



36
20j
UNIVERSIDAD NACIONAL AUTÓNOMA DE MÉXICO

FACULTAD DE CIENCIAS

**"DESARROLLO DE UN PAQUETE ESTADÍSTICO
EN APOYO A LOS CURSOS DE ANÁLISIS
DE REGRESIÓN".**

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

CLAUDIA JIMENEZ VILLASEÑOR

MEXICO, D. F.

SEPTIEMBRE 1991.

**TESIS CON
FALLA DE ORIGEN**



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 :	página
Introducción	1
Capítulo 1.	3
Capítulo 2.	13
Capítulo 3.	26
Capítulo 4.	31
Capítulo 5.	44
Conclusiones	57
Apéndice 1	58
Apéndice 2	60
Bibliografía	64

Introducción.

Las técnicas estadísticas hoy en día recurren a paquetes de cómputo, ya que las operaciones que en ellas intervienen llegan a ser tan extensas que la presencia de errores es muy frecuente.

El Análisis de Regresión, como técnica estadística, no escapa a este tipo de problemas por lo cual se ve en la necesidad de recurrir a paquetes de cómputo, sin embargo el uso de estos tiene ciertas desventajas ya que muchas veces se usan de una manera tan automática que se llegan a olvidar las técnicas a través de las cuales se obtiene la información estadística, por lo que surge la preocupación de crear un paquete de cómputo que proporcione los elementos necesarios para obtener resultados sin perder de vista la manera mediante la cual se llega a ellos.

En este trabajo se elaboró un paquete de cómputo en apoyo a los cursos de Análisis de Regresión que se imparten en la Facultad de Ciencias de la UNAM. Fue adaptado para ser ejecutado en microcomputadoras PC o compatibles con el objeto de que todos los alumnos de la facultad que esten interesados en la materia tengan acceso a él.

El paquete de cómputo, proporciona los cálculos necesarios facilitando de esta manera una rápida interpretación y tratando de que el alumno no pierda de vista las técnicas a través de las cuales llega a los resultados deseados.

Debido a la naturaleza de las cálculos utilizados en esta técnica se utilizó el lenguaje de programación PASCAL (Turbo PASCAL versión 5), el cual es muy apegado a este tipo de operaciones.

El presente trabajo se encuentra estructurado de la siguiente forma :

Los tres primeros capítulos estan dedicados a la teoría en la que se basa el Análisis de Regresión.

El capítulo 1 describe en forma detallada la construcción del modelo de Regresión así como cada uno de los supuestos que lo componen con el objeto de relacionarlo con los elementos estadísticos que se conocen para después tratarlo como una técnica de inferencia estadística.

En el capítulo 2 se presenta el modelo de Regresión Lineal Múltiple el cual servirá de base teórica para la implantación del paquete de cómputo.

El capítulo 3 se refiere a selección de variables. Se presentan algunos métodos muy comentados en la literatura, mismos que podrán ser utilizados a través del paquete de cómputo.

En el capítulo 4 se presenta en forma detallada el diseño de cada uno de los programas que conforman el paquete así como los algoritmos de cada uno de ellos.

Por último en el capítulo 5 se muestra la aplicación del paquete mediante ejemplos.

Es importante mencionar que no se pretende competir con los paquetes estadísticos comerciales y menos aún comparar su utilidad. Como se mencionó anteriormente el objetivo central del presente trabajo es proporcionar una herramienta didáctica a los cursos de Análisis de Regresión.

Capítulo 1. Construcción Del Modelo.

Al conjunto de métodos cuyo objetivo es la formulación de modelos matemáticos que muestran relaciones modeladas con el propósito de predecir y hacer inferencias estadísticas se le conoce como *Análisis de Regresión*.

La aplicación de Análisis de Regresión es muy amplia y variada. Como ejemplo se pueden citar las siguientes situaciones:

Explicar el rendimiento de alfalfa (en toneladas por acre) de una granja experimental a través de la cantidad de agua aplicada (en pulgadas por acre), si las variables guardan cierta relación sería posible tratar de predecir el rendimiento a través de la cantidad de agua.

Una cadena de tiendas esta interesada en abrir una nueva unidad para lo cual considera importante saber cómo se explican las ventas semanales en cada tienda, a través del número de empleados y el tamaño de cada tienda. Si existe una relación razonable entre las variables, sería posible predecir las ventas semanales en cada tienda a través del número de empleados y el tamaño de tienda para poder así evaluar si es conveniente abrir o no una nueva tienda.

La palabra *regresión* la utilizó por primera vez Sir Francis Galton quien analizó la altura de hijos en terminos de la altura promedio de sus padres. De sus observaciones concluyó que hijos de padres muy altos eran generalmente más altos que el promedio pero no tan altos como sus padres y lo mismo para el caso de hijos de padres muy bajos. Este estudio se publicó en 1885 bajo el titulo de "*Regresión hacia la mediocridad en la herencia de estatura*". En este caso, el termino regresión se usó en el sentido de que las alturas de los hijos tendían al promedio más que a valores extremos.

En este trabajo se denotará por "y" a la característica de interes (manifestación del fenómeno) y se le denominará variable de respuesta. A la(s) característica(s) a través de la(s) cual(es) se construirá la relación con "y", se le(s) conocerá como variable(s) explicativa(s) (o factores de interes) y se denotarán por x_1, x_2, \dots, x_p .

Para establecer relaciones funcionales entre variables hay que tomar en cuenta el tipo de relaciones que pueden existir. Por ejemplo, Mendez(1977), menciona las siguientes:

1. *Asociación Simple*. Es cuando un fenómeno o modalidad de fenómeno "A" se presenta frecuentemente acompañado de "B", otro fenómeno o modalidad. Sin embargo puede darse "A" sin que ocurra "B" y "B" sin que ocurra "A". Como ejemplo considérese el hecho de fumar (A) y el cáncer pulmonar (B), los dos aspectos ocurren frecuentemente juntos, sin embargo, hay fumadores a los que no les da cáncer y hay personas con cáncer que no son o han sido fumadores.

Este tipo de asociación puede servir de base para la construcción de hipótesis que ligan los fenómenos en términos de causalidad.

Se llama causa al fenómeno o conjunto de fenómenos que preceden a otro y le dan origen.

Se llama efecto el fenómeno que sigue a otro y es originado por él.

2. *Causalidad Probabilística o Relación producto-productor.* Este es el caso en que un fenómeno "B" es necesario para que ocurra otro "A". Pero "B" no es suficiente para que ocurra "A", es decir dada "B" es probable que ocurra "A"; y si ocurre "A" debe haber ocurrido "B". Un ejemplo de esto es observar que el agua de un recipiente esta hirviendo como fenómeno "B" e inferir que la temperatura de ese recipiente es de 100 grados centígrados o más, fenómeno "A". Sin embargo puede suceder que el agua esté hirviendo sin que se eleve la temperatura, sino únicamente por un efecto de vacío, fenómeno "A1", sobre el recipiente considerado. Se dice que "A" es una entre varias de las causas posibles de "B".

3. *Causalidad Determinística o Relación de causa y efecto.* En este caso un fenómeno "A" es necesario y suficiente para que ocurra otro fenómeno "B". Si "A" ocurre, ocurrirá o ha ocurrido "B" y si "B" ocurre, ocurrirá o ha ocurrido "A". Se habla en este caso de que "A" es la causa única de "B" o viceversa. Como ejemplo de esta relación considérese la ocurrencia de un cromosoma 21 por triplicado en un ser humano (trisomía 21) como fenómeno "A" y el síndrome de down o mongolismo, fenómeno "B".

Se puede retomar el ejemplo del hecho de fumar, "A", y el cáncer pulmonar, "B", podría considerarse "A" como causa y "B" como efecto, aquí existe una probabilidad de que una persona que fuma enferme de cáncer, sin embargo pueden existir otras causas que originen que una persona enferme de cáncer.

De acuerdo al tipo de relación que se tenga y al marco de referencia sobre el fenómeno se puede saber si tiene sentido establecer una relación funcional e intentar el ajuste de funciones entre las variables.

Establecer la relación funcional puede ser un problema complicado sin embargo en ocasiones puede aproximarse por medio de una función matemática lo más simple posible.

Como se vió anteriormente se trata de construir una función que relacione de la mejor manera posible la manifestación del fenómeno bajo estudio con los factores de interés.

Al llevar a cabo un experimento con el fin de producir información existen factores que pueden ser de dos tipos : controlados y no-controlados.

Factores Controlados. Variables de interés bajo control que afectan la manifestación del fenómeno y que toman diferentes modalidades .Dentro de esta categoría están las condiciones experimentales que son los factores que permanecen constantes a lo largo del experimento.

Factores No-controlados. Este tipo de factores, como su nombre lo indica no se controlan debido a diversas causas : puede ser que dichos factores se consideraron no relevantes al fenómeno bajo estudio, o ni siquiera se ha percatado su existencia o bien no se pueden manipular.

Sean

y : Variable de respuesta

x_1, \dots, x_p : variables explicativas, es decir los factores de interés que se controlan y toman diversas modalidades.

A_1, A_2, \dots, A_n : condiciones experimentales.

B, C, D, E, \dots : factores no controlados.

La manifestación y queda descrita por la relación

$$y = f_0(x_1, \dots, x_p, A_1, \dots, A_n, B, C, D, E, \dots)$$

Como se puede observar los factores controlados junto con las condiciones experimentales constituyen el dominio de f_0 .

Dado que A_1, \dots, A_n son factores que toman la misma modalidad a lo largo de todo el experimento, se restringe el dominio a :

$$y = f_1(x_1, \dots, x_p, B, C, D, E, \dots) \quad \dots(1.1)$$

Sin embargo tanto f_0 como f_1 son funciones cuyo dominio es infinito y habrá factores que difícilmente se identificaran por lo que no es posible construir una regla de correspondencia bajo estas condiciones.

Para delimitar el argumento de f_2 es necesario considerar algunos supuestos. Sería deseable inducir una partición de tal forma que

$$P_1 = \{ x_1, \dots, x_p \} \quad \dots \text{factores controlados}$$

$$P_2 = \{ B, C, D, F, \dots \} \quad \dots \text{factores no-controlados}$$

Para estos conjuntos sean f_2 y f_3 tales que

$$y = f_2(x_1, \dots, x_p) + f_3(B, C, D, F, \dots) \dots (1.2)$$

De lo anterior se desprende el supuesto estructural básico de separabilidad aditiva: Los efectos que produce el conjunto P_1 en la variable "y" y los efectos que produce el conjunto P_2

en la misma variable, son separables en forma aditiva.

Así el dominio de f_2 es finito, pero el de f_3 sigue siendo infinito.

Sin embargo hay que reconocer que aunque f_3 no podrá determinarse lo que interesa es su magnitud.

Como f_3 representa la función que agrupa a todos los factores que no se controlan y que afectan a la manifestación del fenómeno en mayor o menor grado, es usual que se le denomine "error" y cuya magnitud se denota por ϵ_1

$$y = f_2(x_1, \dots, x_p) + \epsilon_1 \dots (1.3)$$

Nótese que si ϵ_1 es muy grande podría pensarse que no se ha considerado un factor importante o bien no se tomó en cuenta otro factor como factor experimental o bien la relación entre x y y no es muy fuerte.

De (1.3) se obtiene que

$$\epsilon_1 = y - f_2(x_1, \dots, x_p)$$

Por facilidad supóngase que se cuenta con un sólo factor de interés, x , con $x \in P_1$, entonces

$$\epsilon_1^* = y - f_2(x)$$

ϵ_1^* es desconocido y representa la magnitud del error. Sería deseable que esta cantidad fuese lo más pequeña posible.

Entonces el problema es determinar f tal que ϵ_1^* sea mínimo.

Nótese que la variable x puede ser cualitativa (nominal) o cuantitativa (discreta o continua). Si x es cualitativa se cuenta con una colección de posibles valores (o etiquetas). Por ejemplo si la variable a estudiar es sexo, sólo admite dos posibles valores: hombre y mujer (h,m) o bien "0,1" o bien "1,2", etc. Los valores que se asignan son arbitrarios y por tanto no tiene ningún sentido establecer una función. La solución a este problema se enmarca dentro de las técnicas conocidas como *Diseño de Experimentos*.

Si x es cuantitativa, se tiene que definir una familia de funciones y escoger f_2 en alguna forma.

Sea $f_2 \in L$, con $L = \{ g | g \text{ es función} \}$

Hacer una selección en L resulta sumamente difícil debido a la diversidad de funciones que existen. Sería conveniente delimitar una familia de funciones como por ejemplo :

$C = \{ g | g \text{ es función continua en } [a,b] \}$

Se sabe que toda función continua en un intervalo cerrado y acotado puede ser aproximada uniformemente por un polinomio (teorema de aproximación de Stone-Weierstrass).

Tomando en cuenta este resultado, se puede restringir aún más el espacio de funciones a elegir como :

$F = \{ f : [a,b] \rightarrow [c,d] | f \text{ es polinomio} \}$

Aún restringiéndose a F ¿Cómo elegir $F \in F$?

Sea $h \in F$, lo cual quiere decir que h será una aproximación polinomial de f , y el error de aproximación se sumará a ϵ_1^* de tal forma que

$$y = h(x) + \epsilon \quad \text{con} \quad \epsilon = \epsilon_1^* + (f(x) - h(x))$$

De donde $\epsilon = y - h(x)$

El problema es encontrar una función $h(x)$ tal que minimice a ϵ .

Sean $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ n parejas de datos y

$\epsilon_1, \epsilon_2, \dots, \epsilon_n$ los errores correspondientes.

Se debe determinar un criterio Δ para elegir h tal que minimice los errores.

Sea $\Delta : F \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}$ donde $h(x) \in F$ y $\underline{e} = (e_1, \dots, e_n) \in \mathbb{R}^n$

Una medida de la magnitud de \underline{e} podría ser su norma definida en general como

$$N = \left[\sum_{i=1}^n |e_i| \right]^{1/n}$$

El criterio Δ consistiría en minimizar la norma de \underline{e} .

Restringiéndose a aquellos polinomios de grado k , se tiene que:

$$h(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots + \alpha_k x^k$$

entonces

$$y_i = \alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 + \dots + \alpha_k x_i^k + e_i \quad i=1, \dots, n$$

entonces

$$e_i = y_i - (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 + \dots + \alpha_k x_i^k) \quad i=1, \dots, n$$

Elegir $h(x)$ tal que minimice al error, es equivalente a elegir los coeficientes $(\alpha_0, \alpha_1, \dots, \alpha_k)$ que minimicen a e_i (ya que el valor de x está fijo).

