

2ej.
7



UNIVERSIDAD NACIONAL AUTÓNOMA
DE MÉXICO

FACULTAD DE CIENCIAS

OBSERVACIONES DISCREPANTES
EN REGRESION

T E S I S

Que para obtener el Título de

A C T U A R I O

p r e s e n t a

Jesús Bravo Segovia

México, D. F.

1 9 8 6



Universidad Nacional
Autónoma de México



UNAM – Dirección General de Bibliotecas Tesis Digitales Restricciones de uso

DERECHOS RESERVADOS © PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

Todo el material contenido en esta tesis está protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

CONTENIDO.

1. Introducción. 1

- 1.1. ¿Qué es una Observación Discrepante? 1
- 1.2. ¿Por qué estudiar Observaciones Discrepantes? 6
- 1.3. Causas de Observaciones Discrepantes. 9
- 1.4. Identificación y Acomodo. 11
- 1.5. Breve historia del estudio de Observaciones Discrepantes. 13

2. Enfoque Frequentista 20

- 2.1. El Modelo de Regresión Lineal. 20
 - 2.1.1. Matrices de Proyección. 23
 - 2.1.2. Residuales. 26
- 2.2. Modelos para el Manejo de Observaciones Discrepantes. 29
- 2.3. Métodos para el Manejo de Observaciones Discrepantes. 34
 - 2.3.1. Versión Etiquetada. 34
 - 2.3.2. Versión No Etiquetada. 36
 - 2.3.2.1. Método de Ellenberg. 36
 - 2.3.2.2. Método de Cook y Prescottt. 40
 - 2.3.2.3. Método de Doornbos. 43
 - 2.3.2.4. Método de Weisberg. 46
 - 2.3.2.5. Sobre la Sensibilidad de las Pruebas en Regresión cuando no hay

- normalidad en los errores. 49
- 2.3.2.6. Método de Gentleman y Wilk. 55
- 2.3.2.7. Método de Andrews y Pre-
gibon. 57
- 2.3.2.8. Método de Aitkin y Wilson. 61
- 2.3.2.9. Método de Marasinghe. 63
- 2.4. Comentarios y Conclusiones. 67
- 2.5. Ejemplos. 76

3. Enfoque Bayesiano. 101

3.1. Introducción. 101

3.2. Métodos Bayesianos para el Manejo de Observaciones discrepantes. 105

3.2.1. Método de Box-tiao. 106

3.2.2. Método de Abraham y Box. 112

3.2.3. Método de Guttman, Freeman y Dutter. 116

3.2.4. Método de Dutter y Guttman. 119

3.3. Comentarios. 121

3.4. Ejemplos. 123

3.5. Conclusiones 129

4. Enfoque Robusto. 130

4.1. Introducción 130

4.2. Estimadores Robustos en Regresión 134

4.2.1. Estimadores M. 134

4.2.2. Estimadores R. 141

4.3. Comentarios y Conclusiones 142

4.4. Ejemplos. 146

5. transformaciones. 152

5.1. Introducción 152

5.2. Métodos de transformación de Variables. 153

5.2.1. Método de Máxima Verosimilitud. 153

5.2.2. Método de Carroll 155

5.2.3. Método de Box y Tidwell. 156

5.2.4. Método de Andrews. 158

5.2.5. Método de Atkinson. 158

5.3. Comentarios. 165

5.4. Ejemplos. 167

3.5. Conclusiones. 168

Conclusiones Generales. 171

PROLOGO

El problema de observaciones discrepantes es uno de los mas viejos en estadística. Durante los últimos 150 años, el interés en estas observaciones ha tenido sus altas y sus bajas, algunas veces es área activa de investigación para que años mas tarde sufra de un abandono relativo.

En años recientes se han resuelto un gran número de problemas en la teoría de observaciones discrepantes y se han descubierto otros. Sin embargo, la mayoría de los resultados se encuentran dispersos en revistas especializadas.

El objetivo de este trabajo es presentar una pequeña revisión de la filosofía, teoría y métodos inherentes a tres enfoques de la teoría estadística: Clásico o frecuentista, Bayesiano y Robusto.

Esta tesis se presenta en 5 capítulos. En el primer capítulo, se analizan y proponen algunas definiciones del término "observación discrepante", también se determinan las causas por las cuales éstas se pueden presentar y las razones de su estudio. Por otro lado y dependiendo del interés sobre dichas observaciones, se clasifican a los métodos para el manejo de éstas. Por último, en este capítulo se presenta una reseña histórica de las observaciones discrepantes.

En el segundo capítulo se efectúa el estudio del enfoque frecuentista. En éste, se asientan muchos elementos necesarios (modelos, elementos de diagnóstico, etc) para la mayoría de los métodos. La parte medular de este capítulo, la sección

2.3 referente a los métodos para el manejo de observaciones discrepantes. La intención de presentar en ese orden los métodos es la de diferenciar cuales de estos sirven para una observación discrepante ($K=1$) y cuales para mas de una ($K>1$) haciendo énfasis en si K es conocida o desconocida. Finalmente, en las últimas secciones se efectúan comentarios generales de los métodos, se aplican algunos de éstos a dos problemas "clásicos" y por último se presentan conclusiones referentes a dichos métodos.

Por lo que respecta a los capítulos 3, 4 y 5, podríamos decir que estructuralmente son iguales al capítulo 2, es decir, inicialmente se presenta la filosofía general del enfoque para que en secciones posteriores se presenten los métodos, comentarios, ejemplos y conclusiones. Naturalmente la diferencia radica en el aspecto filosófico con que cada enfoque trata el tema de observaciones discrepantes.

Por último, al final de este trabajo se presentan una serie de conclusiones cuya intención es distinguir que enfoque usar y cual de los métodos presentados podría ser el adecuado en una situación específica.

NOTACION

En este trabajo el modelo de regresión lineal se mejorará regularmente en su forma matricial:

$$\underline{y} = X\underline{\beta} + \underline{\varepsilon}$$

donde $\underline{y} = (y_1, y_2, \dots, y_n)$ es el vector de observaciones de orden $n \times 1$. X es la matriz de dimensiones $n \times p$ y rango $r(X) = p$, cuyos renglones son de la forma $(x_{i1}, x_{i2}, \dots, x_{ip})$ $i = 1, 2, \dots, n$ con $x_{ii} = 1$, $i = 1, 2, \dots, n$ si la recta, plano o hiperplano (en general) de regresión no pasa por el origen. $\underline{\beta} = (\beta_1, \beta_2, \dots, \beta_p)$ es un vector de orden $p \times 1$ de parámetros desconocidos. $\underline{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$ es un vector aleatorio de $n \times 1$ y además se supone $E(\underline{\varepsilon}) = 0$; $\text{Var}(\underline{\varepsilon}) = E[(\underline{\varepsilon} - E(\underline{\varepsilon}))(\underline{\varepsilon} - E(\underline{\varepsilon}))^T] = \sigma^2 I$. I la matriz identidad de dimensión n y adicionalmente $\underline{\varepsilon} \sim N(0, \sigma^2 I)$

Se hará referencia a este modelo como el "modelo de regresión lineal con las hipótesis usuales".

Al estimador de $\underline{\beta}$ y $\underline{\varepsilon}$ se les denotará por $\hat{\underline{\beta}}$ y $\hat{\underline{\varepsilon}}$ respectivamente, y al de σ^2 por S^2 , donde:

$$\begin{aligned}\hat{\underline{\beta}} &= (X^T X)^{-1} X^T \underline{y}, \\ \hat{\underline{\varepsilon}} &= \underline{y} - \hat{\underline{y}} = (I - X(X^T X)^{-1} X^T) \underline{y}, \underline{y} = (I - V) \underline{y} = M \underline{y}; V, y M \text{ matrices de proyección.} \\ S^2 &= \frac{\underline{y}^T M \underline{y}}{n - p},\end{aligned}$$

$k =$ número de observaciones discrepantes, $i = \{i_1, i_2, \dots, i_k / 1 \leq i_1 < i_2 < \dots < i_k \leq n\}$. En este trabajo en ocasiones i , estará como subíndice de alguno de los elementos del modelo de regresión definido anteriormente, esto puede ser de dos formas: i únicamente como subíndice, en el caso de

X_i , indicará a los renglones que están señalados por los elementos de i , es decir X_i representa a la matriz de $k \times n$ cuyos renglones están indicados por los elementos de i , en el caso de vectores, por ejemplo Y_i , representa al vector de $k \times 1$ observaciones cuyas componentes son las k observaciones sospechosas de ser discrepantes, interpretaciones similares tendrán los vectores θ_i , ξ_i . En el caso que el conjunto indicador este entre paréntesis, significará que se eliminan a los elementos, renglones o componentes según sea el caso, que están indicados por los elementos de i . Por ejemplo $X_{(i)}$ representa a la matriz $(n-k) \times n$ que resulta de eliminar a los renglones indicados por i ; $Y_{(i)}$ representa al vector de $(n-k) \times 1$ de observaciones resultantes de eliminar a los indicados por los elementos de i , interpretaciones similares tendrán $\theta_{(i)}$, $\xi_{(i)}$.

Capítulo 1.

INTRODUCCION.

1.1. ¿Que es una observacion discrepante?

El concepto de observacion discrepante ha interesado a los investigadores desde los primeros intentos de la interpretacion de datos. Aun antes del desarrollo del metodo estadistico habia argumentos vehementes sobre si se debian rechazar observaciones de un conjunto de datos, argumentando que no son representativos. Las actitudes variaban de un extremo a otro: Desde el punto de vista que nunca debemos manchar la "santidad" de los datos con "sus pechas para desechar sus propiedades, hasta un pragmatismo "si dudas, rechazalo".

Hoy contamos con puntos de vista más sofisticados, está reconocida una gran variedad de proposiciones en el manejo de observaciones discrepantes, se han propuesto modelos para la generacion de estas observaciones y está disponible una gran variedad de técnicas estadísticas para el procesamiento de datos. A pesar de que existen algunos libros en el tema de observaciones discrepantes, la mayoría de material se encuentra aun disperso en la literatura.

Los primeros trabajos en observaciones discrepantes se caracterizan principalmente por la falta de atención a la modelación del mecanismo de generacion de observaciones discrepantes, sin apoyo de un estudio en términos de las propiedades estadísticas de los modelos propuestos.

Ninguna observación puede garantizarse sea totalmente una manifestación del fenómeno bajo estudio. Un evento

Con una posibilidad en un millón podría ocurrir con la frecuencia apropiada, no importando que tan sorprendidos estemos de que esto nos ocurra. Intuitivamente la confiabilidad de una observación se refleja por su relación con las otras observaciones que se obtienen bajo condiciones similares. Observaciones que en opinión del investigador se apartan del resto de los datos se llaman: observaciones discrepantes, disidentes, espurias, contaminantes o valores sorprendentes, contaminantes, malos o sucios; solo para mencionar algunos términos que se han usado en épocas pasadas.

El tema de observaciones discrepantes es viejo e indudablemente data de los primeros intentos para obtener conclusiones de un conjunto de datos. Comentarios hechos por Bernoulli (1777) acerca de observaciones astronómicas indican que la práctica de descartar observaciones discrepantes era común 200 años atrás. Los primeros intentos por desarrollar métodos estadísticos objetivos fueron reportados por 1850. Hoy en día deberíamos esperar métodos extendidamente aceptados para tratar y manejar observaciones discrepantes, además de una definición universalmente aceptada. Este no es el caso de acuerdo a nuestra interpretación con la literatura. Aunque mucho se ha escrito, la noción de observaciones discrepantes parece vaga hoy como lo fue 200 años atrás. Por ejemplo Edgeworth (1887) escribe:

- Las observaciones discrepantes se pueden definir como aquellas que presentan la apariencia de diferir con respecto a su ley de frecuencias con las otras observaciones con las que están combinadas.—

82 años después Grubbs (1969) establece que:

- Una observación discrepante, es una que parece desviarse marcadamente de los otros miembros de la muestra en que ocurre. —

en los libros recientes sucede algo similar como a continuación se evidencia. Barnett y Lewis (1973) definen observación discrepante de la siguiente manera:

- Definimos observación discrepante en un conjunto de datos, como una observación (o subconjunto de observaciones) que parecen ser inconsistentes con el resto del conjunto de datos. —

mientras que Hawkins (1980) escribe:

- La definición intuitiva de una observación discrepante será: "Una observación que se desvía mucho de las otras observaciones como para despertar sospechas de que fue generada por un mecanismo diferente" —

Estas afirmaciones ilustran que una observación discrepante es un concepto subjetivo después de que los datos se han obtenido. Históricamente, se emplearon métodos "objetivos" para el tratamiento de observaciones discrepantes solamente después de que estos fueron identificados por medio de una inspección visual de los datos.

Para fijar ideas, supongamos que las siguientes 10

observaciones, de Johnson y Hunt (1979) se cree son realizaciones independientes de una población normal:

-1.64, -1.33, -1.10, -0.57, -0.27

1.04, 1.56, 1.34, 2.04, 4.99

¿Cuál de estas observaciones es discrepante? La respuesta será reservada hasta que se construya una prueba o alguna otra medida objetiva en forma natural, pero la decisión para ampliar un criterio objetivo se basa frecuentemente en reacciones subjetivas iniciales de los datos. En este caso, la observación 4.99 parece demasiado grande y también se puede formular algún cuestionamiento para la observación 2.04.

Collet y Lewis (1976) reportan los resultados de un experimento para investigar la naturaleza subjetiva de la decisión para etiquetar una observación como discrepante. Ellos concluyen que la inclinación individual depende del método de presentación (azar, gráfica, orden) experiencia y la escala de los datos; las observaciones extremas tienden a parecer más discrepantes cuando la escala se incrementa.

En conjunto de datos más estructurados, como, cuando el interés se centra en la forma en que observaciones de una variable de interés principal que varía con respecto a otras variables (análisis de regresión) o con respecto al tiempo (series de tiempo) u otras situaciones, debemos esperar encontrar de vez en cuando, datos no representativos en forma de observaciones discrepantes. Aquí también es importante, como en una muestra simple univariada, poder interpretar

tar y acomodar observaciones discrepantes mediante técnicas estadísticas apropiadamente diseñadas. En este caso las observaciones discrepantes pueden ser de interés intrínseco o pueden ser indicadores de especificaciones inapropiadas de la estructura del error o del modelo básico, con implicaciones de consecuencia para el uso de procedimientos apropiados de inferencia. Con datos más estructurados surgen dos complicaciones: observaciones sospechosas tienden a ser menos aparentes, más ocultos en el cuerpo de los datos y además los métodos para su rechazo o su acomodo están menos desarrollados. En regresión por ejemplo, una inspección visual de los datos difícilmente reflejará en un caso es una observación discrepante. La evidencia para esto se refleja regularmente cuando los parámetros son estimados y los residuales tabulados. Debido a que los residuales no son independientes, se hace difícil juzgar a las observaciones discrepantes, además, también se complica la metodología. Mucho de esta metodología se basa en modelos de observaciones discrepantes y son necesarias algunas decisiones acerca de la naturaleza y frecuencia antes de escoger un criterio. La subjetividad se introduce en la etapa del modelaje.

Debido al énfasis en el modelo en años recientes, "observación discrepante" ahora parece usarse por varios autores para indicar cualquier observación que no provenga de la población seleccionada, aunque en publicaciones recientes falta aun una definición formal de observación discrepante. De las 10 observaciones dadas anteriormente 1.04, 2.04 y 4.99 fueron generados de una

población normal con $\mu=3$ y $\sigma^2=1$, mientras que las 7 restantes fueron generadas por una población normal con $\mu=0$ y $\sigma^2=1$. Algunos autores dirían que esta muestra tiene una observación discrepante si utiliza alguna de las definiciones expuestas anteriormente, mientras que otros dirían que este conjunto tiene 3 observaciones discrepantes por pertenecer éstas a otra distribución.

En un esfuerzo por evitar alguna ambigüedad adoptaremos las siguientes definiciones en el resto de este trabajo:

Observación anormal.- cualquier observación que parezca sorprendente o discrepante para el investigador.

Contaminante.- cualquier observación que no sea una realización de la población supuesta.

Observación discrepante.- para referirse a contaminante o anormal.

1.2 ¿Por qué estudiar observaciones discrepantes?

Una respuesta está indicada en la sección anterior. Parece claro que se necesita un método adecuado para manejar tales observaciones o al menos un fuerte entendimiento de los méritos relativos a los métodos disponibles. Aunque esta es la razón históricamente dominante para el estudio de observaciones discrepantes, existen varias razones adicionales e igualmente importantes que se han enfatizado en años recientes.

Interés especial.- puede suceder que la atención principal sobre la(s) observación(es) discrepante(s) sea aun mayor que, por ejemplo la estimación de un parámetro pobla-

cional. En tales situaciones el problema estadístico involucra el hacer inferencias válidas acerca de la(s) observación(es) en cuestión.

Barnett (1978) describe un caso legal que es interesante: 349 días después de que el señor Hadlum marchó al extranjero en servicio militar, la señora Hadlum tuvo un parto. El señor Hadlum juzgó a la observación de 349 días discrepante cuando la comparó con el promedio de gestación de 280 días y por lo tanto hace una petición de divorcio. En este caso la garantía no es cuando descartar la observación discrepante o el efecto en la estimación de algún parámetro, mas bien es como juzgar el peso de la evidencia contra la hipótesis de que la observación discrepante, aunque extrema, es una realización válida de la distribución tiempo de gestación.

Detección de un fenómeno alternativo específico. — Algunos países industrializados para localizar satélites extraños a los suyos han medido los cambios en el nivel de radiación en las regiones en que se sospecha la presencia de estos satélites. Los niveles de radiación que fueran altos relativos al promedio, se toman como indicadores de un fenómeno alternativo. La metodología estadística básica consiste en la aplicación de esquemas de detección de observaciones discrepantes a gran escala. No hay interés en estimar el nivel promedio de radiación en el ambiente, excepto para proporcionar una guía para juzgar observaciones particulares como discrepantes.

Indicadores de diagnóstico. — En análisis más complica-

das, donde una empresa cree que el modelo es erróneo, la presencia de observaciones discrepantes es frecuente indicación de la inconsistencia en el modelo, en los datos o en ambos. Las observaciones discrepantes, pueden conformarse por ejemplo, cuando las respuestas son transformadas a una escala logarítmica. también las observaciones discrepantes pueden ser indicadoras de observaciones altamente influyentes en el modelo o reflejar errores en los renglones de la matriz (en el caso de regresión), no aditividad o heteroscedasticidad.

Los análisis de datos deben incluir la aplicación de técnicas de diagnóstico que puedan suministrar información sobre lo apropiado del análisis y la precisión de las conclusiones. Naturalmente muchas de las técnicas disponibles se basan en la detección de observaciones discrepantes.

Acomodo e influencia.- El estudio de la naturaleza y frecuencia de observaciones discrepantes en un problema particular pueden llevar a modelos y métodos de estimación que acomoden observaciones discrepantes, mejorándose de esta manera la inferencia. Como ejemplo de estos métodos podemos considerar a los métodos robustos.

Las observaciones discrepantes no necesariamente son observaciones influyentes, es decir, los resultados de un análisis pueden permanecer esencialmente sin cambio cuando una observación discrepante se remueve del resto de la muestra. Es útil considerar a las observaciones

influyentes como un tipo especial de observación discrepante. El estudio de las observaciones influyentes hace posible, como veremos más adelante, el obtener un mejor entendimiento de la estructura de un problema, en particular, el problema de regresión. Por ejemplo, tales estudios pueden evidenciar deficiencias inherentes en los datos y llevarnos a experimentaciones adicionales.

1.3 Causas de observaciones discrepantes.

Varios autores (Anscombe 1960; Grubbs 1969) han intentado clasificar las formas diferentes en las que surgen las observaciones discrepantes. En el manejo de observaciones pueden encontrarse diferentes fuentes de variabilidad.

Variabilidad inherente.- Esta es la expresión de la forma en que las observaciones discrepantes varían intrínsecamente en la población; tal variación es una característica incontrolable y natural de la población.

Error de Medición.- Frecuentemente en la toma de mediciones en miembros de una población bajo estudio, los fallos en el instrumental de medición que se usa, introducen un grado adicional de variabilidad, además de redondeos de los valores obtenidos y errores en la recolección entre otros, contribuyen a la variabilidad. En esta situación pueden tomarse algunas precauciones para reducir tal variabilidad.

Error de Ejecución.- también puede surgir una fuente de variabilidad en la colección imperfecta de los datos. Po-

demostramos escoger inadvertidamente una muestra sesgada o que incluye elementos que no son verdaderamente representativos de la población. Se pueden tomar ciertas precauciones para reducir esta variabilidad.

Beckman - Cook (1983), también clasifican las causas por las que surgen las observaciones discrepantes en 3 categorías: a) modelo global de inconsistencia; b) modelo local de inconsistencia; c) variabilidad natural. Pensando en las dos primeras categorías, las observaciones discrepantes deben juzgarse con algún modelo implícito o explícito en mente. Desde un punto de vista estadístico, es útil pensar que la aparición de una observación discrepante se debe a la incapacidad del modelo para proporcionar una explicación estadística o un ajuste adecuado.

Modelo Global de Inconsistencia. - Son todas aquellas causas que llevan al reemplazo del modelo en turno por un modelo revisado o nuevo para toda la muestra. Esta categoría incluye causas como variables de respuesta en una escala equivocada, también nos puede llevar a una transformación o al reemplazo del modelo por un modelo mixto.

Modelo Local de Inconsistencia. - Son aquellas causas que se aplican solamente a las observaciones discrepantes y no al modelo como un total. Tales causas requieren que las observaciones discrepantes sean tratadas con individualidad, ya que regularmente se desconoce algún modelo que incorpore a estas observaciones. Algunos ejemplos incluyen mediciones aisladas, errores en la recolección y observaciones -

altamente influyentes en regresión debido a puntos remotos en el espacio de factores.

Variabilidad Natural = Finalmente una observación discrepante puede ser el resultado de una variación natural en lugar de una inconsistencia del modelo. Neyman y Scott (1971), Green (1974, 1976) introducen el concepto de familias \mathcal{F} de distribuciones "propensas a discrepantes" y "resistentes a discrepantes" en un esfuerzo por caracterizar la extensión para que se produzcan observaciones discrepantes en forma natural. Sean $Y_{(1)}, Y_{(2)}, \dots, Y_{(n)}$ las estadísticas de orden de una muestra cuya distribución $F \in \mathcal{F}$ entonces se dice que \mathcal{F} es "resistente discrepante (k, n) " si

$$P(k, n) = \sup \Pr [Y_{(n)} - Y_{(n-1)} > k (Y_{(n-1)} - Y_{(1)}) / F] < 1$$
de otra manera se dice que \mathcal{F} es "propensa a discrepante (k, n) ". Si $P(k, n) = 1$ para toda $k > 0$ y $n \geq 3$, se dice que \mathcal{F} es completamente propensa discrepante."

1.4. Identificación y Acomodo

Seguendo a Barnett (1978) y a Barnett y Lewis (1978), se distinguen dos grandes métodos para el manejo de las posibles observaciones discrepantes. El primer método es simplemente para la identificación de observaciones discrepantes para un estudio adicional, la identificación de una de estas observaciones nos puede llevar a: a) su subsecuente rechazo; b) información nueva importante contenida en variables concomitantes que de otra manera pasaría inadvertida; c) de su incorporación a través de una revisión del modelo (global de inconsistencia) o del método de

estimación; d) reconocer la inconsistencia inherente en los datos y entonces el reconocimiento también en experimentos posteriores. Posteriormente, son útiles los diagnósticos que se basan en la identificación de observaciones discrepantes principalmente, en la construcción iterativa de un modelo donde sea necesaria una evaluación crítica en las diferentes etapas de desarrollo. En Box (1979, 1980), Cook y Weisberg (1982), podemos encontrar discusiones más generales.

El segundo método es para el acomodo de posibles observaciones discrepantes, para modificaciones deseables del modelo y/o método de análisis. Por ejemplo, podemos usar modelos mixtos para acomodar cierto tipo de contaminantes y usar estimadores M para obtener alguna protección contra observaciones discrepantes cuando el modelo mixto sea simétrico.

De estos dos métodos generales, juzgamos que el de identificación es más importante. Los métodos de acomodo de observaciones discrepantes tienden a requerir mucha información, acerca del proceso de generación de observaciones discrepantes o están diseñados para ser inmunes a la presencia de estas observaciones. Los tipos anteriores a menudo oscurecen la información, mientras que los primeros requieren que mucha de la información esencial se haya obtenido ya usualmente por medio del proceso de identificación. Aunque las filosofías fundamentales de identificación y acomodo parecen distintas, estas regularmente son confundidas, ya que un método de acomodo puede producir un método de identificación

como un producto y viceversa.

Un punto que merece un énfasis especial es que los métodos de identificación y acomodo que se basan en modelos alternativos pueden producir resultados que son específicos al modelo en cuestión; entonces la determinación de una observación discrepante puede depender críticamente de nuestro entendimiento del problema. En particular, es útil distinguir entre situaciones donde se suponen observaciones discrepantes que llevan información acerca de los parámetros en el modelo básico y situaciones donde estas observaciones no llevan información relevante. El modelo de la varianza inflada es un ejemplo de la primera clase, mientras que el modelo de la media desplazada es un ejemplo de la segunda.

1.5 Breve historia del estudio de observaciones discrepantes.

La primera discusión acerca de observaciones discrepantes que se ha podido situar es la realizada por Bernoulli (1777). Bernoulli cuestionaba la suposición de que los errores tenían distribución idéntica y criticaba la práctica general de descartar observaciones como si éstas constituyeran la muestra completa. Un ejemplo de esta situación la observamos cuando Boscovich (1755), para determinar la elipticidad de la tierra en base a un promedio de 10 "seconds pendulums" a diferentes latitudes, descarta las dos observaciones más extremas y calcula la media a partir de la muestra reducida. Durante este periodo y hasta la mitad del siglo XIX, el punto principal de discusión en la literatura, fue cuando se jus.

tificata el rechazo de observaciones sospechosas.

La primer persona que propuso un criterio objetivo para el rechazo de observaciones discrepantes fue Pierce (1859). El procedimiento de Pierce supone que el modelo que genera los datos es una mezcla de distribuciones con parámetro p desconocido. Haciendo la peor suposición acerca de p y que la distribución principal es $N(0, \sigma^2)$ con σ conocida, si el valor absoluto de cualquier residual se excede en $c\sigma$, éste debe ser rechazado, donde c se determina de la siguiente manera: Si $\Phi(\cdot)$ denota la distribución de probabilidad acumulada de la normal estandar y s denota la desviación estandar de la muestra completa y s' la de la muestra que resulta de eliminar k observaciones discrepantes. Si $n' = n - k$, entonces c , es tal que satisface la ecuación $(s'/s)^n \{ \exp \frac{1}{2} (c^2 - 1) \} 2\Phi(-c) \}^n = [k^n n'^n / n^n]^{n/k}$. tres años después Gould (1862) proporciona tablas para la implantación del criterio de Pierce.

Sin embargo en este tiempo hay muchas investigaciones que plantean objeciones teóricas y prácticas al criterio de Pierce, el astrónomo Airy comenta al respecto:

— Pienso que toda la teoría es defectuosa en su fundamentación e ilusoria en sus resultados; no se puede obtener una regla para la exclusión de observaciones, por un proceso fundamentado en una consideración de la discrepancia de esas observaciones. —

Chauvenet (1866) propone un método diferente: Notando que el valor esperado del número de observaciones en una muestra de tamaño n de una $N(0, \sigma^2)$ excediéndose a σ es $n \Phi(-c)$, él propuso que cualquier residual que en valor absoluto exceda a σ debe ser rechazado, donde c satisface

$$n \Phi(-c) = 0.5$$

entonces el criterio de Chauvenet es tal, que rechaza en promedio media observación de datos buenos por muestra, sin tomar en cuenta el valor de n . Este aspecto de la prueba de Chauvenet es especialmente significativo ya que se aboca a la tasa de rechazo del experimento y no a la proporción de datos rechazados, así cuando n crece, c también lo hace y la proporción de observaciones rechazadas decrece.

Un tercer criterio para el rechazo completo de observaciones discrepantes fue creado por Stone (1868). Él encuentra al criterio de Pierce fastidioso, y citando a Airy expresa sus dudas acerca de la validez matemática de dicho procedimiento. Él afirma que cualquier observador tiene una probabilidad de $1/n$ de cometer un error, y propone que c sea tal que

$$n \Phi(-c) = 0.5$$

La aparente similitud entre el criterio de Chauvenet y el de Pierce es engañosa. La regla de Stone rechaza una proporción fija de observaciones buenas y así el número de observaciones discrepantes crece en proporción con

el tamaño de la muestra.

Glaisher (1872-73) fue quizás la primera persona en escribir un procedimiento ponderado. Su método supone que las x_i ($i=1, 2, \dots, n$) tienen distribución $N(\mu, \sigma_i^2)$, con la media μ común a todas ellas, pero con varianzas σ_i^2 desiguales y desconocidas. Se propone un valor inicial del estimador de la media y de los valores de σ_i^2 , con éstos modifica el valor inicial del estimador de la media y de las σ_i^2 y así sucesivamente. Específicamente se calcularon medias ponderadas $m_{(1)}, m_{(2)}, \dots$, siendo $m_{(r)}$ la media muestral ordinaria y los pesos $w_{r,i}$ para la r -ésima media ponderada $m_{(r)} = \sum_i w_{r,i} x_i$ definidos recursivamente por

$$w_{r,i} \propto \text{EXP} \{-2(w_{r-1,1} + \dots + w_{r-1,n})(x_i - m_{(r-1)})^2\}$$

Stone (1873), unos meses después critica este método y propone un procedimiento alternativo que se basa en la maximización de la verosimilitud, que es proporcional a

$$\prod_i [(1/\sigma_i) \text{EXP} \{-\{(x_i - \mu)^2 / 2\sigma_i^2\}],$$

con respecto a μ y a todas las σ_i . Esto lleva a una media ponderada $\tilde{\mu}$ dado por la ecuación de grado $n-1$

$$\sum_{i=1}^n (x_i - \tilde{\mu})^3 = 0$$

Este mismo método fue publicado independientemente por Edgeworth (1883), aunque en 1887 reconoce que Stone fue el primero en enunciarlo.

Otro método ponderado propuesto en esta época fue el de Newcomb (1886), su procedimiento supone que las n observaciones provienen de una mezcla de r distribuciones normales y desarrolla un estimador final de μ , construido como una media ponderada de r medias ponderadas diferentes de las x_i . Tal vez de mayor interés es la referencia de Newcomb en esa publicación al valor "maligno", que resulta ser el error cuadrático medio de un estimador, y es uno de los primeros usos del concepto de función de pérdida. Una de las cosas notables es que los pesos que obtuvo corresponden, en estimación robusta, con una función $\psi = \max[-c, \min(x, c)]$ una de las funciones favoritas de Huber.

Edgeworth (1887) propone 3 modelos para observaciones discrepantes e ilustra cada uno, usando el método de Monte Carlo. Entonces estos modelos fueron usados como una base para un estudio comparativo de los modelos de Airy, Chauvenet, Glaisher, Newcomb, Pierce y Stone, obteniendo los siguientes resultados: a) El "módulo de descuido" de Stone es legítimo; b) Se refiere a la sugerencia de Airy como el "No-método"; c) Reconoce los modelos fundamentales del trabajo de Stone y Glaisher; d) Indica que el criterio de Chauvenet y Pierce es pesimista; e) Expresa una gran consideración por el trabajo de Newcomb.

Daniell (1920) en su trabajo "observaciones ponderadas de acuerdo a su orden" escribe; además de la media que resulta después de la eliminación de observaciones discrepantes, debemos buscar otros estimadores en los que los pesos son múltiplos de funciones que dependen del orden de las observaciones.

Irwin (1925) sienta las bases de una prueba de varias - observaciones discrepantes en la muestra. Considérese un conjunto de datos x_i ($i=1, 2, \dots, n$) con $x_i \sim \text{iid } N(\mu, \sigma^2)$, σ conocida. Irwin consideró las estadísticas de orden $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ de la muestra y derivó la distribución de la k -ésima "brecha" $(X_{(n-k+1)} - X_{(n-k)})/\sigma$, también derivó expresiones que - computacionalmente son tratables para los momentos de esta k -ésima brecha y demostró que para $k=1, 2$, la distribución de la brecha se puede aproximar por medio de las distribuciones normales cuyas medias y varianzas él tabuló para varios valores de n .

Goodwin (1913) propuso que en lugar del criterio $(x_i - \bar{x})/s$ usado por los primeros autores para ciertas pruebas de observaciones discrepantes, uno debería sustituir a \bar{x} y s por la media y desviación estándar que resultarían de eliminar la observación discrepante considerada en la muestra. Algunos años tuvieron que pasar antes de que se demostrara que esta estadística de prueba es una función - monótona de $(x_i - \bar{x})/s$. Este descubrimiento hecho por Thompson (1935) lleva a la derivación de la distribución nula - de un residual studentizado arbitrario. De este resultado, él deduce un procedimiento en contra de observaciones discrepantes y rechaza una proporción fija de datos buenos, produciendo también una tabla de valores críticos.

En 1936 Pearson y Chandra Sekar, usando los resultados de Thompson y algunos lemas matemáticos sobre residuales studentizados, demuestran que para alguna c suficientemente grande, el evento $|z_i - \bar{x}|/s > c$ implica

plica que para toda $j \neq i$, $|x_j - \bar{x}|/s < c$. Entonces si c se escoge de tal manera que sea el $\alpha/2n$ percentil de la distribución de $|x_i - \bar{x}|/s$, entonces una prueba de observaciones discrepantes rechazando cualquier x_i que satisfaga $|x_i - \bar{x}|/s > c$ tiene un nivel de significancia experimental de α . Esta publicación también manifiesta la existencia del efecto de enmascaramiento que se analizará más adelante.

Winsor (1941) propone un procedimiento que hoy en día lleva su nombre y que consiste en que, la observación sospechosa no se rechaza totalmente de la muestra, sino que se sustituya por el valor más cercano a esta observación que no sea considerada discrepante.

Capítulo 2.

ENFOQUE FRECUENTISTA

Este enfoque recibe también el nombre de clásico, estándar, ortodoxo o de teoría de muestreo y se distingue por dos características principales. Considera como única forma de información relevante a aquella proporcionada en forma de datos muestrales y adopta como base para la evaluación y construcción de procedimientos estadísticos, el concepto de probabilidad visto como frecuencia, es decir, para cualquier evento, la probabilidad de éste será una medida de su "regularidad estadística".

2.1 El modelo de regresión lineal.

En innumerables procesos se tiene que el comportamiento de la variable bajo estudio es influenciada por las variables que intervienen en el proceso de obtención de las observaciones. Si se cree que la respuesta podría ser aproximada a partir de una relación funcional, en la cual se consideren únicamente aquellas variables que se juzgue más importantes, el modelo que se propone para esta situación es

$$Y_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \epsilon_i \quad i=1,2,\dots,n \quad (2.1.1)$$

donde Y_i es la i -ésima observación de la variable bajo estudio y es X_{ij} el valor de la j -ésima variable cuando se obtiene la i -ésima observación, ϵ_i $i=1,2,\dots,n$ es el i -ésimo parámetro (desconocido) y

ε_i es el error asociado a la i -ésima observación - causado por aquellos elementos que se considero no eran determinantes para explicar el comportamiento de la variable bajo estudio. La variable Y es - llamada comúnmente variable dependiente mientras que X_1, \dots, X_p reciben el nombre de variables respues- ta o explicativas.

Con el objeto de complementar la definición del modelo, se hacen algunas suposiciones acerca de la va- riable aleatoria ε_i :

$$E(\varepsilon_i) = 0; \quad V(\varepsilon_i) = \sigma^2, \quad \text{cov}(\varepsilon_i; \varepsilon_j) = 0 \quad \forall i \neq j \quad (2.1.2)$$

La notación matricial del modelo (2.1) es

$$\underline{Y} = \underline{X} \underline{\beta} + \underline{\varepsilon} \quad (2.1.3)$$

donde \underline{Y} es un vector n -dimensional cuyas compo- nentes son las y_i , $\underline{\beta}$ es un vector p -dimensio- nal de parámetros desconocidos, \underline{X} es una matriz de $n \times p$ con rango $r(\underline{X}) = p$ cuyo i -ésimo ren- glón es de la forma $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ y $\underline{\varepsilon}$ es un vector aleatorio n dimensional cuyas componen- tes satisfacen (2.1.2).

El problema posterior a la definición del mo- delo, es encontrar el valor de $\underline{\beta}$, que minimice de alguna manera $\underline{\varepsilon}$. El criterio que mayor dife- rencia ha tenido es el de Mínimos Cuadrados, cu- ya metodología esta fuera del objetivo de es-

te trabajo. La solución que se obtiene para $\hat{\beta}$ por medio de este método es

$$\hat{\beta} = (X^T X)^{-1} X^T Y. \quad (2.1.4)$$

Nótese que la matriz X puede pensarse como una transformación lineal, es decir $X: \mathbb{R}^p \rightarrow \mathbb{R}^n$. La imagen de X , $\text{Im}(X)$, es un subespacio de \mathbb{R}^n cuya dimensión es precisamente $r(X) = p$. También es importante hacer notar que la ecuación $Y = X\beta$ no tiene solución, es decir, Y no pertenece a $\text{Im}(X)$.

Sea

$$\begin{aligned} \hat{Y} &= X \hat{\beta} \\ &= X(X^T X)^{-1} X^T Y, \end{aligned} \quad (2.1.5)$$

el vector de valores esperados que resultan de "ajustar" el modelo (2.1.3) por mínimos cuadrados. En la figura (2.1.1) encontramos la interpretación geométrica de la solución al problema de mínimos cuadrados.

