon}$$

donde se observa que el límite existe en z , donde $\delta_z = 1$.

Definiendo $\mathbf{F}_\epsilon = (1 - \epsilon)\mathbf{F} + \epsilon\delta_z$ y escribiendo nuevamente $\psi(z, \mathbf{F}, \mathbf{T})$, así,

$$\psi(z, \mathbf{F}, \mathbf{T}) = \lim_{\epsilon \rightarrow 0} \frac{\mathbf{T}(\mathbf{F}_\epsilon) - \mathbf{T}(\mathbf{F})}{\epsilon} = \frac{d}{d\epsilon} \mathbf{T}(\mathbf{F}_\epsilon) |_{\epsilon=0}$$

Entonces, la curva de influencia es la derivada de la función $\mathbf{T}(\mathbf{F}_\epsilon)$ con respecto a ϵ , evaluada en $\epsilon = 0$.

La curva de influencia tiene varias aplicaciones, una de las cuales, se usa para estudiar las propiedades asintóticas de un estimador, para lo cual, considerando a \mathbf{F}_n como la función de distribución empírica basada en n observaciones, y suponiendo que \mathbf{F}_n converge a \mathbf{F} y que $\mathbf{T}(\mathbf{F}_n)$ converge a $\mathbf{T}(\mathbf{F})$, entonces bajo condiciones de regularidad apropiadas, se tiene que

$$\sqrt{n}(\mathbf{T}(\mathbf{F}_n) - \mathbf{T}(\mathbf{F})) \cong \sqrt{n} \int \psi(z, \mathbf{F}, \mathbf{T}) d\mathbf{F}_n(z) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi(z, \mathbf{F}, \mathbf{T})$$

Por el Teorema Central del Límite, el lado derecho de esta expresión es asintóticamente Normal con media cero y varianza $\Sigma_\psi = \int \psi(z, \mathbf{F}, \mathbf{T}) \psi^T(z, \mathbf{F}, \mathbf{T}) d\mathbf{F}(z)$.

Otra aplicación de la curva de influencia, se utiliza para comparar estimadores y sugerir estimadores robustos como otra alternativa de estimación. A continuación se da otro uso de las curvas de influencia, el cual trata sobre la influencia de las observaciones sobre los estimadores de β y σ^2 .

Curvas de Influencia para $\hat{\beta}$ y $\hat{\sigma}^2$

Si la curva de influencia para \mathbf{T} no es acotada, entonces \mathbf{T} se dice ser un estimador no robusto, puesto que es sensible a observaciones extremas. Los estimadores de interés para β y σ^2 , en este caso, son los que resultan de la estimación mínimo cuadrática, los cuales son obtenidos al solucionar el sistema de las $(k + 1)$ ecuaciones lineales

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i^T \beta) = 0 \quad (3.8)$$

$$\frac{1}{n - k} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 = \sigma^2 \quad (3.9)$$

Suponiendo que el vector $(k + 1)$ -dimensional, dado por (\mathbf{x}^T, y) tiene una función de distribución acumulada conjunta \mathbf{F} , entonces las ecuaciones (3.8) y (3.9), pueden

escribirse como

$$\int_{\mathbf{x}^T, \mathbf{y}} \mathbf{x}(\mathbf{y} - \mathbf{x}^T \boldsymbol{\beta}) dF_n(\mathbf{x}^T, \mathbf{y}) = 0 \quad (3.10)$$

$$\int_{\mathbf{x}^T, \mathbf{y}} (\mathbf{y} - \mathbf{x}^T \boldsymbol{\beta}) dF_n(\mathbf{x}^T, \mathbf{y}) = \sigma^2 \quad (3.11)$$

La solución de (3.10) y (3.11) produce funciones para $\boldsymbol{\beta}$ y σ^2 .

Se supone ahora que,

$$\mathbf{E}_F \left\{ \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} x^T & y \end{bmatrix} \right\} = \begin{bmatrix} \Sigma_{xx}(\mathbf{F}) & \Sigma_{xy}(\mathbf{F}) \\ \Sigma_{xy}^T(\mathbf{F}) & \sigma_{yy}(\mathbf{F}) \end{bmatrix},$$

entonces resulta que

$$\boldsymbol{\beta}(\mathbf{F}) = \Sigma_{xx}^{-1}(\mathbf{F}) \Sigma_{xy}(\mathbf{F})$$

y que

$$\sigma^2(\mathbf{F}) = \sigma_{yy}(\mathbf{F}) - \Sigma_{xy}^T(\mathbf{F}) \Sigma_{xx}^{-1}(\mathbf{F}) \Sigma_{xy}(\mathbf{F})$$

Así, la curva de influencia puede ser obtenida sustituyendo z por el vector $(\mathbf{x}^T, \mathbf{y})$ y reemplazando \mathbf{T} por $\boldsymbol{\beta}(\mathbf{F})$ o $\sigma^2(\mathbf{F})$, entonces se tiene que

$$\psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \mathbf{T}) = \lim_{\epsilon \rightarrow 0} \frac{\mathbf{T} \left((1 - \epsilon) \mathbf{F} + \epsilon \delta_{\mathbf{x}^T, \mathbf{y}} \right) - \mathbf{T}(\mathbf{F})}{\epsilon} \quad (3.12)$$

Se requiere del siguiente lema, para probar el Teorema 13, pues éste último, proporciona formas explícitas de curvas de influencia para $\hat{\boldsymbol{\beta}}$ y $\hat{\sigma}^2$, utilizando la ecuación (3.12).

Lema 5. *Sea \mathbf{A} cualquier matriz, tal que $(\mathbf{I} + \epsilon \mathbf{A})^{-1}$ exista, entonces*

$$(\mathbf{I} + \epsilon \mathbf{A})^{-1} = \mathbf{I} + \sum_{i=1}^{\infty} (-1)^i \epsilon^i \mathbf{A}^i$$

Demostración. Usando repetidamente $(\mathbf{I} + \epsilon \mathbf{A})^{-1} = \mathbf{I} - \epsilon \mathbf{A}(\mathbf{I} + \epsilon \mathbf{A})^{-1}$, se tiene que

$$(\mathbf{I} + \epsilon \mathbf{A})^{-1} = \mathbf{I} - \epsilon \mathbf{A}(\mathbf{I} + \epsilon \mathbf{A})^{-1}$$

$$\begin{aligned}
&= \mathbf{I} - \epsilon \mathbf{A} + \epsilon^2 \mathbf{A}^2 (\mathbf{I} + \epsilon \mathbf{A})^{-1} \\
&= \mathbf{I} - \epsilon \mathbf{A} + \epsilon^2 \mathbf{A}^2 - \epsilon^3 \mathbf{A}^3 (\mathbf{I} + \epsilon \mathbf{A})^{-1} \\
&= \dots \\
&= \dots \\
&= \dots \\
&= \mathbf{I} + \sum_{i=1}^{\infty} (-1)^i \epsilon^i \mathbf{A}^i \blacksquare
\end{aligned}$$

Teorema 13.

1. La curva de influencia para $\hat{\beta}$ es

$$\psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \hat{\beta}(\mathbf{F})) = \Sigma_{xx}^{-1}(\mathbf{F}) \mathbf{x} (\mathbf{y} - \mathbf{x}^T \hat{\beta}(\mathbf{F}))$$

2. La curva de influencia para $\hat{\sigma}^2$ es

$$\psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \hat{\sigma}^2(\mathbf{F})) = (\mathbf{y} - \mathbf{x}^T \hat{\beta}(\mathbf{F}))^2 - \sigma_{yy}(\mathbf{F}) + \Sigma_{xy}^T(\mathbf{F}) \hat{\beta}(\mathbf{F})$$

Demostración.

1. La sustitución de $\beta(\mathbf{F})$ en la ecuación (3.12), produce la curva de influencia para $\hat{\beta}$, es decir,

$$\psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \beta(\mathbf{F})) = \lim_{\epsilon \rightarrow 0} \frac{\hat{\beta} \left((1 - \epsilon) \mathbf{F} + \epsilon \delta_{\mathbf{x}^T \mathbf{y}} \right) - \hat{\beta}(\mathbf{F})}{\epsilon}$$

y luego al utilizar el Lema 5, se prueba el Teorema 13(1).

2. La sustitución de $\sigma^2(\mathbf{F})$ en la ecuación (3.12), produce la curva de influencia para $\hat{\sigma}^2$, entonces se tiene

$$\begin{aligned}
\psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \sigma^2(\mathbf{F})) &= \lim_{\epsilon \rightarrow 0} \frac{\hat{\sigma}^2 \left((1 - \epsilon) \mathbf{F} + \epsilon \delta_{\mathbf{x}^T \mathbf{y}} \right) - \hat{\sigma}^2(\mathbf{F})}{\epsilon} \\
&= \frac{d}{d\epsilon} \text{plim } \hat{\sigma}^2(\mathbf{F}\epsilon) |_{\epsilon=0},
\end{aligned}$$

resolviendo la derivada y evaluando en $\epsilon = 0$, se prueba el Teorema 13(2). \blacksquare

Notando que las curvas de influencia para $\hat{\beta}$ y $\hat{\sigma}^2$, no son acotadas, pues $(\mathbf{y} - \mathbf{x}^T \hat{\beta}(F))$ tampoco lo es, por lo que no existen estimadores robustos para β y σ^2 .

Aproximación de la Curva de Influencia

Las curvas de influencia dadas en el Teorema 13, miden la influencia sobre $\hat{\beta}$ y $\hat{\sigma}^2$, al adicionarse una observación $(\mathbf{x}^T, \mathbf{y})$ a una muestra muy grande. En la práctica no siempre se tiene la posibilidad de contar con muestras muy grandes, por lo que es necesario aproximar la curva de influencia para una muestra finita. A continuación se mencionan cuatro aproximaciones de la curva de influencia para $\hat{\beta}$.

La Curva de Influencia Empírica basada en n Observaciones. Se obtiene a partir de la curva de influencia para $\hat{\beta}$, aproximando F por la función de distribución acumulada empírica \hat{F} , y sustituyendo $(\mathbf{x}_i^T, \mathbf{y}_i)$, $n^{-1}(\mathbf{X}^T \mathbf{X})$ y $\hat{\beta}$ por $(\mathbf{x}^T, \mathbf{y})$, $\Sigma_{xx}(F)$ y $\hat{\beta}(F)$, respectivamente, obteniéndose

$$EIC_i = n(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i (\mathbf{y}_i - \mathbf{x}_i^T \hat{\beta}) = n(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \hat{\epsilon}_i, \quad i = 1, 2, \dots, n$$

La Curva de Influencia Muestral. Se obtiene de la ecuación (3.12), considerando las igualdades $(\mathbf{x}^T, \mathbf{y}) = (\mathbf{x}_i^T, \mathbf{y}_i)$, $F = \hat{F}$, $\epsilon = -(n-1)^{-1}$ y omitiendo el límite, resulta

$$\begin{aligned} SIC_i &= -(n-1) \left(\mathbf{T} \left(\frac{n}{n-1} \hat{F} + \frac{-1}{n-1} \delta_{\mathbf{x}_i^T, \mathbf{y}_i} \right) - \mathbf{T}(\hat{F}) \right) \\ &= (n-1) \left(\mathbf{T}(\hat{F}) - \mathbf{T}(\hat{F}_{(i)}) \right), \end{aligned}$$

donde $\hat{F}_{(i)}$ es la función de distribución acumulada empírica cuando la i -ésima observación es omitida. Tomando ahora las igualdades $\mathbf{T}(\hat{F}) = \hat{\beta}$ y $\mathbf{T}(\hat{F}_{(i)}) = \hat{\beta}_{(i)}$, se tiene que SIC_i se puede escribir simplemente como

$$SIC_i = (n-1) (\hat{\beta} - \hat{\beta}_{(i)}), \quad (3.13)$$

donde $\hat{\beta}_{(i)} = (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)}$ es el estimador de β cuando la i -ésima observación es omitida. Al usar la ecuación (2.4), se obtiene

$$\hat{\beta}_{(i)} = \left((\mathbf{X}^T \mathbf{X})^{-1} + \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1}}{1 - p_{ii}} \right) \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} \quad (3.14)$$

También, notando que

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \begin{bmatrix} \mathbf{X}_{(i)}^T & \mathbf{x}_i \end{bmatrix} \begin{bmatrix} \mathbf{Y}_{(i)} \\ \mathbf{y}_i \end{bmatrix} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{y}_i,\end{aligned}$$

despejando, se obtiene

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} = \hat{\beta} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \mathbf{y}_i, \quad (3.15)$$

sustituyendo la ecuación (3.15) en la ecuación (3.14), y simplificando, se obtiene

$$(\hat{\beta} - \hat{\beta}_{(i)}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{\epsilon}_i}{1 - p_{ii}}, \quad (3.16)$$

si se usa esta identidad, entonces el SIC_i dado en la ecuación (3.13), puede ser escrito como

$$SIC_i = (n - 1)(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{\epsilon}_i}{1 - p_{ii}}, \quad i = 1, 2, \dots, n \quad (3.17)$$

La Curva de Sensibilidad. Considerando $\mathbf{F} = \widehat{\mathbf{F}}_{(i)}$ y $\epsilon = n^{-1}$ y omitiendo el límite en la ecuación (3.12), resulta

$$SC_i = n \left(\mathbf{T} \left(\frac{n-1}{n} \widehat{\mathbf{F}}_{(i)} + \frac{1}{n} \delta_{\mathbf{x}_i^T, \mathbf{y}_i} \right) - \mathbf{T}(\widehat{\mathbf{F}}_{(i)}) \right) = n \left(\mathbf{T}(\widehat{\mathbf{F}}) - \mathbf{T}(\widehat{\mathbf{F}}_{(i)}) \right),$$

si nuevamente se consideran las identidades $\mathbf{T}(\widehat{\mathbf{F}}) = \hat{\beta}$ y $\mathbf{T}(\widehat{\mathbf{F}}_{(i)}) = \hat{\beta}_{(i)}$, se llega a

$$SC_i = n(\hat{\beta} - \hat{\beta}_{(i)}) = n(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{\epsilon}_i}{1 - p_{ii}}, \quad i = 1, 2, \dots, n \quad (3.18)$$

La Curva de Influencia Empírica basada en $(n - 1)$ Observaciones. De la misma manera que las demás, ésta también se obtiene a partir de la curva de influencia para $\hat{\beta}$. Sólo que ahora tomando las siguientes igualdades, $\mathbf{F} = \widehat{\mathbf{F}}_{(i)}$, $(\mathbf{x}^T, \mathbf{y}) = (\mathbf{x}_i^T, \mathbf{y}_i)$, $\Sigma_{xx}(\mathbf{F}) = (n - 1)^{-1}(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})$ y $\hat{\beta}(\mathbf{F}) = \hat{\beta}_{(i)}$, y entonces se produce

$$EIC_{(i)} = (n - 1)(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i (\mathbf{y}_i - \mathbf{x}_i^T \hat{\beta}_{(i)}), \quad i = 1, 2, \dots, n, \quad (3.19)$$

a la cantidad $(\mathbf{y}_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{(i)})$, se le llama *residual predicho*, la cual, utilizando las ecuaciones (2.4) y (3.15), puede escribirse como

$$(\mathbf{y}_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{(i)}) = \frac{\hat{\epsilon}_i}{1 - \mathbf{p}_{ii}}$$

Sustituyendo esta última expresión en la ecuación (3.19) y utilizando nuevamente la ecuación (2.4), se obtiene

$$EIC_{(i)} = (n - 1)(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{\epsilon}_i}{(1 - \mathbf{p}_{ii})^2}, \quad i = 1, 2, \dots, n \quad (3.20)$$

Notando que EIC_i es el menos sensible a puntos palanca, mientras que $EIC_{(i)}$ es el más sensible. También notando que SIC_i y SC_i pueden considerarse equivalentes, pues son proporcionales a la distancia $(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)})$.

La curva de influencia para $\hat{\sigma}^2$ dada en el Teorema 13(2), es un escalar. Sin embargo, la curva de influencia para $\hat{\boldsymbol{\beta}}$ dada en el Teorema 13(1), así como sus aproximaciones, son vectores; por lo tanto, es necesario reducir esta curva de influencia a una cantidad escalar y que las observaciones puedan ser ordenadas de acuerdo a su influencia sobre $\hat{\boldsymbol{\beta}}$ o una función lineal de $\hat{\boldsymbol{\beta}}$. Para realizar tal reducción se puede utilizar,

$$\sup_{\mathbf{x}^T, \mathbf{y}} \|\psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \hat{\boldsymbol{\beta}}(\mathbf{F}))\|,$$

esta expresión es considerada por Hampel (1974), como una medida de la influencia máxima posible de cualquier observación sobre los coeficientes estimados, y la llama el *error de sensibilidad*, pues tiene la desventaja de no ser invariante bajo transformaciones de parámetros de localización y escala.

Una alternativa al error de sensibilidad, es del tipo

$$\sup_{\mathbf{x}^T, \mathbf{y}} \sup_a \frac{|\mathbf{a}^T \psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \hat{\boldsymbol{\beta}}(\mathbf{F}))|}{\sqrt{\mathbf{a}^T \boldsymbol{\Sigma}_\psi \mathbf{a}}},$$

la cual, es invariante bajo transformaciones no singulares de \mathbf{X} , y mide la influencia máxima de cualquier observación sobre cualquier combinación lineal de los

coeficientes relativos al error estándar de la combinación lineal $\sqrt{a^T \Sigma_\psi a}$, donde $\Sigma_\psi = \int \psi(z, \mathbf{F}, \mathbf{T}) \psi^T(z, \mathbf{F}, \mathbf{T}) d\mathbf{F}(z)$. Se ve que si a se restringe a \mathbf{x} , entonces se convierte en una medida de la influencia de cualquier observación sobre los valores predichos.