Sea $z_{ji} = x_i^j$ entonces

$$y_i = \alpha_0 + z_{1i} \alpha_1 + z_{2i} \alpha_2 + \dots + z_{ki} \alpha_k + e_i \quad i=1, \dots, n$$

Que es una combinación lineal en los parámetros (desconocidos) y en las z 's, que son funciones conocidas y totalmente especificadas de la variable explicativa.

Hay que notar que esta modalidad permite escribir en forma de combinación lineal a las potencias sucesivas de x , pero aún más, permite que z_{ji} no sólo sean potencias sino funciones de la x en cuestión e incluso diferentes funciones de varias variables explicativas, siempre y cuando estas estén completamente especificadas.

Por ejemplo :

$$w = \alpha_0 + \alpha_1 f_1(x_1) + \alpha_2 f_2(x_2) + \alpha_3 f_3(x_3) + e$$

Podría reescribirse como :

$$y = \alpha_0 + \alpha_1 z_{11} + \alpha_2 z_{22} + \alpha_3 z_{33} + e_i$$

Bajo esta modalidad incluso funciones no lineales pero susceptibles de linealizarse se tomarán en cuenta.

Supóngase entonces que :

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{p-1i} + \epsilon_i \quad i=1, \dots, n$$

en donde

$\{x_1, x_2, \dots, x_p\}$ es el conjunto de variables explicativas o bien funciones de ellas completamente especificadas.

y es la variable de respuesta y

$\beta_0, \beta_1, \dots, \beta_{p-1}$ son coeficientes desconocidos

entonces

$$\epsilon_i = y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{p-1i})$$

Para minimizar $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)$ es razonable pensar en minimizar la norma de ϵ : La norma p se define como

$$N = \left[\sum_{i=1}^n |\epsilon_i|^p \right]^{1/p}$$

y el criterio Δ consistiría en minimizar N_p , es decir

$$\min_{\{\alpha_0, \dots, \alpha_{p-1}\}} \left[\sum_{i=1}^n |y_i - \alpha_0 - \alpha_1 x_{1i} - \dots - \alpha_{p-1} x_{p-1i}|^p \right]^{1/p}$$

Por facilidad sea $p = 2m$ con lo cual $|\epsilon_i|^p = \epsilon_i^{2m}$

La función $(\sum_{i=1}^n \epsilon_i)^{1/2m}$ es monótona no decreciente por lo que basta minimizar

$$\Delta(\beta_0, \beta_1, \dots, \beta_{p-1}) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{p-1i})^{2m}$$

Derivando

$$\frac{\partial \Delta'}{\partial \beta_0} = \sum_{i=1}^n z^m (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{p-1i})^{2m-1} (-1)$$

$\partial \beta_0$

$$\frac{\partial \Delta'}{\partial \beta_{p-1}} = \sum_{i=1}^n z^m (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{p-1i})^{2m-1} (-x_{p-1i})$$

$\partial \beta_{p-1}$

Igualando a cero :

$$\frac{\partial \Delta'}{\partial \beta_0} = 0$$

$\partial \beta_1$

$$\frac{\partial \Delta'}{\partial \beta_{p-1}} = 0$$

$\partial \beta_{p-1}$

El sistema de p-ecuaciones no es sencillo de resolver ya que queda un sistema de ecuaciones en potencias de (z^{2m-1}) . Sin embargo si $m=1$, se obtiene un sistema de ecuaciones lineales que preserva las condiciones originales :

$$-2 \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_{p-1} x_{p-1i}) = 0$$

$$-2 \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_{p-1} x_{p-1i}) = 0$$

Por tanto basta minimizar

$$\begin{aligned} \Delta'(\beta_0, \beta_1, \dots, \beta_{p-1}) &= \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_{p-1} x_{p-1i})^2 \\ &= \sum_{i=1}^n e_i^2 \end{aligned}$$

Al criterio anterior se le conoce como *Mínimos Cuadrados* .

En el caso de una variable se tendrá :

$$\sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

Derivando con respecto a β_0 y β_1 e igualando a cero se obtienen las ecuaciones :

$$\begin{aligned} \sum_{i=1}^n Y_i &= n \hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n X_i \\ \sum_{i=1}^n Y_i X_i &= \hat{\beta}_0 \sum_{i=1}^n X_i + \hat{\beta}_1 \sum_{i=1}^n X_i^2 \end{aligned}$$

De las cuales se obtiene

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n Y_i X_i}{\sum_{i=1}^n X_i^2} \quad \text{y} \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Como se vió anteriormente el error ϵ es desconocido pues involucra a los factores no controlados. En este sentido existe incertidumbre en torno a él. Es razonable suponer que el error es una variable aleatoria.

Como tal la manera natural de describir su comportamiento será a través de una función de distribución.

Como es natural suponer que $\epsilon \in \mathbb{R}$ entonces la distribución que se considere debe ser continua.

Dada la construcción del modelo, es razonable suponer que :

$$E(\epsilon) = 0 \quad \text{y además que} \quad P(\epsilon > 0) = P(\epsilon < 0)$$

Lo cual lleva a que la distribución será simétrica centrada en cero.

Dado que las fuentes de variación son en principio, las mismas para cada observación (ya que los factores no controlados son los mismos de observación a observación), se puede suponer que

$$V(\epsilon_i) = \sigma^2 \quad i = 1, \dots, n$$

Además los factores controlados junto con las condiciones experimentales generan los datos por lo que adicionalmente se hace el siguiente supuesto :

$$\text{Cov} (\epsilon_i , \epsilon_j) = 0 \quad (i \neq j)$$

Con los supuestos anteriores sobre la esperanza, varianza y covarianza de los errores, se obtienen resultados importantes de los estimadores por minimos cuadrados. Estos resultados se resumen en el llamado *Teorema de Gauss-Markov* el cual garantiza que β_0 y β_1 son los mejores estimadores lineales insesgados en el sentido de varianza minima.

Aún cuando este resultado es poderoso, se requiere de un supuesto adicional para poder explicar tecnicas de inferencia, es decir, se requiere suponer una distribución para los errores.

Bajo las propiedades antes descritas es razonable pensar en la dn Normal para explicar el comportamiento de los ϵ_i , pues esta dn posee características que permiten hacer uso de elementos estadísticos con mayor facilidad que con otra dn .

Ahora bien, tomando el vector $\epsilon = (\epsilon_1, \dots, \epsilon_n)$ se tiene entonces un vector aleatorio tal que

$$E(\epsilon) = \underline{0}$$

Y la matriz de varianzas y covarianzas Σ es de la forma $\sigma^2 I$. Entonces

$$\epsilon \sim N_n (\underline{0} , \sigma^2 I)$$

(En el apéndice 1 se presenta una sección sobre esta distribución).

$$\text{Como } \underline{Y} = X\underline{\beta} + \underline{\epsilon}$$

$$\text{Entonces } \underline{Y} \sim N_n (X\underline{\beta} , \sigma^2 I)$$

CAPITULO 2 : MODELO DE REGRESION LINEAL MULTIPLE.

2.1. El modelo.

Cuando se tiene más de una variable explicativa el modelo de regresión lineal múltiple (RLM) queda representado de la siguiente manera :

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{i,p-1} + \epsilon_i \quad 2.1.1$$

En RLM resulta natural utilizar notación matricial

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1,p-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2,p-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{n,p-1} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

es decir

$$Y_{n \times 1} = X_{n \times p-1} \beta_{p-1} + \epsilon_{n \times 1}$$

en donde

$Y_{n \times 1}$ Es el vector de $n \times 1$ observaciones (variables de respuesta).

$X_{n \times p-1}$: Es una matriz de valores fijos a la cual se le denomina matriz de datos o de diseño, de rango p .

β_{p-1} : Vector de parámetros desconocidos.

$\epsilon_{n \times 1}$ Vector aleatorio que agrupa a los factores no controlados.

De 2.1.1 se tiene que

$$\epsilon = Y - X \beta$$

Generalizando los resultados del capítulo 1, se tiene que

$$i) E(\epsilon) = 0$$

$$ii) V(\epsilon) = \sigma^2 I$$

De acuerdo a lo visto en el capítulo 1, para encontrar β tal que minimice el error es necesario minimizar $\sum e_i^2$, es decir

$$\text{encontrar } \min(e^t e) = \sum_{i=1}^n e_i^2 = \min \{ (Y - X\beta)^t (Y - X\beta) \}$$

Al diferenciar con respecto a β , las ecuaciones normales resultantes son :

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

y se obtiene que :

$$\hat{\beta} = (X^t X)^{-1} X^t Y \quad \text{y } (X^t X)^{-1} \text{ existe ya que el rango de } X = p.$$

El estimador $\hat{\beta}$ por mínimos cuadrados es insesgado y su varianza es $\sigma^2 (X^t X)^{-1}$. Además el teorema de Gauss Markov

asegura que el mejor estimador lineal insesgado de $1^t \beta$ en el sentido de varianza mínima es $1^t \hat{\beta}$ con $\hat{\beta}$ el estimador por mínimos cuadrados y $1 \in \mathbb{R}^p$ vector de constantes conocidas.

Aunque el Teorema de Gauss Markov proporciona una propiedad importante sobre $\hat{\beta}$, el método de Mínimos Cuadrados no proporciona directamente un estimador para σ^2 . Además sería deseable poder hacer uso de técnicas de inferencia. Incorporando el supuesto de Normalidad,

$Y \sim N_n(X\beta, \sigma^2 I)$, los estimadores máximo verosímiles de β y σ^2 se obtienen maximizando la función de verosimilitud:

$$L(\sigma^2, \beta | Y) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} (Y - X\hat{\beta})^t (Y - X\hat{\beta}) \right\}$$

Derivando e igualando a 0 se tiene $\hat{\beta} = (X^t X)^{-1} X^t Y$ que coincide con el obtenido por mínimos cuadrados.

El estimador máximo verosímil de σ^2 es :

$$\hat{\sigma}^2 = \frac{1}{n} (Y - X\hat{\beta})^t (Y - X\hat{\beta}).$$

El vector de observaciones ajustadas es : $\hat{Y} = X\hat{\beta}$.

Si se desea predecir el valor de Y para un vector X_* (con $X_{*p \times 1}$ de valores conocidos) se tiene que:

$$\hat{Y}_* = \hat{\beta}_0 + \hat{\beta}_1 X_{*1} + \hat{\beta}_2 X_{*2} + \dots + \hat{\beta}_{p-1} X_{*p-1} = X_*' \hat{\beta}$$

Por otra parte, el valor esperado de Y cuando $X = X_*$ es :

$$\mu_{X_*} = E (Y_* | X_* = X_*) = \beta_0 + \beta_1 X_{*1} + \dots + \beta_{p-1} X_{*p-1}$$

El estimador máximo verosímil de μ_{X_*} es :

$$\hat{\mu}_{X_*} = \hat{\beta}_0 + \hat{\beta}_1 X_{*1} + \dots + \hat{\beta}_{p-1} X_{*p-1}$$

2.2. Intervalos de confianza y pruebas de hipótesis.

Ahora bien en base a los supuestos distribucionales se tiene que:

$$\hat{\beta} \sim N \left[\beta, \sigma^2 (X'X)^{-1} \right]$$

$$\frac{Y}{\mathbf{1}'} \hat{\beta} \sim N_p \left[\frac{\mathbf{1}'}{\mathbf{1}'} \beta, \frac{\mathbf{1}'}{\mathbf{1}'} \sigma^2 (X'X)^{-1} \mathbf{1} \right]$$

Resulta interesante construir intervalos de confianza para $\frac{\mathbf{1}'}{\mathbf{1}'} \beta$ con $\mathbf{1} \in \mathbb{R}^p$ vector de constantes conocidas.

En particular si $\mathbf{1}_{-p \times 1}$ es el i-esimo vector canónico se estará construyendo un intervalo de confianza para la componente $\beta_{(i)}$ del vector β .

El intervalo para $\mathbf{1}'\beta$ al $(1-\alpha) \times 100\%$ de confianza es :

$$I = \left[\frac{\mathbf{1}'\hat{\beta}}{\mathbf{1}'\hat{\beta}} \pm t_{n-p}^{1-\alpha/2} \sqrt{\frac{\mathbf{1}'(X'X)^{-1}\mathbf{1}}{\mathbf{1}'\hat{\beta}}} \right]$$

Con

$$\hat{\sigma}^2 = \frac{(\underline{Y} - X \hat{\beta})^t (\underline{Y} - X \hat{\beta})}{n - p}$$

el estimador insesgado de σ^2 .

También se construye el intervalo para σ^2 de la forma

$$\sigma^2 \in \left[0, \frac{(n-p) \hat{\sigma}^2}{\chi^2_{n-p, \alpha/2}} \right]$$

Este intervalo proporciona información de que tan grande puede llegar a ser σ^2 .

Si se quiere un intervalo para μ_{x_0} , este es:

$$\left[\hat{\mu}_{x_0} - t_{n-2}^{\alpha/2} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right]$$

El intervalo de predicción para Y_0 es:

$$\hat{\mu}_{x_0} - t_{n-2}^{\alpha/2} \sqrt{k + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Prueba de Hipótesis Lineal General.

Dado que los valores del vector β son desconocidos resulta natural hacer conjeturas sobre sus posibles valores.

En particular, una prueba de hipótesis que proporciona mucha información es la llamada hipótesis de significancia global que propone probar si los elementos $\beta_1, \beta_2, \dots, \beta_{p-1}$ son todos ceros.

Es decir

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0 \quad \text{vs} \quad H_a : \text{alguna } \beta_i = 0 \quad \forall i=1, p-1$$

En caso de rechazar H_0 , se tendría que $Y_0 = \beta + \epsilon$ lo cual

implica que ninguna de las variables explicativas se relacionan con Y , al menos en forma lineal, y permitiría deducir que el modelo no es apropiado. Por el contrario, si se rechaza H_0 se tiene que al menos una variable explicativa guarda relación con Y .

Es posible también estar interesado en probar si una parte del vector β toma ciertos valores ó incluso si uno solo de los componentes, digamos β_j , es igual a un valor predeterminado.

Todas las hipótesis antes mencionadas se engloban en lo que se conoce como *Hipótesis Lineal General* que toma la forma:

$$H_0 : C\beta = \gamma \quad \text{vs} \quad H_a : C\beta \neq \gamma$$

Donde $C_{r \times p}$ de rango r , $r \leq p$, $\gamma_{r \times 1}$ completamente especificadas.

Se tiene $Y \sim N(X\beta, \sigma^2 I)$

La región crítica utilizando Cociente de Verosimilitudes es :

$$C = \left\{ \frac{\sup_{H_0} L(\beta, \sigma^2 | Y)}{\sup_{H_0 \cup H_a} L(\beta, \sigma^2 | Y)} \leq K_1 \right\}$$

$$\sup_{H_0 \cup H_a} L(\beta, \sigma^2 | Y) = L(\hat{\beta}, \hat{\sigma}^2 | Y)$$

Con $\hat{\beta}$ y $\hat{\sigma}^2$ los estimadores máximos verosímiles de β y σ^2 .