Como podemos observar \hat{Y} es la proyección ortogonal de Y en $\text{Im}(X)$. La matriz

$$V = X(X^T X)^{-1} X^T, \quad (2.1.6)$$

cuyas propiedades se estudian más adelante, es la matriz que proyecta a Y en la imagen de X . también debemos notar que la matriz

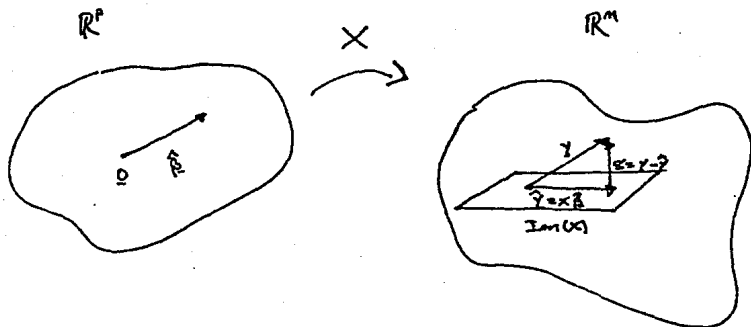


figura 2.1.1. solución al problema de mínimos cuadrados.

$$M = I - V \quad (2.1.7)$$

es la matriz que proyecta ortogonalmente a y en el subespacio de \mathbb{R}^n que es ortogonal a $\text{Im}(X)$, esta proyección ortogonal es precisamente el vector de residuales e , cuyas componentes son

$$y_i - \hat{y}_i = y_i - x_i^t \hat{\beta} ; i = 1, 2, \dots, n \quad (2.1.8)$$

donde x_i^t es el i -ésimo renglón de la matriz X .

2.1.1 La Matriz de proyección

(i) Propiedades. Algunas características están presentes en la sección anterior, además por ser una

matriz de proyección; V es simétrica e idempotente y por lo tanto

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

$$v_{ii} = 0 \Rightarrow v_{ij} = 0$$

$$v_{ii} = 1 \Rightarrow v_{ij} = 0; i \neq j;$$

de lo anterior se desprende que $0 \leq v_{ij} \leq 1$, también se sabe que los valores característicos de una matriz de proyección son cero o uno y que el número de éstos, es igual al rango de la matriz.

En este caso $r(V) = r(X) = p$ y por lo tanto la traza de V es p , es decir $\sum_{i=1}^n v_{ii} = p$. Entonces el promedio de los elementos de la diagonal es p/n .

Una regla empírica señala que un valor v_{ii} es grande si $v_{ii} > 2 \cdot p/n$.

(ii) Importancia de V en el análisis de datos. El vector de residuales ordinarios (2.1.8) está dado por

$$\begin{aligned} e &= Y - \hat{Y} \\ &= (I - V)Y, \end{aligned}$$

la relación entre \underline{e} y $\underline{\varepsilon}$ es

$$\begin{aligned} \underline{e} &= (I - V)(X\beta + \underline{\varepsilon}) \\ &= (I - V)\underline{\varepsilon} \end{aligned}$$

o en forma escalar

$$z_i = \varepsilon_i - \sum_{j=1}^n v_{ij} \varepsilon_j.$$

Esta igualdad, muestra claramente que la relación entre z y ε depende solo a través de V . Si las v_{ij} 's son suficientemente pequeñas, z servirá como un sustituto razonable de ε , de otra manera la utilidad de z , estará limitada. Si $\varepsilon \sim N(0, \sigma^2 I)$ entonces $z \sim N(0, \sigma^2(1-V))$. Casos alejados del espacio de factores tendrán valores de v_{ii} relativamente grandes. Ya que $V(\hat{y}_i) = v_{ii} \sigma^2$ y $\text{Var}(z_i) = (1 - v_{ii}) \sigma^2$, los valores ajustados en puntos remotos tendrán varianzas relativamente grandes y los residuales correspondientes tendrán varianzas relativamente pequeñas. Huber llama a $1/v_{ii}$ "el número efectivo de casos que determinan \hat{y}_i ". Davies y Hutton (1975) señala que si $\max v_{ii}$ no es considerablemente menor que 1, entonces probablemente una observación discrepante será desapercibida cuando se examinen los residuales. Box y Draper (1975) sugieren que para diseñar un experimento que sea insensible a observaciones discrepantes, las v_{ii} 's deben de ser pequeñas y aproximadamente iguales.

La importancia de las v_{ii} 's no se limita al análisis de mínimos cuadrados. Huber (1977) menciona que regresión robusta no funciona bien si $\max v_{ii} = 1$.

Hoaglin y Welsh (1978) sugieren el uso directo de las v_{ii} como un diagnóstico para identificar "puntos altamente influyentes". El motivo detrás de esta sugerencia, se basa en la representación

$$\hat{y}_i = v_{ii} y_i + \sum_{j \neq i} v_{ij} y_j.$$

El valor ajustado \hat{y}_i estara' dominado por $v_{ii} y_i$ si v_{ii} es grande con respecto a los términos restantes. Ellos interpretan a v_{ii} como la cantidad de influencia ejercida en \hat{y}_i por y_i . Es claro, sin embargo, que para cualquier $v_{ii} > 0$, \hat{y}_i estara' dominado por $v_{ii} y_i$ si y_i es suficientemente diferente de los otros elementos de Y (Esto es, una observación discrepante).

2.1.2 Residuales

Cuando el modelo ajustado es incorrecto, cambiara' la distribución de los errores no observables $\underline{\epsilon}$ y por lo tanto de los residuales \underline{e} . El objetivo en el estudio de los residuales \underline{e} es detectar cualquier su posición incorrecta en \underline{e} . Desafortunadamente, la correspondencia entre \underline{e} y $\underline{\epsilon}$ es menos que perfecta. En algunos problemas, las fallas en el modelo no se transmiten a \underline{e} , en otros, los síntomas observados pueden atribuirse a más de una suposición incorrecta.

Para procedimientos de diagnóstico, se han sugerido varias transformaciones de los residuales ordinarios e para superar algunas deficiencias, por ejemplo obtener residuales que sean independientes de los parámetros de escala o para obtener alguna estructura de covarianza específica.

Los residuales ordinarios, tienen una distribución -

que depende de la escala ya que la varianza de cada e_i es una función de s^2 y v_{ii} . Para muchos procedimientos de diagnóstico, es útil definir la versión "studentizada" de los residuales que no dependen de alguna de estas cantidades.

Cuando se usa mínimos cuadrados en regresión, los residuales studentizados internos se definen por

$$r_i = \frac{e_i}{s(1-v_{ii})^{1/2}} \quad i=1, 2, \dots, n \quad (2.1.9)$$

donde $s^2 = \frac{\sum_{i=1}^n e_i^2}{n-p}$. Nos reservaremos el término de residual studentizado al referirnos a (2.1.9).

Los residuales studentizados se usan para reemplazar a los residuales ordinarios en procedimientos gráficos, tal como el diagrama contra valores ajustados (Anscombe y Tukey 1963). También se usan como base para construir la mayoría de las estadísticas que se proponen en este trabajo.

Los residuales studentizados externamente requieren un estimador de s^2 que sea independiente de e_i . Bajo la suposición de normalidad de los errores, sea $S_{(ii)}^2$ el cuadrado medio residual calculado sin el i -ésimo caso, entonces se propone

$$\begin{aligned} S_{(ii)}^2 &= \frac{(n-p)S^2 - e_i^2 / 1-v_{ii}}{n-p-1} \\ &= S^2 \left(\frac{n-p-v_{ii}}{n-p-1} \right). \end{aligned} \quad (2.1.10)$$

Bajo el supuesto de normalidad, $S_{e(i)}^2$ y $e_{(i)}$ son independientes y el residual studentizado externamente se define por

$$t_i = \frac{e_i}{S_{e(i)}(1 - r_i^2)^{1/2}} \quad i=1, 2, \dots, n. \quad (2.1.11)$$

La distribución de t_i es t-student con $n-p-1$ grados de libertad (Beckman y Trussell 1974). La relación entre t_i y r_i se encuentra sustituyendo (2.1.10) en (2.1.11)

$$t_i = r_i \left(\frac{n-p-1}{n-p-r_i^2} \right)^{1/2}, \quad (2.1.12)$$

que demuestra que t_i^2 es una transformación monótona de r_i^2 .

2.2 Modelos para el manejo de observaciones discrepantes.

En un marco de pruebas de hipótesis, Barnett (1973), Barnett y Lewis (1978) plantean los modelos probables que pueden generar observaciones discrepantes: Para una prueba de observaciones discrepantes, la hipótesis nula H expresará el modelo básico de probabilidad que genera todos los datos sin considerar a las observaciones discrepantes, la hipótesis alternativa \bar{H} expresa una manera en la que puede modificarse el modelo para incorporar o explicar las observaciones discrepantes. La hipótesis nula es solamente una afirmación del modelo básico de probabilidad.

Si en base a una prueba de discrepancia juzgamos que existe una o más observaciones discrepantes, rechazamos implícitamente la hipótesis nula en favor de una hipótesis alternativa que naturalmente debemos conocer. Consideremos algunas formas posibles de hipótesis alternativas para pruebas de discrepancia que nos permita considerar que posibilidades existen en el modelo alternativo (generador de observaciones discrepantes) en el caso de una de estas observaciones.

(i) Alternativa Determinística

Este primer tipo cubre los casos de observaciones discrepantes causados por errores gruesos de medición o recolección. Esta hipótesis está totalmente especificada por el

conjunto de datos incluyendo a las observaciones discrepantes. Así si los datos x_1, x_2, \dots, x_n contienen una observación (x_j) que surge de un error de medición o recolección, rechazamos inmediatamente el modelo básico F en favor de un modelo alternativo que diga que todas las x_i ($i \neq j$) provienen de F mientras que x_j es diferente y debe de rechazarse, corregirse o repetirse. No se necesita prueba alguna.

(ii) Alternativa Inherente.

En este caso debemos tomar en cuenta la posibilidad de que las observaciones discrepantes aparecieron en los datos como resultado de la variabilidad inherente

$$H: x_i \in F, (i=1, 2, \dots, n) \text{ VS } \bar{H}: x_i \in G \neq F (i=1, 2, \dots, n).$$

(iii) Alternativa Mixta

En vez de suponer que las observaciones discrepantes reflejan un grado o forma inesperada de variabilidad inherente, debemos admitir la posibilidad de "errores de ejecución" permitiendo la "contaminación" de la muestra por miembros de una población diferente a la representada por el modelo básico

$$H: x_i \in F (i=1, 2, \dots, n) \text{ VS } \bar{H}: x_i \in (1-\alpha)F + \alpha G (0 < \alpha < 1) (i=1, 2, \dots, n).$$

(iv) Alternativa trasladada.

En esta alternativa, se establece que $n-k$ observaciones independientes surgen de un modelo inicial F con parámetros de posición y escala μ y σ^2 respectivamente, mientras que un número "pequeño" k de observaciones independientes surgen de una versión modificada de F en la que μ y σ^2 están desplazados en su valor (μ en cualquier dirección y σ^2 regularmente un valor mayor al inicial)

$$H: x_i \in F \quad (i=1, 2, \dots, n) \quad \text{vs} \quad H: \begin{cases} x_i \in F \quad (i=1, 2, \dots, j-1, j+1, \dots, n) \\ x_j \in F' \end{cases}$$

donde F' es la versión modificada de F . Por ejemplo si

$$F \sim (\mu, \sigma^2),$$

F' puede tomar alguna de estas formas

$$a) F' \sim (\mu + \lambda, \sigma^2); \quad b) F' \sim (\mu, \delta^2 \sigma^2),$$

formas conocidas como modelo de la media "trasladada" (shift) y modelo de la varianza "inflada" respectivamente.

En la revisión que se ha hecho de los métodos, se nota que la literatura sobre observaciones discrepantes en análisis de regresión y en general en modelos lineales, se encuentra dominada por el modelo de la media trasladada

$$y = x\beta + D\lambda + \underline{\epsilon}, \quad (2.2.1)$$

que se detalla a continuación en términos generales.

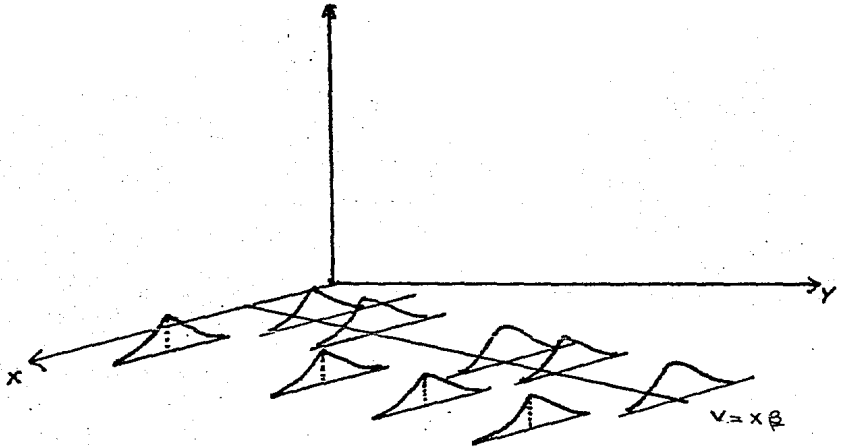
Sea $i = \{i_1, i_2, \dots, i_k / 1 \leq i_1 < i_2 < \dots < i_k\}$ y u_i un vector de $n \times 1$ cuya l -ésima componente ($l \in i$) vale uno y las restantes valen cero. La única diferencia del modelo (2.2.1) con respecto al modelo

$$\underline{Y} = X\underline{\beta} + \underline{\epsilon}, \quad (2.2.2)$$

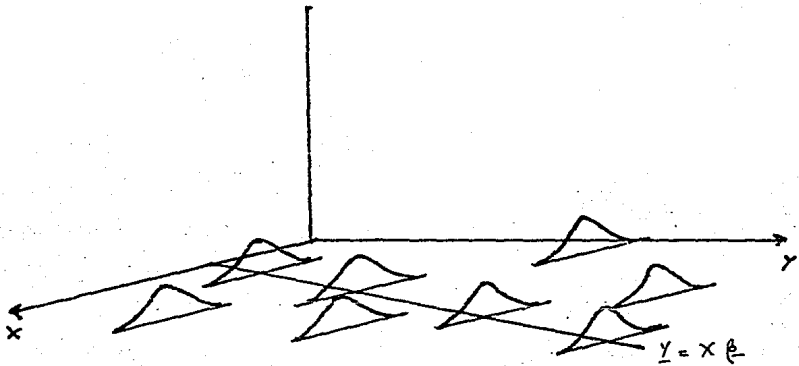
que es el modelo de regresión con las suposiciones de costumbre, es el término $D\underline{\lambda}$ donde D es una matriz de $n \times k$ cuyas columnas son los vectores $u_{i_1}, u_{i_2}, \dots, u_{i_k}$ y $\underline{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_k)$ es un vector de k parámetros desconocidos. Debe notarse que el planteamiento que se está haciendo con el modelo (2.2.1) es que existen k observaciones que no provienen del modelo (2.2.2) y que provienen de un modelo alternativo, a saber el modelo (2.2.1) que especifica que estas k observaciones surgen de una o varias poblaciones, esto dependerá según las λ_i ($i=1, 2, \dots, k$) sean iguales o no, es decir, si $i \in i$ entonces $E(Y_i) = \underline{x}_i^t \underline{\beta} + \lambda_i$ con \underline{x}_i^t el i -ésimo renglón de X . Para una descripción intuitiva, ver figura (2.2.1).

Srikantan (1961) y Ferguson (1961) fueron los primeros en establecer el modelo de la media trasladada como una base para el manejo de observaciones discrepantes, aunque también este enfoque está implícito en el trabajo de Daniel (1960).

El modelo de la media trasladada se considera más apropiado para la identificación de observaciones



- a) Las k observaciones discrepantes provienen de una población alternativa.



- b) Las k observaciones discrepantes provienen de poblaciones diferentes.

Figura 2.2.1

discrepantes que para su acomodo.

Beckman y Cook (1983) distinguen dos situaciones del modelo (2.2.1): La versión etiquetada en la que se supone conocido el conjunto i es decir, se conocen cuántas observaciones discrepantes hay y sus posiciones; y la versión no etiquetada en la que i es desconocido.

2.3 Métodos para el manejo de observaciones discrepantes.

2.3.1 Versión Etiquetada.

Como esta versión es un modelo lineal estándar, el análisis puede llevarse a cabo usando los métodos usuales y los resultados pueden expresarse en términos del modelo (2.2.2) John (1978), Beckman - Cook (1983).

Así, el estimador mínimos cuadrados de λ es:

$$\hat{\lambda} = (I - V_i) e_i \quad (2.3.1)$$

y el estimador mínimos cuadrados de β al suprimir las k observaciones sospechosas resulta

$$\hat{\beta}_{(i)} = (X_{(i)}^t X_{(i)})^{-1} X_{(i)}^t Y_{(i)}. \quad (2.3.2)$$

La reducción en la suma de cuadrados debido al ajuste de λ en el modelo (2.2.1) proporciona una medida del efecto de suponer que en realidad las k observaciones son discrepantes. Esta suma de cuadrados es

$$Q_k(i) = \hat{e}_i^+ \hat{e}_i - \hat{e}_i^+ \hat{e}_i \quad (2.3.3)$$

donde \hat{e} es el vector de residuales que resulta del modelo (2.2.2) y \hat{e}^+ el vector de "residuales revisados" de Gentleman y Wilk (1975):

$$\hat{e}_i^+ = y_i - x_i \hat{\beta}_{(i)} - \hat{\lambda}_i$$

donde si $i \in I$, entonces $\hat{\lambda}_i = y_i - x_i \hat{\beta}_{(i)}$, caso contrario $\hat{\lambda}_i = 0$. John (1978), Draper y John (1981) hacen notar que el modelo (2.2.3) es equivalente al usado cuando se estiman k "valores faltantes" en análisis de covarianza con covariables indicadoras (dummy). Entonces es posible considerar a las observaciones discrepantes como observaciones faltantes y reemplazarlas por sus valores faltantes estimados. Los resultados son los mismos.

La estadística que se propone en Cook y Weisberg (1982), Beckman y Cook (1983) suponiendo normalidad para probar la hipótesis $\lambda = 0$ es

$$F_k(i) = \frac{(n-p-k) Q_k(i)}{(n-p-r_i^2)}$$

Si $k=1$ Pekman y Trussell (1979) demuestran que

$$F_1(i) = \frac{(n-p-1) r_i^2}{(n-p-r_i^2)} \quad (2.3.4)$$

es una función monótona creciente en r_i^2 cuya distribución es t-student con $n-p-1$ grados de libertad y

donde ξ es el i -ésimo residual studentizado.

2.3.2. Versión No Etiquetada.

El análisis de la versión no etiquetada es mucho más complicada ya que tanto i como k son desconocidas. Para una k fija, una de las recomendaciones más comunes es la de examinar todos los $\binom{n}{k}$ subconjuntos y estimar i , usando el conjunto de índices i que maximiza $Q_k(i)$. Desgraciadamente el manejo de esta idea implica un alto costo en el aspecto computacional, razón por la cual aun no se tiene una técnica unificada a pesar de que se han sugerido muchas versiones en un esfuerzo por evitar estos costosos cálculos.

Debido a que en los últimos años ha surgido una gran cantidad de técnicas, en este trabajo se tratará de presentar algunas de las más importantes.

2.3.2.1 Método de Ellenberg

Considere el modelo de regresión lineal con las suposiciones de costumbre y defínase el i -ésimo residual estandarizado por

$$\xi_i = \frac{r_i}{(n-p)^{1/2}}, \quad (2.3.5)$$

con r_i el i -ésimo residual studentizado interno definido en (2.1.9).

La regla de decisión es:

Rechazar H_0 : "no hay observaciones discrepantes" si $\tau^* \geq c_\alpha$ donde c_α es un valor crítico, α es el nivel de significancia deseado y

$$\tau^* = \max_i |\tau_i|. \quad (2.3.6)$$

Como la distribución de τ^* (aun bajo la suposición de normalidad) resulta intratable, se propone el uso de la primera desigualdad de Bonferroni para encontrar cotas para la distribución de τ^* ; así para cualquier constante crítica c_α

$$\begin{aligned} \sum_{i=1}^n P(|\tau_i| > c_\alpha) - \sum_{i>j} P(|\tau_i| > c_\alpha, |\tau_j| > c_\alpha) \\ \leq P(\tau^* > c_\alpha) \leq \sum_{i=1}^n P(|\tau_i| > c_\alpha) \end{aligned} \quad (2.3.7)$$

o equivalentemente

$$\begin{aligned} nP(|\tau_i| > c_\alpha) - \sum_{i>j} P(|\tau_i| > c_\alpha, |\tau_j| > c_\alpha) \\ \leq P(\tau^* > c_\alpha) \leq nP(|\tau_i| > c_\alpha) \end{aligned} \quad (2.3.8)$$

donde los τ_i se distribuyen idénticamente con

$$f(\tau_i) = \frac{\Gamma(\gamma+1)}{\Gamma(\frac{1}{2})\Gamma(\frac{1}{2})} (1 - \tau_i^2)^{\gamma-\frac{1}{2}},$$

para $-1 < \tau_i < 1$, $i = 1, 2, \dots, n$

y

$$f(r_i, r_j) = \frac{\gamma}{\pi(1-e_{ij}^2)^{1/2}} \left[1 - \frac{1}{(1-e_{ij}^2)} (r_i^2 - 2\rho_{ij} r_i r_j + r_j^2) \right]^{\gamma-1} \quad (2.3.9)$$

para $r_i^2 - 2\rho_{ij} r_i r_j + r_j^2 \leq (1-e_{ij}^2)$,

con $\rho_{ij} = -\frac{1-v_{ij}}{[(1-v_{ii})(1-v_{jj})]^{1/2}}$ y $\gamma = \frac{(n-p-2)}{2}$. La función

(2.3.9) es una versión estandarizada de la función student invertida y puede relacionarse implícitamente con la función t-student bivariada.

La solución para la cota inferior de la desigualdad (2.3.8) implica el cálculo de todas las $\binom{n}{2}$ funciones del tipo (2.3.9), dicha carga computacional se puede reducir usando el resultado de Srikantan (1961)

$$P(|r_i| > c_\alpha, |r_j| > c_\alpha) = 0 \text{ si } c_\alpha \geq \left(\frac{1+|\rho_{ij}|}{2} \right)^{1/2} \quad (2.3.10)$$

y determinando c_α , con α dado, tal que en la cota superior de (2.3.8)

$$P(|r_i| > c_\alpha) = \frac{\alpha}{n},$$

entonces todas las parejas (r_i, r_j) donde $c_\alpha \geq \left(\frac{1+|\rho_{ij}|}{2} \right)^{1/2}$ no contribuyen a la cota inferior. Debe de mencionarse que c_α varía inversamente con α , así que teóricamente puede hacerse decrecer α en cualquier problema hasta que c_α alcance las condiciones de la desigualdad (2.3.10) para todas las $\binom{n}{2} \rho_{ij}$, de ese modo pueden producirse re

sultados exactos.

Para aquellos casos donde la distribución bivariada no sea cero, usar la expresión

$$P(|r_i| > c_{\alpha}, |r_j| > c_{\alpha}) = 2P(r_i > c_{\alpha}, r_j > c_{\alpha}) + 2P(r_i > c_{\alpha}, r_j < -c_{\alpha})$$

y hacer la transformación

$$t_i = (r_i - r_j) / [2(1 - e_{ij})]^{1/2},$$

$$t_j = (r_i + r_j) / [2(1 - e_{ij})]^{1/2},$$

entonces la distribución conjunta de t_i, t_j es

$$g(t_i, t_j) = \frac{\delta}{\pi} (1 - t_i^2 - t_j^2)^{\delta-1} \quad \text{para } t_i^2 + t_j^2 \leq 1.$$

Desgraciadamente este procedimiento sólo funciona bien para valores pequeños de $\delta = (n-p-2)/2$, ya que para valores grandes pueden surgir problemas de cálculo y de precisión.

Lund (1975) basándose en la cota superior de la desigualdad (2.3.8) propone tablas para una prueba aproximada en modelos lineales para $\alpha = 0.10, 0.05, 0.01$. La estadística básica de prueba es

$$r^* = \max_i |r_i|, \quad (2.3.11)$$

que es diferente a (2.3.6), pero

$$Pr(\tau^* > \tau_0) = P(\tau^* > \tau_0) \leq nP(|\tau_i| > \tau_0) = 2n \int_{\tau_0}^{\infty} f(\tau) d\tau,$$

con $\tau_0 = \frac{\tau_{\alpha}}{(n-p)}$. Así, para obtener valores críticos τ_0 , es necesario resolver para τ_0 $2n \int_{\tau_0}^{\infty} f(\tau) d\tau = \alpha$

$$\Rightarrow \int_{\tau_0}^{\infty} f(\tau) d\tau = \frac{\alpha}{2n}$$

o equivalentemente

$$\int_0^{\tau_0} f(\tau) d\tau = \frac{1}{2} \left(1 - \frac{\alpha}{n}\right).$$

Nota: Para encontrar los valores críticos c_{α} para la prueba de Eilonberg por medio de las tablas de Lund, es necesario dividir τ_0 entre $(n-p)^{1/2}$

2.3.2.2. Método de Cook y Prescott

Considérese el modelo de regresión lineal con las exposiciones de costumbre y recuerdense las definiciones de r_i y τ_i dadas en (2.1.9) (2.3.5) respectivamente. En este método se parte nuevamente de la desigualdad (2.3.7) con $c_{\alpha} = d$ y $\alpha_i = Pr(|\tau_i| > d)$, $\alpha_{ij} = Pr(|\tau_i| > d, |\tau_j| > d, i \neq j)$, entonces esta desigualdad se transforma a

$$\sum_i \alpha_i - \sum_{i < j} \alpha_{ij} \leq Pr(\tau^* > d) = \sum_i \alpha_i. \quad (2.3.12)$$

Ya que bajo la hipótesis de que no hay observaciones discrepantes, los valores $|\tau_i|$ se distribuyen idéntica-

mente y entonces la cota superior de (2.3.12) puede expresarse convenientemente en términos de una variable aleatoria F cuya distribución es $F(1, n-p-1)$ es decir,

$$\sum_i \alpha_i = n \Pr(\tau_i^2 > d^2) = n \Pr[F > d^2 (n-p-1) / (1-d^2)].$$

La evaluación de la cota inferior de (2.3.12) es, como se vio anteriormente, más difícil ya que las α_{ij} 's deben determinarse por integración (Ellenberg 1976). En su lugar, se usará una aproximación para las α_{ij} 's como siguen. Ya que la distribución de (τ_i, τ_j) es simétrica.

$$\alpha_{ij} = 2 \Pr(\tau_i > d, \tau_j > d) + 2 \Pr(\tau_i > d, -\tau_j > d) \\ \Pr(\tau_i > d, \pm \tau_j > d) \approx \Pr(\tau_i \pm \tau_j > 2d),$$

entonces

$$\alpha_{ij} \leq 2 \Pr(\tau_i + \tau_j > 2d) + 2 \Pr(\tau_i - \tau_j > 2d) \\ = \Pr[(\tau_i + \tau_j)^2 > 4d^2] + \Pr[(\tau_i - \tau_j)^2 > 4d^2] \\ = \beta_{ij}^+ + \beta_{ij}^-; \text{ naturalmente } \beta_{ij}^+ = \Pr[(\tau_i + \tau_j)^2 > 4d^2] \\ \text{y } \beta_{ij}^- = \Pr[(\tau_i - \tau_j)^2 > 4d^2].$$

Este resultado en combinación con (2.3.12) dan como resultado

$$\alpha - \sum_{i,j} (\beta_{ij}^+ + \beta_{ij}^-) \leq \Pr(t^* > d) \leq \alpha, \quad (2.3.13)$$

$$\alpha = \sum_i \alpha_i.$$

Entonces la evaluación de la cota inferior en la desigualdad anterior es sencilla si se reconoce que

$$(\tau_i \pm \tau_j)^2 (n-p-1) / [2(1 \pm \rho_{ij}) - (\tau_i \pm \tau_j)^2] \quad (2.3.14)$$

tiene distribución $F(j; n-p-1)$. Se puede usar (2.3.14) para demostrar que

$$\beta_{ij}^{\pm} = 0 \text{ cuando } zd^2 > (1 \pm \rho_{ij}), \quad (2.3.15)$$

que es la misma condición dada en (2.3.10). Sea $c(\pm) = \{(i, j) / i < j, zd^2 < (1 \pm \rho_{ij})\}$, usando (2.3.13), - (2.3.14) y (2.3.15), la desigualdad de Bonferroni resulta ser

$$\alpha - \beta^+ - \beta^- \leq \Pr(\mathcal{V}^* > d) \leq \alpha \quad (2.3.16)$$

$$\begin{aligned} \text{con: } \alpha &= n \Pr[F > d^2(n-p-1) / (1-d^2)], \\ \beta^+ &= \sum_{c(+)} \Pr[F > d^2(n-p-1) / (\frac{1}{2}(1 + \rho_{ij}) - d^2)], \\ \beta^- &= \sum_{c(-)} \Pr[F > d^2(n-p-1) / (\frac{1}{2}(1 - \rho_{ij}) - d^2)]. \end{aligned}$$

Comentarios finales sobre este procedimiento

i) Se requiere calcular las ρ_{ij} 's ya que la cota inferior en (2.3.16) depende de la distribución conjunta de $(\mathcal{V}_i, \mathcal{V}_j)$ a través de estas.

ii) La cota superior en (2.3.16) es exacta si $c(+)$, $c(-)$ son vacíos o equivalentemente $1 + \max_{i < j} \rho_{ij} < zd^2$

iii) En algunos casos se puede aproximar la cota inferior reemplazando en β^+ y β^- a ρ_{ij} por $\max_{c(+)} \rho_{ij}$ y $\min_{c(-)} \rho_{ij}$ respectivamente.

iv) Para una prueba de una cola, reemplazar en (2.3.5) la cota superior α por $\alpha/2$ y a la cota inferior por $(\alpha - \beta^+)/2$.

2.3.2.3 Método de Doornbos.

Este método tiene la misma finalidad que el anterior, es decir, por medio de la desigualdad (2.3.8) encontrar cotas para la distribución de la estadística de prueba (2.3.6) para probar la hipótesis de que el modelo (2.2.2) es correcto contra la alternativa para alguna i desconocida $E(y_i) = \alpha_i \beta + \lambda$.

teniendo en cuenta $f(\tau_i, \tau_j)$ definida en (2.3.9) y usando la transformación

$$\begin{aligned} t_i &= \tau_i, \\ t_j &= (-c_{ij}\tau_i + \tau_j)(1 - \tau_i^2)^{-\frac{1}{2}}(1 - \tau_j^2)^{-\frac{1}{2}}, \end{aligned}$$

cuya distribución conjunta (Doornbos y Prints (1958)) es

$$\begin{aligned} g(t_i, t_j) &= \frac{\gamma}{\pi} (1 - t_i^2)^{\gamma - \frac{1}{2}} (1 - t_j^2)^{\gamma - \frac{1}{2}} \\ &= \frac{\Gamma(\gamma + 1)}{\Gamma(\frac{1}{2})\Gamma(\gamma + \frac{1}{2})} (1 - t_i^2)^{\gamma - \frac{1}{2}} \cdot \frac{\Gamma(\gamma + \frac{1}{2})}{\Gamma(\frac{1}{2})\Gamma(\gamma)} (1 - t_j^2)^{\gamma - 1}. \end{aligned}$$

En otras palabras t_i, t_j son independientes. Por otro lado

$$\begin{aligned} P(\tau_i \geq c_{\alpha}, |\tau_j| \geq c_{\alpha}) &= 2P(\tau_i \geq c_{\alpha}, \tau_j \geq c_{\alpha}) + 2P(\tau_i \geq c_{\alpha}, \tau_j \leq -c_{\alpha}) \\ &= 2P_1 + 2P_2. \end{aligned}$$

Supongamos por ahora que $\rho_{ij} \geq 0$, entonces

$$P_i = P\{t_i \geq c_\alpha, t_j \geq (-\rho_{ij}t_i + c_\alpha)(1-t_i^2)^{-1/2}(1-\rho_{ij}^2)^{-1/2}\}$$

ahora $\frac{\partial}{\partial t_i} (-\rho_{ij}t_i + c_\alpha)(1-t_i^2)^{-1/2} = (-\rho_{ij} + t_i c_\alpha)(1-t_i^2)^{-3/2}$

que será no negativa para toda $t_i \geq c_\alpha$ si

$$\rho_{ij} \leq c_\alpha^2. \quad (2.3.17)$$

Suponiendo que este es el caso y como t_i y t_j son independientes, entonces podemos escribir

$$\begin{aligned} P_i &\leq P\{t_i \geq c_\alpha, t_j \geq (-\rho_{ij}c_\alpha + c_\alpha)(1-c_\alpha^2)^{-1/2}(1-\rho_{ij}^2)^{-1/2}\} \\ &= P\{t_i \geq c_\alpha\} \cdot P\{t_j \geq (-\rho_{ij}c_\alpha + c_\alpha)(1-c_\alpha^2)^{-1/2}(1-\rho_{ij}^2)^{-1/2}\}. \end{aligned} \quad (2.3.18)$$

Definamos c_α^* tal que $P\{t_i \geq c_\alpha^*\} = \frac{\alpha}{n}$. Debido a que el cociente de las densidades de t_i y t_j decrece con t^2 , los percentiles superiores de t_i serán mayores que los correspondientes a t_j . En otras palabras, $c_\alpha^* > c_\alpha$.

Ahora calculemos los valores de las c_{ij}^* para los que

$$c_\alpha(1-\rho_{ij})(1-c_\alpha^2)^{-1/2}(1-\rho_{ij}^2)^{-1/2} \geq c_\alpha^*$$

o equivalentemente

$$\rho_{ij} \leq \frac{c_\alpha^2 - c_\alpha^{*2} + c_\alpha^4}{c_\alpha^2 + c_\alpha^{*2} - c_\alpha^2} = g_\alpha.$$

Además, usando el hecho de que $c_{\alpha} < c_{\alpha}^*$ se puede demostrar que

$$g_{\alpha} < c_{\alpha}^2,$$

consecuentemente (2.3.19) implica (2.3.17). teniendo en mente que $c_{\alpha}^*(n, k) = c_{\alpha}(n, k+1)$, donde $c_{\alpha}(n, k+1)$ se obtiene de la tabla de Lund al dividir los valores correspondientes por $(n-p)^{1/2}$. Entonces c_{α}^* y g_{α} se pueden obtener de estas tablas o bien, usando las tablas de Doornik (1980)

Regresando a (2.3.18) se observa que para β_{ij} positiva, si (2.3.19) se cumple, entonces

$$P_1 \leq P(t_i \geq c_{\alpha}) P(t_j \geq c_{\alpha}^*) = \left(\frac{1}{2} \alpha/n\right)^2 \quad (2.3.20)$$

Para $\beta_{ij} < 0$, (2.3.20) se verifica fácilmente observando (2.3.17) y (2.3.19). Procediendo de la misma forma.

$$P_2 \leq \left(\frac{1}{2} \alpha/n\right)^2$$

si β_{ij} es positiva, o si es negativa, pero $-\beta_{ij} \leq g_{\alpha}$

Combinando estos resultados obtenemos:

$$Pr(|t_i| \geq c_{\alpha}, |t_j| \geq c_{\alpha}) = 2P_1 + 2P_2 \leq 4\left(\frac{1}{2} \alpha/n\right)^2 = \left(\frac{\alpha}{n}\right)^2 \quad (2.3.21)$$

si

$$|\beta_{ij}| \leq g_{\alpha}, \quad (2.3.22)$$

consecuentemente

$$\sum \Pr(|t_{i1}| \geq c_\alpha, |t_{ij}| \geq c_\alpha) < \binom{p}{2} \left(\frac{\alpha}{n}\right)^2 < \frac{1}{2} \alpha^2, \quad (2.3.23)$$

de donde, la desigualdad (2.3.8) puede expresarse como

$$\alpha - \frac{1}{2} \alpha^2 \leq \Pr(\hat{t}^* \geq c_\alpha) \leq \alpha \quad (2.3.24)$$

si

$$\max_{(i,j)} |S_{ij}| \leq g_\alpha. \quad (2.3.25)$$

Note que si (2.3.25) no se cumple, pero (2.3.22) es válida, excepto para pocas parejas (i, j) , (2.3.21) — sigue siendo útil para proporcionar una cota inferior aceptable para $\Pr(\hat{t}^* \geq c_\alpha)$.

2.3.2.4. Método de Weisberg

Considerese el modelo de regresión lineal con las hipótesis usuales. Una vez que se cuenta con los elementos necesarios para suponer que la i -ésima observación es discrepante, se sugiere seguir los siguientes pasos para su determinación.

- 1) Elimine el i -ésimo caso del conjunto de datos.
- 2) Calcule $\hat{\beta}_{(i)}$ y $S_{(i)}^2$ que son los estimadores de β y σ^2 usando solo $(n-1)$ casos. $S_{(i)}^2$ tiene $(n-1)-p = n-p-1$ grados de libertad.
- 3) Para el caso eliminado i , calcule $\tilde{y}_i = \hat{x}_i^+ \hat{\beta}_{(i)}$ para predecir y_i , debido a que este caso no se usó en la estimación de β . Como el i -ésimo caso fue excluido antes de la estimación, y_i y \tilde{y}_i (y_i y $S_{(i)}^2$) son independientes. La

varianza de \tilde{y}_i está dada por

$$\text{var}(\tilde{y}_i) = \sigma^2 z_i^2 (X_{(i)}^t X_{(i)})^{-1} z_i$$

cuyo estimador está dado por

$$\widehat{\text{var}}(\tilde{y}_i) = S_{(i)}^2 z_i^2 (X_{(i)}^t X_{(i)})^{-1} z_i.$$

4) Ahora, si y_i no es una observación discrepante, $E(y_i - \tilde{y}_i) = 0$, pero si y_i es una observación discrepante $E(y_i - \tilde{y}_i) \neq 0$. La diferencia $y_i - \tilde{y}_i$ tiene varianza $\text{var}(y_i - \tilde{y}_i) = \sigma^2 + (1 - z_i^2 (X_{(i)}^t X_{(i)})^{-1} z_i)$ y como se supone que los errores tienen distribución normal, una prueba t-student de la hipótesis " $E(y_i - \tilde{y}_i) = 0$ " está dada por

$$t_i = \frac{y_i - \tilde{y}_i}{S_{(i)} \sqrt{1 + z_i^2 (X_{(i)}^t X_{(i)})^{-1} z_i}}.$$

Esta prueba tiene $n-p-1$ grados de libertad.