Otra cantidad equivalente a la anterior es

$$\sup_{\mathbf{x}, \mathbf{T}, \mathbf{y}} \frac{\psi^T(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \hat{\boldsymbol{\beta}}(\mathbf{F})) M \psi(\mathbf{x}^T, \mathbf{y}, \mathbf{F}, \hat{\boldsymbol{\beta}}(\mathbf{F}))}{c}$$

En la práctica, lo que interesa es encontrar la influencia máxima posible de cualquier observación sobre los resultados de la regresión, ordenando las n observaciones de acuerdo a su nivel de influencia, para lo cual es necesario hacer una apropiada elección de M y c , por lo tanto, se puede usar

$$D_i(M, c) = \frac{\psi^T(\mathbf{x}_i^T, \mathbf{y}_i, \mathbf{F}, \hat{\boldsymbol{\beta}}(\mathbf{F})) M \psi(\mathbf{x}_i^T, \mathbf{y}_i, \mathbf{F}, \hat{\boldsymbol{\beta}}(\mathbf{F}))}{c}, \quad (3.21)$$

para considerar la influencia de la i -ésima observación sobre los coeficientes de regresión relativos a M y c . De manera que si se registra un valor grande en $D_i(M, c)$, este indica que la i -ésima observación tiene una fuerte influencia sobre $\hat{\boldsymbol{\beta}}$ relativa a M y c .

Se presentan cuatro alternativas de M y c , en $D_i(M, c)$, las cuales son, las distancias de Cook, de Welsch-Kuh, de Welsch y una modificación a la primera.

Distancia de Cook

Para definir esta distancia, es necesario considerar que, bajo normalidad, la región de confianza $100(1 - \alpha)\%$ para $\boldsymbol{\beta}$, se puede obtener a partir de

$$\frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T (\mathbf{X}^T \mathbf{X}) (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}{k \hat{\sigma}^2} \leq F(\alpha; k, n - k),$$

donde $F(\alpha; k, n - k)$ es el punto superior α de la distribución F centrada con $(k, n - k)$ grados de libertad. Esta desigualdad, define una región elipsoidal centrada

en $\widehat{\beta}$, y permite medir la influencia de la i -ésima observación por el cambio en el centro del elipsoide de confianza, cuando la i -ésima observación es omitida.

En forma análoga, Cook (1977), propone la medida

$$C_i = \frac{(\beta - \widehat{\beta}_{(i)})^T (\mathbf{X}^T \mathbf{X}) (\beta - \widehat{\beta}_{(i)})}{k \widehat{\sigma}^2}, \quad i = 1, 2, \dots, n, \quad (3.22)$$

para considerar la influencia de la i -ésima observación sobre el centro del elipsoide de confianza, es decir, sobre $\widehat{\beta}$. Esta medida se conoce como la *distancia de Cook*, y es simplemente, una distancia a escala entre $\widehat{\beta}$ y $\widehat{\beta}_{(i)}$.

Bingham (1977), señala que C_i , también puede ser escrita como

$$C_i = \frac{(\widehat{\mathbf{Y}} - \widehat{\mathbf{Y}}_{(i)})^T (\widehat{\mathbf{Y}} - \widehat{\mathbf{Y}}_{(i)})}{k \widehat{\sigma}^2}, \quad i = 1, 2, \dots, n$$

donde $\widehat{\mathbf{Y}}_{(i)} = \mathbf{X} \widehat{\beta}_{(i)}$ es el vector de valores predichos, cuando se realiza una regresión de $\widehat{\mathbf{Y}}_{(i)}$ con $\widehat{\mathbf{X}}_{(i)}$. Por lo anterior C_i , se puede interpretar como la distancia euclidiana a escala, entre los vectores de valores predichos, incluyendo y excluyendo la i -ésima observación, cuando se hace el ajuste.

A simple vista, pareciera que para encontrar C_i , $i = 1, 2, \dots, n$, se requiere realizar $(n + 1)$ análisis de regresión, uno utilizando todos los datos, y n usando los datos sin cada observación. Sin embargo, si se sustituye la ecuación (3.16) en la ecuación (3.22), se produce

$$C_i = \frac{x_i^T (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} x_i}{k(1 - p_{ii})} \frac{\widehat{\epsilon}_i^2}{\widehat{\sigma}^2(1 - p_{ii})} = \frac{1}{k} \frac{p_{ii}}{(1 - p_{ii})} r_i^2 \quad (3.23)$$

En la ecuación (3.23), se observa que se combinan dos medidas, p_{ii} y r_i , los cuales, como se sabe, proporcionan información sobre puntos palanca y "outliers", respectivamente.

Observando además que $\frac{p_{ii}}{(1 - p_{ii})}$, es la relación dada entre la varianza del i -ésimo valor predicho $\text{Var}(\widehat{y}_i) = \sigma^2 p_{ii}$, y la varianza del i -ésimo residual

$$\text{Var}(\widehat{\epsilon}_i) = \sigma^2(1 - p_{ii}).$$

Cook y Weisberg (1982), se refieren a esta relación, como el potencial de la i -ésima observación en la determinación de $\hat{\beta}$, con $M = \mathbf{X}^T \mathbf{X}$ y $c = (n - 1)^2 k \hat{\sigma}^2$. Otra interpretación de esta relación, es dada por Huber (1981), para la cual se ve que si se utilizan la ecuación (3.16) y el residual predicho antes descrito, $\widehat{\mathbf{y}}_i$ se puede escribir como

$$\begin{aligned} \widehat{\mathbf{y}}_i &= \mathbf{x}_i^T \hat{\beta} \\ &= \mathbf{x}_i^T \left(\frac{\hat{\epsilon}_i^2}{(1 - p_{ii})} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i + \hat{\beta}_{(i)} \right) \\ &= p_{ii} \frac{\hat{\epsilon}_i}{(1 - p_{ii})} + \mathbf{x}_i^T \hat{\beta}_{(i)} \\ &= p_{ii} (\mathbf{y}_i - \mathbf{x}_i^T \hat{\beta}_{(i)}) + \mathbf{x}_i^T \hat{\beta}_{(i)} \\ &= (1 - p_{ii}) \mathbf{x}_i^T \hat{\beta}_{(i)} + p_{ii} \mathbf{y}_i \end{aligned}$$

De esta forma, $\frac{p_{ii}}{(1 - p_{ii})}$, puede también ser considerada como la relación que existe entre la parte de $\widehat{\mathbf{y}}_i$ debida a \mathbf{y}_i y la parte debida al valor predicho $\mathbf{x}_i^T \hat{\beta}_{(i)}$.

Es claro ver en la ecuación (3.23), que \mathbf{C}_i será grande si p_{ii} es grande y r_i^2 también lo es.

Se ve que \mathbf{C}_i es una función monótona de r_i^2 , de manera que si se cumplen las condiciones del Teorema 12(1) y utilizando la expresión

$$r_i^* = r_i \sqrt{\frac{n - k - 1}{n - k - r_i^2}},$$

entonces, resulta que

$$\frac{r_i^2(n - k - 1)}{n - k - r_i^2} \sim F(1, n - k - 1),$$

lo cual indica que \mathbf{C}_i no tiene estrictamente una distribución F .

La distancia de Cook, puede ser obtenida directamente de la curva de influencia para $\hat{\beta}$ dada en el Teorema 13(1), si ésta se aproxima por la curva de influencia muestral \mathbf{SIC}_i dada en la ecuación (3.13) y utilizando para ello la ecuación

(3.21), entonces C_i puede expresarse como

$$C_i = D_i \left(X^T X, k\hat{\sigma}^2(n-1)^{-2} \right)$$

Distancia de Welsch-Kuh

La influencia de la i -ésima observación sobre los valores predichos \hat{y}_i , puede ser medida por el cambio en la predicción en x_i , cuando la i -ésima observación es omitida, relativa al error estándar de \hat{y}_i , esto es,

$$\frac{|\hat{y} - \hat{y}_{i(i)}|}{\sigma\sqrt{p_{ii}}} = \frac{|x_i^T(\hat{\beta} - \hat{\beta}_{(i)})|}{\sigma\sqrt{p_{ii}}} \quad (3.24)$$

Welsch y Kuh (1977), Welsch y Peters (1978) y Belsley *et al.* (1980), prefieren utilizar $\hat{\sigma}_{(i)}$ como un estimador de σ en la ecuación (3.24). Si se usan las ecuaciones (3.16) y (3.5), entonces la ecuación (3.24) puede ser escrita como

$$WK_i = \frac{|\frac{\hat{\epsilon}_i}{(1-p_{ii})} x_i^T (X^T X)^{-1} x_i|}{\hat{\sigma}_{(i)}\sqrt{p_{ii}}} = |r_i^*| \sqrt{\frac{p_{ii}}{1-p_{ii}}} \quad (3.25)$$

Valores grandes de WK_i , indican que la i -ésima observación influye sobre el ajuste del modelo.

Del mismo modo que la distancia de Cook, puede mencionarse que aunque la distancia de Welsch-Kuh sea un estadístico parecido a t , no tiene estrictamente una distribución t . Pero si se dan las condiciones del Teorema 12(1), entonces $r_i^* \sim t(n-k-1)$, lo cual hace que pueda sugerirse como un punto de corte a

$$t\sqrt{\frac{k}{n-k}}$$

Belsley *et al.* (1980), recomiendan usar $2\sqrt{\frac{k}{n}}$, como punto de corte para WK_i , pero parece ser más apropiado utilizar $2\sqrt{\frac{k}{n-k}}$.

La distancia de Welsch-Kuh, también puede ser obtenida directamente de la curva de influencia para $\hat{\beta}$ dada en el Teorema 13(1), si ésta se aproxima por la

curva de influencia muestral SIC_i dada en la ecuación (3.17) y utilizando para ello la ecuación (3.21), entonces WK_i puede expresarse como

$$WK_i = \sqrt{D_i(X^T X, (n-1)\hat{\sigma}_{(i)}^2)}$$

Si en vez de aproximar con la curva de influencia muestral, se aproxima con la curva de sensibilidad, dada en la ecuación (3.18), entonces WK_i se escribe así,

$$WK_i = \sqrt{D_i(X^T X, n\hat{\sigma}_{(i)}^2)}$$

Distancia de Welsch

Si se utiliza la curva de influencia empírica basada en $(n-1)$ observaciones, la cual está dada en la ecuación (3.20), como una aproximación a la curva de influencia para $\hat{\beta}$, y considerando además que, $M = X_{(i)}^T X_{(i)} = (X^T X - x_i x_i^T)$ y $c = (n-1)\hat{\sigma}_{(i)}^2$, entonces la ecuación (3.21) se convierte en

$$\begin{aligned} W_i^2 &= D_i(X_{(i)}^T X_{(i)}, (n-1)\hat{\sigma}_{(i)}^2) \\ &= (n-1) \frac{\hat{\epsilon}_i^2}{\hat{\sigma}_{(i)}^2 (1-p_{ii})^4} x_i^T (X^T X)^{-1} (X^T X - x_i x_i^T) (X^T X)^{-1} x_i \\ &= (n-1) r_i^{*2} \frac{p_{ii}}{(1-p_{ii})^2} \end{aligned} \quad (3.26)$$

Al comparar las ecuaciones (3.25) y (3.26), resulta

$$W_i = WK_i \sqrt{\frac{n-1}{1-p_{ii}}} \quad (3.27)$$

Se observa claramente en la ecuación (3.27), que W_i da un mayor énfasis a los puntos palanca que WK_i , sin embargo, algunos autores prefieren ésta última, por considerarla más fácil de interpretar. Asimismo, en la ecuación (3.27) se puede ver que un punto de corte para W_i , puede obtenerse al multiplicar el punto de corte para WK_i por $\left(\frac{n(n-1)}{(n-k)}\right)^{\frac{1}{2}}$; entonces, si el punto de corte para WK_i es $2\sqrt{\frac{k}{n-k}}$, el correspondiente punto de corte para W_i es $\frac{2}{(n-k)}\sqrt{kn(n-1)}$, sin embargo, si n es grande comparada con k , esta cantidad se convierte en $3\sqrt{k}$.

Distancia de Cook Modificada

Para la detección de observaciones influyentes, Atkinson (1981), propone usar una versión modificada de la distancia de Cook. Esta modificación implica reemplazar $\hat{\sigma}^2$ por $\hat{\sigma}_{(i)}^2$, tomar la raíz cuadrada de C_i y ajustarlo por tamaño de muestra, lo cual produce

$$\begin{aligned}
 C_i^* &= \sqrt{D_i \left(X^T X, \frac{k(n-1)^2}{(n-k)} \hat{\sigma}_{(i)}^2 \right)} \\
 &= |r_i^*| \sqrt{\frac{p_{ii}}{(1-p_{ii})} \frac{(n-k)}{k}} \\
 &= WK_i \sqrt{\frac{n-k}{k}} \tag{3.28}
 \end{aligned}$$

Atkinson (1981), afirma que esta modificación, mejora a C_i , ya que C_i^* proporciona mayor énfasis a puntos extremos.

Si se utilizan diferentes aproximaciones a la curva de influencia para $\hat{\beta}$, en combinación con varias elecciones de M y c de la ecuación (3.21), se pueden obtener algunas mediciones de la influencia de la i -ésima observación sobre los coeficientes de regresión; entonces, la diferencia fundamental en C_i , WK_i , W_i y C_i^* , estriba en la elección de M y c . Así, C_i usa $X^T X$ y $\hat{\sigma}^2$; WK_i y C_i^* usan $X^T X$ y $\hat{\sigma}_{(i)}^2$, y W_i usa $(X_{(i)}^T X_{(i)})$ y $\hat{\sigma}_{(i)}^2$, entonces se puede concluir que C_i mide la influencia de la i -ésima observación solamente sobre $\hat{\beta}$, mientras que WK_i , W_i y C_i^* miden la influencia sobre $\hat{\beta}$ y sobre $\hat{\sigma}^2$.

Medidas Basadas en el Volumen de Elipsoides de Confianza

Hasta el momento, se ha visto que las medidas de diagnóstico basadas en la curva de influencia, se pueden interpretar como medidas que se basan en el cambio en el centro del elipsoide de confianza, cuando la i -ésima observación es omitida; una alternativa, la constituyen, medidas de la influencia de la i -ésima observación, que

se basan en el cambio en el volumen del elipsoide de confianza, cuando la i -ésima observación es omitida. A continuación se describen tres de estas medidas.

Estadístico Andrews-Pregibon

El volumen de la región de confianza para β , como se ha visto, se obtiene de

$$\frac{(\beta - \hat{\beta})^T (\mathbf{X}^T \mathbf{X}) (\beta - \hat{\beta})}{k \hat{\sigma}^2} \leq F(\alpha; k, n - k),$$

la cual, es inversamente proporcional a la raíz cuadrada del $\det(\mathbf{X}^T \mathbf{X})$. Así, la influencia de la i -ésima observación sobre el volumen de los elipsoides de confianza, puede ser medida comparando el $\det(\mathbf{X}^T \mathbf{X})$ con el $\det(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})$. Por otra parte, al omitir una observación con un residual grande, ocasionará una fuerte reducción en la suma de cuadrados de residuales SSE . De tal forma que, la influencia de la i -ésima observación puede ser medida combinando éstas dos ideas, y calculando el cambio en $\hat{\epsilon}^T \hat{\epsilon}$ y en el $\det(\mathbf{X}^T \mathbf{X})$ cuando la i -ésima observación es omitida. Por lo que Andrews y Pregibon (1978), sugieren la relación

$$\frac{SSE_{(i)} \det(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})}{SSE \det(\mathbf{X}^T \mathbf{X})}, \quad i = 1, 2, \dots, n. \quad (3.29)$$

El cociente dado en la ecuación (3.29), se puede simplificar de la siguiente forma, definiendo la matriz aumentada $\mathbf{Z} = (\mathbf{X} : \mathbf{Y})$, y se ve que

$$\mathbf{Z}^T \mathbf{Z} = \begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Y} \\ \mathbf{Y}^T \mathbf{X} & \mathbf{Y}^T \mathbf{Y} \end{bmatrix}$$

al utilizar el Lema 3(1), se tiene que

$$\begin{aligned} \det(\mathbf{Z}^T \mathbf{Z}) &= \det(\mathbf{X}^T \mathbf{X}) \det(\mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) \\ &= \det(\mathbf{X}^T \mathbf{X}) \det(\mathbf{Y}^T (\mathbf{I} - \mathbf{P}) \mathbf{Y}) \\ &= \det(\mathbf{X}^T \mathbf{X}) SSE \end{aligned}$$

Equivalentemente, se tiene

$$\det(\mathbf{Z}_{(i)}^T \mathbf{Z}_{(i)}) = \det(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)}) SSE_{(i)}$$

Entonces, la ecuación dada en (3.29), se convierte en

$$\frac{SSE_{(i)} \det(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})}{SSE \det(\mathbf{X}^T \mathbf{X})} = \frac{\det(\mathbf{Z}_{(i)}^T \mathbf{Z}_{(i)})}{\det(\mathbf{Z}^T \mathbf{Z})},$$

la cual mide el cambio relativo en el $\det(\mathbf{Z}^T \mathbf{Z})$, debido a la eliminación de la i -ésima observación, es decir, la proporción del volumen generado por \mathbf{Z} que no es debida a la i -ésima observación. Si se omite una observación que está lejos del centro de los datos, resultará en una reducción en el determinante y un incremento en el volumen; por lo que, valores pequeños en la ecuación (3.29) requieren una atención especial. Por conveniencia, se ha definido a

$$AP_i = 1 - \frac{\det(\mathbf{Z}_{(i)}^T \mathbf{Z}_{(i)})}{\det(\mathbf{Z}^T \mathbf{Z})} \quad (3.30)$$

por lo que ahora, los que requieren de atención especial, son los valores grandes de la ecuación (3.30). Por otro lado, si en el Lema 4, se considera que $\mathbf{A} = \mathbf{Z}^T \mathbf{Z}$ y $\mathbf{B} = \mathbf{C} = \mathbf{z}_i$, resulta que

$$\begin{aligned} \det(\mathbf{Z}_{(i)}^T \mathbf{Z}_{(i)}) &= \det(\mathbf{Z}^T \mathbf{Z} - \mathbf{z}_i \mathbf{z}_i^T) \\ &= \det(\mathbf{Z}^T \mathbf{Z}) \left(1 - \mathbf{z}_i^T (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{z}_i\right) \\ &= \det(\mathbf{Z}^T \mathbf{Z}) (1 - p_{zii}) \end{aligned}$$

Si este resultado se sustituye en la ecuación (3.30), resulta

$$AP_i = p_{zii},$$

puede concluirse entonces que, AP_i no distingue entre puntos palanca en el espacio \mathbf{X} y "outliers" en el espacio \mathbf{Z} . Asimismo, se ve que AP_i es el i -ésimo elemento de la diagonal de la matriz de predicción \mathbf{P}_Z , se sigue del Teorema 5(1), que $0 \leq AP_i \leq 1$.