Ahora

$$\sup_{H_0} L(\beta, \sigma^2 | Y) = L(\tilde{\beta}, \tilde{\sigma}_0^2 | Y)$$

En donde $\tilde{\beta}$, $\tilde{\sigma}_0^2$ son los estimadores máximo verosímiles bajo la restricción $C\beta = \gamma$ con

$$\tilde{\beta} = \hat{\beta} - (X'X)^{-1}C'(C(X'X)^{-1}C')^{-1}(C\hat{\beta} - \gamma)$$

$$\tilde{\sigma}_0^2 = \frac{1}{n} (Y - X\tilde{\beta})'(Y - X\tilde{\beta})$$

El cociente de verosimilitudes resulta ser :

$$\Delta = \frac{2 \pi \hat{\sigma}_0^2 \exp \{ -n/2 \}}{2 \pi \hat{\sigma}^2 \exp \{ -n/2 \}}$$

$$\Delta = \frac{\hat{\sigma}_0^2}{\hat{\sigma}^2}$$

Se rechaza H_0 si $\Delta \leq K$ es decir

$$\Leftrightarrow \frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} \leq K$$

$$\Leftrightarrow \frac{(\underline{Y} - X\underline{\beta})^t (\underline{Y} - X\underline{\beta})}{(\underline{Y} - X\underline{\hat{\beta}})^t (\underline{Y} - X\underline{\hat{\beta}})} \leq K$$

Sea $SCE_{MC} = (\underline{Y} - X\underline{\beta})^t (\underline{Y} - X\underline{\beta})$

SCE_{MC} representa la variabilidad de Y con respecto a Y cuando se ajuste el modelo completo $Y = X\underline{\beta} + \epsilon$

Puede interpretarse como el error que se comete al ajustar el MC.

Sea $SCE_{MR} = (Y - X\underline{\hat{\beta}})^t (Y - X\underline{\hat{\beta}})$

SCE_{MR} representa la variabilidad de Y con respecto a \hat{Y} cuando el modelo es el modelo reducido, el modelo se reduce dependiendo de la forma de C y γ .

Entonces

$$\Delta = \frac{SCE_{MC}}{SCE_{MR}}$$

Este cociente, por construcción es menor o igual que uno lo que implica que $SCE_{MR} \geq SCE_{MC}$

De hecho puede demostrarse que $SCE_{MR} = SCE_{MC} + SCE_{H_0}$ es decir

$$(\underline{Y} - X\hat{\beta})'(\underline{Y} - X\hat{\beta}) = (\underline{Y} - X\hat{\beta})'(\underline{Y} - X\hat{\beta}) + (C\hat{\beta} - \underline{\gamma})'(C(X'X)C')^{-1}(C\hat{\beta} - \underline{\gamma})$$

La cual se conoce como La Partición Fundamental en Sumas De cuadrados.

Sea $SCE_{H_0} = (C\hat{\beta} - \underline{\gamma})'(C(X'X)C')^{-1}(C\hat{\beta} - \underline{\gamma})$

La SCE_{H_0} es la magnitud en variabilidad, que hace falta para que el modelo completo tenga la misma variabilidad que el reducido.

$$SCE_{H_0} = SCE_{MR} - SCE_{MC}$$

Se rechaza H_0 si $\frac{SCE_{MC}}{SCE_{MR}} \leq K$

$$\Leftrightarrow \frac{SCE_{MR} \geq K_2}{SCE_{MC}}$$

$$\Leftrightarrow \frac{SCE_{H_0} + SCE_{MC} \geq K_2}{SCE_{MC}}$$

$$\Leftrightarrow \frac{SCE_{H_0} + 1 \geq K}{SCE_{MC}}$$

$$\frac{SCE_{H_0}}{SCE_{MC}} \geq K'$$

Entonces $C = \{ \underline{Y} \mid \frac{SCE_{H_0} \geq K}{SCE_{MC}} \}$

$$Pr \left(\frac{SCE_{H_0} \geq K}{SCE_{MC}} \mid H_0 \text{ es cierta} \right) = \alpha$$

Para determinar K es entonces necesario determinar la D_{ig} de $\frac{SCE_{Ho}}{SCE_{Mc}}$

Ahora

$$\frac{SCE_{Mc}}{\sigma^2} \sim X^2_{(n-p)} \quad (\text{ver apendice 2})$$

$$y \quad \frac{SCE_{Ho}}{\sigma^2} \sim X^2_r$$

Bajo Ho.

Las dos variables en cuestion son independientes (ver apéndice 2), y su cociente, afectado por ciertas constantes, resulta ser una variable aleatoria con distribución " F " :

$$\frac{r \text{ SCE}_{MC} / \sigma^2}{n-p \text{ SCE}_{HO} / \sigma^2} \rightarrow F (r, n-p) \text{ bajo } H_0$$

Entonces :

$$\frac{\text{SCE}_{HO} / r}{\text{SCE}_{MC} / n-p} \rightarrow F (r , n-p) \text{ bajo } H_0$$

Se denotan los cuadrados medios del error de la siguiente manera

$$\text{SCE}_{HO} = \text{CME}_{HO} \quad ; \quad \text{SCE}_{MC} = \text{CME}_{MC}$$

entonces la región crítica está dada por

$$C = \left\{ \frac{\text{CME}_{HO}}{\text{CME}_{MC}} \geq F_{r, n-p}^{1-\alpha} \right\}$$

La tabla de análisis de varianza (ANOVA) para probar las hipótesis :

$H_0 : C\beta = \gamma$ VS. $H_a : C\beta \neq \gamma$ en el modelo $\underline{Y} = X\beta + \epsilon$

con $\epsilon \rightarrow N_n (0, \sigma^2 I)$ es

Fuente de variación	g. l.	Suma de Cuad.	Cuad. Medios	F_c F calcula
H_0	r	$(C\hat{\beta} - \gamma)' [C(X'X)^{-1}C'] (C\hat{\beta} - \gamma)$	$\frac{\text{SCE}_{HO}}{\text{CME}_{HO} = r}$	$\frac{\text{CME}_{HO}}{\text{CME}_{MC}}$
Error MC	n - p	$(Y - X\hat{\beta})'(Y - X\hat{\beta})$	$\text{CME}_{MC} = \frac{\text{SCE}_{MC}}{n - p}$	
Total MR	n-p-r	$(Y - X\bar{\beta})'(Y - X\bar{\beta})$		

Como se mencionó anteriormente existen casos particulares para esta prueba entre los que cabe resaltar :

$$1) \quad H_0 : \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0 \quad \text{VS.} \quad H_a : \beta_i \neq 0 \quad i=1, \dots, p-1$$

Es decir

$$H_0 : C\beta = \underline{0} \quad \text{VS.} \quad H_a : C\beta \neq \underline{0}$$

$$C_{(p-1) \times p} = \begin{bmatrix} 0 & 1 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & 1 \end{bmatrix} \quad \beta_{p \times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} \quad \chi_{(p-1) \times 1} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

Esta prueba se le conoce como *Prueba de Significancia Global* lo que indica que existe al menos una $\beta_i \neq 0$ ($i=1, \dots, p-1$) entonces el modelo tiene sentido.

Por lo anterior bajo H_0 cierta el modelo reducido es

$$MR : Y_i = \beta_0 + \epsilon_i \quad (i=1, \dots, n) \quad \text{con} \quad Y_i \sim N(\beta_0, \sigma^2)$$

y la suma de cuadrados del error para el MR queda como

$$SCE_{MR} = \sum_{i=1}^n (y_i - \bar{y})^2 = \underline{Y}'\underline{Y} - n\bar{y}^2 \text{ que es la variabilidad total de las } y\text{'s. El modelo completo es}$$

$$MC : \underline{Y} = X\underline{\beta} + \underline{\epsilon} \text{ y la suma de cuadrados del error para el}$$

MC es

$$SCE_{MC} = \underline{Y}'\underline{Y} - \underline{\beta}'X'\underline{Y}$$

Por la partición fundamental en sumas de cuadrados se tiene

$$SCE_{MR} = SCE_{MC} + SCE_{H_0}$$

De donde

$$SCE_{H_0} = SCE_{MR} - SCE_{MC}$$

$$SCE_{H_0} = \underline{\beta}'X'\underline{Y} - n\bar{y}^2$$

Por tanto la tabla Anova es :

Fuente de variación	g.l.	Suma de Cuad.	Cuad. Medios	F_c F calculada
H_0	$p - 1$	$SCE_{H_0} = \underline{\beta}'X'\underline{Y} - n\bar{y}^2$	$CME_{H_0} = \frac{SCE_{H_0}}{p - 1}$	$\frac{CME_{H_0}}{CME_{MC}}$
Error MC	$n - p$	$SCE_{MC} = \underline{Y}'\underline{Y} - \underline{\beta}'X'\underline{Y}$	$CME_{MC} = \frac{SCE_{MC}}{n - p}$	
Total MR	$n - 1$	$SCE_{MR} = \underline{Y}'\underline{Y} - n\bar{y}^2$		

donde la región crítica es :

$$C = \{ \underline{Y} \mid F_c \geq F_{r, n-p}^{1-\alpha} \}$$

2) $H_0: \beta_i = m$ con $m \in \mathbb{R}$ conocida donde 1 vector de constantes $p \times 1$ conocidas.

Un caso particular de esta situación es cuando 1 es canónico y $m = 0$, es decir 1 es el i -ésimo vector canónico, entonces

$$\beta_i = \beta_{i-1} \quad \text{y la prueba será} \quad H_0: \beta_{i-1} = 0$$

Esta prueba se conoce como prueba *F-parcial*, sirve para probar si una β_i toma el valor cero, lo cual equivale a decir que la x asociada a ese coeficiente no entra al modelo.

A continuación se presentan algunos conceptos importantes

Una herramienta que sirve para medir el grado de asociación lineal entre dos variables x y Y es el coeficiente de correlación lineal que está dado por

$$r_{xy} = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (Y_i - \bar{Y})^2}} \quad \text{con} \quad -1 < r_{xy} < 1$$

si $r_{xy} = 1$ se tiene una asociación lineal directa.

si $r_{xy} = -1$ se tiene una asociación lineal indirecta.

si $r_{xy} = 0$ no se tiene una asociación lineal, lo cual no significa que x y Y no guarden una relación.

Una medida de la bondad del ajuste de un modelo de regresión es la estadística R^2 . Se define como

$$R^2 = \frac{SCE_{H_0}}{SCE_{MR}} \quad 0 < R^2 < 1$$

Con SCE_{MR} de significancia global y se interpreta como la proporción de la variabilidad de la Y que quedó explicada por el modelo.

Residuales. Los residuales se definen como $e_i = Y_i - \hat{Y}_i$. Sirven para hacer un análisis de los supuestos distribucionales.

$$\text{Como vector } \underline{e} = \underline{Y} - \hat{\underline{Y}} = \underline{Y} - X\hat{\beta}$$

$$E(\underline{e}) = \underline{0}$$

$$V(\underline{e}) = \sigma^2 (I - X(X^T X)^{-1} X^T)$$

Residuales Studentizados. Se definen como $r_i = e_i / \sigma \sqrt{1 - h_{ii}}$

donde h_{ii} es el componente de la matriz $(X(X^T X)^{-1} X^T)$

Estos residuales son más utilizados que los e_i estandarizados, pues las pruebas en las que se involucran, como Papel Normal, entre otras, son más precisas.

Para probar el supuesto de Normalidad es usual utilizar estos residuales o los e_i , lo cual consiste en tomar los residuales ordenados de menor a mayor y graficar contra $K_i = (i - 1/2)/n$, en papel Normal, si la gráfica presenta un comportamiento aproximadamente lineal, esto querra decir que el supuesto de Normalidad es correcto.

Otras gráficas de gran utilidad es tomar los residuales contra las y_i estimadas y contra las x_i , las cuales sirven para analizar entre otras cosas, la variabilidad de los e_i .

Una estadística que sirve para analizar la correlación entre los residuales es la propuesta por Durbin y Watson. Se define como:

$$d_i = \frac{\sum_{i=1}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} \quad i=1, \dots, n.$$

Capítulo 3 : Selección de Variables.

En el análisis de regresión surge un problema que consiste en incluir un conjunto de variables de tal manera que para este conjunto el modelo de regresión sea óptimo, esto es, que el modelo refleje al fenómeno bajo estudio. Como se vió anteriormente el modelo es de la forma :

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$$

Donde las x todas las variables independientes que se piensa ayudarán a explicar el fenómeno lo mejor posible. Cuando se plantea el modelo, se trata de involucrar el mayor número de variables explicativas de tal manera que existan menos factores no controlados y así tener un mejor control, sin embargo, es importante notar que el hecho de incluir en el modelo todas las variables explicativas que se tengan, no implica que el modelo sea el mejor pues puede ocurrir que existan variables que no aporten al modelo suficiente información, o que unas estén explicadas a través de otras, además de que el modelo sería más complicado y su interpretación sería menos clara. Por lo anterior surge el problema de seleccionar un subconjunto de q variables ($q < p-1$) a través de las cuales sea posible explicar el fenómeno, donde dicha selección se llevará a cabo mediante técnicas estadísticas que serán discutidas más adelante, con lo que el nuevo modelo sería de la forma :

$$y_i = \beta_0 + \sum_{j=1}^q \beta_j x_{ij} + \epsilon_i$$

Problema de multicolinealidad.

Cuando el modelo presenta variables tales que unas están explicadas a través de otras, lo cual dificultará en un momento determinado precisar la relación que existe entre la manifestación del fenómeno con respecto a alguna variable explicativa, este último hecho es lo que se conoce como multicolinealidad (correlación entre variables explicativas).

El tener efectos de multicolinealidad en un modelo origina algunos problemas. Por ejemplo, si se desea realizar una interpretación de los coeficientes parciales de regresión, el efecto de multicolinealidad puede conducir a contradicciones. Estos problemas se pueden evitar mediante la eliminación de alguna de las variables que presentan tal efecto, proceso que puede llevarse a cabo mediante técnicas de selección de variables.

Selección de Variables.

- Todas las regresiones posibles.
- Backward.
- Forward.
- Stepwise.

Cabe aclarar que ninguno de estos métodos garantiza que el subconjunto de variables explicativas que intervienen en el nuevo modelo sea el mejor.

Todas las regresiones posibles.