Una prueba que ayuda a reducir cálculos es

$$t_i = r_i \left(\frac{n-p-1}{n-p-r_i^2} \right)^{1/2}.$$

Naturalmente, esta estadística es útil tanto en la versión etiquetada como en la no etiquetada. En el primer caso t_i puede compararse directamente con el valor obtenido en la tabla de distribución t-student usual, mientras que en el segundo caso debe de hacerse una modificación.

Usualmente, el investigador no tiene conocimiento de cual observación es discrepante. Si probamos el caso con el valor máximo de t_i , en realidad estamos efectuando n pruebas de significancia, una para cada caso. Supongamos por ejemplo, que no hay observación discrepante y que $n=65$, $p=5$. La probabilidad que una estadística con 60 grados de libertad exceda a 2.00 es 0.05; sin embargo la probabilidad de que la mayor de las 65 pruebas t independientes se exceda a 2.00 es 0.924. Lo anterior sugiere la necesidad de un valor crítico adecuado. La técnica que usamos para encontrar valores críticos se basa en la primera desigualdad de Bonferroni, — que establece que para n pruebas, cada una de tamaño α , la probabilidad de señalar falsamente una observación como discrepante no es mayor que $n\alpha$. Este procedimiento es conservador, ya que la primera desigualdad de Bonferroni debería de especificar solamente que la probabilidad de que el máximo de 65 pruebas que se exceden en 2.00 no es mayor que $65(0.05)$, que es mayor que 1. Sin embargo, escogiendo el valor crítico como el punto $(\alpha/n) \times 100$ de la distribución t dará un nivel de significancia no mayor que $n(\alpha/n) = \alpha$. Entonces deberíamos escoger un nivel de $0.05/65 = 0.00077$ para cada prueba para dar un nivel global no mayor que $65 \times (0.00077) = 0.05$. La tabla B al final de este trabajo es útil para esta prueba.

Esta prueba puede describirse de la siguiente manera: rechazar la hipótesis de que $E(y_i - \bar{y}_i) = 0$ si

$\max_i |t_i| > t(\alpha/2, n-p-1)$ para un nivel nominal α . Cook y Weisberg (1980) sugieren una regla alternativa: rechazar la hipótesis anterior si

$$\max_i |t_i| > t(\sqrt{r_i} \alpha / p; n-p-1)$$

donde afirman que esta regla mantiene el nivel de significancia total, pero proporciona un incremento en la potencia de la prueba en casos donde r_{ii} es "grande".

2.3.2.5. Sobre la sensibilidad de las pruebas en regresión cuando no hay normalidad en los errores.

Muchas de las pruebas aplicadas para observaciones discretas en el modelo de regresión lineal con las suposiciones de costumbre, se basan en los r_i o una función de ellos. Miyashita y Newbold (1983) por medio de simulación y usando como estadística de prueba la expresada en (2.3.6) proporcionan información acerca del comportamiento de las probabilidades en la cola de la distribución de esa estadística cuando la distribución del error no es normal. A continuación se describe este trabajo.

Para evaluar el efecto de no normalidad en el comportamiento de la prueba para observaciones discretas, se considerará una clase de distribuciones conocida como "distribuciones potencia exponencial".

Específicamente, si suponemos que los errores e_i están idénticamente distribuidas con funciones de densidad

$$f(\epsilon_i/\sigma, \theta) = w(\theta) \sigma^{-1} \exp[-c(\theta) \left| \frac{\epsilon_i}{\sigma} \right|^{2/(1+\theta)}], \quad (2.3.26)$$

con

$$c(\theta) = \left[\frac{\Gamma\left[\frac{3}{2}(1+\theta)\right]}{\Gamma\left[\frac{1}{2}(1+\theta)\right]} \right]^{1/(1+\theta)}$$

y

$$w(\theta) = \frac{[\Gamma\left(\frac{3}{2}(1+\theta)\right)]^{1/2}}{(1+\theta)[\Gamma\left(\frac{1}{2}(1+\theta)\right)]^{3/2}}$$

donde σ es la desviación estándar de la distribución de ϵ , mientras que el parámetro θ puede tomar valores entre -1 y 1 y puede considerarse como una medida de curtosis que indica el grado de la "no-normalidad" de la población principal. En particular, cuando $\theta = 0$, la distribución es normal. Cuando $\theta = 1$, la distribución es doble exponencial y finalmente se puede demostrar que cuando θ tiende a -1 , la distribución es rectangular, ver figura (2-3-1) donde las distribuciones tienen la misma varianza. Podemos observar que para $\theta < 0$ y $\theta > 0$ las distribuciones son platocúrticas y leptocúrticas respectivamente.

El proceso de simulación se realiza con un valor de $\sigma = 1$ y manejando primeramente el modelo de regresión simple

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (i = 1, 2, \dots, n),$$

$\theta = -0.75$



$\theta = 0.25$



$\theta = -0.50$



$\theta = 0.50$



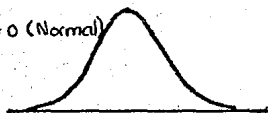
$\theta = -0.25$



$\theta = 0.75$



$\theta = 0$ (Normal)



$\theta = 1.00$
(Doble exponencial)

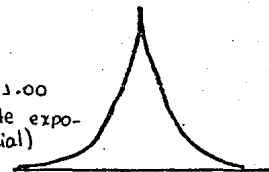


Figura 2.3.1. Distribuciones potencia exponencial con la misma varianza para varios valores de θ

con muestras de tamaño $n=10, 20, 40$ y 2000 repeticiones para los conjuntos

A: $x_1 = x_2 = 1$, con los demás valores $x_i = 0$; $i = 3, 4, \dots, n$

B: $x_1 = 1, x_2 = 0, x_i = 0.5$ para $i = 3, 4, \dots, n$

C: $\frac{n}{2}$ de las $x_i = 1$ y las restantes son valor $x_i = 0$

D: x_i con distribución uniforme en el intervalo $0, 1$

La tabla 2.3.1 muestra la proporción de veces que se rechaza la hipótesis nula "no hay observaciones discrepantes" en base a la estadística (2.3.10) con un nivel de significancia $\alpha = 0.05$. De tal forma que las entradas en el cuerpo de la tabla representan niveles de significancia empíricos basados en el máximo residual studentizado en valor absoluto cuando se supone que los errores tienen distribución normal.

Cualitativamente, los resultados en la tabla (2.3.1) son como se esperaban: para distribuciones platycúrticas, los niveles de significancia empíricos están por abajo del valor supuesto, mientras que para distribuciones leptocúrticas, éstos están por arriba del valor $\alpha = 0.05$.

Más aún, cuando el tamaño de la muestra crece en distribuciones leptocúrticas, también aumenta este nivel de significancia. Esto se esperaba, ya que para ma-

<u>X</u>	<u>n</u>	<u>θ</u>							
		<u>-0.75</u>	<u>-0.50</u>	<u>-0.25</u>	<u>0</u>	<u>0.25</u>	<u>0.50</u>	<u>0.75</u>	<u>1</u>
A	10	.0205	.0185	.0165	.0470	.0690	.0900	.1060	.1310
B	10	.0160	.0210	.0300	.0510	.0630	.0875	.1100	.1330
C	10	.0215	.0285	.0320	.0515	.0675	.0830	.1050	.1250
D	10	.0200	.0305	.0345	.0505	.0650	.0860	.1145	.1305
A	20	.0020	.0065	.0205	.0390	.0890	.1310	.1810	.2150
B	20	.0010	.0070	.0195	.0500	.0880	.1380	.1955	.2270
C	20	.0045	.0085	.0255	.0475	.0990	.1435	.1790	.2210
D	20	.0035	.0070	.0200	.0520	.0920	.1435	.1850	.2265
A	40	.0000	.0015	.0105	.0545	.1120	.1885	.2420	.3015
B	40	.0000	.0015	.0100	.0565	.1150	.1875	.2505	.2955
C	40	.0000	.0010	.0135	.0565	.1135	.1780	.2310	.3055
D	40	.0000	.0015	.0120	.0560	.1175	.1870	.2310	.3070

+110a (2.3.1)

<u>X</u>	<u>n</u>	<u>-0.75</u>	<u>-0.50</u>	<u>-0.25</u>	<u>0</u>	<u>0.25</u>	<u>0.50</u>	<u>0.75</u>	<u>1</u>
A	10	.0205	.0185	.0165	.0470	.0690	.0900	.1060	.1310
B	10	.0160	.0210	.0300	.0510	.0630	.0875	.1100	.1330
C	10	.0215	.0285	.0320	.0515	.0675	.0830	.1050	.1250
D	10	.0200	.0305	.0345	.0505	.0650	.0860	.1145	.1305
A	20	.0020	.0065	.0205	.0390	.0890	.1370	.1810	.2150
B	20	.0010	.0070	.0195	.0500	.0880	.1380	.1755	.2270
C	20	.0045	.0085	.0255	.0475	.0990	.1435	.1790	.2210
D	20	.0035	.0070	.0200	.0520	.0920	.1435	.1850	.2265
A	40	.0000	.0015	.0105	.0545	.1120	.1885	.2420	.3015
B	40	.0000	.0015	.0100	.0565	.1150	.1875	.2505	.2955
C	40	.0000	.0010	.0135	.0565	.1135	.1780	.2310	.3055
D	40	.0000	.0015	.0120	.0560	.1175	.1870	.2370	.3070

+1120 (2.3.1)

por tamaño de muestra, mayor es la probabilidad - de encontrar errores en las colas de distribuciones leptocúrticas. Así, si la suposición de normalidad no es válida estos errores podrían señalarse como discrepantes.

Cuantitativamente los resultados son motivo de alguna preocupación, enfatizando la importancia de la suposición de normalidad, particularmente cuando la muestra es de un tamaño moderadamente grande. En el caso extremo de una distribución de error doble exponencial, los niveles de significancia empíricos están alrededor del 13%, 22% y 30% para muestras de tamaño 10, 20 y 40 respectivamente. Entonces, aún para muestras de tamaño moderado, es bastante probable que una observación sea señalada (erroneamente) como discrepante si la distribución del error es doble exponencial. Quizás sea de mayor importancia la influencia substancial sobre los niveles de significancia de las menores desviaciones de normalidad en los errores cuando la muestra contiene alrededor de 40 observaciones. Aquí, para $\rho = -0.25$ los niveles de significancia son del orden de 1.25% - mientras que para $\rho = 0.25$ están alrededor del 11.5%. Esto, indica que la prueba tiene propiedades robustas muy pobres en muestras moderadamente grandes.

Ahora, consideremos el modelo de regresión lineal múltiple

$$Y_i = X_{i1} \beta_1 + \beta_2 X_{i2} + \beta_3 X_{i3} + \epsilon_i \quad (i = 1, 2, \dots, n).$$

con $X_{ii} = 1$ para los conjuntos A, B, C, D definidos anteriormente. En este modelo, como matrices de diseño se escogieron las 4 combinaciones $(X_2 : X_3) = (A : B)$, $(B : C)$, $(C : D)$, $(D : A)$ y nuevamente se realizaron 2000 repeticiones con muestras de tamaño 10, 20, 40 para cada una.

Aunque los autores no muestran los resultados de este procedimiento de simulación, reportan que los niveles de significancia empíricos son muy similares a los contenidos en la tabla (2.3.1). Esto sugiere que el número de variables independientes tiene muy poco efecto en el comportamiento de la prueba estadística para estas distribuciones del error también afirman que los resultados son muy similares cuando se proponen matrices de diseños diferentes.

Los métodos expuestos anteriormente, se refieren a pruebas para la identificación de una observación ($k=1$). A continuación se presentan algunas pruebas que dominan la literatura para la identificación de más de una observación discrepante ($k > 1$).

2.3.2.6 Método de Gentleman y Wilk.

Usando el modelo de la media trasladada (2.2.1) este método se resume en los siguientes pasos:

- i) Escoger k , número máximo de observaciones discrepantes.

ii) Calcular $G_k(i)$ para cada una de las $\binom{n}{k}$ particiones de los datos.

iii) Sea $G_k^*(i) = \max_{\binom{n}{k}} G_k(i)$. Se dice que el subconjunto

al que le corresponde $G_k^*(i)$ es "el subconjunto de las k observaciones que probablemente son las más discrepantes".

iv) Como la distribución de $G_k^*(i)$ es desconocida, se proponen métodos informales para evaluar la significancia de esta. Una ayuda para juzgar cuando la magnitud de $G_k^*(i)$ es suficientemente grande con respecto a otros valores $G_k(i)$ es una forma de "diagrama de probabilidad" en el que, so 0 o 100 de los valores más grandes de $G_k(i)$ se grafican contra "valores típicos", por ejemplo la mediana, de los correspondientes $G_k(i)$ más grandes que se obtienen por simulación bajo el supuesto de que no hay observaciones discrepantes.

v) Si se concluye que $G_k^*(i)$ no es estadísticamente significativa, entonces el procedimiento anterior se repite para $(k-1)$ y así sucesivamente, caso contrario el método se da por terminado señalando a las k observaciones discrepantes en los datos.

Como se puede constatar este método es muy laborioso, es más, difícilmente se podría efectuar sin la ayuda de una computadora. Con la finalidad de reducir esta carga computacional, Gentleman y Wilk (1975) demuestran que

$$Q_k(i) = \underline{z}_i (I - V_i)^{-1} \underline{z}_i$$

donde \underline{z}_i y V_i son respectivamente, el subvector de \underline{z} y la submatriz de $k \times k$ de V indicados por i . Entonces el cálculo de $Q_k(i)$, bajo la suposición de normalidad, requiere únicamente de los residuales ordinarios y de la matriz de proyección V

Gentleman (1980) con la misma finalidad, propone explotar la estructura de la matriz V , por ejemplo, usa la descomposición Choleski de $I - V_i$; Existe una matriz triangular $R_{k \times k}$ tal que $R^T R = I - V_i$. Entonces si resolvemos para \underline{a} ; $R \underline{a} = \underline{z}_i$ podemos calcular

$$Q_k(i) = \underline{a}^T \underline{a}$$

que es una forma cuyas propiedades computacionales la hace más deseable.

Alternativamente, Cook y Weisberg (1982) usando la descomposición espectral de $V_i = P \Lambda P^T$, donde P es la matriz ortogonal cuyas columnas son los vectores característicos y Λ es la matriz diagonal cuyos elementos son los valores característicos correspondientes, demuestran que

$$Q_k(i) = (P^T \underline{z}_i)^T (I - \Lambda)^{-1} (P^T \underline{z}_i).$$

2.3.2.7 Método de Andrews y Pregibon.

Este método alternativo es capaz de encontrar a

las observaciones discrepantes y/o influyentes más importantes, es decir, identifica varias observaciones discrepantes, con atención directa en aquellas que tienen una influencia substancial en los estimadores resultantes. Para la construcción de esta estadística se combinarán 2 elementos de diagnóstico. Inicialmente, considera los efectos de una observación discrepante en \underline{y} y del renglón correspondiente en forma separada. Primero, al eliminar el caso correspondiente a una observación discrepante en \underline{y} , tiende a reducir marcadamente la suma de cuadrados del residual. La suma de cuadrados del residual es, por lo tanto, un diagnóstico para detectar casos influyentes que surgen debido a una observación discrepante en \underline{y} . Segundo, la influencia de un renglón de X está en parte reflejado por el cambio en $|X^*X|$ cuando se suprime éste. Si $|X^*X|$ cambia substancialmente cuando z_i^2 se suprime, entonces el caso correspondiente (y_i, z_i^2) tendrá una gran influencia en $\hat{\beta}$.

Andrews y Pregibon (1978) sugieren estos diagnósticos separados, que se basan en el cambio en la suma de cuadrados del residual y en $|X^*X|$ que al combinar los nos dan un diagnóstico que se base en el cambio de $(n-p)S^2|X^*X|$ que resulta de la eliminación de uno o más casos. Específicamente se sugiere

$$\hat{r}_i = \frac{|X_{(i)}^* X_{(i)}^*|}{|X^* X^*|} \quad (2.3.30)$$

donde $X^* = (X/\underline{y})$, es decir, la matriz X aumentada -

con Y . Esta cantidad es adimensional. Geométricamente, $R_i^{-1/2} - 1$ corresponde al cambio proporcional en el volumen del elipsoide generado por $X^{*T} X^*$ cuando se eliminan los casos indicados por i . Así, valores "pequeños" de R_i estarán asociados a observaciones discrepantes y/o influyentes. Sin importar que tipo de observaciones se están considerando, es adecuado aislar subconjuntos de observaciones que produzcan valores pequeños de R_i para un análisis adicional.

Draper y John (1981), Cook y Weisberg (1982) demuestran que

$$\begin{aligned}
 R_i &= (1 - Q_k(i) / e_i^T e) |I - V_i| \\
 &= \frac{(n-p-k) \sum_{j=1}^k X_{ji} X_{ji}}{(n-p) S^2 |X^T X|} \\
 &= \frac{(n-p-k) S_{ii}^2}{(n-p) S^2} |I - V_i| \\
 &= \left(1 - \frac{r_i^2}{n-p}\right) |I - V_i|
 \end{aligned} \tag{2.3.31}$$

donde $r_i^2 = \frac{e_i^T (I - V_i)^{-1} e_i}{S^2} = \frac{Q_k(i)}{S^2}$ es el caso múltiple de

los residuales studentizados internos definidos en (2.1.9). Estas expresiones, además de simplificar un poco los cálculos de R_i , permiten las siguientes interpretaciones de sus componentes, (Draper y John 1981): El primer factor será pequeño si $Q_k(i)$ es grande, en este caso puede más identificar conjuntos de observaciones discrepantes - como en Gentleman y Wilk (1975), el segundo factor $(I - V_i)$ proporciona una medida de influencia de las k

observaciones en el ajuste, valores pequeños de $|1 - v_i|$ indicarán los "puntos más remotos" en el espacio de los valores observados. Cuando $k=1$, este factor se reduce a $1 - v_{(1)}$, así, cuando $|1 - v_{(1)}|$ sea "pequeño" se debe a que $v_{(1)}$ es grande, que como se vio en la sección 2.1.1 corresponde a puntos "altamente influyentes".

Bajo el supuesto de normalidad, $(n-p-k) s_{(1)}^2 / (n-p) s^2$ tiene distribución Beta con parámetros $(n-p-k)/2$, $k/2$. Así R_i es proporcional a una variable aleatoria Beta, donde $|1 - v_i|$ es la constante de proporcionalidad, de tal manera que, basados en momentos se pueden obtener valores de referencia.

Se pueden encontrar valores exactos de probabilidad para valores extremos de

$$R_i^* = \min_{\binom{n}{k}} R_i. \quad (2.3.32)$$

El método que se propone es el de Andrews (1971). Esto es, si R_i^* es el mínimo valor observado de las R_i , el conjunto de valores más extremos está dado por $\{R_i / R_i^* \leq R_i^*\}$. El nivel de significancia es la probabilidad de este conjunto, y está dado por

$$\alpha_k = Pr \{ R_i / R_i^* \leq R_i^* \} = Pr \left\{ \sum_{j=1}^k Z_j \leq \xi_{(k)} \right\} \quad (2.3.33)$$

donde Z_j tiene distribución Beta y $\xi_{(k)} = R_i^* / R_i(x)$; $R_i(x) = \frac{|X_{(k)}^+ X_{(k)}^-|}{|x^+ x^-|} \neq 0$.

Nótese que los pesos diferenciales $R_i(x)$ incrementan fuertemente la significancia de observaciones discrepantes influyentes. Para conjuntos pequeños de datos ($n \leq 30$) es deseable calibrar la "pequeñez" de los R_i mínimos, vía pruebas de significancia. Sin embargo para conjuntos grandes de datos, solamente tiene un pequeño efecto, mientras que observaciones discrepantes pueden tener un gran efecto y el número de éstos es probablemente mayor. Por consiguiente, en estos casos se hace más énfasis en la detección, que en la evaluación de un nivel de significancia. En la práctica la expresión

$$R_i = |M_i^*| = |I - V_i^*|$$

que es un determinante de $K \times K$, puede incorporarse a un procedimiento secuencial que se basa en $1, 2, \dots, K_{\max}$ posibles observaciones discrepantes, donde K_{\max} es el número máximo de observaciones discrepantes, este número lo escoge el analista.

2.3.2.8 Método de Aitkin y Wilson

Este método aplica el algoritmo EM (Dempster, Laird y Rubin, 1977) a modelos de la forma.

$$F(y) = \sum_{i=1}^m \pi_i f_i(y) \quad (2.3.34)$$

donde $f_i(y)$ $i=1, 2, \dots, m$, son funciones de densidad normal, donde posiblemente las medias y/o varianzas correspondientes son diferentes. El algoritmo comienza -

con estimadores iniciales máximo verosímiles de los parámetros (basados convenientemente en la especificación de una o más observaciones discrepantes), como el primer paso M, entonces calcular $f_j(y_i)$, la probabilidad de que la i -ésima observación provenga de la j -ésima componente; que se basa en los estimadores iniciales, paso E. Usando estas probabilidades, calcular los nuevos estimadores mínimos cuadrados ponderados de los estimadores para un nuevo paso M y así sucesivamente hasta alcanzar la convergencia. La maximización de la función log verosimilitud puede realizarse y compararse con la maximización log-verosimilitud para una muestra y dar indicación de la existencia de una mezcla.

La programación del algoritmo es muy fácil (Aitkin y Wilson 1980) proporciona los estimadores máximo verosímiles de los parámetros, incluyendo la proporción de cada componente en la mezcla. La matriz de covarianza asintótica y la función log verosimilitud maximizada pueden servir para proporcionar una indicación del número de componentes en la mezcla.

Por ejemplo, considere el caso donde $F(y)$ tiene 2 componentes, es decir

$$F(y_i) = p f_1(y_i) + (1-p) f_2(y_i) \quad (2.3.35)$$

donde $f_1(y_i) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left[- \frac{(y_i - \frac{\sum_{r=1}^p \beta_r X_{ri})^2}{2 \sigma^2}} \right]$

y $f_2(y_i) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left[- \frac{y_i - \mu}{2 \sigma^2} \right]$

Las ecuaciones máximo verosímil son

$$\hat{\beta} = \frac{\sum_i \hat{\beta}(1/y_i)}{n}$$

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

$$\hat{\mu} = \frac{\sum_i y_i \hat{\beta}(2/y_i)}{\sum_i \hat{\beta}(2/y_i)}$$

$$S^2 = \frac{\sum_{j=1}^2 \sum_i [e_{ij}^2 \hat{\beta}(j/y_i)/n]}$$

donde $e_{ij} = y_i - \sum_{r=1}^p \hat{\beta}_r X_{ri}$ y $e_{2i} = y_i - \hat{\mu}$, X y Y elementos del modelo (2.2.1) y $W = \text{diag}(\hat{\beta}(1-y_i))$. Entonces $\hat{\beta}$ es un estimador mínimos cuadrados ponderados cuyos pesos son $\hat{\beta}(1/y_i)$ y

$$-2 \log L = n \log S^2 - 2 \sum_i \log \{ \hat{\beta} \exp(-e_{1i}^2/2S^2) + (1-\hat{\beta}) \exp(-e_{2i}^2/2S^2) \}.$$

2.3.2.9. MÉTODO DE MARASINGHE.

El procedimiento que aquí se propone, se basa en el uso de residuales studentizados para construir una estadística F_k similar a la estadística E_k propuesta por Tietjen y Moore (1972) para muestras simples, modificado por Rosner y para la cual Hawkins (1979) proporciona una tabla de valores críticos para la prueba de k observaciones discrepantes. Para generalizar esta estadística y usarla en la detección de k observaciones discrepantes en regresión, se obtiene un subconjunto de k observaciones donde la primera observación que se incluye, es la correspondiente al máximo residual studentizado en valor absoluto que resulta de ajustar el modelo de

regresión con las hipótesis usuales.

Después, se elimina esta observación y se realiza el ajuste con estas $(n-1)$ observaciones usando el mismo modelo. La segunda observación que ingresa a dicho subconjunto se obtiene de la misma forma - que la primera. Este procedimiento de eliminar observaciones y calcular residuales studentizados con las observaciones restantes se continúa hasta obtener el subconjunto deseado (naturalmente, k se determina de antemano). Sea $Q_k^{**}(i)$, la reducción en la suma de cuadrados que resulta de eliminar las observaciones indicadas por i , obtenidas como se indica líneas arriba. Entonces, la estadística de prueba se define por

$$F_k = \frac{s - Q_k^{**}(i)}{s}, \quad (2.3.36)$$

donde $s = (n-p)s^2$. Se puede demostrar que

$$Q_k^{**}(i) = \sum_{i=1}^k t_i^2 \quad (2.3.37)$$

con $t_i = \frac{e_i}{(1 - v_{ii})^{1/2}}$ $i=1, 2, \dots, n$, el i -ésimo residual "ajustado".

Aunque el orden precedente de las observaciones no asegura que el subconjunto de las k -observaciones escogidas proporcione siempre el valor máximo de $Q_k(i)$ (esto es, que localice las observaciones probablemente más discrepantes según Gentleman y Wilk, (1975)), se puede observar en los resultados de simulación, proporcionado por el

autor, que más del 95% de las veces $Q_k^{**}(i)$ es idéntica a $Q_k(i) = \max_i Q_k(i)$. Además, el método de construcción garantiza que cada observación agregada al subconjunto, produce un incremento máximo en $Q_k^{**}(i)$ para cada $k=1,2,\dots$

Usando F_k , se rechaza la hipótesis nula: no hay observaciones discrepantes, cuando F_k es menor que un valor crítico específico. En la tabla 2.3.2., se proporcionan valores críticos obtenidos por simulación de F_k para el caso de regresión lineal simple para algunos valores de n y k . Una comparación de estos valores críticos con los obtenidos por Hawkins (ver Hawkins 1980; tabla A31) para la estadística E_k de Rosner, indica que éstos pueden usarse como una aproximación razonable para la estadística F_k con la muestra reducida en uno, debido al ajuste extra de un parámetro. También se sugiere usar esta tabla para el caso de regresión múltiple tomando como tamaño de muestra $n-p+1$.

Aunque la estadística F_k puede usarse de la misma forma que se usa E_k para el caso de muestras simples, ésta tiene los mismos problemas asociados con E_k , ver Hawkins (1980), en particular, es susceptible a empantanamiento (declarar más observaciones discrepantes de las existentes) y enmascaramiento (declarar menos observaciones discrepantes de las existentes) dependiendo de si el valor de k es sobrestimado o subestimado. También se ha observado que las pruebas que se basan en estadísticas de bloque como E_k y F_k , tienen bastante potencia cuando k es sobrestimada. Para sacar ventaja de esta propiedad se propone el siguiente método: Seleccione un conjunto de k observaciones para probarse como dis-

tabla 2.3.2, valores críticos de F para $p=2$.

n	k = 2		k = 3		k = 4		k = 5	
	$\alpha = .01$	$\alpha = .05$	$\alpha = .01$	$\alpha = .05$	$\alpha = .01$	$\alpha = .05$	$\alpha = .01$	$\alpha = .05$
6	.0007	.0041						
7	.0006	.0209						
8	.0196	.0487	.0025	.0084				
9	.0403	.0834	.0085	.0219				
10	.0667	.1195	.0203	.0430	.0053	.0125		
11	.0947	.1551	.0358	.0643	.0111	.0239		
12	.1227	.1890	.0511	.0887	.0217	.0390	.0071	.0156
14	.1792	.2522	.0932	.1404	.0471	.0779	.0218	.0398
16	.2292	.3048	.1326	.1906	.0793	.1173	.0444	.0696
18	.2814	.3543	.1725	.2340	.1105	.1570	.0599	.1019
20	.3147	.3939	.2098	.2708	.1399	.1888	.0947	.1318
24	.3828	.4599	.2782	.3451	.2042	.2595	.1524	.1956
28	.4439	.5134	.3406	.4028	.2622	.3173	.1998	.2516
32	.4890	.5529	.3867	.4472	.3065	.3631	.2458	.2994
36	.5242	.5879	.4291	.4859	.3530	.4086	.2937	.3440
40	.5501	.6165	.4632	.5195	.3935	.4439	.3330	.3915
50	.6273	.6771	.5411	.5999	.4714	.5193	.4135	.4596
60	.6742	.7179	.5951	.6388	.5300	.5735	.4748	.5174
70	.7101	.7489	.6371	.6765	.5763	.6160	.5240	.5633
80	.7384	.7733	.6708	.7065	.6138	.6502	.5644	.6007
90	.7614	.7931	.6984	.7311	.6449	.6788	.5889	.6318
100	.7804	.8095	.7214	.7516	.6710	.7022	.6265	.6580

crepantes. Entonces probar con F_k la hipótesis H_1 : no hay observaciones discrepantes. Si H_1 se rechaza, se elimina a la observación que tenga el mayor residual ajustado y se efectúa la hipótesis H_2 : no hay observaciones discrepantes usando F_{k-1} calculada en base a las $(n-1)$ observaciones restantes, donde $F_{k-1} = (S_1 - Q_k^{**}(i^*)) / S_1$ y $Q_k^{**}(i^*) = \sum_{j \neq i^*} t_j^2$, i^* no contiene el índice de la observación eliminada. Este procedimiento continúa hasta que la hipótesis nula no se rechaza, entonces se determina el número de observaciones discrepantes. Los valores críticos para F_k , que se usan en cada etapa están calculados en base a que n y k están disminuidos en uno con respecto a la etapa anterior. Esto supone que la observación removida en cada etapa en realidad es una observación discrepante, en cuyo caso, el error tipo I es el nivel especificado para la prueba.

2.4 COMENTARIOS.

Analizando el caso donde $k=1$, el problema para la detección de observaciones discrepantes en la versión etiquetada es muy simple. Para la versión no etiquetada, la estadística básica de prueba es el máximo en el valor absoluto de los residuales studentizados cuya distribución, como se menciona, es intratable. El problema se resuelve, como se puede ver al analizar los métodos propuestos, de diversas formas. Por ejemplo, de los métodos no expuestos. Mickey, Dunn y Clark (1967) sugieren que una observación será discrepante si su eliminación induce una gran

reducción en la suma de cuadrados de los residuales; método que implica el cálculo de n ecuaciones de regresión, una por cada observación eliminada. La estadística de prueba que usan.

$$F_i(1, n-p-1) = \frac{S^2 - S(i)}{S(i)/(n-p-1)} \quad ; \quad i=1, 2, \dots, n, \quad (2.4.1)$$

donde S^2 es la suma de cuadrados del error si la i -ésima observación es eliminada. Siguiendo la misma idea, Snedecor y Cochran (1968) proponen la estadística

$$t_i(n-p-1) = \frac{y_i - \hat{y}(i)}{\sqrt{(y_i - y(i))^2}} \quad i=1, 2, \dots, n \quad (2.4.2)$$

donde $\hat{y}(i)$ es el estimador mínimos cuadrados de y_i , calculado al eliminar la i -ésima observación. Por otro lado Ellenberg (1976) demuestra que las estadísticas (2.4.1), (2.4.2) y (2.3.5) son equivalentes

Otra forma para obtener la distribución de (2.4.1) o cualquier estadística similar es efectuar simulación a gran escala como lo hicieron Tietjen, Moore y Beckman (1973) y Prescott (1975) para (2.3.11) en el caso de regresión lineal simple.

tal vez la idea más común para resolver dicho problema es el uso de la desigualdad de Bonferroni propuesta por Ellenberg. Dicha desigualdad acota superior e inferiormente la distribución de (2.3.6) y (2.3.11) y los mé-

todos que aquí se presentan para detectar una observación discrepante se basa en ella.

El método de Ellenberg proporciona las dos cotas de la desigualdad de Bonferroni y una de las grandes dificultades que enfrenta es el cálculo de la cota inferior y en particular de $P(|T_{i1}| > c_{\alpha} | T_{j1}| > c_{\alpha})$ ya que ésta requiere del cálculo de algunas de las $\binom{2}{2}$ integraciones, por que $P(|T_{i1}| > c_{\alpha}, |T_{j1}| > c_{\alpha}) = 0$ si $c_{\alpha} \geq \left(\frac{1 + |R_{ij1}|}{2}\right)^{1/2}$. Así, para la aplicación de este método, se requiere aparte de algunas integraciones para el cálculo de probabilidades bivariadas la correlación entre los residuales razón por la cual este proceso no se puede usar como una técnica rápida de diagnóstico.

El cálculo de la cota superior que realiza Lund y la presentación de ésta en tablas para varios niveles de significancia nos permite usar este procedimiento rápidamente. Con respecto al método de Weisberg, él usa la estadística de prueba $\max |t_i|$, donde t_i está definido en (2.1.12) y para la cual, usando la desigualdad de Bonferroni, calcula la cota superior. Debe de mencionarse que no existen diferencias importantes entre usar el método de Lund o el método de Weisberg para la detección de una observación discrepante. Al final de este trabajo se presentan las tablas de Lund y Weisberg.

Uno de los principales inconvenientes es que la potencia de la prueba puede estar influenciada por la posición de la observación discrepante en espacio de diseño;

ver Ellenberg (1976), Tietjen, Moore y Beckman (1973) entre otros.

Por otro lado los métodos de Cook y Prescott, y el de Doornbos son métodos que al igual que el de Ellenberg nos permiten evaluar la precisión de la desigualdad de Bonferroni, solo que éstos no requieren de integración numérica. En cuanto al método de Cook y Prescott, la idea básica para calcular la cota superior es expresada en términos de una distribución F con 1 y $n-p-1$ grados de libertad a partir de la cual se puede calcular el valor crítico deseado tal que dicha desigualdad se cumple. Para el método de Doornbos la idea básica radica en el hecho de si $\max |r_{ij}| \leq g_\alpha$ entonces se garantiza que el tamaño de prueba está entre $\alpha - \frac{1}{2}\alpha^2$ y α . Por cierto que g_α puede obtenerse mediante la tabla de Lund.

Para la cota de los dos últimos métodos su precisión dependerá de las correlaciones, así, en situaciones en donde $r_{ij} = \pm 1$ para alguna $i \neq j$, las cotas nunca serán exactas.

Debe aclararse que, cualquier método de detección de observaciones solo debe considerarse como un elemento en el proceso de análisis de datos y no el único determinante. Una observación con un residual studentizado grande puede ser el resultado de uno o más problemas serios; fallas en las suposiciones del modelo, no normalidad de los errores, etc., puede llevar a la aparición de uno o más datos sospechosos. Un análisis de la estructura de la muestra ayudaría al analista a evaluar la validez de las

suposiciones básicas del modelo, procurando que el tamaño de muestra sea grande. Con muestras pequeñas es difícil distinguir entre un error grueso y fallas en las suposiciones del modelo, y por lo tanto se debe tener cuidado en atribuir la presencia de un residual studentizado a un error grueso.

Por lo expuesto anteriormente, sería aconsejable que antes de pensar en cualquier prueba para la detección de observaciones discrepantes, deben de verificarse todas las posibles causas que podrían ocasionar que una observación tenga la apariencia de discrepante. En particular debe de hacerse la suposición de normalidad que, como se vio en el proceso de simulación realizado por Miyashita y Newbold, es un aspecto fundamental para el buen funcionamiento de la estadística (2.3.11) para detectar observaciones discrepantes y probar la necesidad de una transformación de los datos.

Igualmente, cuando $k > 1$ la versión etiquetada no ofrece dificultad alguna para su prueba. Por otro lado, se mencionó al principio de la sección 2.3.2. que la versión no etiquetada representa un verdadero reto para detectar las observaciones discrepantes ya que hay varios problemas inmersos en los métodos propuestos y principalmente, el que se refiere al aspecto computacional, que hace difícil la aplicación de estas técnicas en forma rutinaria; y los problemas de enmascaramiento y empantamiento.

Cuando uno se enfrenta con la posibilidad de observacio-

nes discrepantes múltiples el problema que surge es como hacer inferencias acerca de $\hat{\mu}$ y $\hat{\sigma}^2$ simultáneamente, (teniendo en mente el modelo de la media trasladada 2.2.1). En esta situación uno también se enfrenta a la especificación del número exacto K de observaciones discrepantes o una cota superior de éste. Mickey, Dunn y Clark (1967) evitan la selección de K_2 usando un procedimiento secuencial que se basa en la eliminación de ciertos casos tal que su eliminación alcanza el máximo decrecimiento en la suma de cuadrados residual. Las reglas usuales de este tipo de procedimiento indican cuando éste debe detenerse. Gentle (1978) propone también un método secuencial que se basa en la aplicación de estadística (2.3.11) en cada paso. El principal inconveniente para este tipo de pruebas es la pérdida de potencia, además de enfrentar problemas de enmascaramiento.

En cuanto a los métodos que se presentan, Gentleman y Wilk (1975) recomiendan que en general las K - observaciones más discrepantes deben determinarse escogiendo un K_2 y encontrar el valor máximo de $Q_K(i)$ para $K=K_2$. Si $\max Q_K(i)$ no es estadísticamente discrepante el procedimiento se repite con $K=K_2-1$, pero en caso contrario el proceso se detiene con la identificación de $\hat{\mu}$ que maximiza $Q_K(i)$. Con este proceso se descarta la posibilidad de observaciones discrepantes cuando $Q_K(i)$ no es estadísticamente significativo para $K=K_2, K_2-1, \dots, 1$. Naturalmente el primer inconveniente de este método es la explosión computacional, principalmente para valores de n grandes (también para valores moderados), a pesar de los esfuerzos de algunos auto-

res para poder calcular $G_k(i)$ en forma simple ayudándose principalmente de la descomposición de la matriz V ; ver Gentleman y Wilk (1975), Gentleman (1980), Cook y Weisberg (1982). Otro inconveniente es el hecho de no existir elementos claros para evaluar estadísticamente la discrepancia de $G_k(i)$ ya que se desconoce la distribución de $G_k(i)$ y se proponen métodos informales para evaluar dicha discrepancia. Estos incluyen diagramas de algún subconjunto grande de las $G_k(i)$ mas grandes - contra "valores típicos" de éstas, obtenidos (aparentemente - por simulación) bajo el modelo nulo o también de diagramas de residuales del modelo (2.2.1) que corresponden al conjunto de las observaciones discrepantes detectados - por la $G_k(i)$ mas grande. El primero de estos diagramas se apoya fuertemente en las suposiciones del modelo y se piensa están influenciados por no normalidad, heteroscedasticidad y no aditividad de los errores; ver Barnett y Lewis (1978). Una ventaja posible del método se da cuando k se escoge suficientemente grande - para que el efecto de enmascaramiento no se presente.