Recordando que

$$p_{zii} = \mathbf{z}_i^T (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{z}_i = p_{ii} + \frac{\hat{\epsilon}_i^2}{\hat{\epsilon}^T \hat{\epsilon}},$$

entonces la relación entre AP_i , r_i y p_{ii} , está dada por

$$AP_i = p_{ii} + \frac{\hat{\epsilon}_i^2}{\hat{\epsilon}^T \hat{\epsilon}} = p_{ii} + (1 - p_{ii}) \frac{r_i^2}{(n - k)},$$

la cual, también puede escribirse como

$$(1 - \mathbf{AP}_i) = (1 - p_{ii}) \left(1 - \frac{r_i^2}{(n-k)} \right)$$

Como señalan Draper y John (1981), $(1 - \mathbf{AP}_i)$ es un producto de dos cantidades, la primera identifica puntos palanca y la segunda detecta “outliers”, por lo que resulta más informativa que \mathbf{AP}_i .

Relación de Varianzas

Constituye una alternativa al estadístico Andrews-Pregibon, para medir la influencia de la i -ésima observación, comparando la varianza estimada de $\hat{\beta}$ dada por $\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})^{-1}$, y la varianza estimada de $\hat{\beta}_{(i)}$ dada por $\hat{\sigma}_{(i)}^2(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1}$. Si el $r(\mathbf{X}_{(i)}) = k$, estas matrices son positivas definidas, las cuales, se comparan mediante la relación de sus trazas o sus determinantes. Belsley *et al.* (1980), prefieren utilizar la relación entre sus determinantes, esto es,

$$\begin{aligned} VR_i &= \frac{\det(\hat{\sigma}_{(i)}^2(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1})}{\det(\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})^{-1})} \\ &= \left(\frac{\hat{\sigma}_{(i)}^2}{\hat{\sigma}^2} \right)^k \frac{\det(\mathbf{X}^T \mathbf{X})}{\det(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})}, \quad i = 1, 2, \dots, n \end{aligned} \quad (3.31)$$

Observando ahora que, esta relación de determinantes tiene que ver con r_i y p_{ii} , para lo cual, definiendo $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ y $\mathbf{B} = \mathbf{C} = \mathbf{x}_i$ y utilizando el Lema 4, resulta que

$$\begin{aligned} \det(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)}) &= \det(\mathbf{X}^T \mathbf{X} - \mathbf{x}_i \mathbf{x}_i^T) \\ &= \det(\mathbf{X}^T \mathbf{X}) \left(1 - \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \right) \\ &= \det(\mathbf{X}^T \mathbf{X}) (1 - p_{ii}). \end{aligned} \quad (3.32)$$

Sustituyendo las ecuaciones (3.7) y (3.32) en la ecuación (3.31), se produce

$$VR_i = \left(\frac{n-k-r_i^2}{n-k-1} \right)^k \frac{1}{(1-p_{ii})} \quad (3.33)$$

Un análisis de la ecuación (3.33), indica que $\mathbf{V}R_i > 1$, cuando r_i^2 es pequeño y p_{ii} es grande, mientras que si r_i^2 es grande y p_{ii} es pequeño, entonces $\mathbf{V}R_i < 1$; pero si r_i^2 y p_{ii} son grandes o pequeños, entonces $\mathbf{V}R_i$ se aproxima a 1. En la práctica, se ha observado que, $\mathbf{V}R_i$ elige exitosamente observaciones influyentes.

Lo ideal ocurre cuando todas las observaciones tienen igual influencia sobre la matriz de covarianzas, en este caso, $\mathbf{V}R_i$ es aproximadamente igual a 1, cualquier desviación de la unidad, indica que la i -ésima observación es potencialmente influyente. Belsley *et al.* (1980), proporcionan dos aproximaciones de puntos de corte para $\mathbf{V}R_i$, que son:

1. Cuando $|r_i| \geq 2$ con $p_{ii} = \frac{1}{n}$, se produce aproximadamente

$$\mathbf{V}R_i \leq 1 - \frac{3k}{n-k},$$

el cual, es útil sólo cuando $n > 4k$.

2. Cuando $r_i = 0$ con $p_{ii} \geq \frac{2k}{n}$, resulta aproximadamente

$$\mathbf{V}R_i \geq 1 + \frac{3k}{n-k}.$$

Estos autores, reemplazan $(n-k)$ por n y ponen estos límites de la forma

$$|\mathbf{V}R_i - 1| \geq \frac{3k}{n}$$

Estadístico Cook-Weisberg

Bajo normalidad, el $100(1 - \alpha)\%$ del elipsoide de confianza conjunta para β es dado por

$$E = \left\{ \beta : \frac{(\beta - \hat{\beta})^T (\mathbf{X}^T \mathbf{X})(\beta - \hat{\beta})}{k\hat{\sigma}^2} \leq F(\alpha; k, n-k) \right\},$$

pero, cuando la i -ésima observación es omitida, E se convierte en

$$E_{(i)} = \left\{ \beta : \frac{(\hat{\beta} - \hat{\beta}_{(i)})^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})(\hat{\beta} - \hat{\beta}_{(i)})}{k\hat{\sigma}_{(i)}^2} \leq F(\alpha; k, n-k-1) \right\}$$

Cook y Weisberg (1980), proponen el logaritmo de la relación entre E y $E_{(i)}$ como una medida de la influencia de la i -ésima observación sobre el volumen del elipsoide de confianza para β , así, se tiene

$$CW_i = \log \frac{\text{Volumen}(E)}{\text{Volumen}(E_{(i)})}, \quad i = 1, 2, \dots, n, \quad (3.34)$$

y como el volumen de un elipsoide es proporcional a la inversa de la raíz cuadrada del determinante de la matriz de forma cuadrática asociada, entonces la ecuación (3.34), se convierte en

$$CW_i = \log \left(\left(\frac{\det(X_{(i)}^T X_{(i)})}{\det(X^T X)} \right)^{\frac{1}{2}} \left(\frac{\hat{\sigma}}{\hat{\sigma}_{(i)}} \right)^k \left(\frac{F(\alpha; k, n-k)}{F(\alpha; k, n-k-1)} \right)^{\frac{k}{2}} \right) \quad (3.35)$$

Sustituyendo las ecuaciones (3.7) y (3.32) en la ecuación (3.35), se obtiene

$$CW_i = \frac{1}{2} \log(1 - p_{ii}) + \frac{k}{2} \log \left(\frac{n-k-1}{n-k-r_i^2} \right) + \frac{k}{2} \log \left(\frac{F(\alpha; k, n-k)}{F(\alpha; k, n-k-1)} \right) \quad (3.36)$$

Cook y Weisberg (1980), señalan que si CW_i es grande y positivo, entonces la eliminación de la i -ésima observación ocasionará una reducción en volumen, y si es grande y negativo, resultará en un incremento en volumen. Un análisis de la ecuación (3.36), indica que CW_i será grande y positivo, cuando r_i^2 es grande y p_{ii} pequeño, y será grande y negativo, cuando r_i^2 es pequeño y p_{ii} grande. Pero si r_i^2 y p_{ii} son grandes o pequeños, entonces CW_i se aproxima a cero.

Se observa en la ecuación (3.33), que se puede relacionar a CW_i con VR_i por medio de

$$CW_i = -\frac{1}{2} \log(VR_i) + \frac{k}{2} \log \left(\frac{F(\alpha; k, n-k)}{F(\alpha; k, n-k-1)} \right), \quad (3.37)$$

en donde se observa que CW_i y VR_i se pueden considerar como cantidades equivalentes.

Medidas Basadas en la Función de Verosimilitud

Sea $\log(\boldsymbol{\beta}, \sigma^2)$ el logaritmo de la función de verosimilitud basada en las n observaciones; sean $\tilde{\boldsymbol{\beta}}$ y $\tilde{\sigma}^2$ los estimadores de máxima verosimilitud (MLE) de $\boldsymbol{\beta}$ y σ^2 , respectivamente; y sea $\log(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2)$ el logaritmo de la función de verosimilitud evaluada en $\tilde{\boldsymbol{\beta}}$ y $\tilde{\sigma}^2$. Se sabe que, el $100(1 - \alpha)\%$ de la región de confianza asintótica para $\boldsymbol{\beta}$ y σ^2 es dada por

$$\{(\boldsymbol{\beta}, \sigma^2) : 2(\log(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) - \log(\boldsymbol{\beta}, \sigma^2)) \leq \chi^2(\alpha; k + 1)\}, \quad (3.38)$$

donde $\chi^2(\alpha; k + 1)$ es el punto superior α de la distribución ji-cuadrada con $(k + 1)$ grados de libertad.

Suponiendo ahora que $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, entonces el logaritmo de la función de verosimilitud para $\boldsymbol{\beta}$ y σ^2 es

$$\log(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2}\log 2\pi - \frac{n}{2}\log \sigma^2 - \frac{(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2} \quad (3.39)$$

de la cual se sigue que los MLE de $\boldsymbol{\beta}$ y σ^2 son

$$\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} \quad y \quad \tilde{\sigma}^2 = \hat{\sigma}^2 \left(\frac{n - k}{n} \right)$$

Y cuando la i -ésima observación es omitida, los MLE de $\boldsymbol{\beta}$ y σ^2 son

$$\tilde{\boldsymbol{\beta}}_{(i)} = \hat{\boldsymbol{\beta}}_{(i)} \quad y \quad \tilde{\sigma}_{(i)}^2 = \hat{\sigma}_{(i)}^2 \left(\frac{n - k - 1}{n - 1} \right) = \hat{\sigma}^2 \left(\frac{n - k - \mathbf{r}_i^2}{n - 1} \right)$$

Tomando $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}$ y $\sigma^2 = \tilde{\sigma}^2$ en la ecuación (3.39) y simplificando, se obtiene

$$\log(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) = -\frac{n}{2}\log 2\pi - \frac{n}{2}\log \tilde{\sigma}^2 \left(\frac{n - k}{n} \right) - \frac{n}{2} \quad (3.40)$$

Similarmente, tomando $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}_{(i)}$ y $\sigma^2 = \tilde{\sigma}_{(i)}^2$ en la ecuación (3.39), se produce

$$\begin{aligned} \log(\tilde{\boldsymbol{\beta}}_{(i)}, \tilde{\sigma}_{(i)}^2) &= -\frac{n}{2}\log 2\pi - \frac{n}{2}\log \tilde{\sigma}_{(i)}^2 - \frac{1}{2\tilde{\sigma}_{(i)}^2} \sum_{r=1}^n (\mathbf{y}_r - \mathbf{x}_r^T \tilde{\boldsymbol{\beta}}_{(i)})^2 \\ &= -\frac{n}{2}\log 2\pi - \frac{n}{2}\log \tilde{\sigma}_{(i)}^2 - \frac{n - 1}{2} - \frac{(\mathbf{y}_i - \mathbf{x}_i^T \tilde{\boldsymbol{\beta}}_{(i)})^2}{2\tilde{\sigma}_{(i)}^2} \end{aligned} \quad (3.41)$$

La influencia de la i -ésima observación sobre la función de verosimilitud, puede ser medida por la distancia que existe entre las funciones de verosimilitud dadas en las ecuaciones (3.40) y (3.41). Por analogía a la ecuación (3.38), Cook y Weisberg (1982), definen la distancia de verosimilitud por

$$LD_i(\beta, \sigma^2) = 2 \left\{ \log(\tilde{\beta}, \tilde{\sigma}^2) - \log(\tilde{\beta}_{(i)}, \tilde{\sigma}_{(i)}^2) \right\}, \quad i = 1, 2, \dots, n \quad (3.42)$$

Sustituyendo los MLE $\tilde{\beta}_{(i)} = \hat{\beta}_{(i)}$ y $\tilde{\sigma}_{(i)}^2 = \hat{\sigma}^2 \left(\frac{n-k-r_i^2}{n-1} \right)$, así como las ecuaciones (3.40) y (3.41) en la ecuación (3.42), se tiene

$$LD_i(\beta, \sigma^2) = n \log \left(\frac{n(n-k-r_i^2)}{(n-1)(n-k)} \right) + \frac{(n-1)(\mathbf{y}_i - \mathbf{x}_i^T \hat{\beta}_{(i)})^2}{\hat{\sigma}^2(n-k-r_i^2)} - 1$$

Usando el residual predicho, y simplificando, $LD_i(\beta, \sigma^2)$ se convierte en

$$LD_i(\beta, \sigma^2) = n \log \left(\frac{n(n-k-r_i^2)}{(n-1)(n-k)} \right) + \frac{(n-1)r_i^2}{(1-p_{ii})(n-k-r_i^2)} - 1$$

La similitud existente entre las ecuaciones (3.38) y (3.42), sugiere que $LD_i(\beta, \sigma^2)$ se pueda comparar con la distribución ji-cuadrada con $(k+1)$ grados de libertad.

Notando que $LD_i(\beta, \sigma^2)$ se utiliza para estimar β y σ^2 , pero si solamente se desea estimar β , entonces el $100(1-\alpha)\%$ de la región de confianza asintótica para β es

$$\left\{ \beta : 2 \left(\log(\tilde{\beta}, \tilde{\sigma}^2) - \max_{\sigma} \log(\beta, \sigma^2) \right) \leq \chi^2(\alpha; k) \right\},$$

y la distancia de verosimilitud se convierte en

$$LD_i(\beta|\sigma^2) = 2 \left\{ \log(\tilde{\beta}, \tilde{\sigma}^2) - \max_{\sigma} \log(\tilde{\beta}_{(i)}, \sigma^2) \right\}, \quad i = 1, 2, \dots, n, \quad (3.43)$$

donde

$$\log(\tilde{\beta}_{(i)}, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{r=1}^n (\mathbf{y}_r - \mathbf{x}_r^T \tilde{\beta}_{(i)})^2$$

El valor de σ^2 que maximiza $(\tilde{\beta}_{(i)}, \sigma^2)$, se encuentra así,

$$\hat{\sigma}^2(\tilde{\beta}_{(i)}) = \frac{1}{n} \sum_{r=1}^n (\mathbf{y}_r - \mathbf{x}_r^T \tilde{\beta}_{(i)})^2 = \hat{\sigma}_{(i)}^2 \left(\frac{n-1}{n} \right) + \frac{\hat{\epsilon}_i^2}{n(1-p_{ii})^2},$$

de manera que

$$\max_{\sigma} \log(\tilde{\boldsymbol{\beta}}_{(i)}, \sigma^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \tilde{\sigma}^2(\tilde{\boldsymbol{\beta}}_{(i)}) - \frac{n}{2} \quad (3.44)$$

Por lo que al sustituir las ecuaciones (3.40) y (3.44) en la ecuación (3.43), se produce

$$\begin{aligned} LD_i(\boldsymbol{\beta}|\sigma^2) &= n \log \left(\frac{\tilde{\sigma}^2(\tilde{\boldsymbol{\beta}}_{(i)})}{\tilde{\sigma}^2} \right) \\ &= n \log \left(\frac{(n-1) \tilde{\sigma}_{(i)}^2}{n \tilde{\sigma}^2} + \frac{\hat{\epsilon}_i^2}{n \tilde{\sigma}^2 (1 - \mathbf{p}_{ii})^2} \right) \\ &= n \log \left(\frac{(n-k-1) \hat{\sigma}_{(i)}^2}{(n-k) \tilde{\sigma}^2} + \frac{\hat{\epsilon}_i^2}{(n-k) \tilde{\sigma}^2 (1 - \mathbf{p}_{ii})^2} \right) \\ &= n \log \left(\frac{(n-k-r_i^2)}{(n-k)} + \frac{r_i^2}{(n-k)(1 - \mathbf{p}_{ii})} \right) \\ &= n \log \left(1 + \frac{r_i^2}{(n-k)} \frac{\mathbf{p}_{ii}}{(1 - \mathbf{p}_{ii})} \right) \\ &= n \log \left(1 + \frac{k}{(n-k)} \mathbf{C}_i \right), \quad i = 1, 2, \dots, n, \end{aligned} \quad (3.45)$$

donde \mathbf{C}_i está dada en la ecuación (3.22), entonces $LD_i(\boldsymbol{\beta}|\sigma^2)$ es equivalente a la distancia de Cook, pero $LD_i(\boldsymbol{\beta}|\sigma^2)$ puede compararse con la distribución $\chi^2(\alpha; k)$.