Bajo este criterio se requiere ajustar todas las posibles combinaciones de las x_i . Para cada x_i existen dos posibilidades de estar o no estar en el modelo de regresión, por lo tanto existen 2^p modelos posibles a examinar.

Cabe mencionar que bajo este criterio es claro que si se tiene un conjunto de variables explicativas grande, se obtendrán un gran número de posibles modelos, lo cual en términos de cálculo y costo tiene grandes limitaciones.

Suponiendo que no existen limitaciones, el análisis para seleccionar el mejor subconjunto de x_i puede hacerse a través de diferentes criterios entre los cuales se puede mencionar a los siguientes :

1. R^2 .
2. R^2 ajustada.
3. Cp-Mallows.

Criterio de la R^2 . R^2 mide la proporción de la variabilidad de la Y que queda explicada por el modelo. R^2 es usada como una medida de a través de la cual se compara la validez de los resultados que se obtienen bajo especificaciones alternativas de las variables explicativas en el modelo.

El uso de la R^2 acarrea varios problemas, en primer lugar es sensible al número de variables explicativas incluidas en el modelo. La inclusión de más variables explicativas nunca hará que R^2 decrezca, por lo que es posible añadir más y más variables, lo que no necesariamente significa que todas esas variables tengan que estar en el modelo.

Bajo este criterio se calculan las R^2 para cada modelo obtenido dentro de todas las posibles combinaciones, así para todos los modelos con q variables se elige a la R^2 mayor y se toma ese modelo como el más significativo. En el caso de existan empates entre variables, se puede calcular su correspondiente coeficiente de correlación parcial, mismos que se definen más adelante, pues puede suceder que una de estas variables quede reflejada por otra, en ese caso se desecha la variable que está siendo redundante para el ajuste del modelo.

R^2 ajustada. La dificultad con R^2 radica en que contiene la variación de Y explicada por el modelo y la no explicada pero no toma en cuenta los grados de libertad del problema. Una solución natural es pensar en las varianzas (muestrales) o cuadrados medios y no en las variaciones (o sumas de cuadrados). Se define como :

$$\bar{R}^2 = 1 - \frac{\text{Var}(e_i)}{\text{Var}(Y)}$$

$$\text{donde } \text{Var}(e_i) = \sigma^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-p}$$

$$\text{Var}(y) = \sigma_0^2 = \frac{\sum_{i=1}^n (y_i - \bar{y}_i)^2}{n-1}$$

con p = número de parámetros.

Aún cuando la SCE_{MC} decrezca, a medida que se añaden más variables, la varianza de los e_i no necesariamente decrece. Hay que notar que tanto el numerador como el denominador de σ^2 cambian al incluir una nueva variable en el modelo. Es fácil checar que

$$\bar{R}^2 = 1 - \frac{(1 - R^2) n-1}{n-p}$$

y resulta inmediato que

- 1) si $p = 1$ entonces $R^2 = \bar{R}^2$
- 2) si $p > 1$ entonces $R^2 \geq \bar{R}^2$
- 3) \bar{R}^2 puede ser negativa

Estas propiedades hacen que \bar{R}^2 resulte una medida más adecuada que R^2 .

Cp-Mallows. Mallows sugiere una estadística para elegir al mejor subconjunto de variables, la cual es comparada con el número de parámetros del modelo y se recurre a una gráfica en la que intervienen la Cp y p, número de parámetros.

$$Cp = \frac{SCE_{MC_p}}{CMEMC(n - 2p)}$$

donde p se refiere a la SCE_{MC} que resulta al tomar todas las variables explicativas.

El valor de cada Cp es muy importante porque es un estimador de la suma de cuadrados global de las discrepancias entre el modelo ajustado y el modelo verdadero (desconocido).

Para un modelo adecuado se tiene que $E(C_p)=p$. Se sigue que una gráfica de C_p vs. p mostrará modelos adecuados en la medida en que los puntos se acerquen a la recta $C_p=p$. Ajustes que adolezcan de carencia de ajuste, esto es, ecuaciones sesgadas, producirán puntos muy por encima de la recta $C_p=p$. Debido a variaciones aleatorias, puntos que representen modelos adecuados pueden estar por debajo de la recta $C_p=p$.

En la gráfica se toma el siguiente criterio, si C_p se aproxima a p o errará decir que el modelo que corresponda a esa C_p será el más significativo.

Backward. (Eliminación hacia atrás). El método consiste en los siguientes pasos :

1. Se ajusta el modelo de regresión con todas las variables.
2. Para cada variable se calcula la F-parcial correspondiente.
3. De todas las F-parciales calculadas se toma la de menor valor, la cual se denotará como F^* , y se compara con una F-tablas seleccionada de antemano.

3.1. Si $F^* < F$ -tablas, entonces la variable que corresponda a esa F-parcial se excluye del modelo de regresión, ajustando un nuevo modelo con las variables restantes.

Es importante aclarar que una vez que sale una variable no puede volver a entrar al modelo, después regresar al paso número 2.

3.2. Si $F^* > F$ -tablas, no sale ninguna variable y se toma el modelo que se ajusto para predecir a Y .

Este criterio de selección en términos de cálculo y de costo resulta mucho más adecuado que el anterior, pues sólo se ajustan los modelos que sugieren las F-parciales escogidas.

Forward. (Eliminación hacia adelante).

En este criterio usualmente se utilizan los coeficientes de correlación de las x_i y Y , siguiendo los siguientes pasos :

1. Se calculan los coeficientes de correlación de todas las x_i con Y . Se selecciona a la x_i más correlacionada con Y , de este modo la x_i seleccionada será la primer variable incluida en el modelo.
2. Se calculan los coeficientes de correlación de Y con las variables que aún no han entrado al modelo de regresión. La variable con coeficiente de correlación parcial más alto es la que entra al modelo.

El coeficiente de Correlación parcial mide la correlación parcial entre la variable y y x_i en presencia de un conjunto de variables explicativas y se define como :

$$r_{y_i \cdot 12 \dots p} = \frac{-R_{yi}}{\sqrt{R_{yy} R_{ii}}}$$

Donde cada R son los cofactores que se determinan mediante los menores que resultan de obtener una submatriz de correlaciones misma que se compone a partir de la eliminación del renglón donde se encuentra la y y columna i en la matriz de correlaciones.

$$R_{ij} = (-1)^{i+j} a_{ij}$$

3. Se calcula la F-parcial correspondiente y se compara con una F-tablas.

Si F-parcial > F-tablas la variable entra al modelo y se regresa a 2.

Si F-parcial < F-tablas la variable no entra al modelo termina el proceso.

Bajo este criterio el cálculo es mucho menor que en los anteriores pues sólo se consideran las variables necesarias.

Stepwise. (Eliminación paso a paso). Este método es una mezcla de los dos anteriores y consiste en lo siguiente :

1. Se seleccionan las variables siguiendo los mismos pasos de Forward. La diferencia es que una vez incluida una variable se calcula la F-parcial y se compara con una F-tablas.

Si F-parcial < F-tablas sale del modelo.

El proceso continúa hasta que ya no puedan entrar ni salir variables.

Es importante notar que en este criterio se cuestiona la permanencia de todas las variables incluidas en cada etapa.

Capítulo 4 : Paquete de cómputo.

4.1. Introducción.

En este capítulo se presenta el paquete de cómputo desarrollado para apoyar los cursos de Análisis de Regresión impartidos en la Facultad de Ciencias.

El paquete ha sido desarrollado de manera que el alumno pueda obtener en forma clara y sencilla la información necesaria para interpretarla bajo el contexto estadístico y no perder de vista los puntos importantes de la técnica que generalmente se pierde al hacer uso de paquetes muy conocidos y ya implantados.

Como todo sistema de cómputo, es necesario analizar la manera a través de la cual será implantado, tratando que esta sea la óptima.

La manera en la que fué desarrollado el paquete consistió básicamente en tres etapas :

- 1) Análisis de Datos.
- 2) Análisis de Manipulación de Datos.
- 3) Análisis de Técnicas de Cálculo Estadístico.

En la primer etapa, se analizó qué tipos de datos serían utilizados así como su forma de entrada, llegando a las siguientes conclusiones :

- a) Datos de entrada de tipo real.
- b) Entrada de Datos a través de la creación de un archivo.

En la segunda etapa, se analizó que una vez creado un archivo de datos, antes de aplicarles cualquier técnica estadística, es necesario verificar que estos sean correctos permitiendo visualizarlos y modificarlos, si es necesario.

En la etapa tres, se analizó la naturaleza de los cálculos estadísticos utilizados en las técnicas del Análisis de Regresión, con lo cual se llegó a la conclusión, junto con el análisis de las dos etapas anteriores, que es necesario emplear un lenguaje de programación lo más estructurado posible y apegado al desarrollo de cálculos matemáticos, por lo cual se decidió aplicar el lenguaje de programación PASCAL.

4.1.1. Estructura del paquete de Cómputo.

La estructura del paquete de cómputo consta de varios módulos :

1. Módulo de Captura de Datos.
2. Módulo de Matrices.
3. Módulo de Cálculos Estadísticos.
4. Módulo de graficación.

5. Módulo de Ayuda.
6. Módulo Principal.

Los módulos 1 a 5 se identificarán como módulos auxiliares al módulo principal.

1. Módulo de Captura de Datos. Se compone de varios submódulos :

- 1.1. Creación de archivo de datos.
- 1.2. Abertura de archivo de datos.
- 1.3. Modificación de archivo de datos.
- 1.4. Visualización de archivo de datos.
- 1.5. Selección del conjunto de columnas donde se encuentran las variables explicativas y variable de respuesta a partir del archivo de datos.
- 1.6. Ayuda.

A continuación se describen brevemente cada uno de los submódulos anteriores.

1.1. Se refiere a la creación del archivo de datos, el cual tendrá máximo 20 columnas y 50 renglones.

1.2. Una vez creado el archivo de datos, siempre será necesario abrirlo para después hacer cualquier manipulación con el mismo, es importante aclarar que si el archivo no es abierto no se podrá proceder a los siguientes submódulos.

1.3. Una vez creado y abierto el archivo de datos será posible hacerle las modificaciones que se deseen para evitar errores en los resultados.

1.4. La visualización del archivo de datos permitirá darse cuenta de posibles errores en el mismo, pudiéndose modificar en el momento que se desee.

1.5. Una vez que se tienen listos los datos, es decir sin errores, se procederá a seleccionar a través de columnas, la(s) variable(s) explicativa(s) y la variable dependiente que intervendrán en el modelo de regresión, este submódulo permitirá hacer tantas regresiones como se deseen, hecho que facilitará el empleo de métodos de selección de variables.

1.6. En este submódulo únicamente se presenta un texto en el cual se describen los pasos que deben seguirse en el módulo de Captura.

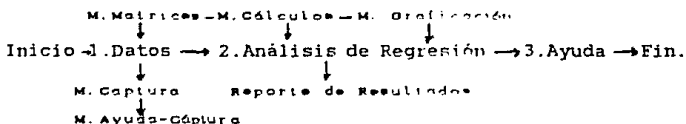
2. Módulo de Matrices. Este módulo consta únicamente de operaciones aplicables a matrices. Como se vió en capítulos anteriores, para la obtención de resultados en el Análisis de Regresión se involucran diversas operaciones con matrices tales como :

Cabe señalar que en este paquete de cómputo se utiliza una unidad creada fuera del mismo por otro programador cuya función es la de leer y verificar números reales, enteros y caracteres.

Para visualizar más claramente la estructura del paquete se presentan a continuación los siguientes diagramas :



El módulo principal está estructurado de la siguiente forma:



4.2. Instrucciones y Especificaciones.

4.2.1. Módulos Auxiliares. A continuación se describen los algoritmos de cada uno de los módulos auxiliares, así como los algoritmos de los submódulos que los componen.

Módulo Captura. Este módulo está formado por varios submódulos, estos a su vez están desarrollados en procedimientos que son instrucciones especiales del lenguaje PASCAL.

Descripción de Procedimientos.

Cada submódulo tiene asociado un procedimiento en PASCAL, a continuación se describen cada uno de ellos.

Submódulo : Creación de Archivo de Datos.

Asociación en PASCAL: *Procedure Crear_Arch.*

Algoritmo :

- ```

Inicio
1. Nombre del Archivo.
2. Lee el nombre del archivo y verifica que
 efectivamente se le da un nombre (cadena
 de caracteres).
3. Abre y prepara el archivo para ser creado.

```

4. Lee las dimensiones del archivo a crear, por renglón y columna.
  5. Captura los datos según las dimensiones leídas.
  6. Graba el archivo de datos en disco.
- Fin.

Submódulo : Abertura de Archivo de Datos.

Asociación en PASCAL : *Procedure Abrio\_Arch.*

Algoritmo :

- Inicio
1. Lee nombre del archivo que se desea abrir.
  2. Verifica que el archivo ya haya sido creado.
  3. Abre el archivo.
- Fin.

Submódulo : Modificación de Archivo de Datos.

Asociación en PASCAL : *Procedure Modifica\_Arch.*

Algoritmo :

- Inicio
1. Lee nombre del archivo a modificar.
  2. Verifica que el archivo exista y esté abierto.
  3. Lee del archivo.
  4. Pregunta tipo de modificación.
    - 4.1. En renglón.
    - 4.2. En Columna.
    - 4.3. En Celda.
  5. Lee tipo de modificación.
  6. Modifica el tipo seleccionado.
  7. Graba la modificación en archivo.
  8. Cierra el archivo de Datos.

Cabe señalar que si se entra en este módulo será necesario volver a abrir el archivo pues como se pudo observar se cerró al final del algoritmo.

Submódulo : Visualización de Archivo de Datos.

Asociación en PASCAL : *Procedure Listar\_Arch.*

Algoritmo :

- Inicio
1. Lee nombre del archivo a listar.
  2. Verifica que el archivo exista y este abierto.
  3. Despliega el archivo de datos.
  4. Cierra el archivo.
- Fin.

Al igual que en *Modifica\_Arch*, si se entra a este submódulo será necesario volver a abrir el archivo para después poder ejecutar cualquier otro submódulo.

Submódulo : Selección del conj. de variables explicativas y variable dependiente.

Asociación en PASCAL : *Procedure Seleccionar*.

Algoritmo :

Inicio

1. Verifica que el archivo este abierto.
2. Lee número de variables explicativas.
3. Asigna l's a la primer columna de la matriz de diseño.
4. Lee columna(s) de la(s) variable(s) explicativa(s).
5. Lee columna de la variable dependiente.

Fin.

Submódulo : Ayuda.

Asociación en PASCAL : *Procedure MenuDatos*(llamado de la unidad Ayuda)

Algoritmo :

Inicio

1. Despliega texto de ayuda para el módulo de Captura.

Fin.