Con respecto al método de Andrews y Pregibon, los comentarios que podemos hacer son de alguna - manera similares a los hechos a Gentleman y Wilk. Observando a R_i en (2.3.30), por ejemplo Cook y Weisberg (1982) comentan: (a) R_i es una medida sin unidad; (b) $R_i^{\frac{1}{2}} - 1$ corresponde al cambio proporcional en el volumen de un elipsoide generado por $x^{*T} x^{**}$ cuando se eliminan los casos influyentes. Finalmente R_i es invariante bajo - permutaciones de las columnas de x^{**} y entonces al vector de

respuestas \underline{y} no se le da reconocimiento especial.

A veces es útil un diagrama de $\log R_i/R_k^0$ para valores pequeños de k , donde R_0 denota el valor mínimo observado de R_i en todas las $\binom{n}{k}$ posibles. En estos diagramas se observa que para cierto valor de k para el cual $\log R_i/R_0$ es el más pequeño, este valor se aisla a la izquierda del resto, declarando este valor de k como el número de observaciones discrepantes $\%/\theta$ influyentes.

Aitkin y Wilson (1980) sugieren que un tratamiento adecuado de observaciones discrepantes se puede apoyar en modelos mezcla de la forma

$$F(\underline{y}) = \sum_{i=1}^m p_i N(\mu_i, \sigma_i^2),$$

que consiste de varias componentes normales donde posible mente las medias y/o las varianzas son diferentes, donde $0 \leq p_i \leq 1$ y $\sum p_i = 1$. Ellos argumentan que el algoritmo EM proporciona una forma relativamente simple para obtener estimadores máxima verosimilitud y que la identificación de observaciones discrepantes se puede basar en la probabilidad a posteriori de que el i -ésimo caso provenga de la j -ésima componente $j=0, 1, 2, \dots, m$. Para ilustrar su metodología, ellos analizan uno de los grupos de datos que proponemos aquí, en el cual (como veremos más adelante), se muestra que los resultados dependen críticamente del número m de componentes usados en el modelo. Pruebas de cociente de verosimilitud para el número de com -

ponentes puede ser útil, pero parece que este enfoque, - básicamente reemplaza el problema de escoger K en la versión no etiquetada por el problema de escoger m . La elección de m estará sujeta al mismo tipo de dificultades e inevitablemente nos llevará a procedimientos seriales cuando m sea desconocida. Este enfoque será más útil cuando el número de contaminantes surjan consistentemente de un número conocido de componentes.

Por último, el procedimiento propuesto por Marasinghe propone una estadística que es una generalización de una que se usa para observaciones discrepantes múltiples en muestras simples, y por lo tanto tiene asociados los mismos problemas a esa estadística. Para superar dichas desventajas, el autor propone un método modificado que aparentemente elimina el problema de enmascaramiento y empantamiento.

Mediante simulación, se obtienen las siguientes conclusiones: (a) parece que el nuevo procedimiento controla el error tipo I cuando se escoge el valor de K adecuadamente. Esto implica que la probabilidad de declarar más observaciones discrepantes de las que hay es α y por lo tanto el problema de empantamiento no se da. (b) la potencia del procedimiento contra una observación discrepante es comparable con R_n (método secuencial, que usa al $\max |r_{i1}|, r_i$ el residual studentizado) cuando $K=2$ y decrece ligeramente para valores grandes de K . La potencia para detectar múltiples observaciones discrepantes mayor o igual al número real de dichas observaciones se mantiene alta. Un hecho

importante es que el proceso de simulación reporta que un porcentaje muy alto de las veces este método coincide con el de Gentleman y Wilk pero con la ventaja de la eliminación del proceso tedioso que se requiere para calcular $\max Q_k(i)$.

2.5. Ejemplos .

A continuación se aplicarán algunas de las técnicas presentadas a dos problemas que dominan la literatura referente a observaciones disrepanantes. El primero de ellos, es

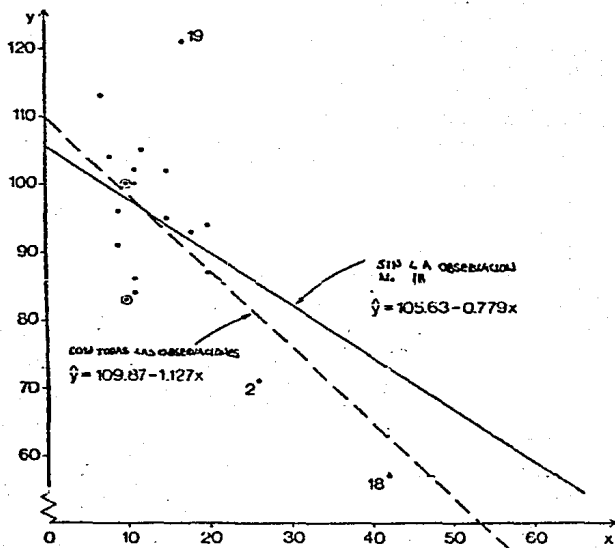
tabla 2.5.1.

Caso.	x	y	Caso	x	y
1	15	95	11	7	113
2	26	71	12	9	96
3	10	83	13	10	83
4	9	91	14	11	84
5	15	102	15	11	102
6	20	87	16	10	100
7	18	93	17	12	105
8	11	100	18	42	57
9	8	104	19	17	121
10	20	94	20	11	86
			21	10	100

analizando inicialmente en Mickey, Dunn y Clark (1967) y se refiere al estudio sobre cierta enfermedad del corazón en niños durante sus primeros meses de vida. Los datos están dados en la tabla 2.5.1. La variable Y es el puntaje adaptable de Gessel y X es la edad del niño (en meses).

La gráfica (2.5.1) representa el diagrama de puntos de las variables de este problema.

Gráfica 2.5.1.



Ahora consideremos el modelo de regresión con las hipótesis usuales para este grupo de datos. A partir de la gráfica, parece posible que dicho modelo se puede establecer a pesar de que las observaciones 19, 18, 2 dominan nuestra percepción. Si los 3 casos fueran removidos, la linealidad es menos pronunciada.

La tabla (2.6.2), contiene los estimadores de los parámetros y la tabla de análisis de varianza al ajustar esos datos por mínimos cuadrados.

tabla 2.6.2.

ESTIMADORES DE LOS PARÁMETROS

CIUDA	VALORES
109.8738	25.8826
-1.1270	0.0932

LA VARIANZA ESTIMADA ES 121.505

TABLA DE ANÁLISIS DE VARIANZA

FV	GL	SC	CM	F
HO : B=0	1	1604.081	1604.081	13.202
ERROR	19	2308.586	121.505	
TOTAL	20	3912.667		

EL COEFICIENTE DE CORRELACION r^2 ES = 0.410

19 es declarada como discrepante, no así, para $\alpha = 0.01$ y $\alpha = 0.001$.

Con la finalidad de obtener el valor de c_α , sin necesidad de efectuar integrales, en algunas ocasiones se podrán usar las tablas de Lund, pero únicamente se podría recomendar su uso cuando se disponga de un método adecuado de interpolación. Para nuestro problema $n=21$, interpolando linealmente, se consigue $c_\alpha = 0.6359$ y como $\tau^* = 0.647$, la conclusión sería la misma.

Si usamos directamente las tablas de Lund y como estadística de prueba a τ^* las conclusiones son las mismas ya que estos métodos son equivalentes, ver Palacios, R. y Castañón V. (1983).

Si además, se requiere la cota inferior de la desigualdad de Bonferroni, entonces, se necesitará la matriz de correlaciones de los residuales y evaluar algunas probabilidades del tipo $P(\tau_i > c_\alpha, |\tau_{ij}| > c_\alpha)$ dependiendo de si $c_\alpha \geq (\frac{1+|\rho_{ij}|}{2})^{1/2}$. Usando la tabla (2.5.4), se encuentra que para $\alpha = 0.05$, la última desigualdad no sirve ya que ninguna correlación residual satisface $0.640 \geq (\frac{1+|\rho_{ij}|}{2})^{1/2}$. Para $\alpha = 0.01$ solo la satisfacen dos: $\rho_{4,2}$ y $\rho_{12,2}$. Sin embargo para $\alpha = 0.01$, solo se requieren calcular 4 probabilidades bivariadas. Los resultados, después de evaluar las probabilidades bivariadas son:

$$\begin{aligned} 0.0499 &\leq P(\max_i |\tau_i| > 0.640) \leq 0.05, \\ 0.0100 &\leq P(\max_i |\tau_i| > 0.708) \leq 0.0100, \\ 0.0010 &\leq P(\max_i |\tau_i| > 0.781) \leq 0.0010. \end{aligned}$$

METODO DE COOK Y PRESCOTT. Nuevamente, si nuestro interés solo radica en probar si una observación es discrepante o no lo es, únicamente se requiere encontrar el valor crítico d que satisface $P(|\hat{b}_i| > d) = P[F > d^2(n-p-1)/(1-d^2)]$ o equivalente - mente $\alpha = nPr [F > d^2(n-p-1)/(1-d)^2]$. Si $\hat{b}_i^* > d$ entonces se rechaza la hipótesis nula: no existen observaciones discrepantes, donde \hat{b}_i^* está definido en (2.3.6). En nuestro ejemplo, se sigue una línea un poco diferente. Tomamos el valor $\hat{b}_i^* = d = 0.6475$ y obtenemos $\alpha = 0.0425$. Entonces no se requiere probar la hipótesis nula con $\alpha = 0.05$ ya que el verdadero valor es menor a 0.05 , por lo tanto, podemos declarar a la observación 19 como discrepante.

Sin embargo, cuando se juzga el peso de la evidencia contra la hipótesis nula, sería deseable conocer que tanto $P(\max_i |\hat{b}_i|)$ es menor que el valor α propuesto. Esto último, puede determinarse calculando la cota inferior de (2.3.16), que suele facilitarse al seguir la recomendación de sustituir en β^+ y β^- a β_{ij} por el $\max_{(i,+)} \beta_{ij}$ y $\min_{(i,-)} \beta_{ij}$ respectivamente. Al observar en la tabla (2.5.4) nos damos cuenta que todas las correlaciones están contenidas en el intervalo $[-0.556, 0.202]$ y aplicar el criterio mencionado, $\alpha = \beta^+ - \beta^-$ es negativo. Para mejorar este resultado, se hace un análisis exhaustivo de la matriz de correlación y obtenemos que $\rho_{18,2} = -0.56$, $\rho_{18,6} = 0.3$, de las restantes 17 caen en el intervalo $[0.002, 0.202]$ y 190 en $[-0.221, -0.016]$. Entonces, para calcular una segunda cota inferior, se usan $\{-0.556, -0.3, -0.016, 0.202\}$ con sus respectivas frecuencias $\{1, 2, 190, 17\}$ para calcular β^+ y $\{-0.556, -0.3, -0.221, 0.002\}$ con las mismas frecuencias para calcular β^- . Como podemos observar este procedi-

miento solo involucra el cálculo de 8 probabilidades para obtener que $|\beta^+ + \beta^-| < 0.0016$. Entonces, la probabilidad del máximo $|T_{ij}|$, tiene como cotas los límites del intervalo $[0.0409, 0.0425]$ para $d = 0.6475$. Igual procedimiento se aplicaría si se quiere encontrar las cotas para un nivel de significancia específico, por ejemplo, para $\alpha = 0.01, 0.05$ y 0.1 , las cotas inferiores resultan $0.00997, 0.0476$ y 0.086 respectivamente.

METODO DE DOORNBOOS. Lo relevante de este método es obtener una cota superior para la probabilidad bivarida de la desigualdad de Bonferroni, expresada en (2.3.21) siempre y cuando $\max_{i,j} |c_{ij}| \leq g_\alpha = \frac{c_1^2 - c_2^2 + c_3^2}{c_1^2 + c_2^2 - c_3^2}$

A continuación, se presentan algunos resultados aproximados para $\alpha = 0.1, \alpha = 0.05, \alpha = 0.01$

$$\alpha = 0.1 ; g_\alpha = 0.1908,$$

$$\alpha = 0.05 ; g_\alpha = 0.225,$$

$$\alpha = 0.01 ; g_\alpha = 0.2969.$$

Como los resultados lo muestran, para ningún valor de α se consigue la condición (2.3.25) y por lo tanto no podemos garantizar que $P(T^* \geq c_\alpha)$ este en el intervalo $[\alpha - \frac{1}{2}\alpha^2, \alpha]$.

METODO DE WEISBERG. La aplicación de este método es muy sencilla. Observando en la tabla (2.5.3) encontramos que $\max |t_{ij}| (t_i = t_i(n-p-1/n-p-r_i^2)^{1/2})$ es $t_{19} = 3.6069$. Nuevamente para $\alpha = 0.05$, la observación 19 es declarada como discrepante ya que el valor de t en la tabla de Weisberg ($n=21, p=2$) es 3.53. Pero, para $\alpha = 0.01$ la conclusión es que no hay observa

ciones discrepantes ya que el valor de t es 4.26.

METODO DE ANDREWS Y PREGIBON. En la tabla (2.5.5) se encuentran resumidos los valores de R_i y z_i para diferentes valores de K , al observarla en $K=1$, nuestra atención se centra en la observación 18, debido principalmente a que es la más alejada en el espacio de factores. Sin embargo las observaciones 2 y 19 también requieren de un análisis adicional. Por lo que respecta a la observación 19, ésta ha sido declarada como discrepante en el uso de los métodos anteriores, con un nivel de confianza del 6%. Este hecho se puede confirmar al observar la segunda y tercer columna de la tabla (2.5.6) que nos proporciona $Q_i(i)$ y la primer componente de la estadística de Andrews y Pregibon (AP_i) que como se ha mencionado anteriormente, son elementos de diagnóstico para detectar observaciones discrepantes. Los valores correspondientes son 969 y .58, valores que confirman que la observación 19 puede declararse como discrepante. A pesar de todo, la estadística de Andrews y Pregibon señala a la observación 18 como influyente y/o discrepante. Debido a su posición en el espacio de factores, la observación 18 es identificada como influyente. Por otro lado, para la observación 2 no hay evidencia de que esta observación sea influyente y/o discrepante.

Para $K=2$, la tabla (2.5.7) proporciona sólo aquellas parejas que potencialmente pudieran importar. Por un lado la estadística A-P selecciona a la pareja (2,18) como más importante con respecto a la pareja (18,19) por pequeño margen. Para esta pareja, (2,18), AP_2 es pequeña indicando que estas observaciones son potencialmente influyentes, mientras que para las

tabla 2.5.5.

k	R_L^0	$\hat{\alpha}_L$	Observaciones
1	-3350940	-3838198	18
2	-1645985	-2310607	18, 2
3	-0897915	-2262575	18, 2, 19
4	-0593024	-7475751	18, 2, 19, 11

Grafica 2.5.2.

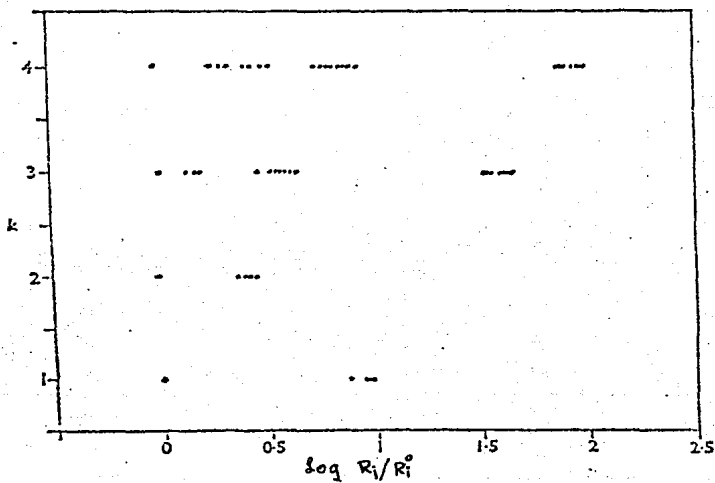


Tabla 2.5.6.

Observación Eliminada	Q_i	AP_1 $100(1-Q_i/ser)$	AP_2 $100(1-v_{ii})$	Estadística AP $100R_i^{(1)}$
1	4	100	95	95
2	108	96	85	81
3	260	89	94	83
4	82	97	93	90
5	86	97	95	92
6	0	100	93	93
7	12	100	94	94
8	7	100	94	94
9	11	100	92	92
10	48	98	93	91
11	133	95	91	86
12	15	99	93	92
13	260	89	94	83
14	193	91	94	86
15	22	99	94	93
16	2	100	94	94
17	79	97	95	92
18	88	95	35	33
19	969	58	95	55
20	140	94	94	89
21	2	100	94	94
Valores usuales	Alto	Bajo	Bajo	Bajo

tabla 2.5.7.

Observación Eliminada	$Q_2(i)$	AP_1	AP_2	Estadística AP
		$100(1 - Q_2^{(i)}/s_n)$	$100 1 - v_i $	$100R_{ij}^{(2)}$
18, 2	442	81	20	16
18, 3	324	86	32	28
18, 11	277	88	30	27
18, 19	983	57	32	18
19, 2	1031	55	80	44
19, 3	1189	48	89	43
19, 11	1128	51	86	44
Valores Usuales	Alto	Bajo	Bajo	Bajo

observaciones (18,19), AP es pequeño y entonces podrían señalarse como discrepantes. Como se ha visto, la observación 19, ha sido anteriormente declarada discrepante y entonces al eliminar dos observaciones en los datos originales que incluya a la observación 19, esta contribuye fuertemente para que $Q_2(i)$ sea un valor grande, así $Q_2(19,3)$ es el valor mas grande, pero mucho de $Q_2(19,3)$ se debe a la observación 19.

De lo expuesto anteriormente, únicamente podremos señalar una observación como discrepante y es, en este caso, la nú

mero 18. Esta misma conclusión puede obtenerse al observar la gráfica (2.5.2) en la que para $K=1$, hay mayor separación entre los valores respectivos de $\log \{R_i/R_i\}$.

Uno de los principales inconvenientes de la estadística AP es que no necesariamente pone atención a puntos que son discrepantes y/o influyentes en términos de la estimación de parámetros ya que, como se puede observar, AP₂ da un peso considerable a conjuntos de puntos que pueden considerarse como remotos en el espacio de factores, para una discusión más amplia acerca de observaciones influyentes. N. R. Draper y J. A. John (1981).

El segundo grupo de datos fue tomado de Brownlee, K.A. (1965), están representados en la tabla (2.5.8) y son el resultado de 21 días de operación de una planta para la oxidación de NH_3 (Amonia) a HNO_3 (ácido nítrico) y las variables son:

Y = Diez veces el porcentaje de amonía que se pierde como óxido nítrico no absorbido (esta es una medida indirecta del rendimiento del ácido nítrico).

X_1 = Flujo de aire en la planta.

X_2 = Temperatura del agua al entrar.

X_3 = Concentración de ácido nítrico en el líquido absorbido.

Al investigar las operaciones de la planta, se puede considerar a las corridas (1,2), (4,5,6), (7,8) y (18,19) como re-

Tabla 2.5.8.

Observacion Numero	Y	x_1	x_2	x_3
1	42	80	27	89
2	37	80	27	88
3	37	75	25	90
4	28	62	24	87
5	18	62	22	87
6	18	62	23	87
7	19	62	24	93
8	20	62	24	93
9	15	58	23	87
10	14	58	18	80
11	14	58	18	89
12	13	58	17	88
13	11	58	18	82
14	12	58	19	93
15	8	50	18	89
16	7	50	18	86
17	8	50	19	72
18	8	50	19	79
19	9	50	20	80
20	15	56	20	82
21	15	70	20	91

plicas.

En la tabla (2.5.9) se presentan los estimadores de los parámetros, la tabla de análisis de varianza y algunos elementos de diagnóstico obtenidos de ajustar estos datos por mínimos cuadrados. El modelo resultante es

$$\hat{Y} = -39.9 + 0.72 X_1 + 1.30 X_2 - 0.15 X_3. \quad (2.5.1)$$

Daniel y Wood (1971) analizan detalladamente este conjunto de datos y la conclusión a la que llegan es que las observaciones 1, 3, 4 y 21 son discrepantes, y que la variable X_3 no es necesaria. Su modelo final es

$$\hat{Y} = -15.4 - 0.07 X_1 + 0.53 X_2 + 0.0068 X_1^2. \quad (2.5.2)$$

Por otro lado si el ajuste (Mínimos cuadrados) se efectúa únicamente eliminando a las observaciones identificadas como discrepantes se obtiene el modelo.

$$\hat{Y} = -37.6 + 0.80 X_1 + 0.58 X_2 - 0.07 X_3. \quad (2.5.3)$$

Como los métodos de Ellenberg, Lund y Weisberg son similares, usaremos este último para el caso $k=1$ (una observación discrepante). Buscando en la tabla (2.5.9), encontramos que $max |t_{21}|$ corresponde a $t_{21} = -3.33$ cuando usamos al modelo (2.5.1). Pero como el valor $t(\alpha/n, n-p) = 3.60$ de la tabla Weisberg no rechazaríamos la hipótesis nula (no existen observaciones discrepantes) para $\alpha = 0.05$, lo mismo ocurre con $\alpha = 0.01$, ya que $t(\alpha/n, n-p) = 4.42$. Como podemos constatar la prueba para detectar una observación discrepante no es útil en esta situación a pesar de que t_1, t_3 ,

tabla 2.5.9.
ESTIMADORES DE LOS PARAMETROS

ESTAD	VARIANZAS
-10.9197	141.5147
0.7156	0.0132
1.2953	0.0132
-0.1521	0.0244

LA VARIANZA ESTIMADA ES

10.519

TABLA DE ANALISIS DE VARIANZA

CV	GL	SC	CM	F
MS	17	1000.400	420.136	59.992
ERROR	17	170.300	10.519	
TOTAL	20	2067.000		

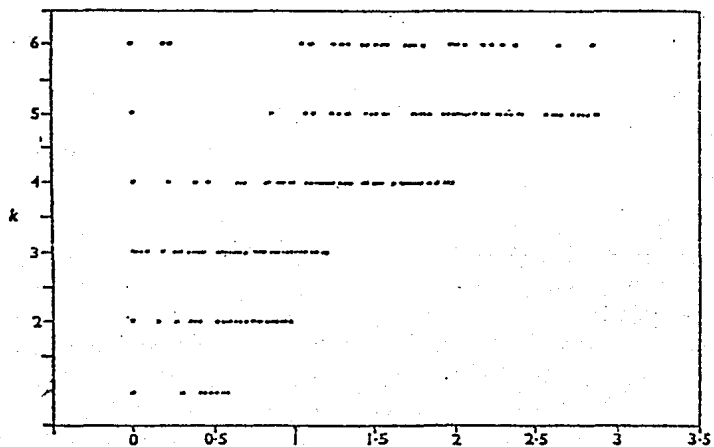
EL COEFICIENTE DE CORRELACION R ES =

0.914

Y
X
V
T
R
S
Q
P
O
N
M
L
K
J
I
H
G
F
E
D
C
B
A

Tabla 2.5.10.

k	R_k^0	α_k	Observaciones
1	·4225375	·0596515	21
2	·2072535	·1093887	21, 4
3	·1147183	·3962458	21, 4, 2
4	·0337735	·0257674	21, 4, 3, 1
5	·0080269	·0008339	21, 4, 3, 1, 2
6	·0039551	·0029340	21, 4, 3, 1, 2, 13



cente al comparar $-2 \log L$ con $\chi^2_{1-\alpha}(4)$ para varios valores de α . (esto es consistente con la prueba de Weisberg aplicada anteriormente). Los parámetros estimados están en la tabla (2.5.11) y la probabilidad a posteriori de que la observación $i=1,2,\dots,21$ sea miembro de la componente j (según el modelo) para los diferentes modelos están en la tabla (2.5.12).

Si inicialmente a la observación 4 (esta tiene el residual positivo más grande en el modelo original) se especifica como discrepante, entonces la convergencia ocurre en 7 iteraciones y los parámetros estimados son idénticos a los obtenidos para el modelo completo de regresión. La observación es "reabsorbida" (reabsorbed) en la primer componente y todas las probabilidades a posteriori de las observaciones de ser miembros de la componente 2 son menores a 10^3 .

Si hay varias observaciones discrepantes, el modelo anterior podría no detectarlas si estas se presentan en el lado opuesto del plano de regresión. Para este propósito será necesario un modelo con al menos 3 componentes.

Consideremos un modelo con 3 componentes (modelo 4), en el cual, la observación 21 y 4 se asignan inicialmente a la segunda y tercera componente respectivamente. La convergencia ocurre en 3 iteraciones y para esta solución $-2 \log L = 58.08$. La reducción comparada con el modelo de 2 componentes es apenas de 3.80. También los parámetros se encuentran en la tabla (2.5.11) y las probabilidades en la (2.5.12). La evidencia para un modelo de 3 componentes no es entonces convincente.

Si las 3 observaciones, 1, 3 y 4 se asignan inicialmente a

tabla 2.5.11.

Parámetros estimados para los modelos del ejemplo no. 2.

Modelo	No. de componentes	Observaciones Discrepancias	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	\hat{u}_2	\hat{u}_3	\hat{u}_4	\hat{u}_5	$\hat{\sigma}^2$ (segado)	-2 log L
1	1	-	-39.9	.716	1.30	-.152	-	-	-	-	8.52	65.98
2	2	(21)	-43.9	.003	.836	-.105	15.0	-	-	-	5.08	61.88
3	2	(4)	-39.9	.716	1.30	-.152	-	.003	-	-	8.51	65.98
4	3	(21) (4)	-42.6	.955	.560	-.108	14.0	28.0	-	-	2.88	58.00
5	4	(21) (4) (1,3)	-38.0	.791	.580	-.006	14.8	28.0	39.5	-	1.54	58.12
6	5	(21) (4) (1) (3)	-38.3	.785	.591	-.005	14.9	28.0	42.0	37.0	0.84	49.96

— 94 —

tabla 2.5.11.

Parámetros estimados para los modelos del ejemplo no. 2.

Modelo	No. de componentes	Observaciones Discrepancias	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	\hat{u}_2	\hat{u}_3	\hat{u}_4	\hat{u}_5	$\hat{\sigma}^2$ (segado)	-2 log L
1	1	-	-39.9	.716	1.30	-.152	-	-	-	-	8.52	65.98
2	2	(21)	-43.9	.883	.836	-.105	15.0	-	-	-	5.00	61.88
3	2	(4)	-39.9	.716	1.30	-.152	-	.003	-	-	8.51	65.98
4	3	(21) (4)	-42.6	.955	.560	-.108	14.0	28.0	-	-	2.80	58.08
5	4	(21) (4) (1,3)	-38.0	.791	.589	-.006	14.0	28.0	39.5	-	1.54	58.12
6	5	(21) (4) (1) (3)	-38.3	.785	.591	-.005	14.9	28.0	42.0	37.0	0.04	49.96

tabla 7.5.12.

Modelo	2		4		5			6			
Componente	2	2	3	2	3	4	2	3	4	5	
Observacion											
1							1			1	
2							.027			.198	
3							1			1	
4				.998		1			1		
5	.051	.022		.004							
6	.076	.032		.016				.001			
7	.027	.005									
8	.008	.001									
9	.133	.101		.124				.193			
10	.072	.068		.078				.091			
11	.078	.072		.100				.138			
12	.057	.043		.035				.019			
13	.036	.032		.006							
14	.042	.028		.013				.003			
15	.001										
16											
17	.001										
18	.001										
19	.003										
20	.092	.117		.260				.073			
21	1	1		1				1			

LOS ESPACIOS EN BLANCO INDICAN UNA PROBABILIDAD MENOR A 10^{-8}

la tercera componente, la convergencia ocurre en 8 iteraciones y los estimadores de los parámetros son idénticos a los obtenidos con el modelo 2. Las tres observaciones son "reabsorbidas" en la primera componente y todas las probabilidades de pertenecer a la componente 3 son menores que 10^{-4} también se consideran modelos con 4, 5 y 6 componentes a pesar de que se requieren sin gran número de parámetros para tales modelos que los hace muy complejos para 21 observaciones. Los modelos 5 y 6 están resumidos en la tabla (2.5.11) y las probabilidades correspondientes en (2.5.12).

Para 6 componentes, las 5 observaciones que se encontraron en Andrews y Pregibon como discrepantes fueron asignadas inicialmente a componentes diferentes. La convergencia ocurre lentamente (23 iteraciones) con una solución de $-2 \log L = 42.11$. Las componentes que contenían a las observaciones 2 y 3 se transforman en una sola y las observaciones 10, 11, 12 y 20 fueron separadas de la observación 21.

Existen varias dificultades con este enfoque: (a) con muestras pequeñas la distribución asintótica (de $-2 \log L$) no se cumple satisfactoriamente, (b) la más seria, χ^2 con el número de grados de libertad igual al número de parámetros por probarse en la hipótesis nula no es válida, para la mezcla de modelos es no-regular, como se puede ver en el modelo de 2 componentes, donde p no es identificable si $\mu_1 = \mu_2$, o equivalentemente, la hipótesis de "no existen observaciones discrepantes".

METODO DE MARRASINGHE. Para la aplicación de este método usaremos el modelo final (2.5.2) obtenido por Daniel y Wood

(1971). En la tabla (2.5.13) se presentan los valores de N_{ij} y t_i para el modelo completo y sin las observaciones 21 (4, 21), (2, 4, 21), (12, 4, 21) y (1, 3, 4, 21) respectivamente. Por otro lado, en la tabla (2.5.14) se presentan los residuales ajustados para el modelo completo - sin la observación 21, (21, 4) y (21, 4, 2).

Analizando las últimas 4 columnas de la tabla (2.5.14) observamos que existen 3 residuales ajustados que difieren fuertemente de los demás, 8.46, 7.68, 4.70 que corresponden a las observaciones 21, 4 y 2 respectivamente.

Por estas razones, podríamos sospechar la presencia de 3 observaciones discrepantes. Para efectos de evitar enmascaramiento, supongamos inicialmente que están presentes cuatro observaciones discrepantes ($k=4$). Las 4 observaciones seleccionadas por el procedimiento son 21, 4, 2, 20 respectivamente, que corresponden a las cuatro observaciones con el residual ajustado (en valor absoluto) más grande. Debe notarse que estas 4 observaciones son las "4 observaciones probablemente más discrepantes" de Gentleman y Wilk seleccionadas por $Q_4^*(k)=159.39$. También debe mencionarse que la prueba para una sola observación discrepante falla en declarar a la observación 21 como discrepante al 5% de significancia. Entonces el procedimiento secuencial basado en R_n (procedimiento que podría usar el criterio de Eilenberg, Lund y Weisberg indistintamente), no detectaría observaciones discrepantes en este modelo.

El procedimiento secuencial descrito anteriormente se aplica a este conjunto de datos, cuyos resultados se encuentran resumidos en la tabla (2.5.15). Como podemos recordar

tabla 2.5.13

Valores de v_{ij}						residuales studentizados						
observaciones eliminadas						observaciones eliminadas						
ninguna	(2,1)	(4,2,1)	(2,4,2,1)	(1,2,4,2,1)	(1,3,4,2,1)	ninguna	(2,1)	(4,2,1)	(2,4,2,1)	(1,2,4,2,1)	(1,3,4,2,1)	
0.409	0.421	0.421	0.727	.	.	1	0.97	0.77	1.08	-1.75	.	.
0.409	0.421	0.421	.	.	0.993	2	-1.06	-1.81	-2.71	.	.	0.57
0.176	0.199	0.201	0.308	0.983	.	3	1.54	1.40	2.30	1.91	1.21	.
0.191	0.192	4	2.27	3.01
0.103	0.108	0.125	0.125	0.125	0.131	5	-0.31	-0.63	-0.31	-0.40	-0.40	-0.12
0.134	0.134	0.164	0.164	0.165	0.169	6	-0.73	-0.97	-0.62	-0.80	-0.82	-0.64
0.191	0.192	0.238	0.239	0.240	0.242	7	-0.84	-0.92	-0.30	-0.35	-0.31	-0.18
0.191	0.192	0.238	0.239	0.240	0.242	8	-0.50	-0.48	0.36	-0.54	0.66	0.84
0.163	0.170	0.206	0.208	0.218	0.208	9	-0.94	-0.89	-0.39	-0.36	-0.18	-0.66
0.139	0.175	0.175	0.176	0.179	0.179	10	0.84	0.38	0.49	0.75	0.94	0.96
0.139	0.175	0.175	0.176	0.179	0.179	11	0.84	0.38	0.49	0.75	0.94	0.96
0.212	0.272	0.275	0.276	0.279	0.280	12	0.96	0.31	0.16	0.30	0.45	0.53
0.139	0.175	0.175	0.176	0.179	0.179	13	-0.17	-0.91	-1.42	-1.82	-1.87	-1.98
0.092	0.110	0.111	0.112	0.116	0.113	14	-0.25	-0.79	-1.04	-1.30	-1.29	-1.41
0.188	0.189	0.191	0.191	0.195	0.193	15	0.17	0.34	0.29	0.29	0.19	0.32
0.188	0.189	0.191	0.191	0.195	0.193	16	-0.17	-0.10	-0.35	-0.57	-0.76	-0.67
0.187	0.195	0.195	0.195	0.198	0.197	17	-0.26	-0.01	-0.01	-0.10	-0.23	-0.21
0.187	0.195	0.195	0.195	0.198	0.197	18	-0.26	-0.01	-0.01	-0.10	-0.23	-0.21
0.212	0.232	0.234	0.234	0.236	0.237	19	-0.35	0.09	0.33	0.37	0.31	0.27
0.064	0.064	0.069	0.070	0.076	0.070	20	0.68	0.75	1.42	2.04	2.38	2.16
0.288	21	-2.63

- 88 -

Tabla 2.514. Residuales Ajustados

Observación	Datos			Observaciones			
	X1	X2	Y	Ninguna	(21)	(21, 4)	(21, 4, 2)
1	80	27	42	3.11	1.96	1.87	-2.25
2	80	27	37	-3.40	-4.61	-4.70	—
3	75	25	37	4.93	3.58	4.00	2.45
4	62	24	28	7.30	7.68	—	—
5	62	22	18	-.98	-1.61	-.54	-.51
6	62	23	13	-2.35	-2.48	-1.07	-1.02
7	62	24	19	-2.71	-2.34	-.52	-.45
8	62	24	20	-1.59	-1.22	.62	.69
9	58	23	15	-3.02	-2.27	-.67	-.46
10	58	18	14	2.69	.98	.64	.96
11	58	18	14	2.69	.98	.84	.96
12	58	17	13	3.09	.79	.28	.39
13	58	18	11	-.54	-2.33	-2.46	-2.34
14	58	19	12	-.79	-2.07	-1.81	-1.68
15	50	18	8	.55	.36	.50	.37
16	50	18	7	-.56	-.25	-.61	-.74
17	50	19	8	-.84	-.01	-.02	-.13
18	50	19	8	-.84	-.01	-.02	-.13
19	50	20	9	-1.13	.22	.58	.43
20	56	20	15	2.17	1.91	2.46	2.62
21	70	20	15	-8.45	—	—	—
Suma de Cuadrados Res.				175.642	104.206	45.240	23.131

tabla (2.5.15)

n	k	observaciones probadas	Q_k^{**}	F_k	valores críticos	observación eliminada
21	4	21, 4, 2, 20	159.34	.0928	.1348(1%)	21
20	3	4, 2, 20	87.93	.1562	.1815(1%)	4
19	2	2, 20	20.96	.36	.3386(5%)	

Q_k^{**} (i) puede obtenerse calculando la reducción en la suma de cuadrados resultante de eliminar k observaciones, o usando (2.3.37). Por ejemplo $Q_4^{**}(21, 4, 2, 20) = (-8.46)^2 + (7.68)^2 + (-4.7)^2 + (2.62)^2 = 159.34$, entonces, $F_4 = (175.642 - 159.34) / 175.642 = 0.0928$. Ya que F es significativo al 1% entonces la observación 21 es eliminada, y se calcula $Q_3^{**}(4, 2, 20) = (7.68)^2 + (-4.7)^2 + (2.62)^2 = 87.93$; $F_3 = (104.206 - 87.93) / 104.206 = 0.1562$ y así sucesivamente. Claramente las observaciones 21 y 4 son declaradas como discrepantes ya que el procedimiento termina en la tercera etapa. Entonces este análisis nos lleva a concluir que existen dos observaciones discrepantes en estos datos con respecto a este modelo.

Capitulo 3.

ENFOQUE BAYESIANO.

3.1. Introducción.

Como podemos recordar, en el enfoque frecuentista, las inferencias se hacen poniendo atención directa a un conjunto hipotético de datos x_1, x_2, \dots, x_n de una muestra aleatoria X_1, X_2, \dots, X_n que se supone son generados por el modelo $p(x/\theta)$ donde $\theta \in \Theta$ (el espacio paramétrico) y además $\theta = (\theta_1, \theta_2, \dots, \theta_r)$ es fijo y representa supuestamente a los valores de los parámetros desconocidos. Posteriormente, se seleccionan estimadores $\hat{\theta}$ que son funciones del vector de datos X y suponiendo que se puede repetir un número "grande" de veces esta situación, se obtienen valores $\hat{\theta}(x_1), \hat{\theta}(x_2), \dots$; ahora se pueden efectuar inferencias comparando los valores de $\hat{\theta}(x)$ con sus "distribuciones de muestreo" generadas al considerar muchas repeticiones. Las funciones $\hat{\theta}(x)$ usualmente se escogen de tal manera que las distribuciones de muestreo de los estimadores $\hat{\theta}(x_i)$ están, en algún sentido, lo más concentradas posible alrededor de los valores verdaderos θ .