Medidas Basadas en un Subconjunto de Coeficientes de Regresión

A continuación se presentan medidas que consideran la influencia que una observación tiene sobre un simple coeficiente de regresión y sobre combinaciones lineales de los coeficientes de regresión.

Influencia sobre un Simple Coeficiente de Regresión

Se supone que la j -ésima variable \mathbf{X}_j es la última columna de \mathbf{X} y particionando \mathbf{X} como $\mathbf{X} = (\mathbf{X}_{(j)} : \mathbf{X}_j)$, donde $\mathbf{X}_{(j)}$ es la matriz \mathbf{X} sin \mathbf{X}_j . El modelo dado en la ecuación (1.1), puede entonces ser escrito como

$$\mathbf{Y} = \mathbf{X}_{(j)}\boldsymbol{\beta}_{(j)} + \mathbf{X}_j\boldsymbol{\beta}_j + \boldsymbol{\epsilon}$$

Utilizando el Teorema 4, se puede descomponer la matriz de predicción P como

$$\begin{aligned} P &= P_{(j)} + \frac{(I - P_{(j)})X_j X_j^T (I - P_{(j)})}{X_j^T (I - P_{(j)})X_j} \\ &= P_{(j)} + \frac{W_j W_j^T}{W_j^T W_j}, \end{aligned} \quad (3.46)$$

donde

$$P_{(j)} = X_{(j)}(X_{(j)}^T X_{(j)})^{-1} X_{(j)}^T, \quad (3.47)$$

es la matriz de predicción para $X_{(j)}$, y

$$W_j = (I - P_{(j)})X_j, \quad (3.48)$$

es el vector de residuales, cuando se hace una regresión de X_j sobre $X_{(j)}$, y $(W_j^T W_j)$ es la suma de cuadrados de residuales.

Definiendo a $\hat{\beta}_j$ y $\hat{\beta}_{j(i)}$, como los estimadores de β_j obtenidos de los datos completos y los datos sin la i -ésima observación, respectivamente. Observando entonces que

$$(\hat{\beta}_j - \hat{\beta}_{j(i)}) = \frac{\hat{\epsilon}_i}{(1 - p_{ii})} \frac{w_{ij}}{W_j^T W_j}, \quad (3.49)$$

donde w_{ij} es la i -ésima componente de W_j .

La influencia de la i -ésima observación sobre el j -ésimo coeficiente de regresión estimado, se obtiene al dividir la ecuación (3.49) por el error estándar de $(\hat{\beta}_j - \hat{\beta}_{j(i)})$, entonces se tiene

$$\text{Var}(\hat{\beta}_j - \hat{\beta}_{j(i)}) = \frac{\sigma^2}{(1 - p_{ii})} \left(\frac{w_{ij}}{W_j^T W_j} \right)^2$$

y así,

$$\frac{(\hat{\beta}_j - \hat{\beta}_{j(i)})}{\sqrt{\text{Var}(\hat{\beta}_j - \hat{\beta}_{j(i)})}} = \frac{\hat{\epsilon}_i}{\sigma \sqrt{1 - p_{ii}}},$$

el cual, dependiendo de la elección del estimador de σ , se puede obtener r_i o r_i^* . Alternativamente, la ecuación (3.49), puede ser dividida por el error estándar de $\hat{\beta}_j$, y así obtener

$$\frac{(\hat{\beta}_j - \hat{\beta}_j(i))}{\sqrt{\text{Var}(\hat{\beta}_j)}} \quad (3.50)$$

Se ve que $\hat{\beta}_j$ puede expresarse como

$$\hat{\beta}_j = \frac{X_j^T (I - P_{(j)}) Y}{X_j^T (I - P_{(j)}) X_j} = \frac{W_j^T Y}{W_j^T W_j} \quad (3.51)$$

Y de la ecuación (3.51), se sigue que

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{W_j^T W_j} \quad (3.52)$$

Sustituyendo ahora, las ecuaciones (3.49) y (3.52) en la ecuación (3.50), se produce

$$\frac{(\hat{\beta}_j - \hat{\beta}_j(i))}{\sqrt{\text{Var}(\hat{\beta}_j)}} = \frac{\hat{\epsilon}_i}{\sigma \sqrt{1 - p_{ii}}} \frac{w_{ij}}{\sqrt{W_j^T W_j}} \frac{1}{\sqrt{1 - p_{ii}}} \quad (3.53)$$

Usando $\hat{\sigma}$ como estimador de σ en la ecuación (3.53), se obtiene

$$D_{ij} = r_i \frac{w_{ij}}{\sqrt{W_j^T W_j}} \frac{1}{\sqrt{1 - p_{ii}}} \quad (3.54)$$

Mientras que, usando $\hat{\sigma}_{(i)}^2$ como estimador de σ en la ecuación (3.53), se llega a

$$D_{ij}^* = r_i^* \frac{w_{ij}}{\sqrt{W_j^T W_j}} \frac{1}{\sqrt{1 - p_{ii}}} \quad (3.55)$$

Belsley *et al.*, (1980), sugieren poner especial atención, a aquellos puntos con valores de $|D_{ij}^*|$, que excedan $\frac{2}{n}$.

Influencia sobre Funciones Lineales de $\hat{\beta}$

Suponiendo ahora, que la i -ésima observación influye un subconjunto de coeficientes de regresión dado, es decir, q combinaciones linealmente independientes de

β , o sea, $L^T \beta$, donde L es una matriz $k \times q$ con rango q . Sea $\Theta = L^T \beta$, el máximo MLE de Θ es $\hat{\Theta} = L^T \hat{\beta}$. La curva de influencia para $\hat{\Theta}$ es $L^T \psi(x^T, y, F, \hat{\beta}(F))$. Si se utiliza la curva de influencia muestral, SIC_i dada en la ecuación (3.13), para aproximar $\psi(\cdot)$, y se usa además, la ecuación (3.21) como una medida de la influencia de la i -ésima observación sobre $\hat{\Theta}$, se obtiene

$$C_i(L) = D_i(M, c) = \frac{(\hat{\beta} - \hat{\beta}_{(i)})^T M (\hat{\beta} - \hat{\beta}_{(i)})}{c} \quad (3.56)$$

donde $c = q(n-1)^2 \hat{\sigma}^2$ y $M = L \left(L^T (X^T X)^{-1} L \right)^{-1} L^T$. Usando la ecuación (3.16), entonces la ecuación (3.56) se reduce a

$$C_i(L) = \frac{r_i^2}{q(1-p_{ii})} x_i^T (X^T X)^{-1} M (X^T X)^{-1} x_i, \quad i = 1, 2, \dots, n \quad (3.57)$$

Se supone que el interés está centrado sólo en q elementos de β , las cuales, son las últimas q componentes de β , y así $L^T = (0 : I)$, donde 0 es una matriz nula de tamaño $q \times (k-q)$. La partición de X es $X = (X_1 : X_2)$ y la de x_i es $x_i^T = (x_{1i}^T : x_{2i}^T)$. Usando el Lema 1(1), se tiene que $(X^T X)^{-1}$ se puede escribir como

$$\begin{bmatrix} (X_1^T X_1)^{-1} + (X_1^T X_1)^{-1} X_1^T X_2 A X_2^T X_1 (X_1^T X_1)^{-1} & -(X_1^T X_1)^{-1} X_1^T X_2 A \\ -A X_2^T X_1 (X_1^T X_1)^{-1} & A \end{bmatrix}$$

donde $A^{-1} = X_2^T (I - X_1 (X_1^T X_1)^{-1} X_1^T) X_2$, y así,

$$\begin{aligned} & (X^T X)^{-1} M (X^T X)^{-1} \\ &= \begin{bmatrix} (X_1^T X_1)^{-1} X_1^T X_2 A X_2^T X_1 (X_1^T X_1)^{-1} & -(X_1^T X_1)^{-1} X_1^T X_2 A \\ -A X_2^T X_1 (X_1^T X_1)^{-1} & A \end{bmatrix} \\ &= (X^T X)^{-1} - \begin{bmatrix} (X_1^T X_1)^{-1} & 0 \\ 0 & 0 \end{bmatrix}, \end{aligned}$$

de lo cual se sigue que

$$x_i^T (X^T X)^{-1} M (X^T X)^{-1} x_i = p_{ii} - x_{1i}^T (X^T X)^{-1} x_{1i},$$

por lo tanto, la ecuación (3.57) puede escribirse como

$$C_i(L) = \frac{r_i^2}{q(1 - p_{ii})} (p_{ii} - x_{1i}^T (X^T X)^{-1} x_{1i}), \quad i = 1, 2, \dots, n \quad (3.58)$$

Así, la ecuación (3.58), mide la influencia de la i -ésima observación sobre q combinaciones linealmente independientes de los coeficientes de regresión especificados para L . Si sólo interesa un coeficiente β_j , entonces la ecuación (3.58) puede simplificarse, esto es, si $q = 1$, $L^T = (0^T : 1)$, $X_1 = X_{(j)}$ y $x_{1i} = x_{i(j)}$, se produce

$$C_i(L) = \frac{r_i^2}{(1 - p_{ii})} (p_{ii} - p_{ii(j)}) \quad (3.59)$$

De la ecuación (3.46), se tiene que

$$p_{ii} = p_{ii(j)} + \frac{w_{ij}^2}{W_j^T W_j},$$

y por consiguiente, la ecuación (3.59) se convierte en

$$D_{ij}^2 = C_i(L) = r_i^2 \frac{w_{ij}^2}{W_j^T W_j} \frac{1}{(1 - p_{ii})}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, k, \quad (3.60)$$

la cual mide la influencia de la i -ésima observación sobre $\hat{\beta}_j$. Se ve que el D_{ij} en (3.60), es el que está dado en la ecuación (3.54).

Método Basado en la Diferenciación

Hasta el momento, se ha visto el efecto que causa una observación individual sobre la ecuación de regresión ajustada, mediante el método de la eliminación de esa observación. Una alternativa, es el método basado en la diferenciación, el cual, en vez de realizar una modificación en los datos, se induce un pequeño cambio en algunos parámetros del modelo y luego se estudian los resultados del análisis de regresión como una función de estos parámetros modificados, Welsch y Kuh (1977), Pregibon (1979 y 1981) y Belsley *et al.* (1980).

$$\nabla \hat{\beta}(0) = \frac{\partial \hat{\beta}(\mathbf{w}_i)}{\partial \mathbf{w}_i} \Big|_{\mathbf{w}_i=0} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{\epsilon}_i}{(1 - \mathbf{p}_{ii})^2} \quad (3.67)$$

Existe equivalencia entre las medidas de influencia basadas en $\nabla \hat{\beta}(\mathbf{w}_i)$ y las basadas en la curva de influencia, esto es, \mathbf{EIC}_i es proporcional a $\nabla \hat{\beta}(1)$, y $\mathbf{EIC}_{(i)}$ es proporcional a $\nabla \hat{\beta}(0)$. Por lo que \mathbf{EIC}_i mide la tasa de cambio en $\hat{\beta}$ con la i -ésima observación ($\mathbf{w}_i = 1$), y $\mathbf{EIC}_{(i)}$ mide la tasa de cambio en $\hat{\beta}$ cuando la i -ésima observación es eliminada ($\mathbf{w}_i = 0$).

Puesto que $\hat{\beta}(\mathbf{w})$ es diferenciable sobre todo el rango de \mathbf{w} , el Teorema del Valor Medio, garantiza la existencia de \mathbf{w}_i^* , tal que $\nabla \hat{\beta}(1) \leq \nabla \hat{\beta}(\mathbf{w}_i^*) \leq \nabla \hat{\beta}(0)$. Entonces

$$\begin{aligned} \nabla \hat{\beta}(\mathbf{w}_i^*) &= \int_0^1 \nabla \hat{\beta}(\mathbf{w}_i) d\mathbf{w}_i \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \hat{\epsilon}_i \int_0^1 \frac{1}{(1 - \mathbf{p}_{ii}(1 - \mathbf{w}_i))^2} d\mathbf{w}_i \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i \frac{\hat{\epsilon}_i}{(1 - \mathbf{p}_{ii})} \end{aligned} \quad (3.68)$$

En las ecuaciones (3.17), (3.18) y (3.68), se ve que

$$\nabla \hat{\beta}(\mathbf{w}_i^*) = \frac{\mathbf{SIC}_i}{(n-1)} = \frac{\mathbf{SC}_i}{n},$$

por lo que, \mathbf{SIC}_i y \mathbf{SC}_i son proporcionales a $\nabla \hat{\beta}(\mathbf{w})$, $0 \leq \mathbf{w} \leq 1$.

EFFECTOS DE MULTIPLES OBSERVACIONES SOBRE LA ECUACIÓN DE REGRESIÓN

Hasta el momento se han descrito algunos métodos para detectar observaciones que individualmente pueden ser consideradas como “outliers”, puntos palanca o puntos influyentes; en este capítulo se extienden los procedimientos de una simple observación al caso general, en el que se consideran los efectos conjuntos de múltiples observaciones, sobre varios resultados del análisis de regresión.

Existen tres problemas inherentes al caso general de observaciones múltiples, que son:

1. La determinación del tamaño del subconjunto de observaciones conjuntamente influyentes. Suponiendo que se desean detectar todos los subconjuntos de tamaño $m = 2, 3, 4, \dots$, de observaciones que se consideran potencialmente influyentes, en donde un método secuencial puede usarse para determinar m , iniciando con $m = 2, m = 3, m = 4$, etc.
2. Se ha visto, en el caso de una simple observación, que para cada medida de diagnóstico, se calculan n cantidades, una para cada observación en el conjunto de datos. En el caso de observaciones múltiples, sin embargo, se tiene la desventaja de que para cada subconjunto de tamaño m , existen $\binom{n}{m}$ posibles subconjuntos para los cuales cada medida de diagnóstico puede ser calculada.
3. Las medidas de influencia, para el caso de una simple observación, pueden ser fácilmente graficadas. Pero en el caso de observaciones múltiples, no es posible la representación gráfica, especialmente para n y m grandes.

Método de la Eliminación de Múltiples Observaciones

Como en el caso de la detección de una observación influyente individual, en el caso general, se tienen las medidas correspondientes que permiten detectar múltiples observaciones, y que son:

1. Medidas basadas en residuales.
2. Medidas basadas en la curva de influencia.
3. Medidas basadas en el volumen de elipsoides de confianza.
4. Medidas basadas en la función de verosimilitud.
5. Medidas basadas en un subconjunto de coeficientes de regresión.

Medidas Basadas en Residuales

Recordando que si la i -ésima observación es omitida, equivale a ajustar el modelo

$$E(Y) = X\beta + u_i\theta, \quad (4.1)$$

donde θ es el coeficiente de regresión del i -ésimo vector unitario u_i .

Esta aproximación en que se adiciona una variable indicadora u_i como una forma para omitir la i -ésima observación, puede ser generalizada al caso de omitir m observaciones y adicionar m variables indicadoras.

Sea $I = \{i_1, i_2, \dots, i_m\}$, $m < (n - k)$, el conjunto que contiene los índices de las m observaciones que serán eliminadas, y sea $U_I = \{u_{i_1}, u_{i_2}, \dots, u_{i_m}\}$, la matriz que contiene las correspondientes variables indicadoras. Se supone ahora que, las m

observaciones que serán omitidas, son las últimas m observaciones, tales que \mathbf{Y} , \mathbf{X} y \mathbf{U}_I , pueden escribirse como

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_{(I)} \\ \mathbf{Y}_I \end{bmatrix} \begin{array}{l} (n-m) \times 1 \\ m \times 1 \end{array} \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_{(I)} \\ \mathbf{X}_I^T \end{bmatrix} \begin{array}{l} (n-m) \times k \\ m \times k \end{array}$$

$$\mathbf{U}_I = \begin{bmatrix} \mathbf{0} \\ \mathbf{I} \end{bmatrix} \begin{array}{l} (n-m) \times m \\ m \times m \end{array}$$

Así, omitiendo las m observaciones indexadas por I , equivale a ajustar el modelo

$$\mathbf{E}(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{U}_I\boldsymbol{\theta} \quad (4.2)$$

donde $\boldsymbol{\theta}$ es un vector $m \times 1$ de los coeficientes de regresión de las variables indicadoras \mathbf{U}_I .

El estadístico para probar,

$$H_0 : \mathbf{E}(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta} \quad \text{vs} \quad H_1 : \mathbf{E}(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{U}_I\boldsymbol{\theta},$$

es

$$F_I = \frac{\frac{SSE(H_0) - SSE(H_1)}{m}}{\frac{SSE(H_1)}{(n-k-m)}} \quad (4.3)$$

Sean \mathbf{P} y \mathbf{P}_{X,U_I} , las matrices de predicción para \mathbf{X} y $(\mathbf{X} : \mathbf{U}_I)$, respectivamente.