El módulo Captura contiene otro submódulo auxiliar, cuya asociación en PASCAL es *Procedure Cap\_Dat*, cuya función es preguntar qué submódulo de los antes descritos se desea ejecutar, permitiendo así ejecutar el submódulo seleccionado.

Algoritmo :

Inicio

1. Pregunta Submódulo a ejecutar.
2. Despliega menú de opciones.
  - 2.1. Captura.
  - 2.2. Abrir.
  - 2.3. Modificar.
  - 2.4. Listar.
  - 2.5. Selección de variable(s) explicativa(s) y var. dependiente.
  - 2.6. Ayuda.
3. Lee submódulo seleccionado.
4. ejecuta submódulo seleccionado.

Fin.

Módulo Matrices. Al igual que Captura, este módulo está formado por varios submódulos, cada uno asociado a un procedimiento definido en PASCAL.

**Descripción de Procedimientos.**

**Submódulo : Inversa.**

**Asociación en PASCAL : *Procedure Inversa.***

**Algoritmo :**

- Inicio
- 1. Inicializa.
  - 1.1. Verifica que la matriz A sea invertible.
- 2. Invierte (Eliminación Gaussiana).
  - 2.1. Calcula pivote (renglón de referencia).
  - 2.2. Cambia renglones.
  - 2.3. Divide renglones en matriz A y  $A^{-1}$ .
  - 2.4. Multiplica y suma renglones en A y  $A^{-1}$ .
  - 2.5. Obtiene inversa.
- Fin.

este submódulo será utilizado para cualquier matriz que sea invertible, la cual será pasada como parámetro.

**Submódulo : Transpuesta.**

**Asociación en PASCAL : *Procedure Trans.***

**Algoritmo :**

- Inicio
- 1. Cambia renglones por columnas.
- 2. Obtiene la transpuesta de una matriz (pasada como parámetro).
- Fin.

**Submódulo : Solución de Sistemas  $Ax = b$ .**

**Asociación en PASCAL : *Bgorro.***

**Algoritmo :**

- Inicio
- 1. Calcula Bgorro utilizando los procedimientos *Inversa* y *Trans.*
- Fin.

Dentro del Módulo de Matrices se desarrolló otro submódulo cuya asociación en PASCAL es *Escribe*, cuyo objeto es mandar a escribir las matrices que se requieran.

**Módulo Cálculos.** Este módulo igual que los anteriores, está formado por varios submódulos asociados a procedimientos y funciones en PASCAL.

**Descripción de Procedimientos.**

**Submódulo : Media y Desviación.**

**Asociación en PASCAL : *Procedure Medio\_Desv.***

Algoritmo :

Inicio

- 1.Verifica Datos.
  - 2.Calcula suma de  $Y_i$ .
  - 3.Calcula suma de  $Y_i^2$ .
  - 4.Calcula suma de  $Y_i$  entre total de datos : Media de las  $Y_i$ .
  - 5.Calcula varianza de las  $Y_i$ .
  - 6.Calcula suma de  $x_i$  para cada  $i$ .
  - 7.Calcula suma de  $x_i^2$  para cada  $i$ .
  - 8.Calcula suma de  $x_i$  entre total de datos para cada  $i$  obteniendo la media de  $x_i$  para cada  $i$ .
  - 9.Calcula varianza de  $x_i$  para cada  $i$ .
- Fin.

Las medias y las varianzas se imprimen a través de los procedimientos *Esc\_medes* y *Esc\_var*, desarrollados exclusivamente para escritura.

Submódulo : Matriz de Correlaciones.

Asociación en PASCAL : *Procedure Matcorr*.

Algoritmo :

Inicio

- 1.Verifica Datos.
  - 2.Calcula coeficiente de Correlación de Y con cada  $x_i$ .
  - 3.Calcula coeficiente de correación de  $x_i$  con  $x_j$ .
- Fin.

La matriz de correlaciones es mandada a escribir en el módulo principal a través del submódulo *Escribe* desarrollado en Matrices.

Submódulo : Residuales, SCEMC, Coeficiente de Determinación.

Asociación en PASCAL : *Procedure Calresid\_est*.

Algoritmo :

Inicio

- 1.Calcula  $Y_i$  estimadas.
  - 2.Calcula residuales de la forma  $Y_i - \hat{Y}_i$ .
  - 3.Calcula SCEMC.
  - 4.Calcula residuales studentizados.
  - 5.Calcula coeficiente de Determinación.
- Fin.

Los residuales estimados y studentizados son escritos a través de un procedimiento de escritura llamado *resid\_est*.

La SCEMC es escrita a través del procedimiento *Esc\_SCEMC*.

Submódulo : Suma de las  $Y_i^2$ .

Asociación en PASCAL : *Function Sc\_de\_Y*.



Algoritmo :

Inicio

1. Calcula suma de cuadrados de las  $y_i$  al cuadrado.

Fin.

Submódulo : Suma de cuadrados de las  $Y_i$  corregidas por la media.

Asociación en PASCAL : *Funcion SC\_de\_Ycorr.*

Algoritmo :

Inicio

1. Calcula suma de las  $Y_i$  al cuadrado corr. por la media.

Fin.

Submódulo : Estadística de Durbin-Watson.

Asociación en PASCAL : *Procedure Durbin\_Watson.*

Algoritmo :

Inicio

1. Calcula sumade  $(e_i - e_{i-1})^2$   $i = 2, \dots, n$ .
2. Calcula sumade  $e_i^2$   $i = 1, \dots, n$ .

3. Divide 1./2. obteniendo a la estadística de Durbin-Watson.  
Fin.

Submódulo : Tabla de Análisis de Varianza. (Para Sig. Global y  $C\beta = \gamma$ ).

Asociación en PASCAL : *Procedure Tanova.*

Algoritmo :

Inicio

1. Genera Matriz  $C_{r \times p}$ .
2. Genera vector  $\gamma_{r \times 1}$ .
3. Calcula  $(X^t X)^{-1} C^t$ .
4. Calcula Inversa de  $C(X^t X)^{-1} C^t$ .
5. Calcula  $C\hat{\beta} - \gamma$ .
6. Calcula  $(C\hat{\beta} - \gamma)^t$ .
7. Calcula  $(C\hat{\beta} - \gamma)^t [C(X^t X)^{-1} C] (C\hat{\beta} - \gamma)$ .
8. Calcula Suma de Cuadrados del Error bajo la Hip. Ho.
9. Calcula Suma de Cuadrados del Error bajo el modelo reducido.
10. Calcula Tabla de Análisis de Varianza TANOVA.

Fin.

Submódulo F's Parciales.

Asociación en PASCAL : *Procedure F-Parcial.*

El algoritmo es el mismo que en *Tanaka*, la única diferencia radica en las dimensiones de la matriz C y el vector  $\gamma$ , es decir  $C_{np}$  y  $\gamma_{1 \times n}$ .

En este módulo se desarrolló el procedimiento *escribe\_vector*, para escribir los vectores que se necesitan posteriormente.

Submódulo : Cofactores.

Asociación en PASCAL : *Procedure Cofactores*.

Algoritmo :

- Inicio.
- 1. Lee  $i, j$  para calcular  $i, j$ .
- 2. Obtiene matriz auxiliar de Matcorre.
- 3. Calcula determinante de Matriz auxiliar.
- 4. reporta cofactor.
- Fin.

Módulo Graficación. Esta formado por varios submódulos asociados a procedimientos en PASCAL.

Descripción de Procedimientos.

Submódulo : Ordenamiento.

Asociación en PASCAL : *Procedure Ordena*.

Este submódulo se encarga de ordenar una serie de datos, lo cual es necesario en el caso de graficación de Papel Normal, el algoritmo utilizado para el ordenamiento fu el Quick Sort, tomado de *Tenenbaum and Augenstein, (1981)*.

Submódulo : Cuantiles.

Asociación en PASCAL : *Procedure Cuantil*.

Algoritmo :

- Inicio
- 1. Calcula cuantiles :  $100 \cdot (j - 0.5) / n$ .
- Fin.

Submódulo : Gráfica de Cp-Mallows.

Asociación en PASCAL : *Procedure Cp\_Mallows*.

Algoritmo :

- Inicio
- 1. Lee número de regresiones (p).
- 2. Lee valores de Cp.
- 3. Lee número de variables explicativas.
- 4. Calcula máximos y mínimos de Cp y p.
- 5. Calcula abscisas y ordenadas.

6. Calcula escalas para los ejes X y Y.
  7. Grafica Cp vs p.
- Fin.

Submódulo : Papel Normal.

Asociación en PASCAL : *Procedure PNormal.*

Algoritmo :

- Inicio
1. Pregunta gráfica que desea.
    - 1.1. Ri vs Ki. (donde Ki son los cuantiles)
    - 1.2. Ri vs Xi.
    - 1.3. Ri vs Yi.
  2. Lee gráfica seleccionada.
  3. Calcula máximos y mínimos de Ri y Ki, o xi o Yi, según sea el caso.
  4. Calcula abscisas y ordenadas.
  5. Calcula escalas para los ejes X y Y.
  6. Grafica Ri vs Ki, xi o Yi según sea el caso.
- Fin.

Submódulo : Graficación.

Asociación en PASCAL : *Procedure GrRiVsKi.*

Este submódulo fué creado con el fin de llamar a cualquiera de los anteriores.

Algoritmo :

- Inicio
1. Pregunta tipo de Gráfica.
    - 1.1. Ki vs Ri.
    - 1.2. Ri vs xi.
    - 1.3 Ri vs Yi
    - 1.4. Cp-Mallows.
  2. Lee tipo de gráfica.
  3. Si tipo = 1 2 ó 3 pregunta :
    - 3.1. Ri Yi - Yi.
    - 3.2. Ri studentizados.
    - 3.3. Lee tipo de Ri.
    - 3.4. Despliega Gráfica.
  - Si tipo = 4  
Despliega gráfica de Cp vs p.

Fin.

Módulo Principal. Este módulo agrupa a todos los anteriores a los cuales llama a ejecutar en el momento que se requiera.

Algoritmo del módulo Principal :

Inicio

1.Verifica Datos.

2.Despliega menu Principal :

1.Datos.

2.Análisis de Regresión.

3.Ayuda.

4.Salir.

3.Lee opción.

3.1.Si 1 : ejecuta *Cap\_dat* (Módulo de Captura).

3.2.Si 2 : Calcula :

medias

varianzas

Transpuestas

Matriz de correlación

Inveras

B-gorro

Residuales

Residuales studentizados

Ejecuta *menu* de Regresión :

Ejecuta *Presenta\_Menu* :

1.Matriz de Datos y vector Y

2.B-gorro

3.Cofactores

4.Matriz de Correlación

5.Hipótesis  $C\beta = \gamma$  tabla Anova

6.F-Parciales

7.Residuales

8.Matriz inversa de  $X^tX$

9.SCE bajo el modelo completo

10.Suma de cuadrados de Y

11.Suma de cuadrados de Y  
corregida por su media

12.Medias

13.Varianzas

14.Graficación

15.Durbin\_Watson

16.Coeficiente de Determinación

17.Fin

Lee opción

Si opción es :

1 : Ejecuta *Escribe* Datos y escribe  
vector Y.

2 : Ejecuta *Esc\_vector* B-gorro.

3 : Ejecuta *Cofactores*.

4 : Ejecuta *Escribe* Matriz de  
Correlación.

5 : Ejecuta *Tanova*.

6 : Ejecuta *F-Parcial*.

7 : Ejecuta *Resid\_est*.

8 : Ejecuta *Escribe* inversa de X'X.  
9 : Ejecuta *Esc\_SCEHC*.  
10 : Escribe *Sc\_y*.  
11: Escribe *Sc\_de\_ycorr*.  
12: Ejecuta *esc\_medes*.  
13: Ejecuta *esc\_var*.  
14: Ejecuta *GrRiusKI*.  
15: Ejecuta *Durbin\_Watson*.  
16: Escribe coef. de Determinación.  
17: Sale de *Ejecuta\_menu*.

3.3. Si 3 :Despliega texto de ayuda en el cual se encuentran las instrucciones del paquete.

3.4. Si 4 : Sale de menu principal y del paquete.

Fin.

En el módulo principal se desarrollaron dos procedimientos, *presenta\_menu* y *ejecuta\_menu*, el primero despliega el menú de 1 a 16 opciones, y el segundo llama a los procedimientos necesarios para la ejecución de cada una de las opciones de *presenta\_menu*.

Capítulo 5 : Aplicación del paquete mediante los datos de HALD.

5.1.Hald.

Para ilustrar el uso del paquete, se escogieron los datos de Hald, presentados en *Industrial and Engineering Chemistry*, 24, 1932, 1207-14, tabla I, por H. Woods, H. H. Stetnour and H. R. Starke.

En estudios hechos sobre el calentamiento generado durante el endurecimiento del cemento de Portland (Oregon, EUA), se supone una función de composición química. Las variables medidas son :

$x_1$  = Cantidad de Tricalcio de Aluminio.

$x_2$  = Cantidad de Tricalcio de Silicato.

$x_3$  = Cantidad de Tetracalcio de Aluminio Férrico.

$x_4$  = Cantidad de Dicalcio de Silicato.

$x_5$  = Y = Calentamiento desarrollado en calorías por gramo de cemento.

| $x_i$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| 1     | 7     | 26    | 6     | 60    | 78.5  |
| 2     | 1     | 29    | 15    | 52    | 74.3  |
| 3     | 11    | 56    | 8     | 20    | 104.3 |
| 4     | 11    | 31    | 8     | 47    | 87.6  |
| 5     | 7     | 52    | 6     | 33    | 95.9  |
| 6     | 11    | 55    | 9     | 22    | 109.2 |
| 7     | 3     | 71    | 17    | 6     | 102.7 |
| 8     | 1     | 31    | 22    | 44    | 72.5  |
| 9     | 2     | 54    | 18    | 22    | 93.1  |
| 10    | 21    | 47    | 4     | 26    | 115.9 |
| 11    | 1     | 40    | 23    | 34    | 83.8  |
| 12    | 11    | 66    | 9     | 12    | 113.3 |
| 13    | 10    | 68    | 8     | 12    | 109.4 |

Datos de Hald.

Las variables  $x_1$  a  $x_4$  son medidas como porcentajes de hulla en el cemento producido.

Se supone que el calor generado durante el endurecimiento es una función lineal de las cuatro variables  $x_1$  a  $x_4$ , el modelo de RLM es

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \beta_4 x_{4i} + \epsilon_i \quad i = 1, \dots, 13$$

Como se vió en capitulos anteriores el error  $\epsilon$  es un vector aleatorio con  $D_n$  Normal  $(0, I\sigma^2)$ .