En estadística Bayesiana, el modelo básico se entiende en otra dirección. Sea X_1, X_2, \dots, X_n una muestra aleatoria de tamaño n que proviene de una densidad $p(x/\theta_0)$ con θ_0 un valor de Θ . Entre las primeras diferencias que se presentan en este enfoque es que θ es considerada como una variable aleatoria cuya función de densidad $g(\theta)$, se supone conocida y que expresa el grado de certeza (degree of belief) o conocimiento inicial acerca de θ ; a $g(\theta)$ se le cono-

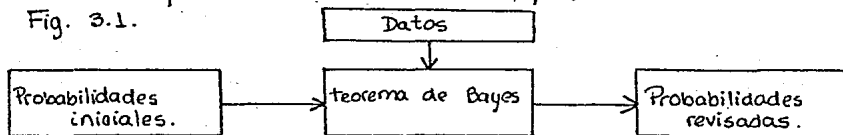
ce con el nombre de distribución a priori para θ o simplemente a priori. Mientras que en enfoque clásico, se considera a la función de verosimilitud como la única expresión que contiene toda la información, en estadística Bayesiana la expresión que nos proporciona la información contenida en la verosimilitud y en la función a priori es precisamente el teorema de Bayes.

$$\begin{aligned}
 p(\theta/x) &= \frac{f(x/\theta) \cdot g(\theta)}{f(x)} \\
 &= \frac{\left[\prod_{i=1}^n f(x_i/\theta) \right] g(\theta)}{\int \left[\prod_{i=1}^n f(x_i/\theta) \right] g(\theta) d\theta} \quad (3.1.1)
 \end{aligned}$$

Así, lo que distingue al esquema Bayesiano de otros enfoques estadísticos es que, antes de la obtención de la muestra, considera su grado de certeza para θ y lo representa en forma de probabilidad. Una vez obtenidos los datos, el teorema de Bayes permite calcular un nuevo conjunto de probabilidades que representan los grados de certeza revisadas en los modelos posibles, tomando en cuenta la nueva información proporcionada por los datos.

El proceso básico bajo la metodología bayesiana se resume esquemáticamente en la figura (3.1).

Fig. 3.1.



Como se ha mostrado inicialmente, las componentes que caracterizan a la estadística Bayesiana son: La información a priori, los datos muestrales, el cálculo de la densidad a posteriori para los parámetros y en algunas ocasiones el cálculo de la distribución predictiva de futuras observaciones. La información a priori está expresada por la densidad $p(\theta)$ del parámetro θ correspondiente al modelo de probabilidad $f(x/\theta)$, donde f es la densidad de la variable aleatoria x . La información en los datos x_1, x_2, \dots, x_n , está contenida en la función de verosimilitud $(\theta/x_1, x_2, \dots, x_n)$ que es la densidad conjunta de los datos muestrales que al combinarse con la densidad a priori de θ por medio del teorema de Bayes, se obtiene la densidad a posteriori de θ

$$P(\theta/x_1, x_2, \dots, x_n) = L(\theta/x_1, x_2, \dots, x_n) \cdot p(\theta)$$

a partir de la cual, se puede resolver el problema de inferencia.

Así la investigación se inicia con un modelo de probabilidad que nos dice como se distribuyen las observaciones $\underline{X} = x_1, x_2, \dots, x_n$; expresamos nuestra información con la densidad $p(\theta)$ (antes de obtener los datos) y una vez que se observa \underline{x} se calcula la distribución a posteriori de θ vía el teorema de Bayes.

La base de la inferencia Bayesiana es la distribución a posteriori, porque los aspectos restantes dependen de esta distribución. A partir de la distribución a posteriori se pueden hacer inferencias directamente, aunque también se pueden construir estimadores puntuales o de intervalos. Si θ es unidimensional, una gráfica de la distribución a posteriori de θ nos

revelara' el comportamiento de este parámetro ; si Θ es de más de una dimensión, se podrían aislar aquellas componentes de Θ que sean de interés.

Información a priori.

De entre todos los aspectos de la teoría bayesiana, éste es el más difícil y controversial ; muchos objetan tratar a los parámetros no observables Θ , como variables aleatorias y darles distribución de probabilidad subjetiva o densidades $f(\Theta)$. Sin embargo no mucha gente objeta darle a las cantidades observables X una distribución de probabilidad de tipo frecuentista ; suponer que la variable aleatoria X tiene distribución de frecuencia, es suponer la existencia de una sucesión interminable de repeticiones del experimento, y por supuesto, ya que tales repeticiones en efecto no ocurren, uno debe suponer que ocurren. Por lo tanto suponer que X tiene una distribución de probabilidad tipo frecuentista podría ser tan subjetivo como suponer que Θ tiene asociada una ley de probabilidad subjetiva, sin embargo, habría que aceptar que la ley de probabilidad de X es más objetiva en el sentido de que X puede observarse y su ley de probabilidad chequearse contra los datos. Apesar de las dificultades para proponer las distribuciones a priori, existen ciertos principios útiles para tal efecto como son : el principio de razón insuficiente, distribuciones a priori medición precisa ; principios que, a pesar de que se usan, no se describen en este trabajo.

Análisis posterior.

Supongamos que se está satisfecho con el modelo de probabilidad $f(\underline{x}/\theta)$ para las observaciones y también con la densidad a priori para los parámetros. Ahora en la etapa de análisis posterior, lo primero que debe de realizarse es encontrar la distribución a posteriori de θ

$$p(\theta/\underline{x}) \propto f(\underline{x}/\theta) \cdot P(\theta)$$

por medio del teorema de Bayes y dependiendo de cual sea el objetivo (estimar θ o alguna de sus componentes, probar hipótesis acerca de los parámetros, producir valores futuros de una observación etc) se procede al trabajo de inferencia tomando como base dicha distribución.

Aunque regularmente son de interés todas las componentes de θ , en algunas ocasiones se requieren solo algunos parámetros; sea θ_1 , un subconjunto de las componentes de θ , entonces para hacer inferencias acerca de θ_1 , se requiere calcular la densidad marginal a posteriori de θ_1

$$P_1(\theta_1/\underline{x}) = \int_{\mathcal{R}_2} P(\theta_1, \theta_2/\underline{x}) d\theta_2, \theta = (\theta_1, \theta_2)$$

naturalmente, se puede usar esta distribución para estimar y probar hipótesis concernientes a θ_1 , pero estos estimadores son características de resumen de $P(\theta_1/\underline{x})$ y solo nos dice parte de la historia acerca de los parámetros y no son un sustituto para la distribución total

3.2 Métodos Bayesianos Para el Manejo de Observaciones Discrepantes.

3.2.1. Método de Box-Tiao

George E. P. Box y George C. Tiao (1968) son de las primeras personas en formular un modelo general para la detección de observaciones discrepantes en modelos lineales. Su modelo es el de la varianza inflada y el objetivo fundamental de ésta es el de obtener inferencias acerca de β .

Consideremos el modelo de regresión lineal con las hipótesis usuales

$$Y = X\beta + \varepsilon,$$

cuyas dimensiones se han mencionado anteriormente. Supongamos que el error asociado con cada una de las observaciones surge de alguna de estas dos fuentes: un modelo central $N(0, \sigma^2)$ o del modelo alternativo $N(\theta, \sigma^2)$ con una probabilidad $1-\alpha$ y α respectivamente, α conocida y también que θ es fija, β y $\log \sigma$ son localmente independientes y uniformes a priori. Usando la misma notación empleada en el enfoque frecuentista, podemos descomponer a los elementos del modelo en $Y = (Y_i)$; $X = (X_{ij})$ y $\varepsilon = (\varepsilon_i)$. Definamos también E_i como el evento; k observaciones son discrepantes donde k se supone conocida; $i = \{i_1, i_2, \dots, i_k\}$. Dadas las observaciones Y , la verosimilitud del (β, σ, E_i) es

$$\begin{aligned} l(\beta, \sigma, E_i) &= f(\varepsilon_i | \sigma) \cdot g(\varepsilon_i | \sigma) \\ &\propto \prod (Y_{ij} - X_{ij}\beta / \sigma) \cdot g(X_i - X_i\beta / \sigma) \end{aligned} \quad (3.2.1)$$

donde f, g representan el producto de las funciones de -

densidad de $\xi_{(i)}$ y ξ_i respectivamente, E_i representa el evento, k de los errores surgen de $g(\xi_i/\sigma)$ y el resto de $f(\xi_{(i)}/\sigma)$ con

$$\left. \begin{aligned} f(\xi/\sigma) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{\xi^2}{2\sigma^2}\right\} \\ g(\xi/\sigma) &= \frac{1}{2\sigma\sqrt{2\pi}} \exp\left\{-\frac{\xi^2}{2\sigma^2}\right\} \end{aligned} \right\} \quad (3.2.2)$$

entonces la verosimilitud completa consiste de 2^n expresiones del tipo (3.2.1) correspondientes a todas las $\binom{n}{k}$ combinaciones posibles para los errores; $k = 0, 1, 2, \dots, n$.

Así, la distribución a posteriori para (β, σ, E_i) es

$$P(\beta, \sigma, E_i / \underline{y}) = \frac{P(E_i) \cdot P(\beta, \sigma) \cdot \prod (y_{(i)} - x_{(i)} \beta / \sigma) \cdot \prod g(y_i - x_i \beta / \sigma)}{\sum P(E_i) h(\underline{y}_i, \underline{y}_{(i)})} \quad (3.2.3)$$

donde $P(E_i) = \alpha^k (1 - \alpha)^{n-k}$ y tomando $P(\beta, \sigma) \propto \frac{1}{\sigma}$ con

$$h(\underline{y}_i, \underline{y}_{(i)}) = \int_{\mathcal{R}} P(\beta, \sigma) \cdot \prod (y_{(i)} - x_{(i)} \beta / \sigma) \cdot \prod g(y_i - x_i \beta / \sigma) \text{ donde}$$

es la distribución marginal de \underline{y} bajo el supuesto de que los errores surgen de las fuentes citadas. Después de integrar (3.2.3) con respecto a σ , la distribución de (β, E_i) puede escribirse como

$$P(\beta, E_i / \underline{y}) = P(E_i / \underline{y}) \cdot P(\beta / E_i, \underline{y})$$

donde

$$P(E_i / \underline{y}) = c \left(\frac{\alpha}{1 - \alpha} \right)^k \cdot \frac{1 \cdot 2^k \cdot \alpha^{1/2}}{1 \cdot 2^k \cdot \alpha - (1 - \frac{\alpha}{2}) \cdot \alpha^{1/2}} \cdot \left\{ \frac{1}{1 - \alpha} \right\}^{-\frac{1}{2}(n-p)} \quad (3.2.4)$$

$$\text{con } J_i^2 = \frac{1}{n-p} S_i(\hat{\beta}_i);$$

$$\hat{\beta}_i = (x_{(i)}^T x_{(i)} + \frac{1}{\alpha^2} x_i^T x_i)^{-1} (x_{(i)}^T y_{(i)} + \frac{1}{\alpha^2} x_i^T y_i),$$

c es una constante que hace que la suma de $P(E_i|Y)$ sea la unidad y

$$P(\beta/E_i, Y) \propto \int_0^\infty \sigma^{-(n+1)} \exp\left\{-\frac{1}{2\sigma^2} S_i(\beta)\right\} d\sigma \quad (3.2.5)$$

$$\propto \{S_i(\beta)\}^{-\frac{1}{2}n}$$

$$\text{donde } S_i(\beta) = (y_{(i)} - x_{(i)}\beta)^T (y_{(i)} - x_{(i)}\beta) + \frac{1}{\alpha^2} (y_i - x_i\beta)^T (y_i - x_i\beta) \quad (3.2.6)$$

$$= S_i(\hat{\beta}_i + (\beta - \hat{\beta}_i)^T (x_{(i)}^T x_{(i)} + \frac{1}{\alpha^2} x_i^T x_i)^{-1} (\beta - \hat{\beta}_i))$$

además si usamos la identidad (3.2.7)

$$(x_{(i)}^T x_{(i)} + \frac{1}{\alpha^2} x_i^T x_i)^{-1} = (x^T x - (1 - \frac{1}{\alpha^2}) x_i^T x_i)^{-1}$$

$$= (x^T x)^{-1} + (1 - \frac{1}{\alpha^2}) (x^T x)^{-1} x_i^T$$

$$\{I - (1 - \frac{1}{\alpha^2}) x_i (x^T x)^{-1} x_i^T\}^{-1} x_i (x^T x)^{-1}$$

entonces

$$\hat{\beta}_i = \hat{\beta} - (1 - \frac{1}{\alpha^2}) (x^T x)^{-1} x_i^T \{I - (1 - \frac{1}{\alpha^2}) x_i (x^T x)^{-1} x_i^T\}^{-1} (y_i - x_i \hat{\beta}) \quad (3.2.8)$$

$$y \quad (n-p) J_i^2 = (n-p) \alpha^2 - (1 - \frac{1}{\alpha^2}) (y_i - x_i \hat{\beta})^T \{I - (1 - \frac{1}{\alpha^2}) x_i (x^T x)^{-1} x_i^T\} (y_i - x_i \hat{\beta})$$

usando de (3.2.4) a (3.2.8) y despues de normalizar, la distribución a posteriori resulta

$$P(\beta/y | E_i) = \frac{P(\frac{n}{2}) / |X^T X - (1 - \frac{1}{a_2}) X_i^T X_i|^{n/2}}{P(\frac{n-p}{2}) (\pi(n-p) \hat{\Delta}_i^2)^{n/2}} \left\{ 1 + \frac{(\beta - \hat{\beta}_i)^T X_i^T X_i - (1 - \frac{1}{a_2}) X_i^T X_i (\beta - \hat{\beta}_i)}{(n-p) \hat{\Delta}_i^2} \right\}^{-\frac{n}{2}} \quad (3.2.9)$$

que es una distribución t multivariada (p -dimensional) con $n-p$ grados de libertad, media $\hat{\beta}_i$ y matriz de dispersión $\hat{\Delta}_i^2 (X_i^T X_i - (1 - \frac{1}{a_2}) X_i^T X_i)^{-1}$. En particular, si $k=0$, la expresión (3.2.9) se reduce a la distribución t ordinaria con $n-p$ g.l. centrada en el estimador mínimos cuadrados $\hat{\beta}$ y matriz de dispersión $\hat{\Delta}^2 (X^T X)^{-1}$. Entonces, la distribución a posteriori para β es un promedio ponderado de distribuciones t-student, es decir

$$P(\beta/y) = \sum_i P(E_i/y) \cdot P(\beta/E_i, y) \quad (3.2.10)$$

donde los pesos $w_i = P(E_i/y)$

Distribución marginal de β_j . Como se mencionó, el objetivo fundamental de este modelo es básicamente lograr inferencias acerca de los parámetros, entonces, si se cuenta con una distribución t-student, (3.2.4), (3.2.9) y (3.2.10); la distribución marginal de la j -ésima componente de β es

$$P(\beta_j/y) = \sum_i w_i P_i(\beta_j/y)$$

$$\text{con } P_i(\beta_j/y) = \hat{\Delta}_i^{-1} (u_i^{jj})^{-\frac{1}{2}} P \left\{ t_{n-p} = \frac{\hat{\beta}_i - \hat{\beta}_{ij}}{\hat{\Delta}_i (u_i^{jj})^{1/2}} \right\},$$

$\hat{\beta}_{ij}$ el j -ésimo elemento de $\hat{\beta}_i$ y u_i^{jj} el j -ésimo elemento de la matriz $(X_i^T X_i - (1 - \frac{1}{a_2}) X_i^T X_i)^{-1}$. Se puede verificar que

$$\bar{\beta}_j = E(\beta_j/y) = \sum_i w_i \hat{\beta}_{ij} \quad y$$

$$\text{Var}(\beta_j | y) = \sum_j w_j \left\{ \frac{\sigma^2}{n-p-2} \Delta_j^2 \sigma_j^{2n} + (\hat{\beta}_j - \bar{\beta}_j)^2 \right\}$$

Mezcla de distribuciones. En el desarrollo del modelo se ha supuesto que el error asociado con cada observación puede surgir de cualquiera de las distribuciones expresadas en (3.2.2) con probabilidad $(1-\alpha)$ y α respectivamente, obteniéndose que la distribución a posteriori para β es un promedio ponderado de 2^n distribuciones t-multivariadas. Dicho resultado, también se puede obtener de la siguiente manera: como

$$f(\xi/\sigma) = (1-\alpha) \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{\xi^2}{2\sigma^2}\right) + \alpha \frac{1}{\alpha \sigma \sqrt{2\pi}} \exp\left(-\frac{\xi^2}{2\alpha^2 \sigma^2}\right)$$

es una mezcla de dos distribuciones normales que contiene 3 parámetros: (α, σ, a) , entonces, la función de verosimilitud de $(\beta, \sigma, a, \alpha)$ es

$$l(\beta, \sigma, a, \alpha) \propto \sigma^{-n} \prod_{i=1}^n \left[(1-\alpha) \exp\left\{-\frac{1}{2\sigma^2} (y_i - x_i^T \beta)^2\right\} + \alpha a^{-1} \exp\left\{-\frac{1}{2\alpha^2 \sigma^2} (y_i - x_i^T \beta)^2\right\} \right] \quad (3.2.10)$$

donde x_i es el i -ésimo renglón de X . Adoptando $P(\beta, \sigma) \propto \frac{1}{\sigma^2}$, entonces, para a, α fijas, la distribución a posteriori para β es

$$P(\beta | a, \alpha, y) = \delta \int_0^\infty \sigma^{-1} l(\beta, \sigma, a, \alpha | y) d\sigma$$

$$\text{con } \delta^{-1} = \int_{\mathbb{R}} \sigma^{-1} l(\beta, \sigma, a, \alpha | y) d\sigma,$$

Esta distribución es matemáticamente equivalente a la distribución (3.2.10). En efecto, al exponer el integrando en (3.2.11) y efectuando la integración término a término, se

obtiene exactamente la expresión (3.2.10)

En aplicaciones prácticas de la distribución a posteriori de θ en (3.2.10), es adecuado, reacomodar la suma de las 2ⁿ posibilidades de la siguiente forma:

$$P(\theta/Y) = w_0 \cdot p_0(\theta/Y) + \sum_{i=1}^n w_i \cdot P_i(\theta/Y) + \sum_{i,j}^c w_{ij} \cdot P_{ij}(\theta/Y) + \dots \quad (3.2.12)$$

La distribución $p_0(\theta/Y)$ en la expresión anterior sería apropiada si todas las observaciones surgen del modelo f de la expresión (3.2.2); La distribución $P_i(\theta/Y)$ corresponde a la posibilidad de que cada observación en turno sea discrepante y así sucesivamente. Similarmente, w_0 es la probabilidad a posteriori de que ninguna observación sea mala; w_i representa a la probabilidad de que la i -ésima observación sea discrepante etc..

3.2.2. Método de Abraham - Box

Al igual que el método de Box-Tiao, B. Abraham y G. Box (1978) desarrollan un método cuyo objetivo fundamental también es el de hacer inferencias acerca de los parámetros en modelos lineales en presencia de observaciones discrepantes. Su modelo, considera una fuente distinta de contaminación

$$Y = X\beta + \lambda z + \epsilon \quad (3.2.13)$$

donde z es un vector de $n \times 1$, cuyos elementos tienen una probabilidad α de ser 1 y $1-\alpha$ de ser 0, la cantidad de contaminación λ se supone la misma para todas las observaciones discrepantes y los demás elementos del modelo están contenidos en el modelo de regresión lineal con las hipótesis usuales. Naturalmente $z^t = \{z_1, z_2, \dots, z_n\}$ es un vector de k unidades y $n-k$ ceros, tal que

$$z_i = \begin{cases} 1 & \text{si la } i\text{-ésima observación es discrepante} \\ 0 & \text{si no lo es.} \end{cases}$$

Podemos observar que el modelo que se maneja en este método, es el modelo de la media trasladada cuya representación geométrica se encuentra reflejada en la figura 2.1.2. a

z conocida. Debe notarse que suponer que z es conocida, implica que $i = \{i_1, i_2, \dots, i_k\}$ es conocida y por lo tanto se conoce k , el número de observaciones discrepantes.

partes. La distribución a posteriori de β dado Y y Z es

$$P(\beta | Y, Z) = \int P(\beta, \lambda, \sigma | Z) \cdot P(Y | \beta, \lambda, \sigma, Z) d\lambda d\sigma. \quad (3.2.14)$$

donde $P(\beta, \lambda, \sigma | Z)$ es la a priori de $(\beta, \lambda, \sigma | Z)$ y $P(Y | \beta, \lambda, \sigma, Z)$ es la distribución conjunta de Y dado $(\beta, \lambda, \sigma, Z)$. Esta distribución conjunta puede escribirse como

$$P(Y | \beta, \lambda, \sigma, Z) = (2\pi\sigma^2)^{-\frac{1}{2}n} \exp\left\{-(2\sigma^2)^{-1} \left\{ (Y - X\beta)^t (Y - X\beta) + k\lambda^2 - 2(Y - X\beta)^t Z\lambda \right\} \right\} \quad (3.2.15)$$

donde $k = Z^t Z$. Suponiendo que $k \neq 0$, n y que $X^t X$ es no singular, entonces (3.2.15) puede reescribirse como

$$P(Y | \beta, \lambda, \sigma, Z) = (2\pi\sigma^2)^{-\frac{1}{2}n} \exp\left[-(2\sigma^2)^{-1} \left\{ (\beta - \beta^*)^t (X^t X) (\beta - \beta^*) + k(\lambda - \lambda_0)^2 + vS_k^2 \right\} \right] \quad (3.2.16)$$

donde $I_1 = I - \frac{1}{k} Z Z^t$; $\lambda_0 = \frac{1}{k} (Y - X\beta)^t Z$; $\beta^* = (X^t X)^{-1} X^t I_1 Y$ y $vS_k^2 = Y^t I_1 (I - X(X^t I_1 X)^{-1} X^t) I_1 Y$.

Supongamos que β , λ , $\log \sigma$ son uniformes a priori y localmente independientes, es decir,

$$P(\beta, \lambda, \sigma | Z) \propto \frac{1}{\sigma}. \quad (3.2.17)$$

Ahora, multiplicando (3.2.16) por (3.2.17) e integrando con respecto a λ y σ obtenemos

$$P(\beta | Y, Z) \propto \left(\frac{1 + (\beta - \beta^*)^t (X^t I_1 X) (\beta - \beta^*)}{v S_k^2} \right)^{-\frac{1}{2}(v+p)}; \text{ donde } v = n - p - 1. \quad (3.2.18)$$

Entonces, si tenemos un conjunto de k -observaciones dis-
crepantes, la distribución a posteriori de β , será una distribu-
ción t -dimensional con media β^* y matriz de dispersión

$S_k^2 (x^t V x)^{-1}$ con $v = n - p - 1$ g. l.

\underline{z} Desconocida. Supongamos que existe una probabilidad α (conocida) de que cualquier elemento de \underline{z} sea la unidad y $1 - \alpha$ de que sea cero. Si V_k denota un vector de $n \times 1$ con k unidades y $(n - k)$ ceros, entonces, a priori tenemos

$$P(\underline{z} = V_k / \alpha) = \alpha^k (1 - \alpha)^{n - k} \quad k = 0, 1, 2, \dots, n.$$

De aquí en adelante se escribirá \underline{z} para indicar $\underline{z} = V_k$ y se evitará la condicional sobre α

Dado \underline{y} la distribución a posteriori de $(\beta, \lambda, \sigma, \underline{z})$ está dada por

$$P(\beta, \lambda, \sigma, \underline{z} / \underline{y}) = \frac{P(\underline{z}) \cdot P(\beta, \lambda, \sigma) \cdot P(\underline{y} / \beta, \lambda, \sigma, \underline{z})}{\sum_{\underline{z}} P(\underline{z}) \cdot P(\underline{y} / \underline{z})} \quad (3.2.19)$$

donde $P(\underline{y} / \beta, \lambda, \sigma, \underline{z})$ está dada en (3.2.16) y

$$P(\underline{y} / \underline{z}) = \int P(\beta, \lambda, \sigma / \underline{z}) P(\underline{y} / \beta, \lambda, \sigma, \underline{z}) d\lambda d\beta d\sigma$$

La suma en (3.2.19) es sobre los 2^n posibles valores de \underline{z} . Integrando (3.2.19) con respecto a λ y σ se obtiene

$$P(\beta, \underline{z} / \underline{y}) = P(\underline{z}) \cdot P(\beta / \underline{y}, \underline{z}),$$

donde $P(\beta / \underline{y}, \underline{z})$ está dada en (3.2.18) y

$$P(\underline{z} / \underline{y}) = P(\underline{z}) \cdot P(\underline{y} / \underline{z}) / \sum_{\underline{z}} P(\underline{z}) \cdot P(\underline{y} / \underline{z})$$

$$= c \{P(\underline{z}) / P(\underline{z}=0)\} \{P(\underline{y} / \underline{z}=0)\} \quad (3.2.20)$$

con $c = P(\underline{z}=0) P(\underline{y} / \underline{z}=0) / \sum P(\underline{z}) P(\underline{y} / \underline{z})$,

por lo tanto $P(\underline{\beta} / \underline{y}) = \sum P(\underline{z} / \underline{y}) P(\underline{\beta} / \underline{y} \underline{z})$

donde $P(\underline{\beta} / \underline{y}, \underline{z})$ esta dada en (3.2.18) y $P(\underline{z} / \underline{y})$ en (3.2.20) que tambien se puede escribir como

$$P(\underline{z} / \underline{y}) = c \left(\frac{\underline{z}}{\underline{1-z}} \right)^k \left(\frac{\underline{z} \underline{1}^T \underline{z}_1 \underline{x}_1}{\underline{z}_0 \underline{z}_1^T \underline{x}_1^T \underline{x}_1} \right)^{-\frac{1}{2}} \left(\frac{\underline{v} \underline{z}_1^T}{\underline{z}_0 \underline{z}_1^T} \right)^{-\frac{1}{2}} \underline{v} B \left(\frac{1}{2}, \frac{1}{2} \underline{v} \right).$$

con c una constante, tal que, $\sum P(\underline{z} / \underline{y}) = 1$. Entonces la distribucion a posteriori de $\underline{\beta}$ dadas las observaciones es un promedio ponderado de distribuciones t .

Inferencias acerca de λ . El modelo que se ha considerado aqui, es una mezcla de dos distribuciones normales, donde solo una de ellas se considera correcta. Este modelo tiene implicaciones en el caso donde se sospecha de observaciones "malas" y se esta interesado en encontrar, cuanto se desvian dichas observaciones de la media de las observaciones buenas.

Como $P(\underline{\beta}, \lambda, \sigma) \propto \frac{1}{\sigma}$ y

$$P(\underline{y} / \underline{\beta}, \lambda, \sigma, \underline{z}) = (2\pi\sigma^2)^{-\frac{1}{2}n} \exp\{(2\sigma^2)^{-1} [\underline{y} - \underline{x}\underline{\beta}] + k\lambda^2 - z(\underline{y} - \underline{x}\underline{\beta})^T \underline{z}\lambda\}$$

y suponiendo que $k \neq 0$ y $\underline{z}^T M \underline{z} \neq 0$; $M = \underline{I} - \underline{x}(\underline{x}^T \underline{x})^{-1} \underline{x}^T$,

$$\text{entonces } P(\underline{\beta}, \lambda, \sigma / \underline{y}, \underline{z}) = (2\pi)^{-\frac{1}{2}n} \sigma^{-(n+1)} \exp\{-(2\sigma^2)^{-1} [(\underline{\beta} - \underline{\beta}_0)^T \underline{x}^T \underline{x} (\underline{\beta} - \underline{\beta}_0) + (\underline{z}^T M \underline{z}) (\lambda - \lambda^*)^2 + \underline{v} \underline{z}_1^T \underline{d}]\},$$

con $\hat{\beta}_n = (x^t x)^{-1} (Y - \lambda z)$, $\lambda^* = (z^t M z)^{-1} z^t M Y$ y

$$v S_d^2 = Y^t M [I - z (z^t M z)^{-1} z^t] M Y.$$

Integrando la expresión (3.2.21) con respecto a β y σ obtenemos

$$P(\lambda | Y, z) d\lambda \propto \left(1 + \frac{z^t M z (\lambda - \lambda^*)^2}{v S_d^2} \right)^{-\frac{v+1}{2}},$$

que es una distribución "escalada" t-student.

3.2.2 Método de Guttman, Freeman y Dutter

I. Guttman, P.R. Freeman y R. Dutter (1978) plantean un modelo más general que los anteriores para la detección de observaciones discrepantes. Ellos suponen que k de las observaciones son generadas por el modelo

$$Y_{ij} = x_{ij}^t \beta + \lambda_j + \epsilon_{ij} \quad ; \quad j = 1, 2, \dots, k \quad (3.2.22)$$

con $i_j \in i = \{i_1, i_2, \dots, i_k\}$ y el número restante $(n-k)$ son generadas por el modelo tradicional

$$Y = X \beta + \epsilon. \quad (3.2.23')$$

Podemos observar que el modelo propuesto, es el modelo de la media trasladada donde las observaciones discrepantes surgen de diferentes poblaciones, situación reflejada geoméricamente en la figura 3.2.1 b.

Supongamos que cualquier subconjunto de k -observaciones - puede ser generado por el modelo (3.2.22') con una probabilidad uniforme a priori de $1/\binom{n}{k}$ y que la información

a priori acerca de β , σ^2 y λ es muy vaga y por lo tanto

$$p(\lambda, \beta, \sigma^2) \propto \frac{1}{\sigma^2}. \quad (3.2.23)$$

entonces la verosimilitud de λ , β y σ^2 dado \underline{y} es

$$l(\lambda, \beta, \sigma^2 | \underline{y}) \propto \sum_i \frac{1}{(n_i)} (\sigma^2)^{-(n_i)/2} \exp\left\{-\frac{1}{2\sigma^2} [(y_{(i)} - X_{(i)}\beta)^T (y_{(i)} - X_{(i)}\beta)]\right\} \times \quad (3.2.24)$$

$$\left\{ (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} [(y_i - x_i\beta - \lambda)(y_i - x_i\beta - \lambda)]\right\} \right\}$$

combinando (3.2.23) con (3.2.24) y usando el teorema de Bayes, encontramos que la distribución a posteriori de $(\lambda, \beta, \sigma^2)$ dado \underline{y} es

$$p(\lambda, \beta, \sigma^2 | \underline{y}) = \sum_i (\sigma^2)^{-\frac{n}{2} + 1} \exp Q, \quad (3.2.25)$$

donde

$$Q = S_{(i)} + (x_{(i)} \bar{x}_{(i)} y_{(i)} - x_{(i)} \bar{x}_{(i)} \beta)^T (x_{(i)} \bar{x}_{(i)} y_{(i)} - x_{(i)} \bar{x}_{(i)} \beta) + [\lambda - (y_i - x_i \beta)]^T [\lambda - (y_i - x_i \beta)] \quad (3.2.26)$$

$$\text{con } S_{(i)} = (y_{(i)} - X_{(i)} \bar{x}_{(i)} y_{(i)})^T (y_{(i)} - X_{(i)} \bar{x}_{(i)} y_{(i)})$$

y $\bar{x}_{(i)}$ satisface

$$\left. \begin{aligned} a) \quad (x_{(i)} \bar{x}_{(i)})^T &= x_{(i)} \bar{x}_{(i)}, \\ b) \quad x_{(i)} \bar{x}_{(i)} x_{(i)} &= x_{(i)}, \end{aligned} \right\} \quad (3.2.27)$$

Si $X_{(i)}$ es el rango de p , entonces $\bar{x}_{(i)} = (X_{(i)}^T X_{(i)})^{-1} X_{(i)}^T$ y $S_{(i)}$ es la suma de cuadrados residual usual

Para encontrar la distribución marginal a posteriori para β , debemos integrar (3.2.25) con respecto a λ y σ^2 .

$$P(\beta, \sigma^2 / y) \propto \prod_i (\sigma^2)^{-\frac{(n_i+1)}{2}} \exp \frac{1}{2\sigma^2} Q \cdot \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \exp \frac{1}{2\sigma^2} (Q - \varepsilon_i)(Q - \varepsilon_i) d\varepsilon_1, \dots, d\varepsilon_n$$

$$\propto \prod_i (\sigma^2)^{-\frac{(n_i+1)}{2}} \exp \left(-\frac{1}{2\sigma^2} Q_{(i)} \right) \quad (3.2.2)$$

con $Q_{(i)} = S_{(i)} + (\beta - \hat{\beta}_{(i)})^t X_{(i)}^t X_{(i)} (\beta - \hat{\beta}_{(i)})$ y $\hat{\beta}_{(i)}$ satisfice:

$$X_{(i)}^t X_{(i)} \hat{\beta}_{(i)} = X_{(i)} Y_{(i)} \quad \beta_{(i)} = X_{(i)}^{-1} Y_{(i)},$$

ahora, integrando (3.2.2) con respecto a σ^2 , se obtiene

$$P(\beta / y) \propto \prod_i [S_{(i)} + (\beta - \hat{\beta}_{(i)})^t X_{(i)}^t X_{(i)} (\beta - \hat{\beta}_{(i)})]^{-\frac{n_i+1}{2}}$$

$$\propto \prod_i \left\{ S_i^{-\frac{n_i+1}{2}} [1 + (\beta - \hat{\beta}_{(i)})^t (S_{(i)}^{-1} (X_{(i)}^t X_{(i)})) (\beta - \hat{\beta}_{(i)})] \right\} \quad (3.2.29)$$

Como se puede observar, $P(\beta / y)$ es un promedio ponderado de 2^n distribuciones t de student.

$$t(w / w_0; B, m; p) = \frac{\Gamma(\frac{m-p}{2}) |B|^{1/2}}{(\Gamma(\frac{1}{2}))^p (\Gamma(\frac{m}{2}))^{p/2} \Gamma(\frac{m}{2})} \left[1 + \frac{(w - w_0)^t B (w - w_0)}{m} \right]^{-\frac{m-p}{2}} \quad (3.2.30)$$

es decir, si suponemos que $X_{(i)}^t X_{(i)}$ es positiva definida para todo $i = i_1, i_2, \dots, i_k$, podemos escribir (3.2.29) como.

$$P(\beta / y) = \sum_i w_i t(\beta / \hat{\beta}_{(i)}; \frac{n_i - k - p}{S_{(i)}} X_{(i)}^t X_{(i)}; n_i - k - p; p)$$

$$\text{donde } w_i = \left[S_{(i)}^{-(n_i-k)} (S_{(i)}^p / (X_{(i)}^t X_{(i)})^{-1}) \right]^{1/2} / \sum_i \left[S_{(i)}^{-(n_i-k)} (S_{(i)}^p / (X_{(i)}^t X_{(i)})^{-1}) \right]^{1/2}$$

son las probabilidades a posteriori de que las observaciones X_1, \dots, X_k sean discrepantes.

3.2.4. Método de Dutter y Guttman

Este trabajo también utiliza el modelo de la media trasladada

$$y = x\beta + \lambda \underline{z} + \underline{\varepsilon},$$

que utilizan Abraham y Box (1978), con la diferencia de que \underline{z} toma el valor uno o el valor cero con la misma probabilidad. Entonces, para este modelo con k observaciones trasladadas de la media (k fija y conocida) por una cantidad λ (desconocida) la verosimilitud de $(\beta, \sigma^2, \lambda)$ dadas las observaciones es

$$l(\beta, \sigma^2, \lambda / \underline{y}) = \sum_i \left\{ \prod_{i \in i} \left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left[-\frac{1}{2\sigma^2} (y_i - z_i \beta)^2 \right] \right\} \left\{ \prod_{i \in i} \left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left[-\frac{1}{2\sigma^2} (y_i - z_i \beta - \lambda)^2 \right] \right\}, \quad (3.2.31)$$

donde i es el conjunto de índices que identifican a las k observaciones discrepantes y que ha sido definido anteriormente. Bajo el supuesto de que k observaciones están trasladadas y que β , $\log \sigma^2$ y λ son localmente uniformes e independientemente distribuidas se tiene que la distribución a priori de $(\beta, \sigma^2, \lambda)$ es proporcional a $\frac{1}{\sigma^2}$. Con esta información se obtiene la distribución a posteriori de $(\beta, \sigma^2, \lambda)$ como:

$$\begin{aligned} p(\beta, \sigma^2, \lambda / \underline{y}) &\propto \sum_{i \in i} \frac{1}{\sigma^{k+2}} \exp \left\{ -\frac{1}{2} \left[\sum_i (y_i - z_i \beta)^2 + \sum_i (y_i - z_i \beta - \lambda)^2 \right] \right\} \\ &= \sum_i \frac{1}{\sigma^{k+2}} \exp \left\{ -\frac{1}{2\sigma^2} [S_i + (\beta - \hat{\beta}_i)' X' I_i X (\beta - \hat{\beta}_i) + k(\lambda - \frac{1}{k} (Y - X\hat{\beta})' \underline{z})^2] \right\} \end{aligned} \quad (3.2.32)$$

con \underline{z}_i' el i -ésimo renglón de la matriz X ,

$$S_i = (Y - X\hat{\beta})' I_i (Y - X\hat{\beta}),$$

$$I_i = I - \frac{1}{k} \mathbf{1} \mathbf{1}^T,$$

$$\hat{\beta}_i = (x^T I_i x)^{-1} x^T I_i y.$$

Se puede demostrar que I_i es simétrica e idempotente, y que $x^T I_i x$ es positiva definida.