Utilizando el Teorema 4, se tiene

$$\begin{aligned} & \mathbf{Y}^T (\mathbf{I} - \mathbf{P}_{X,U_I}) \mathbf{Y} \\ &= \mathbf{Y}^T (\mathbf{I} - \mathbf{P}) \mathbf{Y} - \mathbf{Y}^T (\mathbf{I} - \mathbf{P}) \mathbf{U}_I \left(\mathbf{U}_I^T (\mathbf{I} - \mathbf{P}) \mathbf{U}_I \right)^{-1} \mathbf{U}_I^T (\mathbf{I} - \mathbf{P}) \mathbf{Y} \end{aligned} \quad (4.4)$$

Notando que

$$\left(\mathbf{U}_I^T (\mathbf{I} - \mathbf{P}) \mathbf{U}_I \right)^{-1} = (\mathbf{I} - \mathbf{P}_I)^{-1}, \quad (4.5)$$

donde

$$\mathbf{P}_I = \mathbf{X}_I^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I, \quad (4.6)$$

es la submatriz de \mathbf{P} , cuyas filas y columnas están indexadas por I . También notando que

$$\mathbf{U}_I^T (\mathbf{I} - \mathbf{P}) \mathbf{Y} = \mathbf{U}_I^T \hat{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\epsilon}}_I, \quad (4.7)$$

donde

$$\hat{\boldsymbol{\epsilon}}_I = \mathbf{Y}_I - \mathbf{X}_I^T \hat{\boldsymbol{\beta}}, \quad (4.8)$$

es el subconjunto de los residuales cuyos índices están contenidos en I . Sustituyendo las ecuaciones (4.5) y (4.7) en la ecuación (4.4), y reordenando, se obtiene

$$\mathbf{Y}^T (\mathbf{I} - \mathbf{P}) \mathbf{Y} - \mathbf{Y}^T (\mathbf{I} - \mathbf{P}_{X,U_I}) \mathbf{Y} = \hat{\boldsymbol{\epsilon}}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I, \quad (4.9)$$

entonces, se deduce que

$$SSE(H_0) - SSE(H_1) = \hat{\boldsymbol{\epsilon}}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I \quad (4.10)$$

Sustituyendo la ecuación (4.10) en la ecuación (4.3), se obtiene

$$\mathbf{F}_I = \frac{(n - k - m) \hat{\boldsymbol{\epsilon}}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I}{m \mathbf{Y}^T (\mathbf{I} - \mathbf{P}_{X,U_I}) \mathbf{Y}} = \frac{\hat{\boldsymbol{\epsilon}}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I}{m \hat{\sigma}_I^2}, \quad (4.11)$$

donde

$$\hat{\sigma}_{(I)}^2 = \frac{SSE_{Y(I)X(I)}}{(n - k - m)} = \frac{\mathbf{Y}^T (\mathbf{I} - \mathbf{P}_{X,U_I}) \mathbf{Y}}{(n - k - m)}, \quad (4.12)$$

es el estimador de σ^2 cuando las m observaciones indexadas por I son omitidas. Si se dividen ambos lados de la ecuación (4.10) por $(n - k - m)$, y reordenando, se obtiene

$$\hat{\sigma}_{(I)}^2 = \frac{(n - k)}{(n - k - m)} \hat{\sigma}^2 - \frac{\hat{\boldsymbol{\epsilon}}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I}{(n - k - m)} = \hat{\sigma}^2 \left(\frac{n - k - r_I^2}{n - k - m} \right), \quad (4.13)$$

donde

$$r_I^2 = \frac{\hat{\boldsymbol{\epsilon}}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I}{\hat{\sigma}^2}, \quad (4.14)$$

es el análogo de los r_i^2 .

Suponiendo que existen máximo m "outliers" en el conjunto de datos y que se conocen las etiquetas de esas observaciones, entonces bajo los supuestos de normalidad, el estadístico F_I dado en la ecuación (4.11) se distribuye como $F(m, n - k - m)$. En la práctica, se calcula F_I , para todo $\binom{n}{m}$ posible subconjunto de tamaño m y se busca el F_I más grande. La distribución de $\max(F_I)$ es difícil de obtener, como en el caso de r_{max} .

Medidas Basadas en la Curva de Influencia

Se han visto varias medidas de la influencia de la i -ésima observación sobre $\hat{\beta}$, las cuales, están basadas en la ecuación (3.21). Se presenta ahora, una generalización de $D_i(M, c)$, la cual es dada por

$$D_i(M, c) = \frac{\psi^T \left(X_I^T, Y_I, F, \hat{\beta}(F) \right) M \psi \left(X_I^T, Y_I, F, \hat{\beta}(F) \right)}{c}, \quad (4.15)$$

donde $\psi \left(X_I^T, Y_I, F, \hat{\beta}(F) \right)$ es una generalización de la versión muestral de la curva de influencia para múltiples observaciones.

Al igual que para el caso de una simple observación, se presentan a continuación dos aproximaciones a la curva de influencia para $\hat{\beta}$.

Curva de Influencia Muestral

Una generalización de la curva de influencia muestral SIC_i puede obtenerse de la ecuación (3.12), si se omite el límite y se reemplaza (x^T, y) por (X_I^T, Y_I) y F por \widehat{F} , y se define a $\epsilon = \frac{m}{m-n}$, entonces se obtiene

$$\begin{aligned} SIC_I &= \frac{(n-m)}{-m} \left(T \left(\frac{n}{(n-m)} - \frac{m}{(n-m)} \delta_{X_I^T, Y_I} \right) - T(\widehat{F}) \right) \\ &= \frac{(n-m)}{m} \left(T(\widehat{F}) - T(\widehat{F}_{(I)}) \right), \end{aligned} \quad (4.16)$$

donde $\widehat{\mathbf{F}}_{(I)}$ es la función de distribución acumulada empírica cuando las m observaciones indexadas por I son omitidas. Sustituyendo $\widehat{\beta}$ por $\mathbf{T}(\widehat{\mathbf{F}})$ y $\widehat{\beta}_{(I)}$ por $\mathbf{T}(\widehat{\mathbf{F}}_{(I)})$ en la ecuación (4.16), se obtiene

$$SIC_I = \frac{(n-m)}{m} (\widehat{\beta} - \widehat{\beta}_{(I)}), \quad (4.17)$$

donde

$$\widehat{\beta}_{(I)} = (\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})^{-1} \mathbf{X}_{(I)}^T \mathbf{Y}_{(I)}, \quad (4.18)$$

son los coeficientes de regresión estimados, cuando las m observaciones indexadas por I son eliminadas. Para simplificar la ecuación (4.17), se expresa la diferencia $(\widehat{\beta} - \widehat{\beta}_{(I)})$ en términos de cantidades obtenidas del ajuste del modelo a los datos completos. Reescribiendo $\widehat{\beta}_{(I)}$ como

$$\widehat{\beta}_{(I)} = (\mathbf{X}^T \mathbf{X} - \mathbf{X}_I \mathbf{X}_I^T)^{-1} (\mathbf{X}^T \mathbf{Y} - \mathbf{X}_I \mathbf{Y}_I) \quad (4.19)$$

Y usando el Lema 2 para evaluar la inversa de $(\mathbf{X}^T \mathbf{X} - \mathbf{X}_I \mathbf{X}_I^T)$, se obtiene

$$\begin{aligned} (\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})^{-1} &= (\mathbf{X}^T \mathbf{X} - \mathbf{X}_I \mathbf{X}_I^T)^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{X}_I^T (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned} \quad (4.20)$$

Sustituyendo la ecuación (4.20) en la ecuación (4.19), se obtiene

$$\begin{aligned} \widehat{\beta}_{(I)} &= \left((\mathbf{X}^T \mathbf{X})^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{X}_I^T (\mathbf{X}^T \mathbf{X})^{-1} \right) (\mathbf{X}^T \mathbf{Y} - \mathbf{X}_I \mathbf{Y}_I) \\ &= \widehat{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{X}_I^T \widehat{\beta} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I \mathbf{Y}_I \\ &\quad - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I \mathbf{Y}_I \end{aligned}$$

Adicionando y sustrayendo $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{Y}_I$, y reordenando, se obtiene

$$\begin{aligned} (\widehat{\beta} - \widehat{\beta}_{(I)}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \widehat{\epsilon}_I \\ &\quad + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - (\mathbf{I} - \mathbf{P}_I)^{-1} + (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I) \mathbf{Y}_I \end{aligned}$$

y como $(\mathbf{I} - (\mathbf{I} - \mathbf{P}_I)^{-1} + (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I) = \mathbf{0}$, entonces resulta que

$$(\widehat{\beta} - \widehat{\beta}_{(I)}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \widehat{\epsilon}_I, \quad (4.21)$$

la cual constituye una generalización de la ecuación (3.16). Se puede encontrar una expresión más simple para SIC_I , al sustituir la ecuación (4.21) en la ecuación (4.17), lo cual produce

$$SIC_I = \frac{(n-m)}{m} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I, \quad (4.22)$$

la cual se reduce a la ecuación (3.17), cuando $m=1$.

Curva de Influencia Empírica

La curva de influencia empírica basada en $(n-m)$ observaciones, se encuentra a partir de la curva de influencia para $\hat{\boldsymbol{\beta}}$, definida en el Teorema 13(1). Aproximando \mathbf{F} por $\hat{\mathbf{F}}_{(I)}$, y considerando que $(\mathbf{x}^T, \mathbf{y}) = (\mathbf{X}_I^T, \mathbf{Y}_I)$, $\Sigma_{xx}(\mathbf{F}) = (n-m)^{-1}(\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})$ y $\hat{\boldsymbol{\beta}}(\mathbf{F}) = \hat{\boldsymbol{\beta}}_{(I)}$, se obtiene

$$EIC_{(I)} = (n-m)(\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})^{-1} \mathbf{X}_I (\mathbf{Y}_I - \mathbf{X}_I^T \hat{\boldsymbol{\beta}}_{(I)}) \quad (4.23)$$

De la ecuación (4.20), se tiene que

$$\begin{aligned} (\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})^{-1} \mathbf{X}_I &= \left((\mathbf{X}^T \mathbf{X})^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{X}_I^T (\mathbf{X}^T \mathbf{X})^{-1} \right) \mathbf{X}_I \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} + (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \end{aligned} \quad (4.24)$$

Mientras que de la ecuación (4.21), se tiene

$$\begin{aligned} (\mathbf{Y}_I - \mathbf{X}_I^T \hat{\boldsymbol{\beta}}_{(I)}) &= \mathbf{Y}_I - \mathbf{X}_I^T \left(\hat{\boldsymbol{\beta}} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I \right) \\ &= \hat{\boldsymbol{\epsilon}}_I + \mathbf{P}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I \\ &= (\mathbf{I} + \mathbf{P}_I (\mathbf{I} - \mathbf{P}_I)^{-1}) \hat{\boldsymbol{\epsilon}}_I \\ &= (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\boldsymbol{\epsilon}}_I \end{aligned} \quad (4.25)$$

Sustituyendo las ecuaciones (4.24) y (4.25) en la ecuación (4.23), y simplificando, se obtiene

$$EIC_{(I)} = (n-m)(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-2} \hat{\boldsymbol{\epsilon}}_I, \quad (4.26)$$

la cual se reduce a la ecuación (3.20) cuando $m=1$.

Las medidas de diagnóstico basadas en la curva de influencia pueden también ser generalizadas al caso de observaciones múltiples. A continuación se muestran dos de ellas, las cuales se obtienen utilizando la ecuación (4.15), una aproximación a la curva de influencia para $\hat{\beta}$ (tal como SIC_I y $EIC_{(I)}$) y una elección apropiada para M y c .

Distancia de Cook Generalizada

El análogo a la distancia de Cook C_i , se obtiene sustituyendo SIC_I en la ecuación (4.15), después de una simplificación, se obtiene

$$C_I = D_I \left(\mathbf{X}^T \mathbf{X}, k\hat{\sigma}^2 \left(\frac{m}{n-m} \right)^2 \right) = \frac{\hat{\epsilon}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\epsilon}_I}{k\hat{\sigma}^2} \quad (4.27)$$

Un valor grande de C_I en la ecuación (4.27), indica que las observaciones indexadas por I , son conjuntamente influyentes sobre $\hat{\beta}$. Si se considera el caso específico en que $m = 2$ y $I = \{i, j\}$, entonces la ecuación (4.27) se convierte en

$$k\hat{\sigma}^2 C_I = k\hat{\sigma}^2 (C_i + C_j) \left(1 + \frac{\mathbf{p}_{ij}}{\delta_{ij}} \right)^2 + \frac{\mathbf{p}_{ij}^2 \left(\hat{\epsilon}_i^2 (2 - \mathbf{p}_{jj}) + \hat{\epsilon}_j^2 (2 - \mathbf{p}_{ii}) \right)}{\delta_{ij}^2} + \frac{2\hat{\epsilon}_i \hat{\epsilon}_j \mathbf{p}_{ij} (1 + \mathbf{p}_{ij}^2 - \mathbf{p}_{ii} \mathbf{p}_{jj})}{\delta_{ij}^2}, \quad (4.28)$$

donde $\delta_{ij} = (1 - \mathbf{p}_{ii})(1 - \mathbf{p}_{jj}) - \mathbf{p}_{ij}^2$, y por el Teorema 6(2), $\delta_{ij} \geq 0$. Un análisis de la ecuación (4.5), muestra que

$$\text{si } (\hat{\epsilon}_i \hat{\epsilon}_j \mathbf{p}_{ij}) > 0, \text{ entonces } C_I > (C_i + C_j),$$

se ve además que si δ_{ij} es pequeño, y $(\hat{\epsilon}_i \hat{\epsilon}_j \mathbf{p}_{ij})$ es grande, entonces C_I es grande.

Distancia de Welsch Generalizada

Si se utiliza $EIC_{(I)}$ dado en la ecuación (4.26) como una aproximación a la curva de influencia para $\hat{\beta}$, y se considera además que

$$M = \mathbf{X}_{(I)}^T \mathbf{X}_{(I)} = (\mathbf{X}^T \mathbf{X} - \mathbf{X}_I \mathbf{X}_I^T) \quad \text{y} \quad c = (n - m) \hat{\sigma}_{(I)}^2,$$

entonces, se obtiene

$$\begin{aligned}
 \mathbf{W}_I^2 &= \left(\frac{n-m}{\hat{\sigma}_{(I)}^2} \right) \hat{\mathbf{e}}_I^T (I - \mathbf{P}_I)^{-2} \mathbf{X}_I^T (\mathbf{X}^T \mathbf{X})^{-1} M (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (I - \mathbf{P}_I)^{-2} \hat{\mathbf{e}}_I \\
 &= \left(\frac{n-m}{\hat{\sigma}_{(I)}^2} \right) \hat{\mathbf{e}}_I^T (I - \mathbf{P}_I)^{-2} \mathbf{P}_I (I - \mathbf{P}_I) (I - \mathbf{P}_I)^{-2} \hat{\mathbf{e}}_I \\
 &= \left(\frac{n-m}{\hat{\sigma}_{(I)}^2} \right) \hat{\mathbf{e}}_I^T (I - \mathbf{P}_I)^{-2} \mathbf{P}_I (I - \mathbf{P}_I)^{-1} \hat{\mathbf{e}}_I
 \end{aligned} \tag{4.29}$$

Medidas Basadas en el Volumen de Elipsoides de Confianza

A continuación se presenta una generalización de las medidas \mathbf{AP}_i y \mathbf{VR}_i . No se generaliza \mathbf{CW}_i , pues es una función monótona de \mathbf{VR}_i .

Estadístico Andrews-Pregibon Generalizado

Definiendo la matriz aumentada $\mathbf{Z} = (\mathbf{X} : \mathbf{Y})$, y sean \mathbf{Z}_I^T las m filas de \mathbf{Z} que serán eliminadas y $\mathbf{Z}_{(I)}$ la matriz \mathbf{Z} sin \mathbf{Z}_I^T , entonces, utilizando el Lema 4, se tiene

$$\begin{aligned}
 \det(\mathbf{Z}_{(I)}^T \mathbf{Z}_{(I)}) &= \det(\mathbf{Z}^T \mathbf{Z} - \mathbf{Z}_I^T \mathbf{Z}_I^T) \\
 &= \det(\mathbf{Z}^T \mathbf{Z}) \det\left(I - \mathbf{Z}_I^T (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}_I\right) \\
 &= \det(\mathbf{Z}^T \mathbf{Z}) \det(I - \mathbf{P}_{\mathbf{Z}_I}),
 \end{aligned} \tag{4.30}$$

donde $\mathbf{P}_{\mathbf{Z}}$ es la matriz de predicción para \mathbf{Z} , y $\mathbf{P}_{\mathbf{Z}_I}$ es la submatriz de $\mathbf{P}_{\mathbf{Z}}$, cuyas filas y columnas están indexadas por I . Una generalización obvia de la ecuación (3.30), es

$$\mathbf{AP}_I = 1 - \frac{\det(\mathbf{Z}_{(I)}^T \mathbf{Z}_{(I)})}{\det(\mathbf{Z}^T \mathbf{Z})}$$

De la ecuación (4.30), se sigue que \mathbf{AP}_I puede expresarse como

$$\mathbf{AP}_I = 1 - \det(I - \mathbf{P}_{\mathbf{Z}_I}) \tag{4.31}$$

Del Teorema 6(4), se tiene que

$$P_Z = P + \frac{\widehat{\epsilon}\widehat{\epsilon}^T}{\widehat{\epsilon}^T\widehat{\epsilon}},$$

y así,

$$P_{Z_I} = P_I + \frac{\widehat{\epsilon}_I\widehat{\epsilon}_I^T}{\widehat{\epsilon}_I^T\widehat{\epsilon}_I}$$

Aplicando el Lema 4, se obtiene

$$\begin{aligned} AP_I &= 1 - \det(I - P_I) \left(1 - \frac{\widehat{\epsilon}_I^T (I - P_I)^{-1} \widehat{\epsilon}_I}{\widehat{\epsilon}_I^T \widehat{\epsilon}_I} \right) \\ &= 1 - \det(I - P_I) \left(1 - \frac{r_I^2}{(n - k)} \right), \end{aligned} \quad (4.32)$$

donde r_I^2 se muestra en la ecuación (4.14).