Los métodos de Regresión son usados para estimar  $\beta$  y probar hipótesis acerca de ellas.

## 5.2. Aplicación del paquete a los datos de Hald.

A continuación se exponen los pasos a seguir para ejecutar el paquete mediante los datos de Hald.

1. Creación del archivo de datos. En esta parte se crea el archivo correspondiente a los datos de Hald los cuales podrán ser manipulados de acuerdo a las necesidades del usuario. Los pasos a seguir son los siguientes :

- i) Dentro del menú principal se debe entrar a la opción Datos.
- ii) En el menú de Datos se debe entrar a la opción Crear.
- iii) Dentro de la opción Crear, se leen el nombre del archivo así como el número de renglones y el de columnas, los cuales no deberán ser mayores a 50 y 20 respectivamente, una vez que se proporcionaron estos datos se procede a la captura, la cual terminará automáticamente cuando se llegue al número de renglones leído. En ese momento, se despliega una lista de los datos para que el usuario verifique que están correctos, de no estarlos puede modificarlos en la opción de Modificar que se presenta más adelante, sino hubo errores el archivo queda listo para ser abierto, cabe recordar que como se mencionó en el capítulo 4 una vez que se crea un archivo, ya no es necesario volverlo a crear cuando se vuelva a correr el paquete .

2. Abriendo el archivo. Una vez que se creó el archivo de datos se debe abrir el archivo, para poder pasar a cualquiera de las siguientes opciones, recordando que si se entra a Listar o a Modificar se debe regresar a abrir el archivo, una vez que se abrió el siguiente paso es escoger las columnas que representan las variables explicativas y la variable de respuesta.

3. Listando el archivo. Si se desea ver los datos, basta con dar el nombre del archivo para poderlo ver, recordando que después de entrar a esta opción es necesario regresar a abrir el archivo.

4. Modificando el archivo. Si existe algún error en los datos a través de esta opción se pueden hacer las modificaciones que sean necesarias sólo basta con dar el nombre del archivo, posteriormente aparecerá otro menú el cual se pueden hacer modificaciones por renglón, columna o celda. Una vez que se hicieron los cambios en el archivo se debe regresar a abrirlo nuevamente para asegurar que se tome el archivo ya modificado con el cual se trabajará a lo largo del paquete.

5. Selección de Columnas. Una vez que se tiene el archivo creado y abierto se procede a escoger las columnas donde se encuentran las variables explicativas y la variable dependiente, en el caso de Hald, la variable de respuesta se metió en la columna 5 y las variables explicativas  $X_1$  a  $X_4$  en las columnas 1 a 4, respectivamente, con esta opción se da por terminado lo referente a la opción Datos del menú principal. Cabe señalar que si se presentan dudas de esta opción dentro de la misma existe una de ayuda para facilitar el seguimiento de la misma.

Una vez que se llevó a cabo la opción Datos del menú principal el siguiente paso es pasar a la opción 2. Análisis de Regresión, en la cual antes de entrar al menú presenta algunos avisos con respecto a la impresión de resultados así como la salida de cada una de las opciones del menú de regresión. Inmediatamente después, se presenta en menú de Análisis de Regresión toda la información que se puede obtener a través del mismo, en el caso de Hald se obtuvieron los resultados que a continuación se describen.

### 5.3. Resultados obtenidos a través del paquete de Cómputo.

Opción : Matriz de Diseño y vector de variables de respuesta.

|     |      |      |      |      |       |
|-----|------|------|------|------|-------|
| 1.0 | 7.0  | 25.0 | 6.0  | 60.0 | 78.5  |
| 1.0 | 1.0  | 29.0 | 15.0 | 52.0 | 74.8  |
| 1.0 | 11.0 | 55.0 | 8.0  | 20.0 | 104.8 |
| 1.0 | 11.0 | 31.0 | 8.0  | 47.0 | 87.0  |
| 1.0 | 7.0  | 52.0 | 6.0  | 38.0 | 65.0  |
| 1.0 | 11.0 | 57.0 | 9.0  | 22.0 | 109.2 |
| 1.0 | 9.0  | 71.0 | 17.0 | 6.0  | 107.7 |
| 1.0 | 1.0  | 31.0 | 22.0 | 44.0 | 72.5  |
| 1.0 | 2.0  | 54.0 | 18.0 | 22.0 | 98.1  |
| 1.0 | 21.0 | 47.0 | 4.0  | 26.0 | 115.9 |
| 1.0 | 1.0  | 40.0 | 28.0 | 34.0 | 88.8  |
| 1.0 | 11.0 | 66.0 | 9.0  | 12.0 | 113.8 |
| 1.0 | 10.0 | 68.0 | 8.0  | 12.0 | 100.4 |

Opción : Vector de Parámetros Estimados  $\hat{\beta}$ . De la opción 2.B-gorro, se obtuvo :

$$\hat{\beta} = \begin{bmatrix} 02.4054 \\ 1.9511 \\ 0.1019 \\ -0.1441 \end{bmatrix}$$

De donde

$$\hat{Y} = X\hat{\beta}$$



Opción : Matriz de Correlaciones. ( $\underline{Y}$  con  $x_1, x_2, x_3, y x_4$  )

| Y     | $x_1$   | $x_2$   | $x_3$  | $x_4$   |
|-------|---------|---------|--------|---------|
| Y     | 1.0000  | 0.7907  | 0.8108 | -0.5847 |
| $x_1$ | 0.7907  | 1.0000  | 0.2280 | -0.2454 |
| $x_2$ | 0.8108  | 0.2280  | 1.0000 | 0.1907  |
| $x_3$ | -0.5847 | 0.2454  | 0.1907 | 1.0000  |
| $x_4$ | -0.2454 | -0.2454 | 0.0205 | 1.0000  |

De la matriz de Correlaciones se puede observar que las variables más coorrealcionadas con Y son  $x_2$  y  $x_1$ , mientras que  $x_4$  presenta una correlación negativa con Y, así también  $x_4$  esta altamente coorrealcionada con  $x_2$ , estos resultados son de gran importancia en selección de variables.

Es importante mencionar que la primera columna y el primer renglón de esta matriz es la variable Y.

Opción : Tabla Anova

Prueba de Significancia Global.

Se verá si  $\beta_1, \dots, \beta_4 = 0$  mediante la prueba de significancia global :

$$H_0 : C\beta = 0 \quad \text{vs.} \quad H_a : C\beta \neq 0$$

con

$$C = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Tabla de Análisis de Varianza

| f.var | g.l. | SCE       | CME      | fcalc.   |
|-------|------|-----------|----------|----------|
| $H_0$ | 4    | 2667.9008 | 666.9752 | 111.4823 |
| MC    | 8    | 47.8623   | 5.9828   |          |
| MR    | 12   | 2715.7631 |          |          |

No.de datos : 13    No.de parámetros : 5  
 Coeficiente de Determinación : 0.9824

De la tabla de Análisis de Varianza se obtuvo una F-calculada de 111.29 y comparándola con una F de tablas con 4 y 8 grados de libertad con un  $\alpha$  de 0.05 se obtuvo que :

$$F(4,8) = 3.84 < F\text{-Calculada} = 111.29 \quad \alpha = 0.05$$

Del resultado anterior se puede afirmar que existe suficiente evidencia en contra de alguna  $\beta_i$  sea igual a cero por lo tanto el modelo de regresión es significativo.

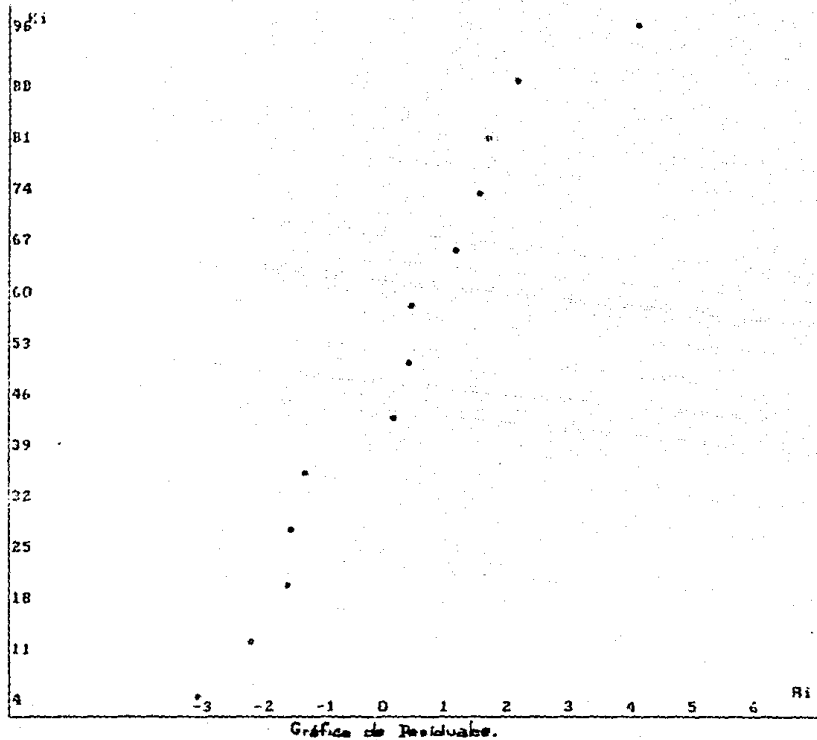
Opción : Coeficiente de Determinación.

$$R^2 = 0.9824$$

El coeficiente de Determinación está indicando que más del 98% de la variabilidad total del as  $Y_i$ 's queda explicada por el modelo.

Opción : Residuales. Aquí se presentan los residuales, de la forma  $Y_i - \hat{Y}_i$  y los studentizados. Posteriormente se puede pasar a la opción de Gráficación, mediante la cual se pueden obtener diferentes gráficas de residuales con las cuales se puede hacer un análisis de supuestos.

| I  | Y        | YESTIMADA | RESIDUAL | RESID.STUD |
|----|----------|-----------|----------|------------|
| 1  | 78.5000  | 78.4952   | 0.0048   | 0.0029     |
| 2  | 74.3000  | 72.7888   | 1.5112   | 0.7566     |
| 3  | 104.3000 | 105.9709  | -1.6709  | -1.0503    |
| 4  | 87.6000  | 89.3271   | -1.7271  | -0.8411    |
| 5  | 95.9000  | 95.4492   | 0.2508   | 0.1279     |
| 6  | 109.2000 | 105.2746  | 3.9254   | 1.7148     |
| 7  | 102.7000 | 104.1487  | -1.4487  | -0.7445    |
| 8  | 72.5000  | 75.6750   | -3.1750  | -1.6878    |
| 9  | 93.1000  | 91.7217   | 1.3783   | 0.6708     |
| 10 | 115.9000 | 115.6185  | 0.2815   | 0.2103     |
| 11 | 83.8000  | 81.8090   | 1.9910   | 1.0739     |
| 12 | 113.3000 | 112.3270  | 0.9730   | 0.4634     |
| 13 | 109.4000 | 111.6943  | -2.2943  | -1.1241    |



## Selección de Variables.

Todas las regresiones posibles.

Para poder aplicar este método mediante el paquete se deben ir eligiendo las columnas que vayan involucrando cada uno de los modelos, como el número de parámetros es igual a 5 se tendrán 2<sup>4</sup> modelos de regresión.

Para los modelos que involucran un cierto número de variables explicativas se tendrá que seleccionar en el módulo de Datos las columnas de la variable de respuesta y la(s) correspondiente(s) a la (s) variable(s) explicativa(s), así por ejemplo si se va a ajustar el modelo con  $x_1$  se elegirá la columna donde se encuentra esta variable junto con la columna de la variable  $y$ , después de esto se regresa al menú principal para ir a la opción de Análisis de Regresión, donde se obtendrá la información requerida para ese modelo.

Cabe aclarar que cada vez que se dese ajustar un modelo será necesario regresar al módulo de Datos para dar el nuevo conjunto de columnas que contienen a las variables de interés.

Tabla 1. Todas las regresiones posibles.

| variables |       |       |       | $\hat{\beta}_0$ | $\hat{\beta}_1$ | $\hat{\beta}_2$ | $\hat{\beta}_3$ | $\hat{\beta}_4$ | $R^2$ | $\bar{R}^2$ |
|-----------|-------|-------|-------|-----------------|-----------------|-----------------|-----------------|-----------------|-------|-------------|
| $x_1$     |       |       |       | 81.48           | 1.86            |                 |                 |                 | 0.583 | 0.401       |
| $x_1$     | $x_2$ |       |       | 57.42           |                 | 0.78            |                 |                 | 0.666 | 0.686       |
| $x_1$     |       | $x_3$ |       | 110.20          |                 |                 | -1.25           |                 | 0.285 | 0.220       |
| $x_1$     |       |       | $x_4$ | 883.9           |                 |                 |                 | -0.78           | 0.674 | 0.644       |
| $x_1$     | $x_2$ |       |       | 52.58           | 1.46            | 0.66            |                 |                 | 0.978 | 0.974       |
| $x_1$     |       | $x_3$ |       | 72.35           | 2.31            |                 | 0.40            |                 | 0.548 | 0.457       |
| $x_1$     |       |       | $x_4$ | 108.10          | 1.44            |                 |                 | -0.61           | 0.972 | 0.967       |
| $x_1$     | $x_2$ | $x_3$ |       | 72.07           |                 | 0.73            | -1.00           |                 | 0.847 | 0.816       |
| $x_1$     | $x_2$ |       | $x_4$ | 94.16           |                 | 0.31            |                 | -0.45           | 0.680 | 0.616       |
| $x_1$     |       | $x_3$ | $x_4$ | 191.28          |                 |                 | -1.70           | -0.72           | 0.995 | 0.922       |
| $x_1$     | $x_2$ | $x_3$ |       | 48.19           | 1.60            | 0.65            | 0.75            |                 | 0.982 | 0.9709      |
| $x_1$     | $x_2$ |       | $x_4$ | 71.05           | 1.45            | 0.41            |                 | -0.23           | 0.982 | 0.9764      |
| $x_1$     |       | $x_3$ | $x_4$ | 111.68          | 1.05            |                 | -0.41           | -0.64           | 0.981 | 0.975       |
| $x_1$     | $x_2$ | $x_3$ | $x_4$ | 208.64          |                 | -0.97           | -1.44           | -1.55           | 0.972 | 0.968       |
| $x_1$     | $x_2$ |       | $x_4$ | 62.38           | 1.55            | 0.71            | 0.10            | -0.14           | 0.982 | 0.978       |

### Criterio de la $R^2$ .

De la tabla 1 se puede observar que del conjunto de modelos que contienen una sola variable explicativa la  $R^2$  mayor es la correspondiente a la variable  $x_4$ , por lo que se escoge el modelo que contiene a esta.