A partir de la expresión (3.2.31) se puede encontrar la distribución marginal de β dadas las observaciones al integrar con respecto a λ y σ^2 para obtener que:

$$\begin{aligned} P(\beta/y) &\propto \sum_i [S_i + (\beta - \hat{\beta}_i)^T x^T I_i x (\beta - \hat{\beta}_i)]^{-\frac{n_i}{2}} \\ &\propto \sum_i S_i^{-\frac{n_i}{2}} \left[1 + \frac{1}{S_i} (\beta - \hat{\beta}_i)^T x^T I_i x (\beta - \hat{\beta}_i) \right]^{-\frac{n_i}{2}} \\ &= \sum_i W_i t_i(\beta/\beta_i; \frac{S_i}{n-p-1} (x^T I_i x)^{-1}, n-p-1) \end{aligned} \quad (3.2.33)$$

$$\text{donde } t_p(\beta, \hat{\beta}, \sigma^2 (x^T x)^{-1}, \nu) = \frac{\Gamma(\frac{n-p}{2}) / (x^T x)^{p/2}}{[\Gamma(\frac{n}{2})]^p \Gamma(\frac{\nu}{2}) (\Delta \sqrt{\nu})^p} \left[1 + \frac{(\beta - \hat{\beta})^T (x^T x)^{-1} (\beta - \hat{\beta})}{\Delta^2 \nu} \right]^{-\frac{\nu+p}{2}},$$

$$\text{y } W_i = \frac{(S_i)^{-\frac{n-p-1}{2}}}{|x^T I_i x|^{1/2}} \bigg/ \sum_i \frac{(S_i)^{-\frac{n-p-1}{2}}}{|x^T I_i x|^{1/2}}$$

Como se puede ver $P(\beta/y)$ es una suma ponderada de $\binom{n}{k}$ distribuciones t multivariadas.

Aunque en algunos casos, se requieren algunas modificaciones se encuentran las distribuciones marginales a posteriori de σ^2 y λ con los siguientes resultados:

$$P(\sigma^2/y) = \sum_i W_i f(\sigma^2/S_i, n-p-1) \quad (3.2.34)$$

donde f denota la densidad de la χ^2 invertida con $(n-p-1)$ grados de libertad (Box-Tiao 1978).

$$f(\sigma^2/\gamma, v) = \left[\Gamma\left(\frac{v}{2}\right) 2^{v/2} \right]^{-1} (\sigma^2)^{-(\frac{v}{2}+1)} (\gamma)^{v/2} \exp\left\{-\frac{\gamma}{2\sigma^2}\right\}$$

y

$$P(\lambda/\underline{y}) = \sum_i W_i t_i(\lambda/\frac{1}{m_i} \underline{y}^T M \underline{z}_i, \frac{V_i}{(n-p-1)}, n-2p) \quad (3.235)$$

$$\text{con } m_i = \underline{z}_i^T M \underline{z}_i,$$

$$M = I - X(X^T X)^{-1} X^T,$$

$$V_i = \underline{y}^T M \left(I - \frac{1}{m_i} \underline{z}_i \underline{z}_i^T \right) M \underline{y}$$

$$y \quad W_i = \frac{[V_i]^{-\frac{n-p-1}{2}}}{[m_i]^{1/2}} \Bigg/ \sum_i \frac{[V_i]^{-\frac{n-p-1}{2}}}{[m_i]^{1/2}}.$$

Es claro que $P(\lambda/\underline{y})$ también es una suma ponderada de distribuciones t - univariadas.

3.3 Comentarios.

Como podemos observar, los métodos Bayesianos que aquí se presentan modelan ligeramente diferente el aspecto de observaciones discrepantes. Box y Tiao (1968), suponen que la distribución del error es del tipo

$$F = (1-\alpha) N(0, \sigma^2) + \alpha N(0, a^2 \sigma^2),$$

$a^2 > 1$; a y a^2 son valores que, dependiendo de sus intereses, el investigador debe especificarlos. Por otro lado, Abraham y Box (1978) aplican estas ideas al modelo de la media trasladada con la restricción que todas las observaciones contaminantes provienen de la misma población. Este modelo es casi idéntico al de Dutter y Guttman con la única diferencia de que en este modelo todas las observaciones tienen la misma probabilidad de ser disorepantes. Esta última característica también es común al modelo de Guttman, Freeman y Dutter (que también utiliza el modelo de la media trasladada) pero la diferencia radica en que el traslado (shift) para cada observación no es el mismo.

Al aplicar el paradigma Bayesiano para cada situación descrita, y dado que el aspecto principal es la estimación de los parámetros, observamos que en todos los casos, la distribución de θ dado Y es una suma ponderada de distribuciones t multivariadas. ver ecuaciones (3.2.10), (3.2.18), (3.2.29), (3.2.33).

Cada uno de estos enfoques puede considerarse como un método de acomodo, cuya base es la distribución a posteriori respectiva de θ o un método de identificación basado en la probabilidad a posteriori W_i , donde i identifica a las observaciones contaminantes. Por ejemplo, Box y tián sugieren solo incluir términos seleccionados en la expresión (3.2.12) y "renormalizar" los pesos: Si no se sospecha de observaciones disorepantes, entonces solo el primer término de dicha expresión es relevante y $P(\theta/Y) = P_0(\theta/Y)$; si solo se sospecha de a lo más una observación disorepante, $P(\theta/Y) = W_0 P_0(\theta/Y) + \sum_{i=1}^n W_i P_i(\theta/Y)$; si se sospecha de dos contaminantes, se agrega el término $\sum_{i,j} W_{ij} P_{ij}(\theta/Y)$ y así sucesivamente. En este enfoque, el problema de especificar el número de observaciones disorepantes se transporta en escoger el

número de términos para $P(B|Y)$.

Por otro lado, el uso de cualquiera de los modelos expuestos dependerá de los intereses del investigador.

3.4 Ejemplos.

Con la finalidad de mostrar algunos de los hechos mencionados anteriormente con respecto a los métodos Bayesianos, en esta sección únicamente se desarrollará un ejemplo en el cual se usará el Método de Ditter y Guttman.

Los datos que manejaremos en este ejemplo, están contenidos en la tabla (3.4.1) y se obtuvieron por medio de un generador normal pseudo-aleatorio tal que la distribución de las primeras 8 observaciones es $N(0, 1)$, mientras que las dos últimas están trasladadas con $\mu=2$.

tabla 3.4.1

i	y_i
1	-0.44168
2	0.80814
3	-0.01549
4	0.01265
5	-1.08758
6	-1.38324
7	0.48106
8	-0.59508
9	2.32219
10	2.05387

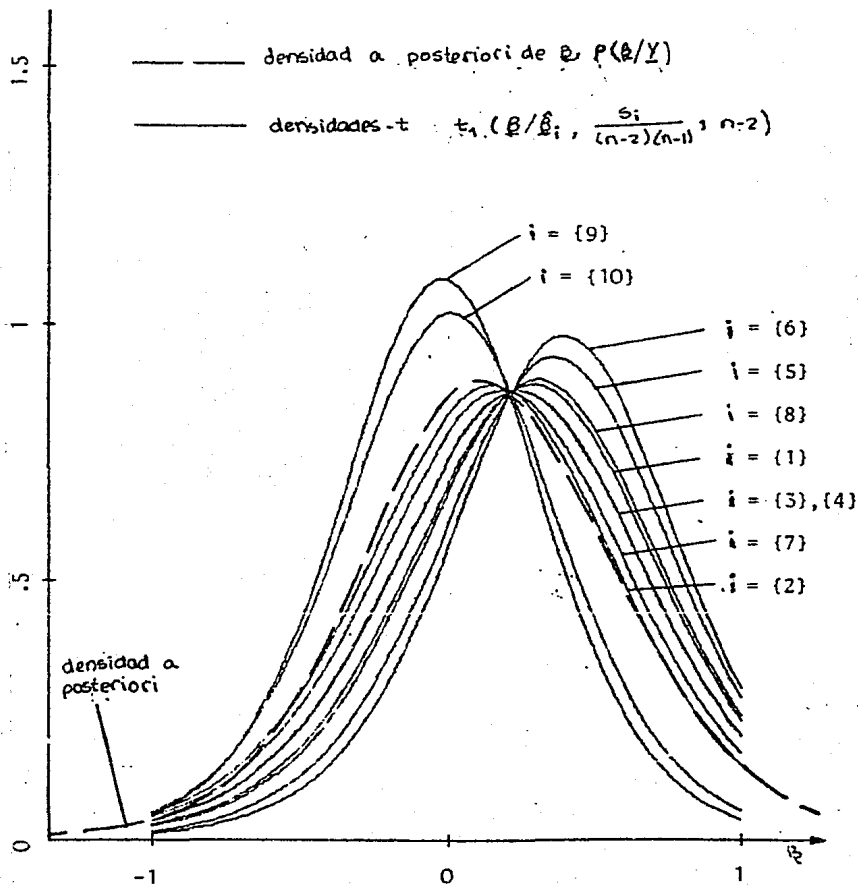
Sin conocimiento de la forma en que se generan los datos y si suponemos que hay una observación discrepante, entonces, al concentrar nuestra atención en la figura (3.4.1) observamos que precisamente las densidades t correspondientes a estas dos últimas observaciones están separadas claramente del resto de las densidades, mas aún, $P(\hat{\theta}/Y)$ también se encuentra dominada por estas dos densidades, de tal forma que, podríamos pensar que estas observaciones son discrepantes. Este hecho también puede constatarse si observamos los valores de β_{ci} , W_i y S_i que se encuentran en la tabla (3.4.2). Como podemos notar, las observaciones con mayor peso también lo son la 9 y 10 para las cuales se obtiene una mayor reducción en la suma de cuadrados al eliminar dichas observaciones.

tabla 3.4.2.

$k = 1$			
$i=i$	β_{ci}	S_i	W_i
1	.2887	13.19	.0562
2	.1499	13.28	.0547
3	.2414	13.61	.0496
4	.2382	13.63	.0494
5	.3605	11.79	.0883
6	.3933	10.83	.1237
7	.1862	13.60	.0499
8	.3056	12.95	.0606
9	-.0184	8.74	.2916*
10	.0114	9.92	.1760*

* pesos dominantes

figura 3.4.1.



Igualmente, para $k=2$, la tabla (3.4.3) nos proporciona los valores $\hat{e}(i)$, \hat{e}_i , S_i y w_i para las (10^2) posibles parejas. Los pesos dominantes son para $i = \{9, 10\}$; $w_i = 0.6489$ y $i = \{5, 6\}$; $w_i = 0.032$ que es más de 20 veces menor que $w(9, 10)$. Por estas razones, a las observaciones 9 y 10 se les declara disrepan-
tes.

tabla 3.4.3.

		k = 2			
	i	$\hat{e}(i)$	\hat{e}_i	S_i	w_i
1	1 2	.2258	0.183	13.67	.0045
2	1 3	.3268	-0.229	13.18	.0052
3	1 4	.3232	-0.215	13.21	.0052
4	1 5	.4608	-0.765	11.27	.0098
5	1 6	.4977	-0.912	10.49	.0130
6	1 7	.2647	0.020	13.58	.0046
7	1 8	.5990	-0.517	12.33	.0068
8	1 9	.0345	0.940	12.36	.0068
9	1 10	.0681	0.806	12.80	.0059
10	2 3	.1705	0.396	13.59	.0046
11	2 4	.1670	0.410	13.58	.0046
12	2 5	.3045	-0.140	13.36	.0050
13	2 6	.3415	-0.288	13.04	.0055
14	2 7	.1085	0.645	13.21	.0052
15	2 8	.2427	0.108	13.64	.0046
16	2 9	-.1217	1.565	9.12	.0228
17	2 10	-.0882	1.431	9.98	.0159
18	3 4	.2700	-0.001	13.56	.0047
19	3 5	.4075	-0.552	12.20	.0071
20	3 6	.4445	-0.699	11.58	.0088
21	3 7	.2114	0.233	13.67	.0045
22	3 8	.3457	-0.304	13.00	.0055
23	3 9	-.0187	1.153	11.48	.0091
24	3 10	.0148	1.019	12.06	.0075
25	4 5	.4040	-0.537	12.26	.0070
26	4 6	.4409	-0.685	11.64	.0086
27	4 7	.2079	0.247	13.67	.0045
28	4 8	.3422	-0.290	13.03	.0055
29	4 9	-.0223	1.167	11.41	.0093
30	4 10	.0113	1.033	12.00	.0076

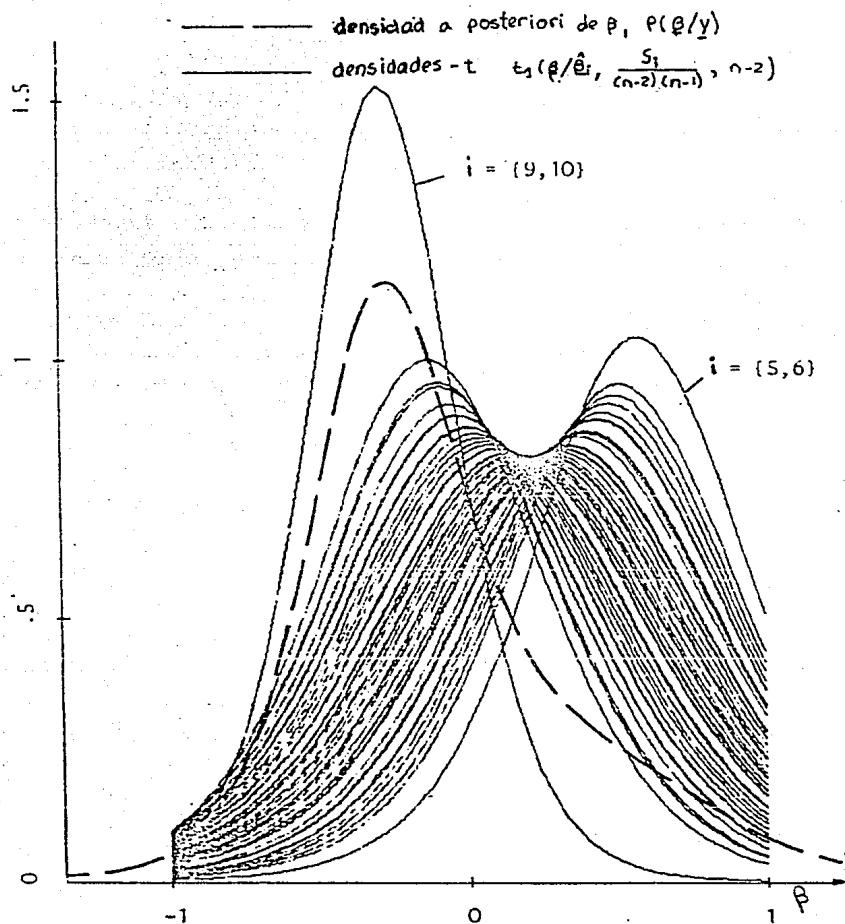
Tabla 3.4.3. (continuación)

31	5	6	.5785	-1.235	8.41	.0315
32	5	7	.3454	-0.303	13.00	.0055
33	5	8	.4797	-0.840	10.89	.0112
34	5	9	.1153	0.617	13.27	.0051
35	5	10	.1488	0.483	13.50	.0048
36	6	7	.3824	-0.451	12.56	.0063
37	6	8	.5167	-0.988	10.05	.0155
38	6	9	.1522	0.469	13.51	.0047
39	6	10	.1858	0.355	13.64	.0046
40	7	8	.2836	-0.056	13.49	.0048
41	7	9	-.0808	1.402	10.16	.0148
42	7	10	-.0473	1.267	10.91	.0111
43	8	9	.0535	0.865	12.62	.0062
44	8	10	.0870	0.730	13.01	.0055
45	9	10	-.2774	2.188	3.95	.6489*

* pesos dominantes

Al igual que para $k=1$, esta conclusión también puede obtenerse al observar la figura (3.4.2) en donde claramente la densidad para la pareja (9,10) domina a las restantes incluyendo a $P(B/Y)$ que es la más cercana a esta densidad.

figura 3.4.2.



3.5 Conclusiones

Una de las primeras cosas que resultan al exponer los métodos Bayesianos es que, al menos teóricamente, se puede modelar cualquier fuente de discrepancia concebible (modelo de la media trasladada, de la varianza inflada, etc), obteniéndose para cada situación un modelo. Naturalmente, en la medida en que se agreguen más condiciones a éste, la dificultad y complejidad en la interpretación de los resultados aumentará.

Por otro lado, parece ser que los métodos bayesianos únicamente contemplan a las observaciones discrepantes en el sentido de contaminantes y no en el sentido de sorprendentes o discrepantes, ver L.J. Pettit y A.F.M. Smith (1983).

Existen varios problemas asociados a los métodos bayesianos, el primero en el aspecto teórico: ya se mencionó que en la medida en que se aumenten el número de condiciones a modelarse, el modelo resultante será mucho más complejo. El segundo, en el aspecto práctico: como cada situación requiere un modelo, entonces, es difícil usar los métodos en forma rutinaria porque cada uno requiere un programa especial. Otro problema práctico es la explosión computacional que surge al tratar de detectar el número de observaciones discrepantes.

Capitulo 4

ENFOQUE ROBUSTO.

4.1. Introducción

La palabra robusto tiene varias connotaciones (algunas de ellas inconsistentes), en este trabajo usaremos este concepto en un sentido relativamente restringido: para nuestros propósitos, robustes significa insensibilidad a pequeñas desviaciones de las suposiciones inherentes al modelo.

Como se sabe, las inferencias estadísticas se basan solo en parte en las observaciones. Otro aspecto igualmente importante está formado por las suposiciones a priori acerca de la distribución principal. Aún en los casos más simples, existen suposiciones implícitas o explícitas acerca de la aleatoriedad e independencia de las observaciones, acerca de los modelos de distribución y en algunas ocasiones para las distribuciones a priori para algunos parámetros desconocidos. Naturalmente, dichas suposiciones no están impuestas por ser exactas, más bien, son supuestos simplificadores usados en la mayoría de las ramas más temáticas que expresan un conocimiento o certeza del fenómeno estudiado y su uso queda justificado por el principio de esta bilidad: "El menor error en el modelo matemático solo debe de causar un error pequeño en las condiciones finales. Desafortunadamente, no siempre se presenta esta situación.

Supongamos que tenemos un modelo paramétrico, con la esperanza de que sea una buena aproximación a la situación verdadera, pero no se debe ni se puede suponer que sea exactamente el correcto, es por eso que cualquier procedimiento estadístico debe de tener las siguientes caracte-

riáticas deseables.

- (a) Debe tener una eficiencia (óptima o cerca del óptimo) razonablemente buena para el modelo propuesto.
- (b) Debe de ser robusto en el sentido de que pequeñas desviaciones de las suposiciones del modelo deben de afectar sólo ligeramente su funcionamiento.
- (c) Algunas desviaciones grandes del modelo no deben de causar una catástrofe.

¿Que es un procedimiento robusto? . Peter J. Huber (1972), plantea esta pregunta y comenta al respecto lo siguiente: Al principio "robustez" era un concepto vago, por ejemplo, Box y Anderson (1955) introdujeron la noción diciendo que "se requieren procedimientos que sean robustos" (insensibles a cambios en factores extraños no contemplados en la prueba) y poderosos (sensibles a factores específicos bajo la prueba)

Pero, si se desea escoger en forma racional entre diferentes competidores robustos para un procedimiento, se deben precisar los objetivos que se quieren alcanzar. Desafortunadamente, no se ha llegado a un consenso: aunque los objetivos raramente se establecen en forma explícita, se pueden discutir cinco o seis situaciones, de las cuales pienso que no todas deberían de llevar el nombre de robustas.

Para fijar ideas, consideremos el problema de estimar un parámetro de posición para un número grande de observaciones independientes x_1, x_2, \dots, x_n , distribuidas de acuerdo con $P(x_i = x) = F\left(\frac{x-\theta}{\sigma}\right)$ donde la forma de F no se conoce exactamente. Uno juzgará a los estimadores en términos de

- (a) varianza asintótica $\sigma_F^2(T)$.
- (b) eficiencia absoluta $1/I(F)\sigma_F^2(T)$; con $I(F)$ la información de Fisher.
- (c) Eficiencia relativa $\sigma_F^2(T')/\sigma_F^2(T)$.

De acuerdo al primer objetivo, un estimador robusto debe de poseer

(i) Una eficiencia absoluta alta para todas las formas suaves de F que se deseen. Mientras que este objetivo puede alcanzarse asintóticamente para muestras grandes, la convergencia parece ser muy lenta para propósitos prácticos. Entonces, modificamos los requerimientos a uno de los siguientes.

(ii) Una eficiencia relativa alta para la media muestral para toda F .

(iii) Una eficiencia absoluta alta sobre un conjunto $\{F_i\}$ de formas seleccionadas estratégicamente (p.e. normal, doble exponencial, cauchy, etc).

Una variante de (iii) es

(iii') Una eficiencia absoluta alta sobre una familia paramétrica seleccionadas estratégicamente.

(iv) Una varianza asintótica pequeña sobre una vecindad

de una familia, en particular la de la normal.

Ninguno de estos objetivos garantiza la estabilidad cualitativa requerida (la convergencia en (i) y en (iv) no necesariamente es uniforme):

(v) La distribución del estimador cambiará poco bajo variaciones arbitrarias pequeñas de la distribución principal F , y esto uniformemente con el tamaño muestral n .

Huber considera a los objetivos locales (iv) y (v) como los más importantes y comenta que al igual que Anscombe (1960), se inclina a ver robustez como un problema de Seguros: "Estoy dispuesto a pagar una prima (Una pérdida de entre el 5% y 10% del modelo ideal) para protegerme contra malos efectos causados por pequeñas desviaciones de él; aunque estoy feliz si el procedimiento funciona bien bajo grandes desviaciones."

Primera mente se estaría interesado con robustez distribucional: la forma de la verdadera distribución principal se desvía ligeramente del modelo supuesto (usualmente el modelo normal). Este es el caso más importante y más difundido.

Hay otros aspectos de robustez acerca de los cuales se conoce poco, por ejemplo, insensibilidad respecto a desviaciones de independencia, distribución idéntica etc. mucho menos es conocido de que pasa cuando las otras suposiciones estándar de estadística no se satisfacen y acerca de la protección en estos otros casos.

Debe enfatizarse una vez más que, la ocurrencia de errores gruesos en una fracción pequeña de observaciones es con

siderada como una desviación pequeña, y que, en vista de la sensibilidad extrema de algunos problemas clásicos un objetivo fundamental de los procedimientos robustos es protegerse contra errores gruesos.

4.2. Estimadores robustos en regresión

4.2.1. Estimadores M. Reciben este nombre todos los estimadores que minimizan una función ρ de los residuales

$$\min_{\beta} \sum_{i=1}^n \rho(e_i) = \min_{\beta} \sum_{i=1}^n \rho(y_i - x_i^T \beta) \quad (4.2.1)$$

donde x_i^T denota el i -ésimo renglón de X . Un estimador de este tipo recibe el nombre de estimador M, debido a que puede pensarse como un estimador semejante al obtenido por máxima verosimilitud, esto es, la función ρ , está referida a la función de verosimilitud para una elección apropiada de la distribución del error, por ejemplo, si se usa el método de mínimos cuadrados (se supone que la distribución del error es normal), entonces $\rho(e) = \frac{1}{2} e^2$; $-\infty < e < \infty$.

Uno de los inconvenientes de los estimadores M, es que no siempre son invariantes bajo reescalamiento, es decir, si los residuales $y_i - x_i^T \hat{\beta}$ se multiplican por una constante, la solución a (4.2.1) no será la misma que se obtuvo anteriormente. Para obtener una versión invariante bajo escala de este estimador, usualmente se resuelve

$$\min_{\beta} \sum_{i=1}^n \rho\left(\frac{e_i}{d}\right) = \min_{\beta} \sum_{i=1}^n \rho[(y_i - x_i^T \beta)/d] \quad (4.2.2)$$

por los métodos: a) H ; b) Newton e) Mínimos cuadrados ponderados

ver Hogg (1979b).

En la expresión (4.2.2) d es un parámetro robusto de escala cuya elección más común es

$$d = \frac{\text{mediana } |e_i - \text{mediana}(e_i)|}{0.6745}$$

La constante 0.6745 hace a d un estimador casi insesgado de σ cuando n es grande y la distribución del error es normal.

Para minimizar (4.2.2), encuentre las primeras derivadas parciales de ρ con respecto a β_j ($j=1, 2, \dots, p$) (suponiendo que ρ es una función derivable y convexa) e igualando a cero; esto da un sistema de p ecuaciones

$$\sum_{i=1}^n x_{ij} \Psi[(y_i - x_i^T \beta) / d] = 0 ; \quad j=1, 2, \dots, p \quad (4.2.3)$$

donde $\Psi = \rho'$ y x_{ij} es la i -ésima observación de la variable j . En general la función Ψ no es una función lineal y (4.2.3) debe de resolverse por algún método iterativo mencionado anteriormente. Por ejemplo, si usamos el método de mínimos cuadrados ponderados y además contamos con un estimador inicial de β y σ ; $\hat{\beta}_0$, entonces, (4.2.3) puede reescribirse como

$$\begin{aligned} \sum_{i=1}^n x_{ij} \Psi[(y_i - x_i^T \beta) / d] &= \sum_{i=1}^n x_{ij} \frac{\Psi[(y_i - x_i^T \hat{\beta}_0) / d]}{(y_i - x_i^T \hat{\beta}_0) / d} \cdot (y_i - x_i^T \beta) / d = 0 \\ &= \sum_{i=1}^n x_{ij} w_i (y_i - x_i^T \beta) = 0 ; \quad j=1, 2, \dots, p \quad (4.2.4) \end{aligned}$$

$$\text{donde } w_{i_0} = \begin{cases} \frac{\varphi[(y_i - x_i^T \hat{\beta}_0)/d]}{(y_i - x_i^T \hat{\beta}_0)/d} & \text{si } y_i \neq x_i^T \hat{\beta}_0 \\ 1 & \text{si } y_i = x_i^T \hat{\beta}_0 \end{cases} \quad (4.2.5)$$

La expresión (4.2.4) se puede escribir en notación matricial como sigue:

$$X^T W_0 X \beta = X^T W_0 Y, \quad (4.2.6)$$

donde W_0 es la matriz diagonal de pesos, cuyos elementos son $w_{10}, w_{20}, w_{30}, \dots$, dados en (4.2.5). La expresión (4.2.6), son las ecuaciones usuales de mínimos cuadrados ponderados. Consecuentemente el estimador en la primera iteración es

$$\hat{\beta}_1 = (X^T W_0 X)^{-1} X^T W_0 Y \quad (4.2.7)$$

En la siguiente iteración, se vuelven a calcular los pesos en (4.2.5) usando $\hat{\beta}_1$ en lugar de $\hat{\beta}_0$. Para alcanzar la convergencia del método, regularmente se requieren pocas iteraciones, aunque claro, esto depende de que tan buenos son los estimadores iniciales de β y σ . Para obtener estos estimadores, algunos autores sugieren usar los estimadores de norma L_1 (más adelante se hace referencia a ellos), mientras que Hogg (1979. b), sugiere el uso del algoritmo desarrollado por Jutter (1977) que proporciona simultáneamente buenos estimadores de β y σ .

En la tabla (4.2.1) se muestran las funciones ρ más usuales y sus correspondientes funciones φ . La función ρ controla los pesos que se asignan a cada residual y es por ello que algu-

tabla 4.2.1.

Estimador	$\rho(e)$	$\psi(e)$	$V(e)$	RANGO
MINIMOS CUADRADOS	$\frac{1}{2} e^2$	e	1.0	$ e < \infty$
HUBER *t=2	$\frac{1}{2} e^2$ $ e t - \frac{1}{2} t^2$	e $t \text{ signo}(e)$	1.0 $t/ e $	$ e \leq t$ $ e > t$
ANDREWS *a=1.389	$a[1 - \cos(e/a)]$ $2a$	$\text{sen}(e/a)$ 0	$\frac{\text{sen}(e/a)}{e/a}$ 0	$ e \leq a\pi$ $ e > a\pi$
RAMSAY *a=0.3	$a^2[1 - \exp(-a e) \cdot (1+a e)]$	$e \exp(-a e)$	$\exp(-a e)$	$ e < \infty$
HAMPEL *a=1.7 b=2.4 c=8.6	$\frac{1}{2} e^2$ $a e - \frac{1}{2} a^2$ $\frac{a(e e - \frac{1}{2} e^2)}{c-b} - (\frac{7}{8}) a^2$ $a(b+e-a)$	e $a \text{ signo}(e)$ $\frac{a \text{ signo}(e)(e- e)}{c-b}$ 0	1.0 $a/(e)$ $\frac{a(e- e)}{ e (c-b)}$ 0	$ e \leq a$ $a < e \leq b$ $b < e \leq c$ $ e > c$
TUKEY	$\frac{e^2}{2} [1 - (\frac{e}{k})^2 + \frac{1}{3} (\frac{e}{k})^4]$ $\frac{k^2}{2}$	$e[1 - (e/k)^2]^2$ 0	$[1 - (e/k)^2]^2$ 0	$ e \leq k$ $ e > k$
	donde $e = \frac{y_i - \bar{y}}{s}$ y t, a, b, c pueden tomar otros valores, esto dependerá de la robustez deseada.			

nas veces se les llama funciones de influencia, por ejemplo, la función mínimos cuadrados no está acotada y entonces, el estimador resultante, no es robusto si los datos surgen de una distribución con colas pesadas.

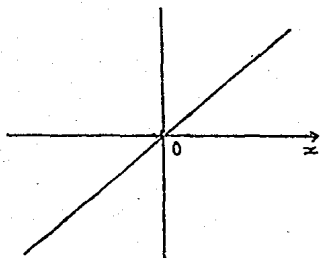
La función de Huber, es idéntica a la de mínimos cuadrados, en el intervalo $(-t, t)$, ver (b) en figura (4.2.1), pero fuera de este la función de influencia toma valores $\psi(x) = t \text{ signo}(x)$, dando así un peso menor a los residuales en la medida en que estos crecen.

La función de Hampel, (e) de (4.2.1), puede pensarse como una versión discreta de la función de Andrews o la de Tukey. también puede observarse que estas últimas funciones son muy similares y por lo tanto una sería sustituto de la otra.

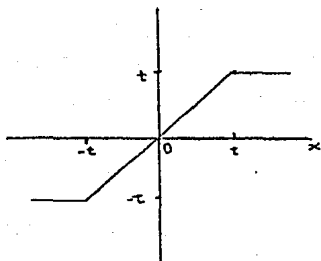
Ahora, si una observación es disrepante, el residual correspondiente es grande. Si se utiliza como estimador el propuesto por Andrews, Tukey o Hampel, entonces estas observaciones tendrán poco peso en el ajuste. Pero si se utiliza el estimador de Huber, el residual tendría que ser muy grande para que el peso $\frac{t}{x}$, fuera pequeño y esta observación no afectará al ajuste.

Debido a que las funciones ρ asociadas con las funciones de influencia no son convexas, pueden surgir ciertos problemas de convergencia en el proceso iterativo, a pesar de esto, se cuenta con buenos procedimientos que pueden usarse con cierto cuidado. En caso de que se utilice el método de mínimos cuadrados, no hay manera de protegerse de las observaciones -

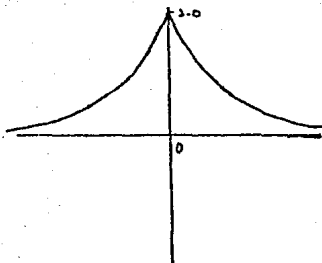
figura 4.2.1.



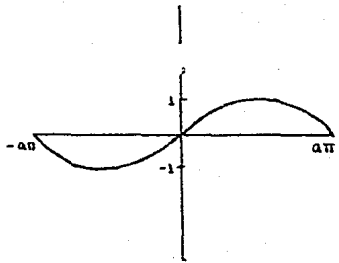
(a) MINIMOS CUADRADOS.



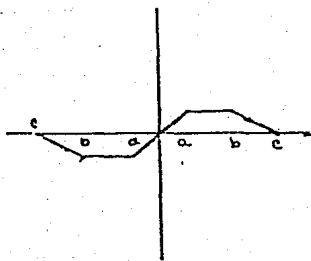
(b) HUBER



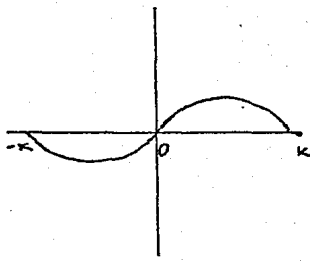
(c) RAMSAY



(d) ANDREWS



(e) HAMPPEL



(f) TUKEY

discrepantes.

Entonces, para que el estimador $\hat{\beta}$ del parámetro no se vea influenciado por algunas observaciones lejanas al plano de regresión, la función $\Psi(x)$ debe de ser acotada y tender a una constante es decir

$$\lim_{|x| \rightarrow \infty} \Psi(x) = 0$$

Como resultado de esta condición se pueden tener varios mínimos locales, por lo que la solución dada por el método iterativo de optimización dependerá del punto $\hat{\beta}_0$ que se usa para iniciar el proceso.

4.2.2 Estimadores R. Los métodos conocidos como estimación R, son procedimientos que se basan en los rangos. El procedimiento general reemplaza un factor de cada término en la función objetivo de mínimos cuadrados.

$S_0(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2$ por su rango correspondiente, así, si R_i es el rango de $y_i - x_i^T \beta$, se desea minimizar

$$\sum_{i=1}^n (y_i - x_i^T \beta) \cdot R_i$$

Más en general, se reemplazan los rangos por una función de puntaje $a_n(i)$, $i=1, 2, \dots, n$, entonces la función objetivo que se debe de minimizar es

$$S_1(y - X\beta) = \sum_{i=1}^n (y_i - x_i^T \beta) \cdot a_n(R_i) \quad (4.2.8)$$

Como S_1 es una función no negativa de β , continua y convexa (Jaekel 1972) y por estar definida para toda β , además de ser derivable en casi todas partes, el número puede hallarse por algún método iterativo, como el de "steepest descent" de Luenberger (1973).

Para puntajes $a_n(i)$ tales que

$$\sum a_n(i) = 0,$$

S_1 es invariante en posición, ya que: $S_1(y - X\beta + e) = S_1(y - X\beta)$ y por lo tanto la ordenada al origen, para la ecuación de regresión, no podrá ser estimada por medio de (4.2.8). Para superar este problema, se hace lo siguiente: con los coeficientes que se puedan obtener de (4.2.8), calcular los residuales

$e_i = (y_i - \sum_{j=1}^k \beta_j x_{ij})$ y entonces, minimizar $\sum_{i=1}^n a_n(R_i) \text{signo}(e_i - \alpha)$ con respecto a α ; donde R_i es el rango de $|e_i - \alpha|$, a_n^+ es una función de puntaje con signo.

Algunas de las funciones de puntaje mayormente utilizadas:

Wilcoxon: $a_n(i) = i \quad 1 \leq i \leq n$

Mediana: $a_n(i) = \begin{cases} 1 & \text{si } i > (n+1)/2 \\ -1 & \text{si } i \leq (n+1)/2 \end{cases}$;

Normal: $a_n(i) = \Phi^{-1} \left(\frac{i - 1/2}{n} \right)$; $1 \leq i \leq n$.

Vander Warden: $a_n(i) = \Phi^{-1} \left(\frac{i}{n+1} \right) \quad 1 \leq i \leq n$

con Φ la función de distribución normal $(0, 1)$

4.3 Comentarios acerca de los estimadores Robustos.

Naturalmente, los estimadores robustos discutidos en las secciones anteriores no son los únicos, se cuenta también con los estimadores L (combinaciones lineales de estadísticas de orden), de norma L_p (minimizar la suma de las desviaciones absolutas del hiperplano de regresión elevadas a la potencia p), método de Andrews, etc; entre otros, ver Gracia - Medrano - (1984).

En este trabajo, se ha hecho énfasis en los estimadores M debido principalmente a que estos pueden hacer frente a casi cualquier tipo de distribución (fuerte contaminación de la distribución normal, distribuciones con colas pesadas, etc.). Además de que

éstos pueden obtenerse fácilmente de un programa estándar de M.C. ponderados y los pesos en la iteración final identifica puntos discrepantes. Los estimadores L no se prestan para una generalización de regresión múltiple. Además bajo ciertas condiciones, los estimadores R son asintóticamente equivalentes a los estimadores M.

Uno de los inconvenientes de los estimadores M, es la estimación de la matriz covarianza y en general para hacer inferencias.

Algunos factores para poder elegir algún estimador pueden encontrarse en Gracia - Medrano V. (1974) y son

- a) Conocimiento acerca de la distribución de la variable respuesta.
- b) tamaño de muestra.
- c) Recursos con los que se cuenta para el cómputo.
- d) La situación, ya sea de observación o experimentación.
- e) Número de variables explicativas.

En general, la mayoría de los estimadores robustos requieren de programas más elaborados para su obtención que los de M.C., así como mayor tiempo de cómputo. En el caso de los estimadores M que se obtienen por métodos iterativos, el estimador inicial debe elegirse con mucho cuidado ya que juega un papel importante. Los factores principales para decidir que tipo de estimador inicial usar, son el tamaño

de muestra y recursos con que se cuenta para el cómputo.

Los métodos robustos de regresión se ven fuertemente afectados por el mal condicionamiento de la matriz X (multicolinealidad). Si el modelo lineal no es el adecuado, usar un método robusto no resolverá este problema. En algunos casos será necesario transformar las variables y después aplicar un método robusto.

El uso de los métodos robustos no se ha generalizado debido a varias razones.

a) Su cómputo es más complicado que el de mínimos cuadrados.

b) Los criterios que se siguen para la obtención de estimadores robustos en la mayoría de los casos son más difíciles de interpretar en comparación con el criterio de minimizar una suma de residuales al cuadrado.

c) Los estimadores robustos no son invariantes bajo escala.

d) La obtención de resultados teóricos, distribuciones asintóticas, estimación de matrices de covarianza, pruebas asociadas, es mucho más difícil y complicada que para el caso de mínimos cuadrados.