Relación de Varianzas Generalizada

En forma análoga a la ecuación (3.31), la influencia de las m observaciones indexadas por I sobre la varianza de los coeficientes de regresión estimados es medida por

$$VR_I = \frac{\det\left(\widehat{\sigma}_{(I)}^2 (\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})^{-1}\right)}{\det\left(\widehat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}\right)} = \left(\frac{\widehat{\sigma}_{(I)}^2}{\widehat{\sigma}^2}\right)^k \frac{\det(\mathbf{X}^T \mathbf{X})}{\det(\mathbf{X}_{(I)}^T \mathbf{X}_{(I)})} \quad (4.33)$$

Usando la demostración del Teorema 10(3) y sustituyendo la ecuación (4.14) en la ecuación (4.33), se obtiene

$$VR_I = \left(\frac{n - k - r_I^2}{n - k - m}\right)^k (\det(I - P_I))^{-1} \quad (4.34)$$

donde r_I^2 es dada en la ecuación (4.14).

Medidas Basadas en la Función de Verosimilitud

Cuando las m observaciones indexadas por I son omitidas, los MLE de β y σ^2 son dados por

$$\tilde{\beta}_{(I)} = \widehat{\beta}_{(I)}, \quad (4.35)$$

y

$$\tilde{\sigma}_{(I)}^2 = \hat{\sigma}_{(I)}^2 \left(\frac{n-k-m}{n-m} \right) = \hat{\sigma}^2 \left(\frac{n-k-r_I^2}{n-m} \right) \quad (4.36)$$

Sustituyendo estos estimadores, por sus respectivos parámetros en la ecuación (3.39), se obtiene

$$\log(\tilde{\beta}_{(I)}, \tilde{\sigma}_{(I)}^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \tilde{\sigma}_{(I)}^2 - \frac{(\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)})^T (\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)})}{2\tilde{\sigma}_{(I)}^2} \quad (4.37)$$

Usando las ecuaciones (4.21) y (4.35), se tiene

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)}) &= (\mathbf{Y} - \mathbf{X}\hat{\beta}_{(I)}) \\ &= \mathbf{Y} - \mathbf{X}\hat{\beta} + \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\epsilon}_I \\ &= \hat{\epsilon} + \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\epsilon}_I, \end{aligned}$$

y así,

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)})^T (\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)}) &= \hat{\epsilon}^T \hat{\epsilon} + 2\hat{\epsilon}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\epsilon}_I \\ &\quad + \hat{\epsilon}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\epsilon}_I \end{aligned} \quad (4.38)$$

Puesto que $\hat{\epsilon}^T \mathbf{X} = 0$, entonces la ecuación (4.38) se reduce a

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)})^T (\mathbf{Y} - \mathbf{X}\tilde{\beta}_{(I)}) &= \hat{\epsilon}^T \hat{\epsilon} + \hat{\epsilon}_I^T (\mathbf{I} - \mathbf{P}_I)^{-1} \mathbf{P}_I (\mathbf{I} - \mathbf{P}_I)^{-1} \hat{\epsilon}_I \\ &= \hat{\sigma}^2 ((n-k) + k\mathbf{C}_I), \end{aligned} \quad (4.39)$$

donde \mathbf{C}_I es definido en la ecuación (4.27).

Sustituyendo las ecuaciones (4.36) y (4.39) en la ecuación (4.37), y simplificando se obtiene

$$\log(\tilde{\beta}_{(I)}, \tilde{\sigma}_{(I)}^2) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \hat{\sigma}^2 \left(\frac{n-k-r_I^2}{n-m} \right) - \frac{(n-m)((n-k) + k\mathbf{C}_I)}{2(n-k-r_I^2)} \quad (4.40)$$

Finalmente, la distancia de verosimilitud para β y σ^2 , cuando m observaciones indexadas por I son eliminadas, es definida por

$$\begin{aligned} LD_I(\beta, \sigma^2) &= 2 \left(\log(\tilde{\beta}, \tilde{\sigma}^2) - \log(\tilde{\beta}_{(I)}, \tilde{\sigma}_{(I)}^2) \right) \\ &= n \log \left(\frac{n(n-k-r_I^2)}{(n-m)(n-k)} \right) + \frac{(n-m)(n-k+k\mathbf{C}_I)}{(n-k-r_I^2)} - n \end{aligned} \quad (4.41)$$

Medidas Basadas en un Subconjunto de Coeficientes de Regresión

La influencia de las m observaciones indexadas por I pueden medirse, análogamente a la ecuación (3.53), por

$$\frac{(\hat{\beta}_j - \hat{\beta}_{j(I)})}{\sqrt{\text{Var}(\hat{\beta}_j)}}, \quad (4.42)$$

el numerador de la ecuación (4.42), representa el cambio en el j -ésimo coeficiente de regresión debido a la eliminación de las m observaciones indexadas por I , y es dado por

$$(\hat{\beta}_j - \hat{\beta}_{j(I)}) = \hat{\epsilon}_I^T (I - P_I)^{-1} W_{jI} (W_j^T W_j)^{-1}, \quad (4.43)$$

donde W_{jI} es el subconjunto de W_j indexado por I , y W_j es el vector de residuales obtenido de la regresión de X_j sobre los demás componentes de X , para esto, se ve la ecuación (3.48).

De las ecuaciones (3.52) y (4.43), se sigue que

$$\frac{(\hat{\beta}_j - \hat{\beta}_{j(I)})}{\sqrt{\text{Var}(\hat{\beta}_j)}} = \sigma \hat{\epsilon}_I^T (I - P_I)^{-1} W_{jI} (W_j^T W_j)^{-\frac{1}{2}}, \quad (4.44)$$

la cual, es el análogo de la ecuación (3.53). Para $m = 1$, la ecuación (4.44) se reduce a las ecuaciones (3.54) o (3.55), dependiendo de cual cantidad se utilice para estimar σ .

Una modificación de la distancia de Cook

José A. Díaz-García

Departamento de Estadística y Cálculo
Universidad Autónoma Agraria Antonio Narro
25350 Buenavista, Saltillo, Coahuila, MÉXICO.

jadiaz@narro.uaaan.mx

y

Oscar Alejandro Martínez-Jaime

Instituto de Ciencias Agrícolas
Universidad de Guanajuato
Ex-Hda. "El Copal", Irapuato, Guanajuato, MÉXICO.

PALABRAS Y FRASES CLAVES: Medidas de diagnóstico, Distancia de Mahalanobis Generalizada, Puntos influyentes.

ABSTRACT

A modification of the classical Cook's distance is proposed, providing us with a generalized Mahalanobis distance in the context of multivariate normal linear regression models. We establish the exact distribution of a pivotal type statistics based on this generalized Mahalanobis distance, which provides critical points for the identification of data points. We illustrate the procedure with an example, in the context of multiple and multivariate linear regression.

RESUMEN

En el presente artículo se propone una modificación de la distancia de Cook, basándose en la distancia de Mahalanobis generalizada, en el contexto del modelo de regresión lineal multivariado con distribución normal. Se establece además, la

distribución exacta del estadístico basado en esta distancia de Mahalanobis generalizada, la cual proporciona puntos críticos para identificar “outliers” en un conjunto de datos. Este procedimiento, se ilustra con un ejemplo, en el caso de la regresión lineal múltiple y multivariada.

1. INTRODUCCIÓN

El problema de la identificación de “outliers” o puntos influyentes, en el caso de la regresión lineal univariada o multivariada y bajo el supuesto de que los errores se distribuyen normales, ha sido estudiado por varios autores, tales como Cook (1977), Besley et al. (1980), Cook y Weisberg (1982) y Chatterjee y Hadi (1988), sólo por mencionar algunos. Muchos de estos resultados han sido extendidos al caso de las distribuciones de contorno elíptico, ver por ejemplo Galea et al. (1997), Liu (2000) y Díaz-García et al. (2001), entre otros. En todos estos trabajos, la idea es utilizar la distancia de Cook como una medida de diagnóstico, para identificar observaciones influyentes, individualmente o en conjunto. Sin embargo, cuando se usa este criterio, se cuenta solamente con puntos críticos, los cuales son proporcionados por una aproximación a la distribución \mathcal{F} centrada, tal como lo propuso Cook (1977). En este artículo, el propósito es realizar una modificación a esta distancia y derivar su distribución exacta.

Suponiendo que $\mathbf{Y} \in \mathbb{R}^{n \times p}$ tiene una distribución normal con media $\mu \in \mathbb{R}^{n \times p}$ y matriz de covarianzas $\Sigma \otimes \Theta \in \mathbb{R}^{np \times np}$ con $\Sigma \in \mathbb{R}^{p \times p}$, $\Sigma > 0$ y $\Theta \in \mathbb{R}^{n \times n}$, $\Theta > 0$, entonces la función de densidad es dada por

$$\mathbf{f}_{\mathbf{Y}}(\mathbf{Y}) = (2\pi)^{-pn/2} |\Sigma|^{-p/2} |\Theta|^{-n/2} \text{etr}(\Sigma^{-1}(\mathbf{Y} - \mu)^T \Theta^{-1}(\mathbf{Y} - \mu))$$

Este hecho, se denota también, como $\mathbf{Y} \sim \mathcal{N}_{n \times p}(\mu, \Sigma \otimes \Theta)$.

Considerando el modelo de regresión lineal multivariado:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \tag{5.1}$$

donde $\mathbf{Y} \in \mathbb{R}^{n \times p}$ es la matriz respuesta, $\mathbf{X} \in \mathbb{R}^{n \times q}$, con $r(\mathbf{X}) = q$, $\boldsymbol{\beta} \in \mathbb{R}^{q \times p}$

es la matriz de parámetros desconocidos y $\epsilon \in \mathbb{R}^{n \times p}$ es una matriz de errores, tal que $\epsilon \sim \mathcal{N}_{n \times p}(0, \Sigma \otimes \mathbf{I}_n)$. Este modelo es conocido como modelo de regresión lineal normal multivariado. Los estimadores de máxima verosimilitud para β y Σ están dados por, ver Muirhead (pp. 83, 1982),

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{X}^- \mathbf{Y} \quad \text{y} \quad \hat{\Sigma} = \frac{1}{n} (\mathbf{Y} - \mathbf{X} \hat{\beta})^T (\mathbf{Y} - \mathbf{X} \hat{\beta}),$$

donde \mathbf{X}^- es la inversa de Moore-Penrose de \mathbf{X} .

Se considera entonces la regresión lineal con distribución normal multivariada, y se propone una extensión y modificación a la distancia de Cook. Esto permite derivar la distribución exacta para la nueva distancia, la cual, a su vez, proporciona un punto crítico para decidir si una observación en particular (o un conjunto de observaciones) se comportan como un "outlier".

2. DISTANCIA MODIFICADA : UNA OBSERVACIÓN

Considerando el modelo de regresión lineal normal multivariado con la siguiente modificación,

$$\mathbf{Y}_{(i)} = \mathbf{X}_{(i)} \beta_{(i)} + \epsilon_{(i)}, \quad \epsilon_{(i)} \sim \mathcal{N}_{(n-1) \times p}(0, \Sigma_{(i)} \otimes \mathbf{I}_n), \quad (5.2)$$

el cual se obtiene del modelo dado en (5.1), eliminando la i -ésima fila de \mathbf{Y} , \mathbf{X} y ϵ , esto es, eliminando la i -ésima observación.

Para el modelo modificado, se tiene que:

$$\hat{\beta}_{(i)} = (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} = \mathbf{X}_{(i)}^- \mathbf{Y}_{(i)} \quad \text{y}$$

$$\hat{\Sigma}_{(i)} = \frac{1}{(n-1)} (\mathbf{Y}_{(i)} - \mathbf{X}_{(i)} \hat{\beta}_{(i)})^T (\mathbf{Y}_{(i)} - \mathbf{X}_{(i)} \hat{\beta}_{(i)})$$

Primero, se requiere una representación simple para $\hat{\beta} - \hat{\beta}_{(i)}$. Para esto, considerando la siguiente partición en las matrices:

$$\mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1^T \\ \mathbf{Y}_2^T \\ \vdots \\ \mathbf{Y}_n^T \end{pmatrix}, \quad \mathbf{Y}_i \in \mathbb{R}^p \quad \epsilon = \begin{pmatrix} \epsilon_1^T \\ \epsilon_2^T \\ \vdots \\ \epsilon_n^T \end{pmatrix}, \quad \epsilon_i \in \mathbb{R}^p \quad \mathbf{X} = \begin{pmatrix} \mathbf{X}_1^T \\ \mathbf{X}_2^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix}, \quad \mathbf{X}_i \in \mathbb{R}^q.$$

por consiguiente

$$\begin{aligned} \mathbf{X}^T \mathbf{X} &= (\mathbf{X}_1 \mathbf{X}_2 \cdots \mathbf{X}_n) \begin{pmatrix} \mathbf{X}_1^T \\ \mathbf{X}_2^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix} = \sum_{k=1}^n \mathbf{X}_k \mathbf{X}_k^T = \sum_{k \neq i}^n \mathbf{X}_k \mathbf{X}_k^T + \mathbf{X}_i \mathbf{X}_i^T \\ &= \mathbf{X}_{(i)}^T \mathbf{X}_{(i)} + \mathbf{X}_i \mathbf{X}_i^T, \end{aligned}$$

y

$$\begin{aligned} \mathbf{X}^T \mathbf{Y} &= (\mathbf{X}_1 \mathbf{X}_2 \cdots \mathbf{X}_n) \begin{pmatrix} \mathbf{Y}_1^T \\ \mathbf{Y}_2^T \\ \vdots \\ \mathbf{Y}_n^T \end{pmatrix} = \sum_{k=1}^n \mathbf{X}_k \mathbf{Y}_k^T = \sum_{k \neq i}^n \mathbf{X}_k \mathbf{Y}_k^T + \mathbf{X}_i \mathbf{Y}_i^T \\ &= \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} + \mathbf{X}_i \mathbf{Y}_i^T. \end{aligned}$$

Notando que \mathbf{e}_i^n es el i -ésimo vector de la base canónica en \mathbb{R}^n , esto es, el vector unitario dado por $\mathbf{e}_i^n = (0 \cdots 0 \ 1 \ 0 \cdots 0)^T$, entonces $\mathbf{e}_i^{nT} \mathbf{Y} = \mathbf{Y}_i^T$, $\mathbf{e}_i^{nT} \mathbf{X} = \mathbf{X}_i^T$ y $\mathbf{e}_i^{nT} \boldsymbol{\epsilon} = \boldsymbol{\epsilon}_i^T$.

Rao (pp. 33, 1973), señala que si \mathbf{A} es no singular, \mathbf{v} y \mathbf{u} son dos vectores arbitrarios, entonces

$$(\mathbf{A} - \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} + \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^T\mathbf{A}^{-1}}{1 - \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u}}$$

por lo tanto, si se define $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ y $\mathbf{u} = \mathbf{v} = \mathbf{X}_i$, se obtiene

$$\begin{aligned} (\mathbf{X}^T \mathbf{X} - \mathbf{X}_i \mathbf{X}_i^T)^{-1} &= (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} + \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1}}{(1 - \mathbf{p}_{ii})}, \quad (5.3) \end{aligned}$$

con $\mathbf{p}_{ii} = \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i$.

De (5.3), se obtiene que

$$\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} - (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)}$$

$$\begin{aligned}
&= \left((\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} - \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1}}{(1 - p_{ii})} \right) \mathbf{X}^T \mathbf{Y} \\
&\quad - (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} \\
&= (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} (\mathbf{X}^T \mathbf{Y} - \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)}) - \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}}{(1 - p_{ii})} \\
&= (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_i \mathbf{Y}_i^T - \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}}{(1 - p_{ii})}. \tag{5.4}
\end{aligned}$$

Usando (5.3) en la primera parte de (5.4), se obtiene

$$\begin{aligned}
(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_i \mathbf{Y}_i^T &= \left((\mathbf{X}^T \mathbf{X})^{-1} + \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1}}{(1 - p_{ii})} \right) \mathbf{X}_i \mathbf{Y}_i^T \\
&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{Y}_i^T + \frac{p_{ii} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{Y}_i^T}{(1 - p_{ii})} \\
&= \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{Y}_i^T}{(1 - p_{ii})}. \tag{5.5}
\end{aligned}$$

Sustituyendo (5.5) en (5.4), resulta

$$\begin{aligned}
\hat{\beta} - \hat{\beta}_{(i)} &= \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{Y}_i^T - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}}{(1 - p_{ii})} \\
&= \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i}{(1 - p_{ii})} (\mathbf{Y}_i^T - \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) \tag{5.6}
\end{aligned}$$

Ahora, puesto que $\hat{\epsilon} = (\mathbf{Y} - \mathbf{X}\hat{\beta}) = (\mathbf{I} - \mathbf{X}\mathbf{X}^{-})\mathbf{Y} = (\mathbf{I} - \mathbf{P})\mathbf{Y}$, donde \mathbf{P} es el proyector ortogonal sobre la imagen de \mathbf{X} . Entonces

$$e_i^{n^T} \hat{\epsilon} = \hat{\epsilon}_i^T = e_i^{n^T} (\mathbf{Y} - \mathbf{X}\hat{\beta}) = \mathbf{Y}_i^T - \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y},$$

así, se obtiene:

$$\hat{\beta} - \hat{\beta}_{(i)} = \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \hat{\epsilon}_i^T}{(1 - p_{ii})} \tag{5.7}$$

Bajo el supuesto de la matriz de distribución normal, se propone la siguiente modificación a la distancia de Cook, denotada como \mathcal{D}_m :

$$\mathcal{D}_m = \text{vec}(\hat{\beta} - \hat{\beta}_{(i)})^T \widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))^{-1} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) \tag{5.8}$$

Nota 1. La expresión dada en (5.8), es una extensión del cuadrado de la distancia de Mahalanobis generalizada, ver Rao y Mitra (203-206, 1971).