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_4 x_4 \quad R^2 = .674$$

Del conjunto de modelos que contienen dos variables explicativas, existe un empate para los modelos que contienen a las variables  $x_1, x_2, x_3$  y  $x_1, x_2$  y  $x_4$  ya que presentan una  $R^2$  de .982, en este caso habría que analizar la correlación entre esta variables.

A partir de la  $R^2$  se pueden ir calculando las  $R^2$  ajustadas y tomar el modelo que tenga la  $R^2$  ajustada mayor.

### Opción :Graficación : Cp-Mallows.

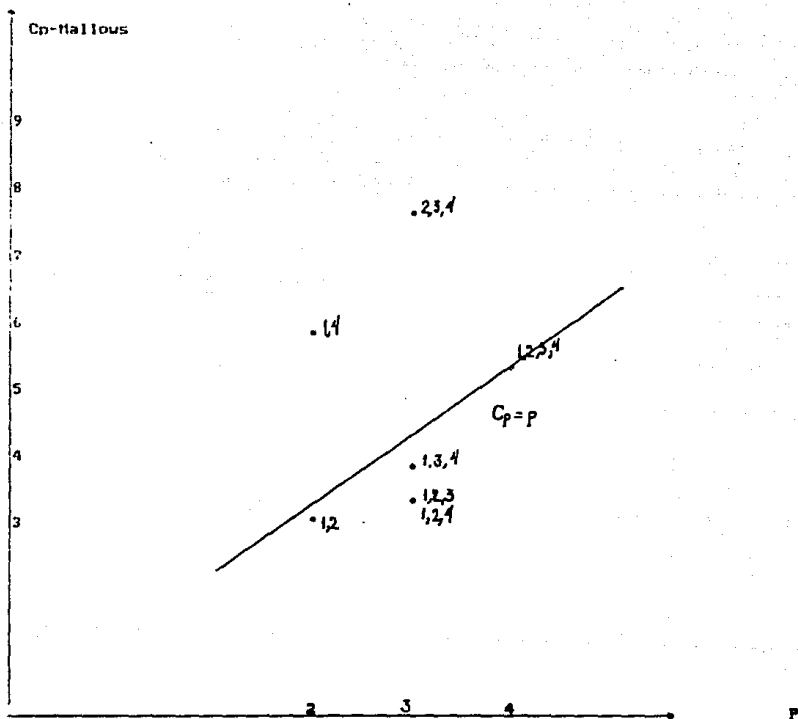
De acuerdo a la definición de la Cp vista en el capítulo 3, se tomaron algunos modelos para describir el uso del paquete, obteniendo los siguientes valores :

$$SCE_{MC} = 5.983$$

|                 |                     |            |
|-----------------|---------------------|------------|
| Para $x_1, x_4$ | $SCE_{MC} = 1265.6$ | $Cp = 5.5$ |
| $x_1, x_2$      | $SCE_{MC} = 57.9$   | $Cp = 2.7$ |
| $x_2, x_3, x_4$ | $SCE_{MC} = 73.8$   | $Cp = 7.3$ |
| $x_1, x_3, x_4$ | $SCE_{MC} = 50.8$   | $Cp = 3.5$ |
| $x_1, x_2, x_4$ | $SCE_{MC} = 48.0$   | $Cp = 3.0$ |
| $x_1, x_2, x_3$ | $SCE_{MC} = 48.1$   | $Cp = 3.0$ |

Para calcular cada uno de los valores de las Cp, se debe ajustar el modelo de acuerdo a las variables que se estén involucrando, es decir se debe ir al módulo de Datos y dar las columnas adecuadas de acuerdo a las variables que contenga cada modelo pues para cada uno se reportan en el paquete las  $SCE_{MC}$  correspondientes.

En esta opción el alumno deberá calcular las Cp para cada uno de los modelos, pues el paquete aporta la información necesaria para hacerlo. Bajo este criterio se llega a lo siguiente :



De la gráfica (opción Graficación-Cp Mallows) se observa que los mejores subconjuntos son :

$(x_1, x_2)$  y el  $(x_1, x_2, x_3, y x_4)$

Para aplicar métodos de selección de variables es necesario ir haciendo regresiones con diferentes subconjuntos de variables, mediante el paquete se pueden obtener estas regresiones para lo cual se deben seguir los siguientes pasos, primero se regresa al menú principal dentro del cual se deberá entrar de nuevo a la opción Datos, una vez que se encuentra en dicha opción, se procede a hacer una nueva selección de las columnas donde se encuentran las variables explicativas según el subconjunto que se haya seleccionado, la variable de respuesta, en este caso  $x_5$ , sigue estando en la columna 5 del archivo, si por ejemplo se desea hacer una regresión con las variables  $x_1, x_2$  y  $x_3$  se deberán dar las columnas 1, 2 y 3 respectivamente.

Backward. :

Opción : F-parcial. Con esta opción será posible aplicar métodos de selección de variables, pues según el modelo que se tenga se reportan las respectivas F's parciales facilitando así la aplicación del método de selección.

En el caso de Backward, se debe empezar con el modelo completo, para lo cual se deben de proporcionar las columnas correspondientes que contiene dicho modelo para después pasar al menú de Análisis de Regresión y en este obtener las F's parciales adecuadas.

Paso 1.

$$\text{Modelo Completo : } Y = 02.41 + 1.55 X_1 + 0.51 X_2 + 0.10 X_3 - 0.14 X_4$$

| Variable | F-parcial ( $\beta_j = 0$ ) |
|----------|-----------------------------|
| $x_1$    | 4.33                        |
| $x_2$    | 0.50                        |
| $x_3$    | 0.02                        |
| $x_4$    | 0.04                        |

La F-parcial de menor valor es  $F = 0.02$   
Comparando con una F-Tablas

$$F_{(4,8)}^{0.05} = 9.84 > 0.02 \quad \therefore$$

Se elimina la variable  $x_3$

En este momento se debe regresar al módulo de Datos para dar el nuevo subconjunto de variables ( $x_1$ ,  $x_2$  y  $x_4$ ).

Paso 2.

$$\text{Nuevo Modelo : } \hat{Y} = 71.64 + 1.45 X_1 + 0.41 X_2 - 0.28 X_4$$

| Variable | F-parcial ( $\beta_j = 0$ ) |
|----------|-----------------------------|
| $x_1$    | 154.01                      |
| $x_2$    | 5.03                        |
| $x_4$    | 1.86                        |

La F-parcial de menor valor es  $F = 1.86$   
Comparando con una F-tablas

$$F_{(3,9)}^{0.05} = 3.86 > 1.86 \quad \therefore$$

Se elimina la variable  $x_4$ .

$$R^2 = 0.982$$

Después de haber dado las columnas del nuevo subconjunto de variables se procede con el siguiente paso.

Paso 3.

$$\text{Nuevo Modelo : } \hat{Y} = 52.57 + 1.46 X_1 + 0.66 X_2$$

| Variable | F-parcial ( $\beta_j = 0$ ) |
|----------|-----------------------------|
| $x_1$    | 146.52                      |
| $x_2$    | 208.58                      |

La f-parcial de menor valor es  $F = 146.52$   
Comparando con una F-Tablas se tiene que

$$F_{(2,10)}^{0.05} = 4.96 < 146.52 \quad \therefore$$

No sale ninguna variable.

$$R^2 = .97$$



Para este subconjunto de variables :

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$$

El vector de parámetros estimados es :

$$\hat{\beta} = \begin{bmatrix} 22.97 \\ 1.40 \\ 0.00 \end{bmatrix}$$

Obteniendo el nuevo modelo :

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}$$

Este modelo se tomará como el más significativo.

Cabe mencionar que para aplicar otros métodos de Selección de Variables, como Stepwise y Forward, el paquete proporciona los cofactores a través de los cuales se calculan los coeficientes de correlación parcial, mismos que son utilizados bajo estos metodos, por ejemplo si se desea calcular  $r_{y x_1 x_2}$  el paquete proporciona

los cofactores correspondientes, recordando que para este conjunto de variables la matriz de correlaciones será :

$$\begin{bmatrix} 1.0 & 0.78 & -0.82 \\ 0.78 & 1.0 & -0.24 \\ -0.82 & -0.24 & 1.0 \end{bmatrix}$$

$R_{yx_1} = (-1)^{1+1} a_{12}$  eliminando el renglón 1 y columna 2

$$\text{donde } a_{12} = \begin{vmatrix} 0.78 & -0.24 \\ -0.82 & 1.0 \end{vmatrix} = -0.5291$$

Se llega a que

$$R_{yx_1} = R_{12} = -(-0.5291) = 0.5291$$

De la misma forma se obtiene

$$R_{22} = 0.32 \quad , \quad R_{33} = 0.93$$

$$r_{y x_1 x_2} = -R_{12} / \sqrt{R_{11} R_{22}} = 0.9566$$

## Conclusiones.

En este paquete se intentó proporcionar elementos a través de los cuales se llegara de una manera más rápida a los cálculos necesarios para el Análisis de Regresión, tratando de no perder de vista la teoría en la que se basa esta técnica estadística.

El paquete de cómputo se elaboró abarcando lo mínimo que se utiliza en un curso de Análisis de Regresión aunque cabe señalar que aún puede ampliarse mucho más. Sin embargo con lo tratado en este trabajo se puede llegar a realizar numerosas aplicaciones.

Este trabajo puede tomarse como material didáctico por quienes imparten los cursos de Análisis de Regresión así como también puede ser de utilidad para quienes tienen interés en la materia.

A través de este paquete los alumnos de Análisis de Regresión e incluso de otras materias relacionadas con la materia podrán obtener información, interpretarla y no perder de vista la forma en la que se llega a la misma, ya que proporciona las herramientas necesarias para calcular las estadísticas relevantes permitiendo así, no perder de vista los métodos de cálculo que frecuentemente se olvidan cuando se utilizan paquetes reconocidos, que más bien son recomendables para quienes dominan la materia y no para los alumnos pertenecientes a un curso de la misma.

En resumen se cubrieron tres objetivos principalmente :

- La presentación muy general de la teoría del Análisis de Regresión, lo cual puede servir en momento dado como material de consulta para quienes estén interesados en la materia. Si se desea profundizar más en la misma, la bibliografía de este trabajo puede servir de referencia.

- La aplicación de un paquete de cómputo que permita a los alumnos de los cursos de Análisis de regresión, obtener resultados perdiendo menos tiempo que si tuvieran que hacer todos los cálculos manualmente y sobre todo sin perder de vista la teoría tanto matemática como estadística, lo cual es una de las inquietudes de los profesores que imparten dichos cursos.

- La aplicación de un paquete de cómputo a una de las técnicas estadísticas más importantes que se conocen, a través del cual se obtiene la información suficiente para hacer una inferencia acerca del fenómeno que se esté estudiando en un momento dado.

Cabe señalar que hoy en día la teoría de la computación crece a pasos agigantados por lo cual este trabajo puede ser retomado en el futuro y ser mejorado con las nuevas innovaciones que esta teoría computacional ofrezca.

Por último, no se pretende competir con ningún paquete estadístico conocido y mucho menos cubrir todos los casos que podrían enumerarse en Regresión.

**A P E N D I C E I.**

## Apéndice 1.

### 1. La Normal Multivariada.

La función de densidad de probabilidad Normal, en el caso multivariado, se expresa como sigue :

$$f(\underline{Y}, \underline{\theta}, \Sigma) = \frac{1}{(2\pi)^{p/2}} \exp\left[-\frac{1}{2}(\underline{Y}-\underline{\theta})' \Sigma^{-1} (\underline{Y}-\underline{\theta})\right]$$

$\forall \underline{Y} \in \mathbb{R}^p$ ,  $\underline{\theta} \in \mathbb{R}^p$  y  $\Sigma$  definida positiva.

Resultados respecto a la Normal Multivariada.

1.  $E(\underline{Y}) = \underline{\theta}$  y  $V(\underline{Y}) = \Sigma$
2.  $(\underline{Y} - \underline{\theta})' \Sigma^{-1} (\underline{Y} - \underline{\theta}) \sim \chi^2_{(p)}$
3. Si  $\underline{Y} \sim N_p(\underline{\theta}, \Sigma)$  entonces  $(\underline{Y} - \underline{\theta}) \sim N_p(0, \Sigma)$

### 2. Formas Cuadráticas.

Si  $\underline{Y}$  es un vector de  $n \times 1$  y  $A$  es una matriz de  $n \times n$  simétrica entonces la forma cuadrática

$$\underline{Y}' A \underline{Y} = \sum_{i=1}^n \sum_{j=1}^n Y_i Y_j a_{ij}$$

- 2.1. Se dice que  $\underline{Y}' A \underline{Y}$  es positiva definida (p.d.)  $\Leftrightarrow \underline{Y}' A \underline{Y} > 0 \quad \forall \underline{Y} \neq 0$
- 2.2. Se dice que  $\underline{Y}' A \underline{Y}$  es positiva semidefinida  $\Leftrightarrow \underline{Y}' A \underline{Y} \geq 0 \quad \forall \underline{Y} \in \mathbb{R}^n$  tal que  $\underline{Y}' A \underline{Y} = 0$
- 2.3.  $A$  es positiva definida si  $\underline{Y}' A \underline{Y} > 0 \quad \forall \underline{Y} \neq 0$   
 $A$  es positiva semidefinida si  $\underline{Y}' A \underline{Y} \geq 0 \quad \forall \underline{Y}$  y  $\exists \underline{Y} \neq 0$  tal que  $\underline{Y}' A \underline{Y} = 0$
- 2.4. Si  $P$  es no singular y  $A$  es positiva definida entonces  $P' A P$  es p.d.
- 2.5. Una condición necesaria y suficiente para que  $A$  (simétrica) sea p.d. es que exista  $P$ , una matriz no singular tal que  $A = P' P$ .
- 2.6. Una condición necesaria y suficiente para que  $A$  sea p.d. es que los menores principales sean positivos ( $\text{Det } A > 0$ ).
- 2.7.  $A$  es definida negativa si los menores principales alternan el signo.
- 2.8. Si  $A_{n \times m}$  es una matriz de rango  $n \times m$  entonces  $A' A$  es positiva definida y  $A A'$  es positiva semidefinida.
- 2.9. Si  $A$  es positiva semidefinida entonces  $\text{Tr}(A) \geq 0$ .
- 2.10. Si  $A$  es positiva definida  $A^{-1}$  también lo es.

2.11. Los elementos diagonales de una matriz positiva definida son todos positivos.

### 3. Distribución de Formas Cuadráticas.

3.1. Si  $Y \sim N_n(0, I)$  una condición necesaria y suficiente para que la forma cuadrática  $Y^t A Y \sim X_{(k)}^2$  es que  $A$  sea una matriz idempotente de rango  $k$ .