En cuanto al cómputo de los procedimientos robustos, actualmente, gracias a los avances en la computación, se ha desarrollado una gran cantidad de rutinas para obtenerlos. Ge

neralmente éstas son más complicadas y requieren de mayor tiempo que el procedimiento de mínimos cuadrados. Desgraciadamente éstas rutinas, en su mayoría no están todavía implantadas en los paquetes usuales de regresión y esto hace difícil su acceso para la mayoría de investigadores interesados en aplicar métodos robustos.

Los métodos robustos de regresión tienen mucho que ofrecer, éstos son extremadamente útiles en la localización de observaciones discrepantes e influyentes. Así, cuando se efectue un análisis mínimos cuadrados sería útil ejecutar también un ajuste robusto. Si los resultados de los dos procedimientos substancialmente concuerdan, entonces usar los obtenidos por mínimos cuadrados debido principalmente a que en la actualidad son los más comprendidos. Sin embargo, si los resultados de los dos análisis difieren, entonces, deben de identificarse las razones de esta(s) diferencia(s). Las observaciones que están ponderadas en el ajuste robusto, deben examinarse cuidadosamente. Algunas de las causas por las cuales los dos procedimientos no concuerdan pueden ser:

- a) Errores gruesos, ya sea en la variable respuesta, ó en las explicativas.
- b) Un modelo lineal inadecuado.
- c) El análisis debió hacerse en otra escala
- d) El error tiene una distribución con colas más pesadas que la normal.

4.4. Ejemplos.

Por tener solamente disponibilidad de los estimadores M con la inclusión de la función de Huber, se llevo a cabo el análisis de los dos ejemplos presentados en la sección 2.5. En las tablas 4.4.1 y 4.4.2 se encuentran resumidos los resultados para diferentes valores de t , del primer y segundo ejemplo respectivamente. Para el primero de ellos, observamos que para los cuatro valores de t propuestos, la observación que conserva el menor peso es precisamente la número 19, que por los métodos anteriores había sido declarada como discrepante. Por otro lado, en la medida que t es pequeño ($t < 2$), nos damos cuenta que los estimadores (de los parámetros) robustos difieren marcadamente de los obtenidos por mínimos cuadrados, y para valores grandes de t los valores son más parecidos. Esto es natural, ya que mientras más grande sea t , la función $\rho(e)$ de Huber se parece más a la de mínimos cuadrados. Por lo que respecta a la suma de cuadrados residual, esta entre 2263.9607 para $t=1$ y 2127.6259 para $t=4$ con lo que podemos observar que se obtiene una reducción en la suma de cuadrados residual con respecto a mínimos cuadrados a pesar de tener valores extremos para t .

En términos generales podríamos decir que los resultados obtenidos por mínimos cuadrados y métodos robustos no difieren cualitativamente ni cuantitativamente para este ejemplo, y siguiendo nuestras recomendaciones podríamos manejar con confianza los resultados obtenidos por mínimos cuadrados.

Un análisis similar en el grupo de datos nos lleva -

tabla 4.4.1

ESTIMACIÓN

EL VALOR DE t FS 1.000000

EL VALOR DE t ES 1.200000

PARAMETROS ESTIMADOS

PARAMETROS ESTIMADOS

RETA

RETA

110.955
-1.209

110.022
-1.181

-LAT -

VECTOR Y	VECTOR Y ESTIMADO	RESIDUALES	PESOS	DATO VECTOR Y ESTIMADO	RESIDUALES	PESOS
95.000000	97.445686	-2.445686	0.408883	1 97.700198	2.700198	0.444412
71.000000	62.741460	-8.258540	0.221096	0.62.496101	-1.103699	0.144514
83.000000	67.402179	-15.597821	0.064112	0.67.592878	-15.207122	0.078910
91.000000	82.193478	-8.806522	0.113552	4.62.511414	-8.338588	0.143652
102.000000	111.445686	9.445686	0.185368	5.5111.700198	5.700198	0.123709
87.000000	87.449197	0.449197	1.000000	6.67.607519	6.607519	1.000000
93.000000	97.071790	4.071790	0.245592	6.67.644590	4.244590	0.282713
160.000000	102.610881	2.610881	0.383013	8.8102.607518	2.974342	0.403451
104.000000	106.984777	2.984777	0.385033	9.9107.429550	3.429550	0.349459
94.000000	101.489193	7.489193	0.133526	10.101.607518	3.607518	0.157739
113.000000	123.776076	10.776076	0.092798	11.1124.245486	3.245486	0.106681
96.000000	92.193478	-3.806522	0.0262767	12.1252.751414	1.248448	0.035430
83.000000	67.402179	-15.597821	0.0064112	13.1367.612873	-1.389583	0.078910
84.000000	70.610881	-13.389119	0.0074688	14.1470.774342	-1.207122	0.092126
102.000000	106.610881	4.610881	0.0216878	15.15106.774342	4.974342	0.241238
100.000000	111.402179	1.402179	0.0713175	16.16101.155106	1.792887	0.0669315
109.000000	111.819582	8.819582	0.113384	17.17114.155106	9.15590	0.131064
105.000000	54.080021	-2.919376	0.042535	18.1853.599724	-3.400227	0.552913
57.000000	151.863089	30.863089	0.032401*	19.19152.663126	1.063126	0.038631
121.000000	74.610881	-11.389119	0.087803	20.2074.774342	-1.207122	0.108837
86.000000	106.402179	1.402179	0.0713175	21.21101.792878	1.792878	0.0669315

también a obtener las conclusiones de Brownlee, es decir, declarar a las observaciones 1, 3, 4, 21 como disorepantes. Como en este caso si existen ciertas diferencias en los estimadores finales con respecto a los obtenidos por mínimos cuadrados, podríamos concluir que se requiere un análisis adicional, éste se hará en el capítulo 5.

Capitulo 5.

TRANSFORMACIONES.

5.1. Introducción

Como hemos observado, las observaciones discrepantes son casos para los cuales, y por diversas razones, el modelo hipotético no funciona bien, es decir, el modelo no ajusta adecuadamente dichas observaciones. Los candidatos para considerarse como observaciones discrepantes, son casos cuyos residuales estudentizados resultan extremos en comparación con los restantes, el valor de v_{ii} o de la distancia D_i de Cook. Si D_i de Cook. Si D_i es pequeña, entonces la observación discrepante (si existe) tiene poco efecto en los estimadores de los parámetros y por lo tanto no es importante excluir el caso del análisis, a pesar de que realmente sea una observación discrepante. Por otro lado, los casos con valores v_{ii} grandes requieren de atención especial, aún cuando los correspondientes r_i 's ya que si D_i es grande, el caso puede tener una influencia excesiva en los estimadores de los parámetros.

Un paso generalmente útil, es efectuar nuevamente la regresión sin el (los) caso(s) sospechoso(s) y evaluar los cambios en el análisis. Este es un enfoque razonable, especialmente si la forma del modelo (por ejemplo la escala de todas las x 's es conocida). Entonces, casos con valores grandes de v_{ii} , D_i o r_i pueden ser indicadores de la necesidad de una transformación y no de una observación discrepante.

Una vez detectadas las observaciones discrepantes, se podrían efectuar cualquiera de las siguientes cosas, Atkinson - (1982):

- a) Efectuar nuevamente la regresión sin los casos discrepantes y evaluar los cambios registrales.
- b) Considerar transformaciones que hagan a los casos influyentes menos importantes o mejoren el ajuste de casos que se consideran en principio discrepantes.

Los métodos para escoger transformaciones requieren de un conocimiento específico acerca de las relaciones entre las variables y hacen uso de técnicas de diagnóstico para sugerir transformaciones posibles, éstas generalmente escogen una transformación para maximizar alguna función de interés.

5.2. Métodos de transformación de Variables.

5.2.1. Método de máxima verosimilitud. Una clase importante de transformaciones puede obtenerse cuando Y toma únicamente valores positivos. Si solo consideramos transformaciones de Y , esto es equivalente a ajustar el modelo

$$Y^\lambda = X\beta + \underline{\varepsilon}, \quad \text{Var}(\underline{\varepsilon}) = \sigma^2 I, \quad (5.2.1.)$$

donde λ es la potencia de la transformación.

Cuando λ es desconocido, puede verse a (5.2.1) como un modelo con un parámetro desconocido adicional a β y σ^2 . Entonces se puede estimar λ de la misma manera que β y σ^2 .

Box y Cox (1964) sugieren, para estimar λ , encontrar el estimador conjunto $(\hat{\beta}, \hat{\sigma}^2, \hat{\lambda})$ máximo verosímil efectuando suposiciones acerca de ε ; en este trabajo se supone que ε tiene distribución normal.

Por un momento, supongamos que λ es conocida, entonces la estimación de $\hat{\beta}$ y el cálculo de suma de cuadrados del residual es inmediata

$$\hat{\beta}_\lambda = (x^t x)^{-1} x^t y^\lambda \quad (5.2.2)$$

$$SRE_\lambda = (y^\lambda)^t (I - V) y^\lambda \quad (5.2.3)$$

que podrá obtenerse fácilmente con un programa de regresión lineal simple.

Ya que λ es desconocido, (5.2.3) puede calcularse para un rango de valores de λ . Para cada $\lambda \neq 0$ calcular

$$-L(\lambda) = \frac{n}{2} \ln(\lambda^2) - \frac{n}{2} \ln(SRE_\lambda) + (\lambda - 1) \sum_{i=1}^n \ln(y_i) \quad (5.2.4)$$

y si $\lambda = 0$

$$-L(\lambda) = -\frac{n}{2} \ln(SRE_\lambda) - \sum_{i=1}^n \ln y_i \quad (5.2.5)$$

Se puede demostrar que la λ que maximiza (5.2.4) y (5.2.5) es el estimador máximo verosímil (suponiendo normalidad de los errores). Debe notarse que no tiene sentido escoger la λ que maximiza a SRE_λ porque para cada λ la SRE_λ es medida en una escala diferente, mientras que las ecuaciones (5.2.4) y (5.2.5) convierte cada SRE_λ

a una escala común.

5.2.2. Método de Carroll

En analogía con los estimadores M y en particular con lo desarrollado por Huber, R.J. Carroll (1980), propone una prueba para obtener una transformación en y . Así, usando

$$c(x) = \begin{cases} \frac{1}{2} x^2, & |x| \leq k, \\ k(|x| - k/2), & |x| > k, \end{cases}$$

La verosimilitud resulta ser

$$L(\beta, \sigma, \lambda) = \sigma^{-n} \prod_{i=1}^n \exp \left\{ -\rho \left(\frac{y_i^{(\lambda)} - x_i^T \beta}{\sigma} \right) + (\lambda - 1) \log y_i \right\} \quad (5.2.6)$$

Para λ fija, tomar un valor inicial de σ y maximizar (5.2.6) con respecto a β tal que

$$\sum_{i=1}^n \psi \left\{ (y_i^{\lambda} - x_i^T \hat{\beta}) / \sigma \right\} x_i = 0 \quad (5.2.7)$$

entonces se actualiza σ , el procedimiento iterativo

$$(n-p)^{-1} \sum_{i=1}^n \psi^2 \left\{ (y_i^{(\lambda)} - x_i^T \hat{\beta}) / \sigma \right\} = E_{\Phi} \psi^2(Z),$$

donde Φ es la función de distribución normal estándar.

Para una λ dada, denotamos a los estimadores como $\hat{\beta}_R(\lambda)$, $\hat{\sigma}(\lambda)$. Para encontrar λ , entonces maximizamos $L(\hat{\beta}_R(\lambda), \hat{\sigma}(\lambda))$.

λ), la solución la denotaremos por λ_2 . Note que cuando $k = \infty$, $P(\infty) = \frac{1}{2}x^2$ y el resultado obtenido se reduce al método de Box-Cox.

Para probar la hipótesis $H_0: \lambda = \lambda_0$, la estadística que se usará bajo el modelo (5.2.6) es

$$\Lambda_{\lambda} = -2 \log \left(\frac{L(\hat{\beta}(\lambda_0), \hat{\sigma}^2(\lambda_0), \lambda_0)}{L(\hat{\beta}(\lambda_{\lambda}), \hat{\sigma}^2(\lambda_{\lambda}), \lambda_{\lambda})} \right),$$

que, bajo condiciones apropiadas, es asintóticamente χ^2 con un grado de libertad.

5.2.3. Método de Box y Tidwell.

Supongamos que tenemos p variables independientes x_1, x_2, \dots, x_p , y se quiere transformar cada una por una transformación de potencia para obtener w_1, w_2, \dots, w_p , donde para $j = 1, 2, \dots, p$

$$w_j = \begin{cases} x_j^{\alpha_j} & \text{si } \alpha_j \neq 0 \\ \ln(x_j) & \text{si } \alpha_j = 0 \end{cases} \quad (5.2.8)$$

Entonces ajustaremos el modelo lineal

$$Y = \beta_0 + \sum \beta_j w_j + \epsilon \quad (5.2.9)$$

En (5.2.8) y (5.2.9) no se requiere que las α_j sean iguales. Ya que para toda $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_p)$, la variable dependiente en (5.2.9) es la misma y podemos encontrar una $\hat{\alpha}$ y una $\hat{\epsilon}$ que sean estimadores mínimos cuadrados (no lineales).

Esto sugiere que para cada α podríamos ajustar el modelo (5.2.9); calcular la suma de cuadrados del residual SCR_{α} y encontrar $\hat{\alpha}$ tal que minimice SCR_{α} .

Un procedimiento iterativo para encontrar $\hat{\alpha}$ los proponen Box y Tidwell (1962). Ellos inician con un valor arbitrario de α que se cree es el más adecuado, usualmente $\hat{\alpha}_1 = \hat{\alpha}_2 = \dots = \hat{\alpha}_p = 1$. Entonces, ajustando un modelo lineal aumentado, obtiene un estimador mejorado de α ; repitiendo este proceso hasta que se alcanza un nivel deseado de convergencia.

Especificando, iniciar con $\hat{\alpha}_1 = \dots = \hat{\alpha}_p = 1$, así $w_j = x_j^{\hat{\alpha}_j} = x_j$ y ajustar el modelo de regresión (5.2.9) para obtener

$$Y = \hat{\beta}_0 + \sum \hat{\beta}_j w_j \quad (5.2.10)$$

Después, construir p variables nuevas Z_1, Z_2, \dots, Z_p definidas por

$$Z_j = w_j \ln(w_j) \quad j = 1, 2, \dots, p \quad (5.2.11)$$

Las Z_j 's se definen de tal manera que, en el modelo lineal aumentado

$$\hat{Y} = \hat{\beta}_0^* + \sum \beta_j^* w_j + \sum \hat{\gamma}_j Z_j \quad (5.2.12)$$

entonces, si se requiere una transformación, cada $\hat{\gamma}_j$ será grande (sin considerar signo) y pequeña en caso contrario. Las β_j^* 's no son iguales a las $\hat{\beta}_j$ del modelo (5.2.10). Box y Tidwell muestran que un estimador apropiado de α_j es

$$\hat{\alpha}_j = \left(\frac{\hat{\alpha}_j}{\hat{\sigma}_j} + 1 \right) (\text{valor corriente de } \hat{\alpha}_j). \quad (5.2.13)$$

El procedimiento anterior, con las $\hat{\alpha}_j$ obtenidas de (5.2.13), puede repetirse hasta hacer decrecer la SER suficientemente. Sin embargo, frecuentemente un ciclo es adecuado.

5.2.4. Método de Andrews

Este método se basa en la prueba de hipótesis $\lambda = \lambda_0$. La estadística de prueba se construye expandiendo $\underline{y}^{(\lambda)}$ con λ_0 .

$$\underline{y}^{(\lambda_0)} = \underline{y}^{(\lambda)} + (\lambda_0 - \lambda) \underline{G}_y^{(\lambda_0)} \quad \text{donde } \underline{G}_y^{(\lambda_0)} = \frac{\partial \underline{y}^{(\lambda)}}{\partial \lambda} \Big|_{\lambda = \lambda_0}.$$

$$\text{Ya que } \underline{y}^{(\lambda)} = X\beta + \underline{\varepsilon}; \quad \underline{y}^{(\lambda_0)} \approx X\beta + (\lambda_0 - \lambda) \underline{G}_y^{(\lambda_0)} + \underline{\varepsilon}.$$

La estadística para la prueba de Andrews es igual a la estadística t para la hipótesis $\lambda_0 - \lambda = 0$ en el modelo

$$\underline{y}^{(\lambda_0)} = X\beta + (\lambda_0 - \lambda) \hat{\underline{G}}_y^{(\lambda_0)} + \underline{\varepsilon}$$

donde $\hat{\underline{G}}_y^{(\lambda_0)}$ es igual a $\underline{G}_y^{(\lambda_0)}$ evaluada en los valores ajustados del modelo nulo $\underline{y}^{(\lambda_0)} = X\beta + \underline{\varepsilon}$, se sigue inmediatamente, a partir del trabajo de Milliken y Graybill (1970) que la estadística t, tiene una distribución t-student con $n-p-1$ grados de libertad, es decir, la prueba de Andrews es exacta.

5.2.5. Método de Atkinson.

El enfoque que Atkinson (1982) usa es similar al usado en Box y Cox (1964). Se supone que existe una transfor-

mación paramétrica de los datos $\underline{z} = \underline{z}^\lambda$ tal que, para alguna λ , las observaciones satisfacen las suposiciones de segundo orden para el modelo lineal $E(\underline{z}^\lambda) = X\beta$ y además que, al menos aproximadamente, las observaciones transformadas pueden tratarse como una muestra de una distribución normal. Para un análisis de verosimilitud $\hat{\lambda}$. Como una forma alternativa para determinar la evidencia para una transformación, Atkinson usa la estadística t_0 para probar la hipótesis $\lambda = \lambda_0$.

$$t_0 = \frac{-\underline{z}^T M \underline{w}}{\sqrt{\underline{w}^T M \underline{w}}} \quad (5.2.14)$$

donde $\underline{w} = \frac{\partial \underline{z}}{\partial \lambda}$; $M = I - X(X^T X)^{-1} X^T$, y Δ_z es un estimador de la varianza de \underline{z} derivada de la suma de cuadrados después de la regresión en X y \underline{w} , entonces

$$(n-p-1)\Delta_z^2 = \underline{z}^T M \underline{z} - \frac{(\underline{z}^T M \underline{w})^2}{\underline{w}^T M \underline{w}}, \quad (5.2.15)$$

además, todas las cantidades se calculan usando $\lambda = \lambda_0$. Asíntoticamente, la distribución de $t_0(\lambda_0)$ es normal estándar, pero su distribución en muestras pequeñas es inestable ya que \underline{z} y \underline{w} son variables aleatorias con distribuciones no estándar. La transformación paramétrica estándar que se usa regularmente es

$$\underline{z}^\lambda = \begin{cases} \frac{y^\lambda - 1}{\lambda y^{\lambda-1}}, & \lambda \neq 0 \\ \dot{y} \log y, & \lambda = 0 \end{cases} \quad (5.2.16)$$

donde \dot{y} es la media geométrica de las observaciones

En realidad, w también se puede manejar como una variable adicional al modelo de regresión con las suposiciones de costumbre, es decir

$$E(y) = x\beta + w\gamma \quad (5.2.17)$$

El estimador mínimos cuadrados de γ en (5.2.17) puede escribirse en términos de los residuales de la regresión de w y x :

$$e = MY = (I - V)Y$$

$$w^* = Mw = (I - V)w$$

donde M y V están definidos en la sección 2.1 de este trabajo. Las fórmulas estándar de análisis de covarianza muestran que el estimador mínimos cuadrados de γ está dado por

$$\hat{\gamma} = \frac{w^{*T}(I-H)Y}{w^{*T}(I-H)w} = \frac{w^{*T}MY}{w^{*T}Mw} = \frac{w^{*T}e}{w^{*T}w^*}$$

esta expresión demuestra que $\hat{\gamma}$ es el estimador mínimos cuadrados de la pendiente de la regresión de los residuales e en los residuales w . Los diagramas de estos dos conjuntos de residuales, llamados diagrama de variable agregada, se usan frecuentemente para evaluar la evidencia de términos adicionales en el modelo de regresión y observaciones influyentes ver Cook y Weisberg (1982, fig 2.3-8, 2.3.10), Belsley, Kuh y Welsch (1980 fig. 2.9-2.13)

En realidad para el estimador $\hat{\gamma}$, requerimos la prueba

ba t para la significancia de la regresión, la varianza de $\hat{\delta}$; $\sigma_{\hat{\delta}}^2 = \sigma^2 / \underline{w}^T M \underline{w}$, así que, la estadística para la hipótesis de no regresión es

$$t_w = \frac{\hat{\delta}}{\sqrt{\hat{\sigma}_{\hat{\delta}}^2}} = \frac{\underline{w}^T M \underline{Y}}{S(\underline{w} M \underline{w})^{1/2}} = \frac{S(\underline{w}, \underline{Y})}{S(\underline{w}, \underline{w})^{1/2}} \quad (5.2.18)$$

donde $S(\underline{a}, \underline{b}) = \underline{a}^T (I - V) \underline{b} = \underline{a}^T M \underline{b}$ y con

$$(n-p-1) s_w^2 = \underline{Y}^T M \underline{Y} - (\underline{Y}^T M \underline{w})^2 / \underline{w}^T M \underline{w} \\ = S(\underline{Y}, \underline{Y}) - S^2(\underline{Y}, \underline{w}) / S(\underline{w}, \underline{w}) \quad (5.2.19)$$

Entonces, como quedaría indicado por una tabla de análisis de covarianza, el efecto de agregar una variable \underline{w} depende solamente de la suma de cuadrados y productos de \underline{w} y \underline{Y} . Nótese que excepto por un signo sin importancia la expresión (5.2.18) y (5.2.19) son iguales.

Al escribir la transformación (5.2.16) en series de Taylor evaluada en λ_0 obtenemos

$$\underline{Z}(\lambda) \approx \underline{Z}(\lambda_0) + (\lambda - \lambda_0) \underline{w}(\lambda_0) \quad (5.2.20)$$

$$\text{con } \underline{w}(\lambda_0) = \frac{\partial \underline{Z}(\lambda)}{\partial \lambda} \Big|_{\lambda=\lambda_0}$$

y se encuentra que el modelo lineal aproximado es

$$\underline{Z}(\lambda_0) = \underline{x} \beta - (\lambda - \lambda_0) \underline{w}(\lambda_0) + \underline{\varepsilon} \quad (5.2.21)$$

Al comparar esta última expresión con (5.2.17), se muestra que podemos obtener un estimador rápido (aproximado) para

λ como

$$\begin{aligned}\tilde{\lambda} &= \tilde{\lambda}(\lambda_0) = \lambda_0 - \underline{w}^t(\lambda_0) M_{\underline{z}}(\lambda_0) / \underline{w}^t(\lambda_0) M_{\underline{w}}(\lambda_0) \\ &= \lambda_0 - S(\underline{w}, \underline{z}; \lambda_0) / S(\underline{w}, \underline{w}; \lambda_0)\end{aligned}\quad (5.2.22)$$

similarmente la estadística aproximada de prueba para la transformación es

$$\begin{aligned}t_{\underline{w}}(\lambda_0) &= -\underline{z}^t(\lambda_0) M_{\underline{w}}(\lambda_0) / s_{\underline{z}} \{ \underline{w}^t(\lambda_0) M_{\underline{w}}(\lambda_0) \}^{1/2} \\ &= -S(\underline{w}, \underline{z}; \lambda_0) / \{ s_{\underline{z}}^2 S(\underline{w}, \underline{w}; \lambda_0) \}^{1/2}\end{aligned}\quad (5.2.23)$$

En (5.2.23) el estimador de $s_{\underline{z}}^2$, por analogía con (5.2.19) es

$$(n-p-1) s_{\underline{z}}^2 = S(\underline{z}, \underline{z}; \lambda_0) - \frac{S^2(\underline{w}, \underline{z}; \lambda_0)}{S(\underline{w}, \underline{w}; \lambda_0)}\quad (5.2.24)$$

Un valor significativo de $t_{\underline{w}}(\lambda_0)$ indicaría que se requiere una transformación de un valor diferente a λ_0 . Debido a que la log-verosimilitud maximizada parcialmente es raramente cuadrática, la estadística $t_{\underline{w}}(\lambda_0)$ no tiene exactamente una distribución T con $n-p-1$ grados de libertad. Similarmente $\tilde{\lambda}$, el estimador rápido de la transformación, puede estar un poco distante al estimador máxima verosimilitud de λ .

Ahora desarrollaremos unas aproximaciones del efecto en $\tilde{\lambda}$ y en $t_{\underline{w}}(\lambda_0)$ al eliminar algunas observaciones. En las muestras se está interesado regularmente en una observación. La teoría sin embargo, se desarrolla fácilmente para el caso más general en la eliminación de k observaciones señaladas por i .

Empecemos con la suma de cuadrados del residual del modelo usual de regresión, que es

$$(n-p)\hat{\Delta}^2 = \mathbf{y}^t \mathbf{M} \mathbf{y} = S(\mathbf{y}, \mathbf{y})$$

Después de eliminar k observaciones, la suma de cuadrados residual está dado por

$$(n-p-k)\hat{\Delta}_{(i)}^2 = S(\mathbf{y}, \mathbf{y}) - \mathbf{e}_i^t (\mathbf{I} - \mathbf{V}_i)^{-1} \mathbf{e}_i \quad (5.2.25)$$

donde (i) y i como subíndice tiene la misma conotación que en el enfoque frecuentista. Así, \mathbf{e}_i es el conjunto de k residuales indicados por i y \mathbf{V}_i es la submatriz de $k \times k$ de \mathbf{V} . Por conveniencia, escribiremos una notación general para la suma de cuadrados como los presentados en (5.2.25). Para vectores \mathbf{a} , \mathbf{b} si los residuales, definidos por $\mathbf{a}^* = (\mathbf{I} - \mathbf{V})\mathbf{a} = \mathbf{M}\mathbf{a}$, $\mathbf{b}^* = (\mathbf{I} - \mathbf{V})\mathbf{b} = \mathbf{M}\mathbf{b}$. Entonces podemos escribir

$$S_{(i)}(\mathbf{a}, \mathbf{b}) = S(\mathbf{a}, \mathbf{b}) - \mathbf{a}_i^{*t} (\mathbf{I} - \mathbf{V}_i)^{-1} \mathbf{b}_i^* \quad (5.2.26)$$

donde \mathbf{a}_i^* y \mathbf{b}_i^* son los k miembros de los vectores de residuales. Para el modelo aumentado (5.2.19), el coeficiente estimado de regresión en w después de eliminar k observaciones está dado por

$$\hat{\gamma} = S_{(i)}(w, \mathbf{y}) / S_{(i)}(w, w) \quad (5.2.27)$$

Usando esta expresión con w reemplazada por $w(\lambda_0)$ para obtener un estimador rápido del parámetro de la transformación después de eliminar k observaciones. De (5.2.26) y (5.2.27).

$$\begin{aligned}\tilde{\lambda}_{(i)}(\lambda_0) &= \lambda_0 - S_{(i)}(\underline{w}, \underline{z}; \lambda_0) / S_{(i)}(\underline{w}, \underline{w}; \lambda_0) \\ &= \lambda_0 - \frac{\underline{w}_i^* M \underline{z} - \underline{w}_i^{*'} (\mathbf{I} - \mathbf{V}_i)^{-1} \underline{z}_i^*}{\underline{w}_i^* M \underline{w} - \underline{w}_i^{*'} (\mathbf{I} - \mathbf{V}_i)^{-1} \underline{w}_i^*}\end{aligned}\quad (5.2.28)$$

donde los residuales \underline{w}_i^* y \underline{z}_i^* se evalúan en λ_0 .

La derivación de un estimador rápido de $\lambda(\lambda_0)$ requiere aproximación de series de Taylor para la transformación del modelo. La derivación de $\tilde{\lambda}_{(i)}(\lambda_0)$ vía (5.2.27) incluye una aproximación adicional, ya que las variables $\underline{z}(\lambda_0)$ son funciones de la media geométrica de todas las observaciones. Cuando k observaciones son eliminadas, la media geométrica de las $n-k$ observaciones restantes no tendrá el mismo valor. En general podemos esperar que el cambio en \hat{y} sea despreciable si k es pequeña relativa a n , excepto si una de las observaciones eliminadas fuera muy cercana a cero como lo son las otras. Como el propósito de esta técnica es proporcionar cantidades de diagnóstico que se calculen rápidamente, esta dificultad será ignorada en el desarrollo teórico.

El efecto de la eliminación de observaciones en la estadística (5.2.23) puede derivarse en forma similar. Se sigue de (5.2.25) que el efecto de eliminación para producir una estadística aproximada

$$\tilde{T}_{w(i)}(\lambda_0) = -S_{(i)}(\underline{w}, \underline{z}; \lambda_0) / \{S_{z(i)}^2 S_{(i)}(\underline{w}, \underline{w}; \lambda_0)\}^{1/2}\quad (5.2.29)$$

donde, por analogía con (5.2.17)

$$(n-p-k-1)A_{\underline{w}}^2 = S_{\underline{w}}(\underline{z}, \underline{z}; \lambda_0) / S_{\underline{w}}(\underline{w}, \underline{w}; \lambda_0) \quad (5.2.30)$$

En particular, si solo eliminamos la observación

$$\lambda_{(i)}(\lambda_0) = \lambda_0 - \frac{w_i^* M z_i - w_i^* z_i / (1 - v_i)}{w_i^* M w - w_i^* z_i / (1 - v_i)} \quad (5.2.31)$$

también se puede obtener una simplificación similar para la estadística $\tilde{T}_{\underline{w}_{(i)}}(\lambda_0)$ definida en (5.2.29) y (5.2.30). Adicionalmente, se puede obtener una simplificación para la variable $w(\lambda)$ definida en (5.2.20) para la transformación de potencia (5.2.16). Para cualquier λ

$$w(\lambda) = \frac{\partial \tilde{z}(\lambda)}{\partial \lambda} / \lambda = \lambda_0 = \frac{y^\lambda \log y}{\lambda y^{\lambda-1}} - \frac{y^\lambda - 1}{\lambda y^{\lambda-1}} (1/\lambda + \log y) \quad (5.2.32)$$

Debido a que las cantidades de diagnóstico dependen solo de los residuales de $w(\lambda)$, pueden encontrarse formas simples para $w(\lambda)$ para ciertos valores de λ . En particular

$$w(1) = y \{ \log(y/y) - 1 \}$$

$$w(0) = y \log y (\log y/2 - \log y) \quad (5.2.33)$$

5.3 COMENTARIOS

Uno de los principales inconvenientes que presentan los

métodos de máxima verosimilitud para estimar el parámetro de la transformación λ , es que, éstos se ven afectados por observaciones discrepantes, confirmando el hecho de que máxima verosimilitud en el modelo lineal normal tiende a no ser robusto: Andrews (1971). Además, estos métodos requieren de cálculos repetidos al realizar varias transformaciones de los datos originales, resultando por esto, un método bastante tedioso cuando se usa. Adicionalmente, las pruebas y los límites de confianza que se basan en ellos solo tienen validez asintótica y el número de parámetros debe ser chico comparado con el número de observaciones.

Por otro lado, el método de Andrews presentará las siguientes ventajas sobre el de máxima verosimilitud.

- a) La prueba de significancia es exacta.
- b) La cantidad de cálculos involucrados se reduce cuando se prueba una o pocas transformaciones.
- c) La precisión de la transformación puede evaluarse teóricamente.

Una de sus principales desventajas es que no se pueden obtener conclusiones por medio de gráficas como las que se obtienen con diagramas de verosimilitud y cuando la distribución de los errores es normal, la prueba de Andrews es menos potente que la prueba del cociente de verosimilitud: Atkinson (1973). En cuanto al método de Carroll (1980), éste demuestra por medio de un experimento

de Monte Carlo que su procedimiento es preferible al de verosimilitud en términos de inferencias acerca de λ y β cuando se consideran modelos para el error diferentes al normal. también el método parece preferible al de Andrews debido a la falta de potencia de éste último; esta ventaja se ve disminuida un poco, por un pequeño incremento en el nivel de la prueba de una transformación particular. Además el método de Carroll parece ligeramente menos robusto contra observaciones discrepantes que el método de Andrews.

5.4. Ejemplos.

En esta sección únicamente se aplican algunos métodos de transformaciones al grupo de datos Brownlee, K. A. (1965) Para determinar si es necesaria una transformación, se calcula la estadística T_D de Atkinson (1973) cuyo valor es -3.29 . Dado que este resultado indica la necesidad de una transformación se calcula $\hat{\lambda}$ (estimador de máxima verosimilitud); donde $\hat{\lambda} = 0.30$ y un intervalo de confianza del 95% ($-0.19, 0.74$) que da evidencia de una transformación logarítmica. Cuando se efectúa dicha transformación y calcular T_D , ésta no proporciona elementos para una nueva transformación, ver Atkinson (1981), (1982) para mayores detalles.

Como hemos visto la observación 21 es la más discrepante, si la eliminamos $T_D = -2.53$ y $\hat{\lambda} = 0.48$ con $(0.05, 0.87)$ un intervalo de confianza del 95%. De lo anterior podemos concluir que en la observación más discrepante no está contenida toda la información necesaria para una transformación.

Por otro lado, al aplicar procedimientos estándar de regresión

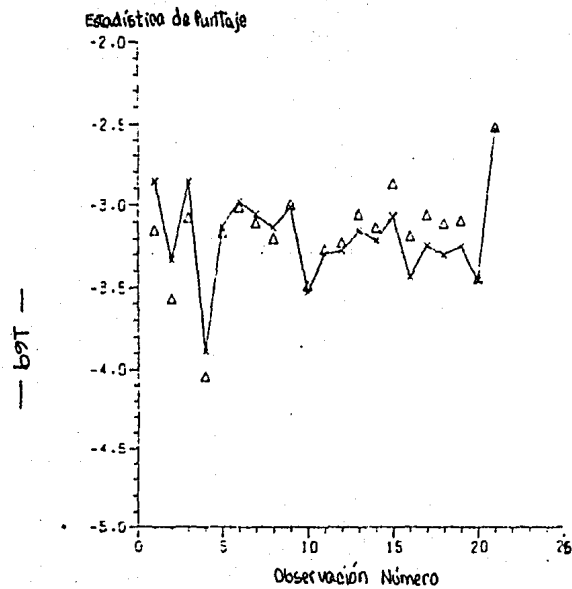
para la adición y eliminación de variables nos conduce (Atkinson (1981)) al modelo $\log Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_2 + \beta_4 X_1 X_2 + \epsilon$ para el cual la suma de cuadrados residual es 64.4 que es todavía mayor que 20.40, valor que se obtiene con el modelo de primer orden (2.5.3); pero a pesar de todo, no existe con este modelo evidencia de observaciones discrepantes.

En la gráfica (5.4.1) se muestra un diagrama de la estadística $\tilde{T}_{w(i)}(\lambda_0)$ y $T_{w(i)}(\lambda_0)$ definido en (5.2.23), $T_{w(i)}(\lambda_0)$ la estadística (5.2.19) cuando se eliminan del modelo (2.5.1) las k observaciones indicadas por i ; para este caso $i = \{i\}$ y $\lambda_0 = 1$. Aunque la eliminación de la observación 21 causa solo una pequeña reducción en la evidencia para una transformación, las dos estadísticas, $\tilde{T}_{w(i)}(\lambda_0)$ y $T_{w(i)}(\lambda_0)$ convierten en que los datos deben de ser transformados. Por último, la gráfica (5.4.2) muestra los valores de $\tilde{\lambda}_{(i)}^{(1)}$, $\hat{\lambda}_{(i)}^{(1)}$ y $\hat{\lambda}_{(i)}$ y de la cual podemos concluir que la eliminación de la observación 21 causa que la transformación estimada que era de aproximadamente de $1/3$ con todos los datos a un valor muy cercano a $1/2$.

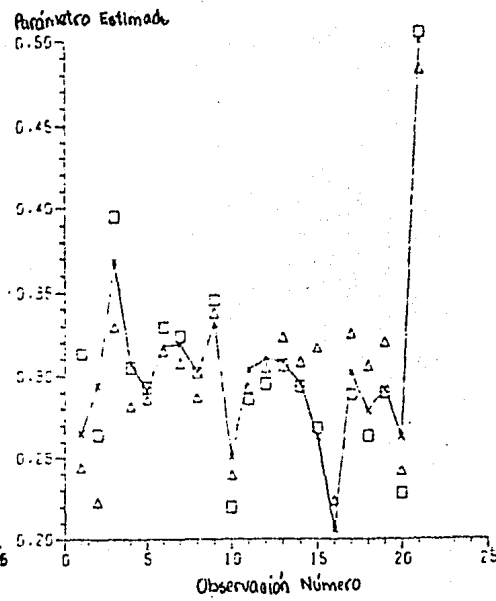
5.5. Conclusiones

Como ya se mencionó anteriormente, antes de aplicar algún método para el manejo de observaciones discrepantes, sería adecuado verificar todas las hipótesis subyacentes a los modelos, si en esta etapa no existe algo anormal, entonces buscar alguna transformación que permita "unir" los datos. De entre estos métodos, tal vez el que mayor difusión a tenido es el de mínimos cuadrados, pero debido a que es un proceso tedioso se buscan nuevas alternativas. Aunque siempre es conveniente tener un método robusto, de los métodos presentados y que más desarrollo ha tenido es el de Atkinson. Este

Grafica 5.4.1.



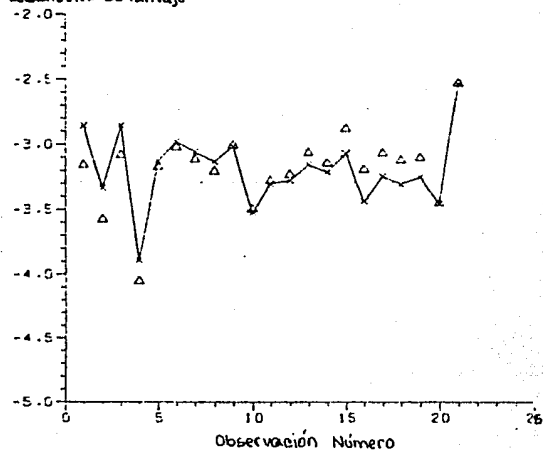
Grafica 5.4.2.