El segundo paso, es encontrar una expresión simple para la matriz de varianzas y covarianzas $\text{Cov}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)}))$. Para lo cual, recordando que $\mathbf{Y} \sim \mathcal{N}_{n \times p}(\mu, \Sigma \otimes \Theta)$, entonces $\mathbf{E}(\mathbf{Y}) = \mu$, $\text{Cov}(\text{vec}(\mathbf{Y})) = (\Sigma \otimes \Theta)$, ver Muirhead (pp. 79, 1982).

Puesto que $\widehat{\epsilon}_i^T = e_i^{nT}(\mathbf{Y} - \mathbf{X}\widehat{\beta}) = e_i^{nT}(\mathbf{I} - \mathbf{P})\mathbf{Y}$, es claro que

$$\text{vec } \widehat{\epsilon}_i^T = (\mathbf{I}_p \otimes e_i^{nT}(\mathbf{I} - \mathbf{P})) \text{vec } \mathbf{Y},$$

por consiguiente

$$\begin{aligned} \text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)}) &= \frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i)}{(1 - p_{ii})} (\mathbf{I}_p \otimes e_i^{nT}(\mathbf{I} - \mathbf{P})) \text{vec } \mathbf{Y} \\ &= \frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i e_i^{nT}(\mathbf{I} - \mathbf{P}))}{(1 - p_{ii})} \text{vec } \mathbf{Y} \\ &= \frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{H}_i^T)}{(1 - p_{ii})} \text{vec } \mathbf{Y} \end{aligned} \quad (5.9)$$

donde $\mathbf{H}_i^T = e_i^{nT}(\mathbf{I} - \mathbf{P})$ es la i -ésima fila de la matriz $(\mathbf{I} - \mathbf{P})$. Entonces

$$\begin{aligned} \text{Cov}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)})) &= \frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{H}_i^T)}{(1 - p_{ii})} \text{Cov}(\text{vec } \mathbf{Y}) \frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{H}_i^T)^T}{(1 - p_{ii})} \\ &= \frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{H}_i^T)}{(1 - p_{ii})^2} (\Sigma \otimes \mathbf{I})(\mathbf{I}_p \otimes \mathbf{H}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1}) \\ &= \frac{\|\mathbf{H}_i\|^2 (\Sigma \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1})}{(1 - p_{ii})^2} \end{aligned} \quad (5.10)$$

Notando que

$$\begin{aligned} \|\mathbf{H}_i\|^2 &= e_i^{nT}(\mathbf{I} - \mathbf{P})(\mathbf{I} - \mathbf{P})e_i^n \\ &= e_i^{nT}(\mathbf{I} - \mathbf{P})e_i^n \\ &= e_i^{nT}e_i^n - e_i^{nT}\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}\mathbf{X}^T e_i^n \\ &= 1 - \mathbf{X}_i(\mathbf{X}^T \mathbf{X})^{-1}\mathbf{X}_i^T \\ &= 1 - p_{ii} \end{aligned} \quad (5.11)$$

Sustituyendo (5.11) en (5.10), se obtiene

$$\text{Cov}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)})) = \frac{(\Sigma \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1})}{(1 - p_{ii})} \quad (5.12)$$

Sea $S_1 = n\widehat{\Sigma}/(n - q)$ y observando que $E(S_1) = \Sigma$, ver Muirhead (pp. 84, 1982), entonces

$$\widehat{\text{Cov}}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)})) = \frac{(S_1 \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1})}{(1 - p_{ii})} \quad (5.13)$$

Sea $\mathbf{r}_i = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i$. Basándose en los siguientes resultados:

1. Para $a \in \mathbb{R}^n$, $a^- = a^T / \|a\|^2$,
2. Dada $A \in \mathbb{R}^{p \times q}$, entonces $(AA^T)^- = A^{T-} A^-$ con $A^{-1} = A^-$ si A es no singular,
3. Dadas las matrices A y B , $(A \otimes B)^- = A^- \otimes B^-$,

se obtiene

$$\begin{aligned} (\widehat{\text{Cov}}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)})))^- &= \widehat{\text{Cov}}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)}))^- \\ &= \left(\frac{(S_1 \otimes \mathbf{r}_i \mathbf{r}_i^T)}{(1 - p_{ii})} \right)^- \\ &= \frac{(1 - p_{ii})}{\|\mathbf{r}_i\|^4} (S_1^{-1} \otimes \mathbf{r}_i \mathbf{r}_i^T) \end{aligned}$$

Por lo tanto, la distancia de Cook modificada puede ser re-escrita como:

$$\begin{aligned} \mathcal{D}_m &= \text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)})^T \widehat{\text{Cov}}(\text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)}))^- \text{vec}(\widehat{\beta} - \widehat{\beta}_{(i)}) \\ &= \left(\frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{H}_i^T) \text{vec } \mathbf{Y}}{(1 - p_{ii})} \right)^T \frac{(1 - p_{ii})(S_1^{-1} \otimes \mathbf{r}_i \mathbf{r}_i^T)}{\|\mathbf{r}_i\|^4} \\ &= \left(\frac{(\mathbf{I}_p \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \mathbf{H}_i^T) \text{vec } \mathbf{Y}}{(1 - p_{ii})} \right) \\ &= \frac{(1 - p_{ii})^{-1}}{\|\mathbf{r}_i\|^4} \text{vec}^T \mathbf{Y} (S_1^{-1} \otimes \mathbf{H}_i \mathbf{r}_i^T \mathbf{r}_i \mathbf{r}_i^T \mathbf{r}_i \mathbf{H}_i^T) \text{vec } \mathbf{Y} \\ &= (1 - p_{ii})^{-1} \text{vec}^T \mathbf{Y} (S_1^{-1} \otimes \mathbf{H}_i \mathbf{H}_i^T) \text{vec } \mathbf{Y}. \end{aligned} \quad (5.14)$$

Alternativamente, ya que

$$\text{tr}(\mathbf{B} \mathbf{X}^T \mathbf{C} \mathbf{X} \mathbf{D}) = \text{vec}^T \mathbf{X} (\mathbf{B}^T \mathbf{D}^T \otimes \mathbf{C}) \text{vec } \mathbf{X} = \text{vec}^T \mathbf{X} (\mathbf{D} \mathbf{B} \otimes \mathbf{C}^T) \text{vec } \mathbf{X},$$

para matrices de órdenes adecuados, se puede escribir \mathcal{D}_m como

$$\mathcal{D}_m = (1 - \mathbf{p}_{ii})^{-1} \text{tr}(S_1^{-1} \mathbf{Y}^T \mathbf{H}_i \mathbf{H}_i^T \mathbf{Y}).$$

Por otro lado, puesto que $\hat{\boldsymbol{\epsilon}}_i^T = \mathbf{e}_i^{n^T} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}) = \mathbf{H}_i \mathbf{Y}$, entonces

$$\begin{aligned} \mathcal{D}_m &= (1 - \mathbf{p}_{ii})^{-1} \text{tr}(S_1^{-1} \hat{\boldsymbol{\epsilon}}_i \hat{\boldsymbol{\epsilon}}_i^T) \\ &= (1 - \mathbf{p}_{ii})^{-1} \text{tr}(\hat{\boldsymbol{\epsilon}}_i^T S_1^{-1} \hat{\boldsymbol{\epsilon}}_i) \\ &= (1 - \mathbf{p}_{ii})^{-1} \hat{\boldsymbol{\epsilon}}_i^T S_1^{-1} \hat{\boldsymbol{\epsilon}}_i \end{aligned}$$

Por lo que se tienen las siguientes expresiones como alternativa para el cuadrado de la distancia de Cook modificada:

$$\mathcal{D}_m = \begin{cases} \text{vec}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)})^T \widehat{\text{Cov}}(\text{vec}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}))^{-1} \text{vec}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}) \\ (1 - \mathbf{p}_{ii})^{-1} \text{vec}^T \mathbf{Y} (S_1^{-1} \otimes \mathbf{H}_i \mathbf{H}_i^T) \text{vec} \mathbf{Y} \\ (1 - \mathbf{p}_{ii})^{-1} \text{tr}(S_1^{-1} \mathbf{Y}^T \mathbf{H}_i \mathbf{P}_i^T \mathbf{Y}) \\ (1 - \mathbf{p}_{ii})^{-1} \hat{\boldsymbol{\epsilon}}_i^T S_1^{-1} \hat{\boldsymbol{\epsilon}}_i \end{cases} \quad (5.15)$$

Nota 2. De acuerdo con Chatterjee y Hadi (pp. 124, 1988), se puede reemplazar la matriz S_1 por otra obtenida usando la muestra reducida ($n - 1$), denotada por $S_1(i)$.

Nota 3. Cook (1977), Chatterjee y Hadi (pp. 117, 1988), Díaz-García et. al (2001), y muchos otros autores, utilizan la matriz de varianzas y covarianzas del $\text{vec}(\hat{\boldsymbol{\beta}})$ para construir la medida de distancia. La reformulación que se propone, está basada en el reemplazo de ésta matriz, por la matriz de varianzas y covarianzas del $\text{vec}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)})$. Se puede encontrar esta idea en Chatterjee y Hadi (pp. 150, 1988), para el caso univariado, pero para la evaluación de observaciones influyentes sobre un particular coeficiente de regresión, solamente se utiliza la varianza de un coeficiente, en vez de la varianza de la diferencia. El problema que se presenta, cuando esta idea es extendida al caso multivariado, es que tal matriz es singular,

por lo que se necesita considerar la inversa de Moore-Penrose para la matriz de varianzas y covarianzas del $\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})$.

Teorema 1. *Considerando el modelo de regresión normal multivariado dado en (5.1). Entonces, el cuadrado de la distancia de Cook modificada para detectar un "outlier", puede ser escrita como:*

$$\mathcal{D}_m = \begin{cases} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)})^T \left(\frac{(S_1 \otimes (X^T X)^{-1} X_i X_i^T (X^T X)^{-1})}{(1 - p_{ii})} \right) \text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) \\ (1 - p_{ii})^{-1} \text{vec}^T Y (S_1^{-1} \otimes H_i H_i^T) \text{vec} Y \\ (1 - p_{ii})^{-1} \text{tr} (S_1^{-1} Y^T H_i P_i^T Y) \\ (1 - p_{ii})^{-1} \hat{\epsilon}_i^T S_1^{-1} \hat{\epsilon}_i \end{cases} \quad (5.16)$$

Nota 4. En (5.16) es fácil ver que si se quiere implementar esta medida para todo el conjunto de datos, es suficiente ajustar el modelo una sola vez, y de la forma usual se puede construir la distancia modificada para cada punto. Notando que la expresión \mathcal{D}_m , en el caso normal univariado, coincide con el análisis de residuales estudentizados, ver Besley et al. (p. 201, 1980) y Chatterjee y Hadi (p. 78, 1988).

3. DISTANCIA MODIFICADA : MULTIPLES OBSERVACIONES

Sea $I = \{i_1, i_2, \dots, i_k\}$ un subconjunto de tamaño k de $\{1, 2, \dots, n\}$, de forma tal que $(n - k) \geq q$. Ahora, bajo el modelo (5.1), denote por $\mathbf{X}_{(I)}$, $\mathbf{Y}_{(I)}$ y $\hat{\epsilon}_{(I)}$, las matrices de regresión, de datos y de errores, respectivamente, después de eliminar las correspondientes observaciones de acuerdo con los subíndices en I . Sea $\hat{\beta}_{(I)}$, y $\hat{\Sigma}_{(I)}$ los correspondientes estimadores de máxima verosimilitud en el modelo

$$\mathbf{Y}_{(I)} = \mathbf{X}_{(I)} \beta_{(I)} + \epsilon_{(I)}, \quad \epsilon_{(I)} \sim \mathcal{N}_{(n-k) \times p}(0, \Sigma_{(I)} \otimes \mathbf{I}_n).$$

Basándose en el Lema 2, y usando procedimientos similares a los de la sección 2, es fácil verificar que

$$\hat{\beta} - \hat{\beta}_{(I)} = (X^T X)^{-1} X_I (I - P_I)^{-1} \hat{\epsilon}_I,$$

con $(I - P_I) = (I_k - X_I^T (X^T X)^{-1} X_I)$ y X_I es la matriz con las correspondientes filas de X de acuerdo con I . Observando que $\hat{\epsilon}_I = U_I^T \hat{\epsilon} = U_I^T (I - P) Y$, donde

$$U_I^T = \begin{pmatrix} e_{i_1}^{n^T} \\ e_{i_2}^{n^T} \\ \vdots \\ e_{i_k}^{n^T} \end{pmatrix}$$

Se obtiene

$$\text{vec}(\hat{\beta} - \hat{\beta}_{(I)}) = (I_p \otimes (X^T X)^{-1} X_I (I - P_I)^{-1} H_I) \text{vec} Y,$$

con $H_I = U_I^T (I - P)$.

De donde

$$\text{Cov}(\text{vec}(\hat{\beta} - \hat{\beta}_{(I)})) = (\Sigma \otimes (X^T X)^{-1} X_I (I - P_I)^{-1} X_I^T (X^T X)^{-1})$$

y

$$\widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(I)})) = (S_1 \otimes (X^T X)^{-1} X_I (I - P_I)^{-1} X_I^T (X^T X)^{-1})$$

Entonces se tiene,

Teorema 2. *Considerando el modelo de regresión normal multivariado dado en (5.1). Entonces, el cuadrado de la distancia de Cook modificada para detectar k observaciones influyentes, puede ser escrita como:*

$$\mathcal{D}_{m_I} = \begin{cases} \text{vec}(\hat{\beta} - \hat{\beta}_{(I)})^T (S_1 \otimes (X^T X)^{-1} X_I (I - P_I)^{-1} X_I^T (X^T X)^{-1})^{-1} \text{vec}(\hat{\beta} - \hat{\beta}_{(I)}) \\ \text{vec}^T Y (S_1^{-1} \otimes H_I^T (I - P_I)^{-1} H_I) \text{vec} Y \\ \text{tr}(S_1^{-1} Y^T H_I^T (I - P_I)^{-1} H_I Y) \\ \text{tr}(S_1^{-1} \hat{\epsilon}_I^T (I - P_I)^{-1} \hat{\epsilon}_I) \end{cases} \quad (5.17)$$

4. FUNCIONES DE DISTRIBUCIÓN ASOCIADAS CON LAS DISTANCIAS MODIFICADAS

La razón principal para estudiar las modificaciones dadas en las secciones 2 y 3 para el cuadrado de la distancia de Cook es que, en vez de utilizar una aproximación a la distribución \mathcal{F} , se derivará la distribución exacta para \mathcal{D}_m . Análogamente se encuentra la distribución exacta para \mathcal{D}_{m_i} , es decir, para el caso de la detección de varias observaciones influyentes, simultáneamente.

Se derivará la distribución del estadístico para el caso de una observación influyente a la vez y para el caso múltiples observaciones.

Teorema 3. *Bajo los supuestos del Teorema 1, se tiene*

$$\frac{(n - q - p + 1)}{p(n - q)} \mathcal{D}_m \sim \mathcal{F}_{p, (n - q - p + 1)} \quad (5.18)$$

donde $\mathcal{F}_{p, (n - q - p + 1)}$ denota una distribución \mathcal{F} centrada, con p y $(n - q - p + 1)$ grados de libertad (gl).

Prueba: Se sigue directamente del Teorema 5.2.2 en Anderson (p. 163, 1984). ■

Del Teorema 3, dando un nivel de significancia α , se puede escribir la siguiente regla de decisión:

Y_i , $i = 1, 2, \dots, n$, es un "outlier", si

$$\frac{(n - q - p + 1)}{p(n - q)} \mathcal{D}_m \geq \mathcal{F}_{\alpha; p, (n - q - p + 1)} \quad (5.19)$$

donde $\mathcal{F}_{\alpha; p, (n - q - p + 1)}$ es el correspondiente α -percentil superior de una distribución \mathcal{F} con p y $(n - q - p + 1)$ gl.

Nota 5. Para el caso univariado, $p = 1$, la regla de decisión se convierte en: Y_i , $i = 1, 2, \dots, n$, es un "outlier" si

$$\mathcal{D}_m \geq \mathcal{F}_{\alpha; 1, (n - q)} \quad (5.20)$$

donde $\mathcal{F}_{\alpha;1,(n-q)}$ es el α - percentil de una distribución \mathcal{F} con 1 y $(n - q)$ gl, o equivalentemente:

$$\mathcal{D}_m^{1/2} \geq t_{\alpha/2;(n-q)} \quad (5.21)$$

donde $t_{\alpha/2;(n-q)}$ es el $\alpha/2$ - percentil superior de una distribución t con $(n - q)$ gl.

Similarmente, para el caso de múltiples observaciones:

Teorema 4. *Bajo los supuestos del Teorema 2, se tiene que*

$$\mathcal{D}_{m_I} \sim \mathcal{LH}_{s,m,h} \quad (5.22)$$

donde $\mathcal{LH}_{s,m,h}$ denota la distribución centrada para el estadístico Lawley-Hotelling con parámetros $s = \min(p, k)$, $m = (|p - k| - 1)/2$ y $h = (n - q - p + 1)/2$.

Prueba: Se sigue directamente del Teorema 10.6.2, Corolario 10.6.3, en Muirhead [pp. 468-471 y p. 471, 1982]. ■

5. UNA APLICACION

Se ilustra la utilización de la prueba exacta dada en la sección 4, bajo una regresión simple y bajo una regresión múltiple multivariada.