3.2. Si  $Y \sim N_n(\theta, I)$  entonces  $Y^t A Y \sim X_{(k, \lambda)}^2$  con  $\lambda = \theta^t A \theta / 2 \iff A$  es una matriz idempotente de rango  $k$ .

3.3. Si  $Y \sim N_n(\theta, \sigma^2 I)$  entonces  $Y^t A Y / \sigma^2 \sim X_{(k, \lambda)}^2$  con  $\lambda = \theta^t A \theta \iff A$  es una matriz idempotente de rango  $k$ .

3.4. Si  $Y \sim N_n(0, \Sigma)$ , entonces  $Y^t B Y \sim X_{(k)}^2 \iff B \Sigma^{-1}$  es idempotente de rango  $k$ .

3.5. Si  $Y \sim N_n(\theta, \Sigma)$  entonces  $Y^t B Y \sim X_{(k, \lambda)}^2$  con  $\lambda = \theta^t B \theta / 2$  y  $B$  de rango  $k \iff B \Sigma^{-1}$  es idempotente.

3.6. Distribución de Ji-cuadrada no central ( $X^2$ )

3.6.1. Si  $Y \sim N_n(0, I)$  entonces  $Y^t Y = \sum_{i=1}^n Y_i^2 \sim X_{(n)}^2$

3.6.2. Si  $Y \sim N_n(\theta, I)$  entonces  $Y^t Y \sim X_{(n, \lambda)}^2$ , es decir,  $Y^t Y$  es una  $X^2$  no central con  $n$  grados de libertad y parámetro de no centralidad  $\lambda = \theta^t \theta / 2$

3.6.3. Si  $w_1, w_2, \dots, w_k$  son independientes y si  $w_i \sim X_{(p_i, \lambda_i)}^2$  ( $i=1, \dots, k$ ), entonces  $W = \sum_{i=1}^k w_i \sim X_{(p, \lambda)}^2$  con  $p = \sum_{i=1}^k p_i$  y  $\lambda = \sum_{i=1}^k \lambda_i$

3.6.4. Si  $Y \sim N_n(\theta, \sigma^2 I)$ , entonces  $Y^t Y / \sigma^2 \sim X_{(n, \lambda)}^2$  con  $\lambda = \theta^t \theta / 2\sigma^2$

3.6.5. Si  $Y \sim N_n(\theta, D)$  con  $D$  diagonal, entonces  $Y^t D^{-1} Y \sim X_{(n, \lambda)}^2$  con  $\lambda = \theta^t D^{-1} \theta / 2$

### 4. Independencia de Formas Cuadráticas.

4.1. Si  $Y \sim N_p(\theta, \Sigma)$  entonces  $Y^t A Y$  y  $Y^t B Y$  son independientes  $\iff A \Sigma^{-1} B = 0$

4.2. Si  $Y \sim N(\theta, \sigma^2 I)$  la forma cuadrática positiva semidefinida  $Y^t A Y$  y  $Y^t B Y$  son independientes si  $\text{Tr}(AB) = 0$ .

ESTA TESIS NO DEBE  
SALIR DE LA BIBLIOTECA

APENDICE 2.

Lo cual significa que la variable  $y$  y  $x_4$  dado  $x_4$  están altamente correlacionadas.

Aplicando los métodos de selección antes mencionados se toma la variable más correlacionada con  $Y$ , y se continúa con el proceso según lo indiquen dichos métodos .

En el mesú de Analisis de Regresión también se reporta la estadística de Durbin y Watson la cual toma el valor de 2.0526.

## Apéndice 2.

Independencia entre las variables aleatorias  $SCE_{MC}$  y  $SCE_{HO}$ .

Nótese que

$$SCE_{HO} = (\bar{C}\bar{\beta} - \gamma)' [C(X^t X)^{-1} C^t] (\bar{C}\bar{\beta} - \gamma) \quad \text{con } \bar{\beta} = (X^t X)^{-1} X^t Y$$

entonces

$$SCE_{HO} = [C(X^t X)^{-1} X^t Y - \gamma]' [C(X^t X)^{-1} C^t]^{-1} [C(X^t X)^{-1} X^t Y - \gamma]$$

ahora

$$[C(X^t X)^{-1} X^t] X C^t = C C^t \Rightarrow [C(X^t X)^{-1} X^t] X C^t (C C^t)^{-1} = I$$

en donde  $(C C^t)^{-1}$  existe pues  $C$  es de rango completo.

La  $SCE_{HO}$  puede entonces expresarse como :

$$SCE_{HO} = [C(X^t X)^{-1} X^t (Y - X C^t (C C^t)^{-1} \gamma)]' [C(X^t X)^{-1} C^t]^{-1} [C(X^t X)^{-1} X^t (Y - X C^t (C C^t)^{-1} \gamma)] \\ = [Y - X C^t (C C^t)^{-1} \gamma]' X (X^t X)^{-1} C^t [C(X^t X)^{-1} C^t]^{-1} C (X^t X)^{-1} X^t (Y - X C^t (C C^t)^{-1} \gamma)$$

$$\text{Sea } \underline{Y}_* = \underline{Y} - X C^t (C C^t)^{-1} \gamma$$

Dado que  $\underline{Y} \Rightarrow N_n(X\beta, \sigma^2 I)$  se sigue que  $\underline{Y}_* \sim N_n(X\beta - X C^t (C C^t)^{-1} \gamma, \sigma^2 I)$

Sea  $\underline{\theta} = X\beta - X C^t (C C^t)^{-1} \gamma$  entonces  $\underline{Y}_* \sim N_n(\underline{\theta}, \sigma^2 I)$

$$SCE_{HO} = \underline{Y}_*^t X (X^t X)^{-1} C^t [C(X^t X)^{-1} C^t]^{-1} C (X^t X)^{-1} X^t \underline{Y}_*$$

$$\text{Sea } A = X (X^t X)^{-1} C^t [C(X^t X)^{-1} C^t]^{-1} C (X^t X)^{-1} X^t$$

Con respecto a  $SCE_{MC}$  se tiene que

$$SCE_{MC} = Y^t [I - X (X^t X)^{-1} X^t] Y$$



$$\text{además } X [I - X(X^t X)^{-1} X^t] = [I - X(X^t X)^{-1} X^t] X = 0$$

La SCE<sub>MC</sub> se puede expresar tambien como

$$\begin{aligned} SCE_{MC} &= (Y - X C^t (C C^t)^{-1} \gamma) [I - X(X^t X)^{-1} X^t] (Y - X C^t (C C^t)^{-1} \gamma) \\ &= Y^t [I - X(X^t X)^{-1} X^t] Y \quad \text{con } Y \text{ definida anteriormente} \end{aligned}$$

$$\text{Sea } B = [I - X(X^t X)^{-1} X^t]$$

Entonces se quiere demostrar que  $Y^t A Y$  (SCE<sub>MO</sub>) es independiente de  $Y^t B Y$  (SCE<sub>MC</sub>) donde  $Y \sim N_n(\theta, \sigma^2 I)$ , i.e.  $\Sigma = \sigma^2 I$

Se debe demostrar de acuerdo al resultado 4.1 del apendice 1 que  $A B = 0$

Sustituyendo se tiene que

$$\sigma^2 X(X^t X)^{-1} C^t [C(X^t X)^{-1} C^t]^{-1} C(X^t X)^{-1} X^t [I - X(X^t X)^{-1} X^t] = 0$$

$$\text{De donde } X^t [I - X(X^t X)^{-1} X^t] = 0$$

Por lo tanto SCE<sub>MC</sub> y SCE<sub>MO</sub> son independientes.

Distribución de SCE<sub>MC</sub>

$$SCE_{MC} = (Y - X \hat{\beta})^t (Y - X \hat{\beta}) = \frac{(Y - X(X^t X)^{-1} X^t Y)^t (Y - X(X^t X)^{-1} X^t Y)}{\frac{1}{Y^t [I - X(X^t X)^{-1} X^t] Y} \frac{1}{Y^t [I - X(X^t X)^{-1} X^t] Y}}$$

forma cuadrática con matriz asociada \*

Ahora sea  $P = X(X^t X)^{-1} X^t$  entonces la forma cuadrática se puede escribir como

$$Y^t [I - P] [I - P] Y$$

$[I - P]$  es idempotente ya que

$$[I - P][I - P] = [I - P] \quad \text{lo cual es facil de checar sustituyendo } P$$

P.D. que  $(n-p)\sigma^2$

$$\frac{\quad}{\sigma^2} X_{(n-p)}$$

es decir

$$\frac{Y^t [I - P] Y}{\sigma^2} X_{(n-p)}^t \quad (I - P) \text{ es idempotente}$$

De acuerdo al resultado 3.5 del apéndice 1, en este caso se tiene

$Y \sim N(X\beta, \sigma^2 I)$  entonces  $\hat{\beta} = X\beta$  y  $\Sigma = \sigma^2 I$  en este caso  $\frac{Y^t B Y}{2} \sim \chi^2_{(n, \lambda)}$   
 con  $\lambda = \hat{\beta}^t B \hat{\beta} / 2$  y  $B$  de rango  $k \Leftrightarrow B\Sigma$  es idempotente

Sea  $B = (I - P) / \sigma^2$

$B\Sigma = \frac{(I - P)}{\sigma^2} \sigma^2 I = (I - P)$  que es idempotente

Recordar que  $X$  de rango  $p$ ,  $(X^t X)_{p \times p}$  de rango  $p$ .

$(X^t X)_{p \times p}^{-1}$  de rango  $p$ ,  $X(X^t X)_{n \times p}^{-1}$  de rango  $p$ ,  $X(X^t X)^{-1} X^t$  de rango  $p$ .

$(I - P)$  es simétrica e idempotente, es decir es una matriz proyección

entonces  $\text{Tr}(A) = \text{rango}(A)$

$\text{rango}(I - P) = \text{Tr}(I - P) = \text{Tr}(I) - \text{Tr}(P)$

pero  $\text{Tr}(P) = \text{Tr}(X(X^t X)^{-1} X^t) = \text{Tr}(X^t X(X^t X)^{-1} X^t) = \text{Tr}(I) = p$

por lo tanto el rango de  $(I - P) = n - p$

entonces el rango de  $B$  es  $n - p$  por lo que los grados de libertad de la  $X$  son  $n - p$

Por otro lado  $\lambda = \hat{\beta}^t B \hat{\beta} / 2$

$\hat{\beta} = X\beta$  ,  $B = \frac{(I - P)}{\sigma^2}$

$$\lambda = \frac{1}{2} \frac{(X\beta)^t (I - P) (X\beta)}{\sigma^2}$$

$$= \frac{1}{2\sigma^2} (\hat{\beta}^t X^t) (I - P) (X\beta)$$

$$= \frac{1}{2\sigma^2} (\hat{\beta}^t X^t - \hat{\beta}^t X^t P) (X\beta)$$

$$= \frac{1}{2\sigma^2} (\hat{\beta}^t X^t - \hat{\beta}^t X^t X (X^t X)^{-1} X^t) X\beta$$

$$= \frac{1}{2\sigma^2} (\hat{\beta}^t X^t - \hat{\beta}^t X^t) X\beta = 0 \text{ por lo tanto } \lambda = 0$$

por lo tanto  $\frac{Y^t (I - P) Y}{\sigma^2} \sim \chi^2_{(n-p)}$

### Distribución de $SCE_{Ho}$

Del resultado 3.6.5 del apéndice 1, se tiene que

$$SCE_{Ho} = (C\hat{\beta} - \gamma)' \frac{C(X^t X)^{-1} C'}{\sigma^2} (C\hat{\beta} - \gamma) = Q$$

como  $(C\hat{\beta} - \gamma) \sim N_n(C\hat{\beta} - \gamma)$ ,  $C(X^t X)^{-1} C'$  y según el resultado 3.6.5.  $Y = (C\hat{\beta} - \gamma)$  basta probar que  $B \frac{1}{\sigma^2} C(X^t X)^{-1} C'$  es idempotente

$$\underbrace{\frac{1}{\sigma^2} C(X^t X)^{-1} C'}_B \cdot \underbrace{\sigma^2 C(X^t X)^{-1} C'}_{\Sigma} = I \quad \text{idempotente}$$

entonces

$$Q = SCE_{Ho} \sim \chi_{(r)}^2 \quad \text{Se tienen } r \text{ grados de libertad ya que } B = \frac{1}{\sigma^2} C(X^t X)^{-1} C' \text{ es de rango } r.$$

Dado que  $SCE_{Ho}/\sigma^2 \sim \chi_{(r)}^2$  bajo  $H_0$ .

$$SCE_{Mc}/\sigma^2 \sim \chi_{(n-p)}^2$$

Se tiene que

$$\frac{SCE_{Ho} / r \sigma^2}{SCE_{Mc} / (n-p) \sigma^2} \sim F(r, n-p) \quad \text{bajo } H_0$$

$$\frac{SCE_{Ho} / r}{SCE_{Mc} / (n-p)} \sim F(r, n-p) \quad \text{bajo } H_0$$

## Bibliografía.

Bernhard, F. Riedwil, H. (1988). *Multivariate Statistics "A practical Approach"*. Chapman & hall. London.

Chatterjee, S. and Price, B. (1977). *Regression Analysis by Example*. Wiley & Sons. New York.

Dongarra, J.J. (1979). *Linpac Users' Guide*. Siam. Philadelphia.

Draper, N.R, and Smith, H. (1981). *Applied Regression Analysis*. Wiley & Sons. New York.

Freund, R.J. and Minton, P.M., (1979). *Regression Methods*. Dekker, Inc. New York.

Hennefeld, J. (1989). *Turbo Pascal con aplicaciones 3.0, 4.0 y 5.0*. Grupo Editorial Iberoamérica. México, D.F.

Koffman, E.B. (1985). *Introducción al lenguaje y resolución de problemas con programación estructurada*. Fondo Educativo Interamericano. México, D.F.

Lang, S. (1976). *Algebra lineal*. Fondo Educativo Interamericano, S.A. México, D.F.

Mendez Ramírez Ignacio. (1981). *Modelos Estadísticos Lineales*. CONACYT. México, D.F.

Pindyck, R.S. and Rubinfeld, D.L. (1981). *Econometric Models and Economic Forecasts*. Mc.Graw. Hill. Tokio.

Tenenbaum, A.M. and Augenstein, M.J. (1981). *Estructura de Datos en Pascal*. Prentice-Hall. México, D.F.

Younger, M.S. (1979). *A handbook for linear Regression*. Duxbury Press. Massachusetts.