— 607 —

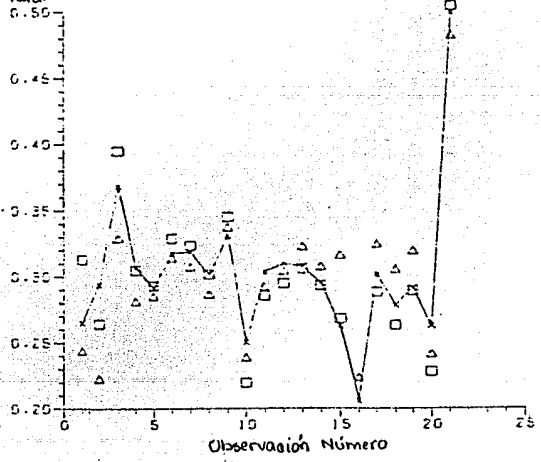
Gráfica 5.4.1.

Estadística de Furtaje



Gráfica 5.4.2.

Argumento Estimado



método presenta aparentemente mayores ventajas debido principalmente a que su aplicación puede efectuarse casi directamente a partir de los programas usuales de regresión, además pueden calcularse valores rápidos para el valor λ de la transformación.

Conclusiones Generales.

En el desarrollo de este trabajo, hemos visto una gran variedad de métodos para el manejo de observaciones discrepantes y notado, al menos en modelos lineales, no existe un método óptimo en todos sentidos, y además, en la mayoría de los casos éstos proporcionan esencialmente los mismos resultados. Por otro lado, también se observa que dichos métodos, para el caso de $K=1$, se basan principalmente en residuales y matrices de proyección, y para $K>1$ en cantidades equivalentes como son $A_k(i)$, $V(i)$.

Naturalmente la elección de un método no es fácil y esto dependerá de la filosofía que cada uno de nosotros considere tiene mayores ventajas, aunque claro está, uno podría mezclar varios enfoques para un análisis más profundo, también de la información a priori disponible, y por supuesto de los objetivos específicos del análisis. La verdad es que muchos de los métodos que hemos analizado pueden tener un lugar en un problema específico, dependiendo de los requisitos del investigador: desde métodos cuya aplicación no representa problema alguno (gráficas, resultados tradicionales de mínimos cuadrados, etc.); métodos ligeramente sofisticados (métodos robustos principalmente) hasta métodos para los cuales, casi en toda situación, se tenga que construir un modelo particular (métodos Bayesianos).

Pero la cuestión de cuando una observación es mala o informativa acerca de lo inapropiado del modelo y cuál es la mejor técnica para distinguir tales observaciones, no puede decidirse en base a un conjunto de datos analizados a través de los años, para los cuales cada autor de un método propone su modelo. lo

que se requiere hacer es aplicar éstos métodos a "datos reales" donde se puedan checar las mediciones sospechosas y si es necesario, repetirlos.

¿Debemos sujetar a rutina todas las muestras a algún tipo de observaciones discrepantes? La respuesta es sí, en modelos lineales, por ejemplo, una inspección de los residuales studentizados, la diagonal de V y una medida de influencia por lo menos para un caso deberían de incluirse como parte de la rutina en fase de diagnóstico de cada análisis. Generalmente, dichos análisis son una herramienta importante para guiar el análisis subsecuente y puede ser un indicador de la necesidad de acomodado, transformación o revisión del modelo.

¿Debemos buscar en forma rutinaria rechazar observaciones discrepantes en base a pruebas formales? La respuesta es no. Durante la fase de construcción del modelo, por ejemplo, las pruebas formales no deben de usarse ya que requieren que el modelo nulo sea exacto en términos de la esperanza y distribución de los errores. Deben de buscarse pruebas formales solo en situaciones en donde se sabe que el modelo nulo es exacto en ausencia de contaminantes, y el interés se centra en el estudio del fenómeno alternativo. Existe poca información sobre el funcionamiento de los estimadores mínimos cuadrados usuales en combinación con rechazo vía pruebas formales.

En modelos lineales normales, quizás el problema más difícil y que no se ha resuelto a satisfacción, es el concerniente a cómo tratar con la posibilidad de más de una observación discrepante ($K > 1$); cuando hay poca información a priori importante sobre el número y tipos de observaciones discrepantes cuando, identificación

es lo más importante. Algunos métodos recomendados de identificación surgen como derivado de métodos de acomodo. Los pesos de regresión robusta y métodos bayesianos, por ejemplo, caen en esta categoría. Sin embargo, los métodos de acomodo se basan en información a priori específica como simetría, por ejemplo y se desarrolla con las propiedades de los estimadores resultantes en mente.

Table de Lund.

($\alpha = .10$)

r	1	2	3	4	5	6	7	10	15	25
5	1.57									
6	2.00	1.89	-							
7	2.19	2.02	1.93							
8	2.18	2.12	2.03	1.91						
9	2.24	2.20	2.13	2.05	1.92					
10	2.30	2.26	2.21	2.15	2.06	1.92				
12	2.39	2.37	2.33	2.29	2.24	2.17	1.93			
14	2.47	2.45	2.42	2.39	2.36	2.32	2.19	1.94		
16	2.53	2.51	2.50	2.47	2.45	2.42	2.34	2.20		
18	2.53	2.57	2.56	2.54	2.52	2.50	2.44	2.35		
20	2.63	2.62	2.61	2.59	2.58	2.56	2.52	2.46	2.11	
25	2.72	2.72	2.71	2.70	2.69	2.64	2.66	2.63	2.50	
30	2.80	2.79	2.79	2.78	2.77	2.77	2.75	2.73	2.66	2.13
35	2.86	2.85	2.85	2.85	2.84	2.84	2.82	2.81	2.77	2.55
40	2.91	2.91	2.90	2.90	2.90	2.89	2.87	2.87	2.84	2.72
45	2.95	2.95	2.95	2.95	2.94	2.94	2.93	2.93	2.90	2.82
50	2.99	2.99	2.99	2.99	2.98	2.98	2.98	2.97	2.95	2.89
60	3.06	3.06	3.05	3.05	3.05	3.05	3.05	3.04	3.03	3.03
70	3.11	3.11	3.11	3.11	3.11	3.11	3.10	3.10	3.09	3.07
80	3.16	3.16	3.16	3.15	3.15	3.15	3.15	3.15	3.14	3.12
90	3.20	3.20	3.19	3.19	3.19	3.19	3.19	3.19	3.18	3.17
100	3.23	3.23	3.23	3.23	3.23	3.23	3.23	3.22	3.22	3.21

($\alpha = .05$)

r	1	2	3	4	5	6	7	10	15	25
5	1.92									
6	2.07	1.93								
7	2.19	2.03	1.94							
8	2.24	2.23	2.13	1.94						
9	2.35	2.29	2.21	2.10	1.95					
10	2.42	2.37	2.31	2.22	2.11	1.95				
12	2.52	2.49	2.45	2.39	2.33	2.24	1.96			
14	2.61	2.54	2.55	2.51	2.47	2.41	2.25	1.96		
16	2.63	2.65	2.63	2.62	2.57	2.53	2.43	2.26		
18	2.73	2.72	2.70	2.68	2.65	2.62	2.55	2.44		
20	2.75	2.77	2.76	2.74	2.72	2.70	2.64	2.57	2.15	
25	2.89	2.83	2.87	2.86	2.84	2.83	2.80	2.76	2.60	
30	2.96	2.96	2.95	2.94	2.93	2.93	2.93	2.84	2.79	2.17
35	3.03	3.02	3.02	3.01	3.00	3.00	2.98	2.97	2.91	2.64
40	3.04	3.04	3.07	3.07	3.06	3.06	3.05	3.00	3.00	2.84
45	3.13	3.12	3.12	3.12	3.11	3.11	3.10	3.09	3.06	2.96
50	3.17	3.16	3.16	3.16	3.15	3.15	3.14	3.14	3.11	3.04
60	3.23	3.23	3.23	3.23	3.22	3.22	3.22	3.21	3.20	3.15
70	3.29	3.29	3.23	3.23	3.23	3.23	3.27	3.27	3.26	3.23
80	3.33	3.33	3.33	3.33	3.33	3.33	3.32	3.32	3.31	3.29
90	3.37	3.37	3.37	3.37	3.37	3.37	3.36	3.36	3.36	3.34
100	3.41	3.41	3.40	3.40	3.40	3.40	3.40	3.40	3.39	3.38

Continuación

($\alpha = .01$)

n	p									
	1	2	3	4	5	6	8	10	15	25
5	1.93									
6	2.17	1.98								
7	2.32	2.17	1.98							
8	2.44	2.32	2.18	1.98						
9	2.54	2.44	2.33	2.18	1.99					
10	2.62	2.55	2.45	2.33	2.18	1.99				
12	2.76	2.70	2.64	2.56	2.46	2.34	1.99			
14	2.86	2.82	2.78	2.72	2.65	2.57	2.35	1.99		
16	2.95	2.92	2.88	2.84	2.79	2.73	2.58	2.35		
18	3.02	3.00	2.97	2.94	2.90	2.85	2.75	2.59		
20	3.08	3.06	3.04	3.01	2.98	2.95	2.87	2.76	2.20	
25	3.21	3.19	3.18	3.16	3.14	3.12	3.07	3.01	2.73	
30	3.30	3.29	3.28	3.26	3.25	3.24	3.21	3.17	3.04	2.21
35	3.37	3.36	3.35	3.34	3.34	3.33	3.30	3.23	3.19	2.81
40	3.43	3.42	3.42	3.41	3.40	3.40	3.33	3.36	3.30	3.08
45	3.48	3.47	3.47	3.46	3.46	3.45	3.44	3.43	3.38	3.23
50	3.52	3.52	3.51	3.51	3.51	3.50	3.49	3.48	3.45	3.34
60	3.60	3.59	3.59	3.59	3.58	3.58	3.57	3.56	3.54	3.48
70	3.65	3.65	3.65	3.65	3.64	3.64	3.64	3.63	3.61	3.57
80	3.70	3.70	3.70	3.70	3.69	3.69	3.69	3.68	3.67	3.64
90	3.74	3.74	3.74	3.74	3.74	3.74	3.73	3.73	3.72	3.70
100	3.78	3.78	3.78	3.77	3.77	3.77	3.77	3.77	3.76	3.74

n = número de observaciones.

p = número de variables independientes.

Tabla de Weisberg para la estadística t_i

$\alpha = .05$

n / P	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
6	4.85	6.23	10.89	76.39											
7	4.38	5.67	6.58	11.77	89.12										
8	4.12	4.53	5.26	6.90	12.57	101.9									
9	3.95	4.22	4.64	5.44	7.18	13.36	114.6								
10	3.83	4.03	4.32	4.77	5.60	7.45	14.09	127.3							
11	3.75	3.90	4.10	4.40	4.88	5.75	7.70	14.78	140.1						
12	3.69	3.81	3.96	4.17	4.49	4.98	5.89	7.94	15.44	152.8					
13	3.65	3.74	3.84	4.02	4.24	4.56	5.08	6.02	8.16	16.08	165.5				
14	3.61	3.69	3.79	3.91	4.07	4.30	4.63	5.16	6.14	0.37	16.69	178.2			
15	3.58	3.65	3.73	3.83	3.95	4.12	4.36	4.70	5.25	6.25	8.58	17.28	191.0		
16	3.56	3.62	3.68	3.77	3.87	4.00	4.17	4.41	4.76	5.33	6.36	9.77	17.85	203.7	
17	3.54	3.59	3.65	3.72	3.80	3.90	4.04	4.21	4.46	4.82	5.40	6.47	8.95	18.40	216.4
18	3.53	3.57	3.62	3.68	3.75	3.83	3.94	4.08	4.26	4.51	4.88	5.47	6.57	7.13	18.93
19	3.52	3.56	3.60	3.65	3.71	3.78	3.86	3.97	4.11	4.30	4.55	4.93	5.54	6.67	9.30
20	3.51	3.54	3.58	3.62	3.67	3.73	3.81	3.89	4.00	4.15	4.33	4.59	4.98	5.60	6.76
21	3.50	3.53	3.57	3.60	3.65	3.70	3.76	3.83	3.92	4.03	4.18	4.37	4.64	5.03	5.67
22	3.50	3.52	3.55	3.59	3.63	3.67	3.72	3.78	3.86	3.95	4.06	4.21	4.40	4.68	5.09
23	3.49	3.52	3.54	3.57	3.61	3.65	3.69	3.75	3.81	3.88	3.98	4.09	4.24	4.44	4.71
24	3.49	3.51	3.53	3.56	3.59	3.63	3.67	3.71	3.77	3.83	3.91	4.00	4.12	4.27	4.47
25	3.48	3.50	3.53	3.55	3.58	3.61	3.65	3.69	3.73	3.79	3.85	3.93	4.02	4.14	4.30
26	3.48	3.50	3.52	3.54	3.57	3.60	3.63	3.66	3.70	3.75	3.81	3.87	3.95	4.05	4.17
27	3.48	3.50	3.52	3.54	3.56	3.58	3.61	3.65	3.68	3.72	3.77	3.83	3.89	3.97	4.07
28	3.48	3.50	3.51	3.53	3.55	3.58	3.60	3.63	3.66	3.70	3.74	3.79	3.84	3.91	3.99
29	3.48	3.49	3.51	3.53	3.55	3.57	3.59	3.62	3.64	3.68	3.71	3.76	3.81	3.86	3.93
30	3.48	3.49	3.51	3.52	3.54	3.56	3.58	3.60	3.63	3.66	3.69	3.73	3.77	3.82	3.88
31	3.48	3.49	3.50	3.52	3.54	3.55	3.57	3.59	3.62	3.64	3.67	3.71	3.74	3.79	3.84
32	3.48	3.49	3.50	3.52	3.53	3.55	3.57	3.59	3.61	3.63	3.66	3.67	3.72	3.76	3.80
33	3.48	3.49	3.50	3.52	3.53	3.54	3.56	3.58	3.60	3.62	3.64	3.67	3.70	3.74	3.77
34	3.48	3.49	3.50	3.51	3.53	3.54	3.56	3.57	3.59	3.61	3.63	3.66	3.68	3.71	3.75
35	3.48	3.49	3.50	3.51	3.52	3.54	3.55	3.57	3.58	3.60	3.62	3.64	3.67	3.70	3.73
36	3.48	3.49	3.50	3.51	3.52	3.54	3.55	3.56	3.58	3.60	3.61	3.63	3.66	3.68	3.71
37	3.48	3.49	3.50	3.51	3.52	3.53	3.55	3.56	3.57	3.59	3.61	3.62	3.65	3.67	3.69
38	3.48	3.49	3.50	3.51	3.52	3.53	3.54	3.56	3.57	3.58	3.60	3.62	3.64	3.66	3.68
39	3.49	3.49	3.50	3.51	3.52	3.53	3.54	3.55	3.57	3.58	3.59	3.61	3.63	3.65	3.67
40	3.47	3.49	3.50	3.51	3.52	3.53	3.54	3.55	3.56	3.58	3.59	3.60	3.62	3.64	3.66
50	3.51	3.51	3.51	3.52	3.53	3.53	3.54	3.54	3.54	3.54	3.57	3.57	3.58	3.59	3.60
60	3.53	3.53	3.53	3.54	3.54	3.54	3.55	3.55	3.56	3.56	3.57	3.57	3.58	3.58	3.59
70	3.55	3.55	3.55	3.55	3.56	3.56	3.56	3.56	3.57	3.57	3.57	3.58	3.58	3.59	3.59
80	3.57	3.57	3.57	3.57	3.57	3.58	3.58	3.58	3.58	3.58	3.59	3.59	3.59	3.60	3.60
90	3.58	3.59	3.59	3.59	3.59	3.59	3.59	3.60	3.60	3.60	3.60	3.60	3.60	3.61	3.61
100	3.60	3.60	3.60	3.60	3.61	3.61	3.61	3.61	3.61	3.61	3.61	3.62	3.62	3.62	3.62
200	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.73	3.74
300	3.81	3.81	3.81	3.81	3.81	3.81	3.81	3.81	3.81	3.81	3.82	3.82	3.82	3.82	3.82
400	3.87	3.87	3.87	3.87	3.87	3.87	3.87	3.87	3.88	3.88	3.88	3.88	3.88	3.88	3.88
500	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92	3.92

tabla de Weisberg - Continuación

$\alpha = .01$

n / p	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
4	7.53	10.87	24.46	382.0											
7	6.35	7.84	11.45	26.43	445.6										
8	5.71	6.54	8.12	11.93	29.26	509.3									
9	5.31	5.84	6.71	8.38	12.47	29.92	573.0								
10	5.04	5.41	5.96	6.87	8.41	12.92	31.60	636.6							
11	4.85	5.12	5.50	6.07	7.01	8.83	13.35	33.14	700.3						
12	4.71	4.91	5.19	5.58	6.17	7.15	9.03	13.75	34.62	763.9					
13	4.60	4.76	4.97	5.25	5.66	6.26	7.27	9.22	14.12	36.03	827.6				
14	4.51	4.64	4.81	5.02	5.32	5.73	6.33	7.39	9.40	14.48	37.40	891.3			
15	4.44	4.55	4.68	4.85	5.05	5.37	5.80	6.43	7.50	9.57	14.82	38.71	954.9		
16	4.38	4.48	4.59	4.72	4.90	5.12	5.43	5.96	6.51	7.60	9.73	15.15	39.98		
17	4.34	4.41	4.51	4.62	4.76	4.94	5.17	5.48	5.92	6.59	7.70	9.83	15.46	41.21	
18	4.30	4.36	4.44	4.54	4.66	4.80	4.98	5.21	5.53	5.98	6.64	7.80	10.03	15.76	42.41
19	4.26	4.32	4.39	4.47	4.57	4.69	4.83	5.01	5.25	5.57	6.03	6.72	7.89	10.17	16.05
20	4.23	4.29	4.35	4.42	4.50	4.60	4.72	4.86	5.05	5.29	5.62	6.08	6.79	7.98	10.31
21	4.21	4.25	4.31	4.37	4.44	4.52	4.62	4.74	4.89	5.08	5.33	5.66	6.13	6.85	8.05
22	4.19	4.23	4.29	4.33	4.39	4.46	4.55	4.65	4.77	4.92	5.11	5.36	5.70	6.18	6.91
23	4.17	4.21	4.25	4.30	4.35	4.41	4.49	4.57	4.67	4.80	4.95	5.14	5.40	5.74	6.22
24	4.15	4.19	4.22	4.27	4.32	4.37	4.43	4.51	4.59	4.70	4.82	4.98	5.17	5.43	5.78
25	4.14	4.17	4.20	4.24	4.28	4.33	4.39	4.45	4.53	4.62	4.72	4.85	5.00	5.20	5.46
26	4.12	4.15	4.18	4.22	4.26	4.30	4.35	4.41	4.47	4.55	4.64	4.74	4.87	5.03	5.23
27	4.11	4.14	4.17	4.20	4.24	4.27	4.32	4.37	4.43	4.49	4.57	4.66	4.76	4.89	5.05
28	4.10	4.13	4.15	4.18	4.21	4.25	4.29	4.33	4.38	4.44	4.51	4.59	4.68	4.78	4.91
29	4.09	4.12	4.14	4.17	4.20	4.23	4.26	4.30	4.35	4.40	4.46	4.53	4.60	4.69	4.80
30	4.09	4.11	4.13	4.15	4.18	4.21	4.24	4.28	4.32	4.36	4.42	4.47	4.54	4.62	4.71
31	4.08	4.10	4.12	4.14	4.17	4.19	4.22	4.26	4.29	4.33	4.38	4.43	4.49	4.56	4.64
32	4.07	4.09	4.11	4.13	4.15	4.18	4.21	4.24	4.27	4.31	4.35	4.39	4.45	4.50	4.57
33	4.07	4.08	4.10	4.12	4.14	4.17	4.19	4.22	4.25	4.28	4.32	4.36	4.41	4.46	4.52
34	4.06	4.08	4.09	4.11	4.13	4.15	4.18	4.20	4.23	4.26	4.29	4.33	4.37	4.42	4.47
35	4.06	4.07	4.09	4.11	4.12	4.14	4.16	4.19	4.21	4.24	4.27	4.31	4.34	4.37	4.43
36	4.05	4.07	4.08	4.10	4.12	4.13	4.15	4.18	4.20	4.22	4.25	4.28	4.32	4.36	4.40
37	4.05	4.06	4.08	4.09	4.11	4.13	4.14	4.16	4.19	4.21	4.24	4.26	4.29	4.33	4.37
38	4.05	4.06	4.07	4.09	4.10	4.12	4.13	4.15	4.17	4.20	4.22	4.25	4.27	4.31	4.34
39	4.04	4.05	4.07	4.08	4.10	4.11	4.13	4.14	4.16	4.18	4.21	4.23	4.26	4.28	4.32
40	4.04	4.05	4.06	4.08	4.09	4.10	4.12	4.14	4.15	4.17	4.19	4.22	4.24	4.27	4.29
50	4.03	4.03	4.04	4.05	4.06	4.07	4.08	4.09	4.10	4.12	4.14	4.16	4.18	4.21	4.12
60	4.03	4.03	4.04	4.04	4.05	4.05	4.06	4.06	4.07	4.08	4.08	4.09	4.10	4.11	4.12
70	4.03	4.03	4.04	4.04	4.05	4.05	4.06	4.06	4.07	4.08	4.08	4.09	4.10	4.11	4.12
80	4.04	4.04	4.04	4.04	4.05	4.05	4.06	4.06	4.06	4.07	4.07	4.07	4.08	4.08	4.09
90	4.05	4.05	4.05	4.05	4.06	4.06	4.06	4.06	4.07	4.07	4.07	4.07	4.08	4.08	4.09
100	4.06	4.06	4.06	4.06	4.06	4.07	4.07	4.07	4.07	4.08	4.08	4.08	4.08	4.09	4.09
200	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15	4.15
300	4.21	4.21	4.21	4.21	4.21	4.21	4.22	4.22	4.22	4.22	4.22	4.22	4.22	4.22	4.22
400	4.26	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27	4.27
500	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31	4.31

BIBLIOGRAFIA.

Abraham, B., and Box, G.E.P (1978) "Linear Models and Spurious Observations", Applied Statistics, 27, 131-138.

Aitkin, M., and Wilson, G.T. (1980) "Mixture Models, Outliers and the E.M. Algorithm", Technometrics 22, 325-332.

Andrews, D.F. (1971), "A note on the Selection of data transformations", Biometrika, 58, 249-254.

———, Bickel, P.J., Hampel, F.R., Huber, P.J., Rogers, W.H., and Tukey, J.W. (1972), "Robust Estimates of Location", Princeton, N.J.: Princeton University Press.

———, and Pregibon, D. (1978) "Finding the Outliers that Matter", Journal of the Royal Statistical Society Ser. B, 40, 37-43.

Ansecombe, F.J. (1960) "Rejection of Outliers" Technometrics 2, 123-147.

Atkinson, A.C. (1981), "Robustness transformations and two Graphical Displays for Outlying and Influential Observations in Regression", Biometrika, 68, 13-20.

———, (1982), "Robust and Diagnostic Regression Analysis" Commun. Statist. Theor. Meth., 11(22), 2559-2571.

———, (1982), "Regression Diagnostics Transformations and Constructed -

Variables" (with discussion), Journal of the Royal Statistics Society, Ser. B, 44, 1-35.

_____, (1986), "Diagnostics for transformations", *Technometrics*, 28, 29-34.

Barnett, V. (1973), "Comparative Statistical Inference", John Wiley and Sons, Inc. New York.

_____, (1978), "The Study of Outliers: Purpose and Model", *Applied Statistics* 27, 242-250.

_____, and Lewis, T. (1978), "Outliers in Statistical Data", New York: John Wiley.

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.

_____, and Cook, R. D. (1983) "Outlier....s", *Technometrics* 25, 119-149.

Bernoulli, D. (1777), "The Most Probable Choice Between Several Discrepant Observations and the Formation there from of the Most Likely Induction", in C. G. Allen (1961), *Biometrika* 48, 3-13.

Besley, D. A., Kuh, E. and Welsch, R. E. (1980), "Regression Diagnostics", New York: John Wiley.

Bikel, P. J. and Doksum, K. A. (1981), "An Analysis of Transformations Revisited", *Journal of the Royal Statistical Association*, 76, 296-311.

Box, G.E.P. (1979), "Robustness in Strategy of Scientific Model - Building" eds. R.L. Launer and G.N. Wilkinson, New York, Academic Press.

_____, (1980), "Sampling and Bayes Inference in Scientific Modelling and Robustness (with discussion)." J. R. Statist. Soc. A., 143, 383-480.

_____, and Anderson S.L. (1955), "Permutation theory in the Derivation of Robust Criteria and the Study of Departures from Assumption," J. Roy Statist. Soc. Ser: B, 17, 1-34.

_____, and Cox, D.R. (1964), "An Analysis of transformations (With discussion)," J.R. statist Soc. A, 143, 383-430.

_____, and tiao, G.C. (1968), "A Bayesian Approach to Some Outlier Problems", Biometrika, 55, 119-129.

_____, and tiao, G.C. (1973), "Bayesian Inference in Statistical Analysis. Addison-Wesley. Reading, Massachusetts.

_____, and Tidwell, P.W. (1962), "transformations of the Independent Variables", Technometrics 4, 531-550.

Broemeling, L.D. (1985) "Bayesian Analysis of Linear Models", Oklahoma State University Marcel Dekker, Inc: New York and Basel.

Brownlee, K. A. (1965). Statistical theory and Methodology in Science and Engineering, New York: John Wiley.

- Carroll, R.J. (1980 a) "A Robust method for testing transformations to Achieve Aproximate normality", J.R. Statist. Soc. B, 42, 71-78.
- , (1980b), "Robust Methods for factorial Experiments with Outliers" Applied Statist 29, 246-251.
- , (1982), "two Examples of transformations when there are possible Outliers". Appl. Statist. 31, 149-152.
- Chambers, R.L. and Heathcote, C.E. (1981), "On the estimation of slope and the Identification of Outliers in Linear Regression", Biometrika, 68, 21-33.
- Collet, D. and Lewis, T (1976), "the subjective Nature of Outlier Rejection Procedures", Applied Statistics, 25, 228-237.
- Cook, R.D. (1977), "Detection of Influential Observations in Linear Regression" Technometrics, 19, 15-18.
- , (1979), "Influential Observations in Linear Regression", Journal of the American Statistical Association 74, 169-174.
- , Hasehuh, N., and Weisberg, S. (1982), "A note on an Alternative Outlier Model", Journal of the Royal Statistical Societus, Ser B. 44, 370-376.
- , and Prescott, P. (1981), "On the Accuracy of Bonferroni Significance Levels for Detecting Outliers in Linear Models", Technometrics 24, 59-63.

- _____, and Wang, P.C. (1983), "transformations and Influential Cases in Regression", *Technometrics* 25, 337-343.
- _____, and Weisberg, S. (1980), "characterizations of an Empirical Influence Function for Detecting Influential Cases in Regression", *Technometrics*, 22, 495-508
- _____, (1982), "Residuals and Influence in Regression", New York: Chapman Hall.
- Daniel, C., and Woods, F.S. (1971), "Fitting Equations to Data", New York: John Wiley.
- Daniel, P.J. (1920), "Observations Weighted According to Order", *American Journal of Mathematics*, 42, 222-236.
- Davies, R.B., and Hutton, B. (1975) "the Effect of Errors in the Independent Variables in Linear Regression" *Biometrika* 62, 383-391.
- Demster, A.P. and Gasko-Green, M. (1981), "New tools for Residual Analysis", *Annals of Statistics*, 9, 945-959.
- _____, Laird, N.M. and Rubin, D.B. (1977), "Maximum Likelihood from Incomplete Data via the E.M. Algorithm (with discussion) *J. Roy. Statist. Soc. B*, 39, 1-38.
- Doornik, R. (1981), "testing for a Single Outlier in a Linear Model", *Biometrika* 67, 705-712.

Draper, N. R., and John J. A. (1980), "Testing for three or fewer outliers in two-way tables", *Technometrics* 22, 9-15.

———, and John J. A. (1981), "Influential Observations and Outliers in Regression" *Technometrics* 23, 21-26.

———, and Smith, H. (1981) "Applied Regression Analysis, Wiley, New York.

Dutter, R. (1977), "Numerical Solution of Robust Regression Problems: Computational Aspects, a Comparison," *Journal of Statistical Computation and Simulation*, 5, 207-238.

———, and Guttman (1979), "On Estimation in the linear Model when Spurious Observations are Present a Bayesian Approach" *Commun. Statist. theor. Meth.*, A8 (7), 61 - 635).

Edgeworth, F. Y. (1887), "On Discordant Observations", *Philosophical Magazine*, 23 Ser. 5, 364-375

Efron, B. (1965), "the Convex Hull of a Random Set of Points", *Biometrika*, 52, 331-343

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, 941-963.

——— (1976) "Testing for a Single Outlier From a General linear Regression" *Biometrics* 32, 637-645.

Freeman, P.R. (1981) "On the number of Outliers in data from a linear Model"; In Bayesian Statistics (eds. Bernardo, J.M.), 349 - 365. University Press. Valencia.

Furnival, G.M. and Wilson, R.W. Jr. (1974) Regressions by leaps and bounds. *Technometrics*. 16, 499-551.

Geisser, S. (1980). In Discussion of Box (1980)

Gentle, J.E. (1978), "Testing for Outliers in linear Regression" in Contributions to Surety Sampling and Applied Statistics in Honor of H.O. Hartley, ed. H.A. David, New York: Academic Press.

Gentleman, J.F. (1980), "Finding the k Most Likely Outliers in two-Way tables" *Technometrics* 22, 591-600.

———, and Wilk, M.B. (1975a), "Detecting Outliers in a two-Way table: I statistical Behavior of Residuals" *Technometrics*, 17, 1-14.

———, "Detecting Outliers. II. Supplementing the Direct Analysis of Residuals" *Biometrics* 31, 387-410.

Glaisher, J.W.L. (1873), "On the Rejection of Discordant Observations", *Monthly Notices of the Royal Astronomical Society*, 23, 141-160.

Gould, B.A:Jr (1855) "On Peirces criterion for the Rejection of Doubtful Observations with tables for facilitating its

Application", *Astronomical Journal* 6, 81-83.

Graedà-Medrano, L.E. (1984), "Aplicación de técnicas de Regresión Robusta", tesis para obtener el título de Actuario, Facultad de Ciencias: UNAM

Green, P.J. and Silverman, B.W. (1979), "Constructing the Convex Hull of a Set of Points, in the Plane. the *Computer Journal* 22, 262, 266.

Grubbs, F.E. (1969), "Procedures for Detecting Outlying Observations in Samples", *Technometrics*, 15, 385-404.

Guttman, I. (1973a), "Care and Handling of Univariate or Multivariate Outliers in Detecting Spuriousity. A Bayesian Approach", *Technometrics*, 15, 723-738.

_____, and Dutter, R., and Freeman, P.R. (1978). Care and Handling of Univariate Outliers in the General Linear Model to detect Spuriousity: A Bayesian Approach", *Technometrics* 20, 187-193.

_____, and Dutter, R., (1976), "Procedures for Investigating Outliers when Estimating in the General Univariate Linear Situation Nonfull Rank Case"; *Communications in Statistics A5*, 819-835.

_____, and Kahatri, C.G., (1975), "A Bayesian Approach to Some Problems Involving the detection of Spuriousity", in *Applied Statistics Symposium*, ed. R.P. Gupta, Amsterdam: North. Holland, 111-146.

_____, and Smith, D.F. (1969), "Investigation of Rules for Dealing with Outliers in Small Samples from the Normal Distribution: I: Estimation of the Mean" *technometrics*, 11, 527 - 550.

Hawkins, D.M. (1979), "Fractiles of an Extended Multiple Outliers test", *Journal of Statistical Computation and Simulation*, 8, 227 - 236.

_____, (1980), "Identification of Outliers", London: Chapman and Hall.

Hewett, J. and Bulgren, W.G. (1971), "Inequalities for Some Multivariate F-Distributions with Applications", *technometrics*, 13, 397 - 402.

Hoaglin, D.C., and Welsh, R. (1978), "The hat Matrix in Regression and ANOVA", *the American Statistician* 32, 17-22.

Hocking, R.R. (1983), "Developments in Linear Regression, Methodology: 1859-1982", *technometrics*, 25, 219-249.

Hogg, R.V. (1979), "Statistical Robustness: One View of its Use in Applications today", *the American Statistician* 33, 108-115.

_____, (1979b) "An Introduction to Robust Estimation" *Robustness in Statistics*, Academic Press.

Holland, P.W. and Welsh, R.E. (1977), "Robust Regression using iteratively reweighted least squares" *Communications in Statistics*, A6,

013-028.

Huber, P.J. (1972), "the 1972 Wald Lecture Robust Statistics: A Review" *Annals of Mathematical Statistics*, 35, 73-101.

———, (1981), *Robust Statistics*, New York: John Wiley.

Irwin, J.D. (1925), "On a Criterion for the Rejection of Outlying Observations", *Biometrika*, 17, 238-250.

Jaekel, L.A. (1972) "Estimating Regression Coefficients by minimizing the Dispersion of the Residual", *Ann. Math. Statistics*, 43, 1449 - 1468.

John, J.A. (1978), "Outliers in Factorial Experiments", *Applied Statistics*, 27, 111-119.

———, and Draper, N.R. (1978), "On testing for two Outliers One Outlier in two-way tables", *Technometrics* 20, 69-78.

———, and Prescott (1975), "Critical values of a test Detect Outliers in Factorial Experiments", *Applied Statistics* 24, 56-59.

Johnson, B.A. and Hunt, H.H. (1979), "Performance characteristics for certain tests to Detect Outliers" *Proceedings of the Statistical Computing Section, American Statistical Association* 247-249.

Joshi, P.C. (1972) "Some Slippage tests of Mean for a Simple Outlier in Linear Regression", *Biometrika*, 59, 109-120.

- Kitagawa, G. (1979), "On the use of AIC. for Detection of Outliers", *Technometrics* 21, 193-199.
- Krasker, W.S., and Welch, R.E. (1982), "Efficient Bounded - Influence Regression Estimation", *Journal of the American Statistical Association*, 695-699.
- Lindley, D.V. and Smith, A.J.M. (1972), "A Bayes Estimates for the Linear Model (With discussion)", *J. R. Statist Soc. B.*, 34, 1-41.
- Luendberger, D.G. (1973) *Introduction to Linear and non linear Programming*. Addison - Wesley. Menlo Park California. 148-155.
- Lund, R.E. (1975) "Tables for an Approximate test for Outliers in Linear Regression", *Technometrics*, 17, 473-476.
- Mickey, M. R., Dunn, O. J., and Clark, V. (1967), "Note on the Use of Stepwise Regression, in Detecting Outliers", *Computers and Biomedical Research*, 1, 105-111.
- Miyashita, H. and Newbold, P. (1983), "On the Sensitivity to Non-Normality of a test of Outliers in Linear Models", *Commun Statist. - theor. Meth.*, 12(12), 1413-1419.
- Montgomery, D.C. and Peck, E.A. (1982), "Introduction to Linear Regression Analysis", New York: John Wiley.
- Neave, H. R. (1978), "Statistics tables", Allen and Unwin - London.
- Newcomb, S. (1886) "A Generalized theory of the Combination of

Observations to obtain the best result", American Journal of Mathematics, 8, 343-366.

Neyman, J. and Scott, E.L. (1971), "Outlier Proneness of phenomena and of Related Distributions", in Optimizing Methods in Statistics, ed. J.S. Rustagi New York: Academic Press.

Palacios, R.B. y Castañón, V.I., "Observaciones Disordantes en Análisis de Regresión: Detección y tratamiento", tesis para obtener título de Actuario, Facultad de Ciencias, UNAM.

Pearson, E.S., and Chandra Sekar C. (1936), "The Efficiency of Statistical tools and Criterion for the Rejection of Outlying Observations", Biometrika 28, 308-320.

Pierce, B. (1882) "Criterion for the Rejection of Doubtful Observations", Astronomical Journal. 2, 161-163.

Pettit, L.I., and Smith, A.F.M. (1983), "Outliers and Influential Observations in Linear Models", Second Valencia International Meeting on Bayesian Statistics.

Plackett, R.L. (1952), "Some theorems in Least Squares", Biometrika 37, 149-157.

Prescott, P. (1975), "An Approximate test for Outliers in Linear Models" Technometrics 17, 129-132.

Rosner, B. (1975), "On the Detection of Many Outliers" Technometrics 17, 221-227.

Schweder, t. (1976), "Some Optimal Methods to Detect Structural Shift or Outliers in Regression" *Journal of the American Statistical Association*, 68, 872-879.

Snedecor, G.W. and Cochran, W.G. (1967), "Statistical Methods" (6th ed.), Iowa State University Press.

Srikantan, K.S. (1961) "Testing for a single Outlier in a Regression Model", *Sankya, Ser. A*, 251-260.

Stefansky, W. (1971), "Rejecting Outliers by Maximum Normed Residual", *the Annals of Mathematical of Statistics* 42, 35-45

———, (1972). "Rejection Outliers in Factorial Designs", *technometrics*, 14, 469-479.

Stigler, S.M. (1973), "Simon Newcomb; Percy Dannel and the History of Robust Estimation (1885-1920)", *Journal of the American Statistical Association*, 68, 872-879.

Stone, E.J. (1868) "On the Rejection of Disoordant Observations" *Monthly Notices of the Royal Astronomical Society*, 34, 9-15.

tiao, G.C. and Box, G.E.P. (1974) "Some comments on Bayes Estimators", in S.E. Fienberg and A. Zellner (eds.) *Studies in Bayesian Econometrics and Statistics*: North-Holland

tietjen, G.L., Moore, R.L. and Beckman, R.J. (1973), "testing for a Outlier in Simple Linear Regression?" *technometrics* 15, 717-721.

Weisberg, S. (1980), *Applied Linear Regression* New York: John Wiley.