El primer conjunto de datos fue presentado por Cook y Weisberg (pp. 204-207, 1994). Este es un conjunto de datos pequeño, con observaciones sobre 21 niños, dando la edad (AGE) en meses, en la que dijeron su primera palabra, y un $SCORE$, el cual es una medida del desarrollo del niño. En la figura 1(a), se presenta la gráfica ($AGE, SCORE$). Es claro que las observaciones 18, 1 y 17, muestran un comportamiento diferente. Si se considera un ajuste lineal usando cuadrados mínimos ordinarios, el $SCORE$ parece decrecer con AGE . El caso 18 parece tener un pobre ajuste hacia una tendencia lineal, comparado con los demás datos. Los casos 1 y 17, tienen relativamente valores grandes de AGE . La figura 1(b), muestra el gráfico Q-Q de los residuales provenientes del ajuste lineal, usando cuadrados mínimos or-

dinarios. Es claro que, solamente la observación 18 parece ser un candidato para ser considerado “outlier”, tal como es definido en Chatterjee y Hadi (pp.94-95, 1988). La figura 2 muestra la identificación y detección de “outliers” y puntos influyentes, usando diferentes técnicas. Este análisis enfatiza el hecho mencionado antes, es decir, los residuales estudentizados coinciden con la distancia modificada para el caso univariado. Además, los residuales estudentizados son la base para la distancia de Draper y John, y se discuten en Draper y Smith (pp. 169-175, 1981).

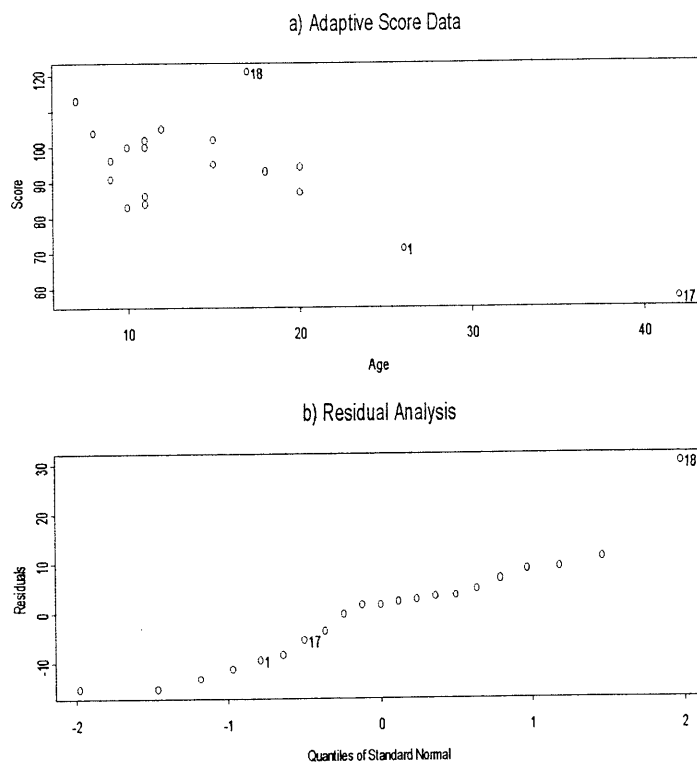


Figura 1. Datos originales para el ajuste del *SCORE* y un gráfico Q-Q para los residuales en la regresión simple de *SCORE* sobre *AGE*.

La figura 2(a), muestra cómo la distancia de Cook detecta la observación 17, la cual es un punto palanca, tal como se describe en Cook y Weisberg (1994). Considerando la figura 1(b) y de acuerdo con las figuras 2(b) y 2(c), la observación 18, es un candidato para ser considerado “outlier”, y las pruebas están a favor para declararla

como tal. Notando que el valor crítico en la figura 2(c), es el valor aproximado de una distribución \mathcal{F} , multiplicado por $s^2(i)$ para $i = 18$, de acuerdo al gráfico de la distancia original de Draper y John.

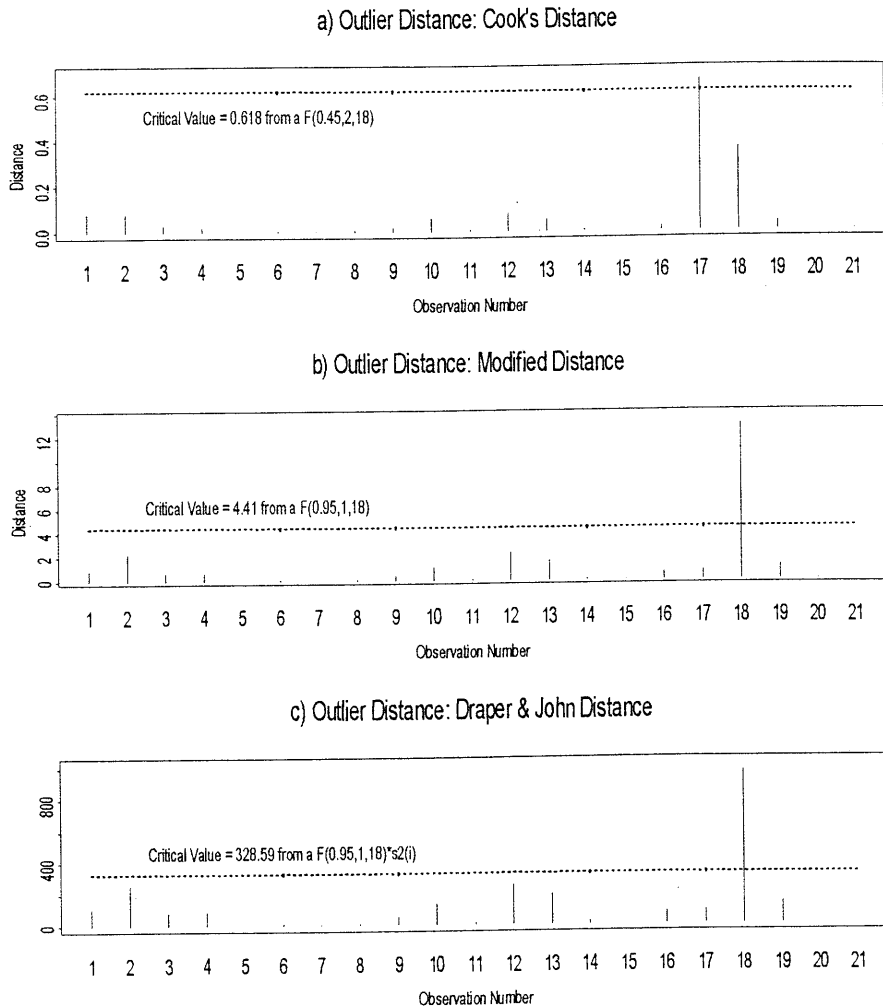


Figure 2. Identificación de la influencia y “outliers”, basados en a) Distancia de Cook, b) Distancia de Cook Modificada de acuerdo con (5.20), c) Distancia de Draper y John. $n = 21$, $q = 2$, $p = 1$ y $s^2(i)$. La varianza de los residuales sin la i -ésima observación se usa en todos los casos, de acuerdo con la Nota 2.

Para el ejemplo de la regresión múltiple multivariada, se generan 19 observaciones de un modelo $y = X\beta + \epsilon$, con errores normales, $p = 2$, y $q = 4$.

Se ajusta el modelo, la matriz de la suma de cuadrados de residuales, entonces aplicando el cuadrado de la distancia de Mahalanobis, como se muestra en Seber (pp. 152-153, 1984). La figura 3(a) presenta el cuadrado de la distancia de Mahalanobis y sugiere a las observaciones 10 y 11 como posibles “outliers”. Aplicando (5.22), con $k = 2$ y $Dm = 4.77$, comparado con un valor crítico de 3.015, el percentil es aproximado por una distribución \mathcal{F} , tal como sugiere Seber (pp. 38-39, 563-564, 1984). La prueba está a favor de considerar las observaciones 10 y 11 como “outliers”.

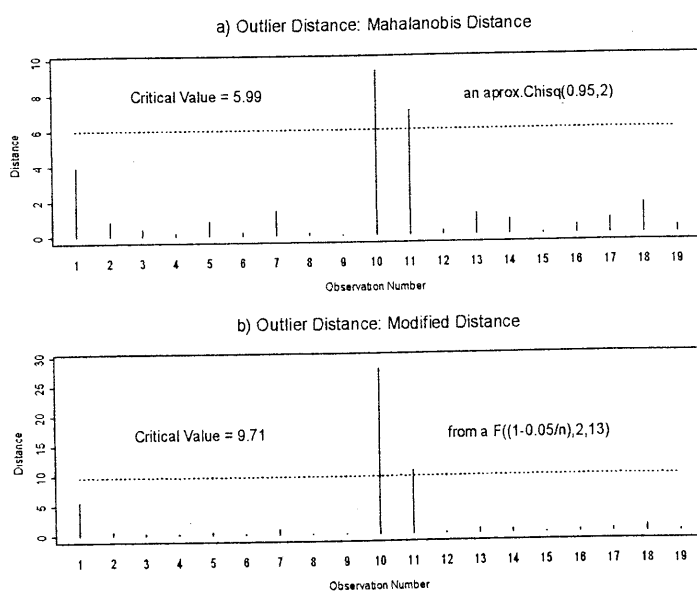


Figure 3. Identificación de “outliers” basada en la Distancia de Mahalanobis sobre la matriz de residuales, y detección de “outliers” basada en la Distancia de Cook Modificada dada en el Teorema 3.

La figura 3(b), muestra el mismo análisis tomando una observación a la vez. Se usa una prueba \mathcal{F} , basada en $(1 - \alpha/n)$ en vez de $(1 - \alpha)$, para obtener una prueba simultánea con un nivel mínimo α . En este caso, se obtiene la misma conclusión, como con la prueba basada en k observaciones.

Se recomienda el uso de la prueba basada en k observaciones, dada en (5.22), y se

seleccionan los k puntos, usando un método gráfico, como el mostrado en la figura 3(a).

REFERENCIAS

- Anderson, T. W. (1984). *An Introduction to Multivariate Statistical Analysis*, John Wiley & Sons, New York.
- Besley, D.A., Kuh, E., and Welsch, R.E. (1980). *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*, John Wiley & Sons, New York.
- Chatterjee, S., and Hadi, A. S. (1988). *Sensitivity Analysis in Linear Regression*, John Wiley & Sons. New York.
- Cook, R.D. (1977). "Detection of influential Observations in Linear Regression," *Technometrics*, 19, 15-18.
- Cook, R. D., and Weisberg, S. (1982). *Residual and Influence in Regression*, Chapman and Hall, London.
- Cook, R. D., and Weisberg, S. (1994). *An Introduction to Regression Graphics*, John Wiley & Sons. New York.
- Díaz-García, J.A., Galea, M., and Leiva- Sánchez, V. (2001). "Influence Diagnostics for Elliptical Regression Linear Models", Submitted for Publication.
- Draper, N., and Smith, H. (1981). *Applied Regression Analysis*, (2nd ed.), John Wiley & Sons, New York.
- Galea, M., Paula, G., and Bolfarine, H. (1997). "Local Influence in Elliptical Linear Regression Models", *The Statistician* 46, 71-79.
- Liu, S.Z. (2000). "On Local Influence for Elliptical Linear Models", *Statistical Papers*, 41, 211-224.
- Muirhead, R. J. (1982). *Aspects of Multivariate Statistical Theory*, John Wiley & Sons, New York.
- Rao, C.R. (1973). *Linear Statistical Inference and its Applications* (2nd ed.), John

Wiley & Sons, New York.

Rao, C. R. and Mitra, S. K.(1971). *Generalized Inverse of Matrices and its Applications* (2nd ed.), John Wiley & Sons, New York.

Seber, G.A.F.(1984). *Multivariate Observations* , John Wiley & Sons, New York.

LITERATURA CITADA

- Anderson, T. W. (1984). *An Introduction to Multivariate Statistical Analysis*, John Wiley & Sons, New York.
- Andrews, D.F., and Pregibon, D. (1978). "Finding Outliers That Matter," *Journal of the Royal Statistical Society*, (B), 40, 85- 93.
- Atkinson, A.C. (1981). "Two Graphical Displays for Outlying and Influential Observations in Regression," *Biometrika*, 68, 13- 20.
- Atkinson, A.C. (1982). "Regression Diagnostics, Transformations, and Constructed Variables (With Discussion)," *Journal of the Royal Statistical Society*, (B), 44, 1-36.
- Beckman, R.J., and Trussell, H.J. (1974). "The Distribution of an Arbitrary Studentized Residual and the Effects of Updating in Multiple Regression," *Journal of the American Statistical Association*, 69, 199-201.
- Behnken, D.W., and Draper, N.R. (1972). "Residuals and their Variance," *Technometrics*, 11, No. 1, 101-111.
- Besley, D.A., Kuh, E., and Welsch, R.E. (1980). *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*, John Wiley & Sons, New York.
- Bingham, C. (1977). "Some Identities Useful in the Analysis of Residuals from Linear Regression," Technical Report 300, School of Statistics, University of Minnesota.
- Chatterjee, S., and Hadi, A.S. (1986). "Influential Observation, High Leverage Points, and Outliers in Linear Regression," *Statistical Science*, 1, No. 3, 379-416.
- Chatterjee, S., and Hadi, A.S. (1988). *Sensitivity Analysis in Linear Regression*, John Wiley & Sons. New York.
- Chatterjee, S., and Price, B. (1977). *Regression Analysis by Example*, John Wiley & Sons, New York.
- Cook, R.D. (1977). "Detection of influential Observations in Linear Regression," *Technometrics*, 19, 15-18.
- Cook, R.D., and Weisberg, S. (1980). "Characterization of an Empirical Influence Function for Detecting Influential Cases in Regression," *Technometrics*, 22, 495-508.

- Cook, R.D., and Weisberg, S. (1982). *Residuals and Influence in Regression*, Chapman and Hall, London.
- Cook, R. D., and Weisberg, S. (1994). *An Introduction to Regression Graphics*, John Wiley & Sons. New York.
- Díaz-García, J.A., Galea, M., and Leiva- Sánchez, V. (2001). "Influence Diagnostics for Elliptical Regression Linear Models", Submitted for Publication.
- Draper, N.R., and John, J.A. (1981). "Influential Observations and Outliers in Regression," *Technometrics*, 23, 21-26.
- Draper, N.R., and Smith, H. (1981). *Applied Regression Analysis*, 2nd ed., John Wiley & Sons, New York.
- Ellenberg, J.H. (1973). "The Joint Distribution of the Standardized Least Squares Residuals From a General Linear Regression," *Journal of the American Statistical Association*, 68, 941-943.
- Ellenberg, J.H. (1976). "Testing for Single Outlier From a General Linear Regression Model," *Biometrics*, 32, 637-645.
- Galea, M., Paula, G., and Bolfarine, H. (1997). "Local Influence in Elliptical Linear Regression Models", *The Statistician* 46, 71-79.
- Hampel, F.R. (1968). "Contributions to the Theory of Robust Estimation," Ph.D. thesis, University of California, Berkeley.
- Hampel, F.R. (1974). "The Influence Curve and Its Role in Robust Estimation," *Journal of the American Statistical Association*, 62, 1179-1186.
- Henderson, H.V., and Searle, S.R. (1981). "On Deriving the Inverse of a Sum of Matrices," *SIAM Review*, 23, 53-60.
- Hoaglin, D.C., and Welsh, R.E. (1978). "The Hat Matrix in Regression and ANOVA," *The American Statistician*, 32, 17-22.
- Huber, P. (1977). *Robust Statistical Procedures*, No. 27, Regional Conference Series in Applied Mathematics, Philadelphia: SIAM.
- Huber, P. (1981). *Robust Statistics*, John Wiley & Sons, New York.
- Liu, S.Z. (2000). "On Local Influence for Elliptical Linear Models", *Statistical Papers*, 41, 211-224.
- Lund, R.E. (1975). "Tables for an Approximate Test for Outliers in Linear Models," *Technometrics*, 17, 473-476.

- Muirhead, R. J. (1982). *Aspects of Multivariate Statistical Theory*, John Wiley & Sons, New York.
- Pregibon, D. (1979). "Data Analytic Methods for Generalized Linear Models," Unpublished Ph.D. Thesis, University of Toronto.
- Pregibon, D. (1981). "Logistic Regression Diagnostics," *The Annals of Statistics*, 9, No. 4, 705-724.
- Rao, C.R. (1973). *Linear Statistical Inference and its Applications* (2nd ed.), John Wiley & Sons, New York.
- Rao, C. R. and Mitra, S. K. (1971). *Generalized Inverse of Matrices and its Applications* (2nd ed.), John Wiley & Sons, New York.
- Searle, S.R. (1982). *Matrix Algebra Useful for Statistics*, John Wiley & Sons, New York.
- Seber, G.A.F. (1984). *Multivariate Observations*, John Wiley & Sons, New York.
- Tietjen, G.L., Moore, R.H., and Beckman, R.J. (1973). "Testing for a Single Outlier in Simple Linear Regression," *Technometrics*, 15, 717-721.
- Velleman, P.F., and Welsch, R.E. (1981). "Efficient Computing of Regression Diagnostics," *The American Statistician*, 35, 234-242.
- Welsch, R.E., and Kuh, E. (1977). "Linear Regression Diagnostics," *Technical Report 923-77*, Sloan School of Management, Massachusetts Institute of Technology.
- Welsch, R.E., and Peters, S.C. (1978). "Finding Influential Subsets of Data in Regression Models," *Proceedings of the Eleventh Interface Symposium on Computer Science and Statistics*, (A.R. Gallant and T.M. Gerig, eds.), Raleigh: Institute of Statistics, North Carolina State University.